

Escaping from the Barren Plateau via Gaussian Initializations in Deep Variational Quantum Circuits

Kaining Zhang*, Liu Liu*, Min-Hsiu Hsieh[†], and Dacheng Tao^{*,‡}

*School of Computer Science, Faculty of Engineering, The University of Sydney, Sydney, Australia

[†]Hon Hai Quantum Computing Research Center, Taipei, Taiwan

[‡]JD Explore Academy, Beijing, China

Abstract

Variational quantum circuits have been widely employed in quantum simulation and quantum machine learning in recent years. However, quantum circuits with random structures have poor trainability due to the exponentially vanishing gradient with respect to the circuit depth and the qubit number. This result leads to a general standpoint that deep quantum circuits would not be feasible for practical tasks. In this work, we propose an initialization strategy with theoretical guarantees for the vanishing gradient problem in general deep quantum circuits. Specifically, we prove that under proper Gaussian initialized parameters, the norm of the gradient decays at most polynomially when the qubit number and the circuit depth increase. Our theoretical results hold for both the local and the global observable cases, where the latter was believed to have vanishing gradients even for very shallow circuits. Experimental results verify our theoretical findings in the quantum simulation and quantum chemistry.

1 Introduction

Quantum computing has attracted great attention in recent years, especially since the realization of quantum supremacy [1, 2] with noisy intermediate-scale quantum (NISQ) devices [3]. Due to mild requirements on the gate noise and the circuit connectivity, variational quantum algorithms (VQAs) [4] become one of the most promising frameworks for achieving practical quantum advantages on NISQ devices. Specifically, different VQAs have been proposed for many topics, e.g., quantum chemistry [5–13], quantum simulations [14–23], machine learning [24–31], numerical analysis [32–36], and linear algebra problems [37–39]. Recently, various small-scale VQAs have been implemented on real quantum computers for tasks such as finding the ground state of molecules [10–12] and exploring applications in supervised learning [25], generative learning [30] and reinforcement learning [29].

Typical variational quantum algorithms is a trainable quantum-classical hybrid framework based on parameterized quantum circuits (PQCs) [40]. Similar to classical counterparts such as neural networks [41], first-order methods including the gradient descent [42] and its variants [43] are widely employed in optimizing the loss function of VQAs. However, VQAs may face the trainability barrier when scaling up the size of quantum circuits (i.e., the number of involved qubits or the circuit depth), which is known as the barren plateau problem [44].

Roughly speaking, the barren plateau describes the phenomenon that the value of the loss function and its gradients concentrate around their expectation values with exponentially small variances. We remark that gradient-based methods could hardly handle trainings with the barren plateau phenomenon [45]. Both the machine noise of the quantum channel and the statistical noise induced by measurements could severely degrade the estimation of gradients. Moreover, the optimization of the loss with a flat surface takes much more time using inaccurate gradients than ideal cases. Thus,

solving the barren plateau problem is imperative for achieving practical quantum advantages with VQAs. In this paper, we propose Gaussian initializations for VQAs which have theoretical guarantees on the trainability. We prove that for Gaussian initialized parameters with certain variances, the expectation of the gradient norm is lower bounded by the inverse of the polynomial term of the qubit number and the circuit depth. Technically, we consider various cases regarding VQAs in practice, which include local or global observables, independently or jointly employed parameters, and noisy optimizations induced by finite measurements. To summarize, our contributions are fourfold:

- We propose a Gaussian initialization strategy for deep variational quantum circuits. By setting the variance $\gamma^2 = \mathcal{O}(\frac{1}{L})$ for N -qubit L -depth circuits with independent parameters and local observables, we lower bound the expectation of the gradient norm by $\text{poly}(N, L)^{-1}$ as provided in Theorem 4.1, which outperforms previous $2^{-\mathcal{O}(L)}$ results.
- We extend the gradient norm result to the global observable case in Theorem 4.2, which was believed to have the barren plateau problem even for very shallow circuits. Moreover, our bound holds for correlated parameterized gates, which are widely employed in practical tasks like quantum chemistry and quantum simulations.
- We provide further analysis on the number of necessary measurements for estimating the gradient, where the noisy case differs from the ideal case with a Gaussian noise. The result is presented in Corollary 4.3, which proves that $\mathcal{O}(\frac{L}{\epsilon})$ times of measurement is sufficient to guarantee a large gradient.
- We conduct various numerical experiments including finding the ground energy and the ground state of the Heisenberg model and the LiH molecule, which belong to quantum simulation and quantum chemistry, respectively. Experiment results show that Gaussian initializations outperform uniform initializations, which verify proposed theorems.

1.1 Related work

The barren plateau phenomenon was first noticed in [44], which proves that if the circuit distribution forms unitary 2-designs [46], the variance of the gradient of the circuit vanishes to zero with the rate exponential in the qubit number. Subsequently, several positive results are proved for shallow quantum circuits such as the alternating-layered circuit [45, 47] and the quantum convolutional neural network [48] when the observable is constrained in small number of qubits (local observable). For shallow circuits with N qubits and $\mathcal{O}(\log N)$ depth, the variance of the gradient has the order $\text{poly}(N)^{-1}$ if gate blocks in the circuit are sampled from local 2-design distributions. Later, several works prove an inherent relationship between the barren plateau phenomenon and the complexity of states generated from the circuit. Specifically, circuit states that satisfy the volume law could lead to the barren plateau problem [49]. Expressive quantum circuits, which is measured by the distance between the Haar distribution and the distribution of circuit states, could have vanishing gradients [50]. Since random circuits form approximately 2-designs when they achieve linear depths [46], deep quantum circuits were believed to suffer the barren plateau problem generally.

The parameterization of quantum circuits is achieved by tuning the time of Hamiltonian simulations, so the gradient of the circuit satisfies the parameter-shift rule [51]. Thus, the variance of the loss in VQAs and that of its gradient have similar behaviors for uniform distributions [44, 52]. One corollary of the parameter-shift rule is that the gradient of depolarized noisy quantum circuits vanishes exponentially with increasing circuit depth [53], since the loss itself vanishes in the same rate. Another corollary is that both gradient-free [54] and higher-order methods [55] could not solve the barren plateau problem. Although most existing theoretical and practical results imply the barren plateau phenomenon in deep circuits, VQAs with deep circuits do have impressive advantages from other aspects. For example, the loss of VQAs is highly non-convex, which is hard to find the global minima [56] for both shallow and deep circuits. Meanwhile, for VQAs with shallow circuits, local minima and global minima have considerable gaps [57], which could severely influence the training performance of gradient-based methods. Contrary to shallow cases, deep VQAs have vanishing gaps between local minima and global minima [58]. In practice, experiments show that overparameterized VQAs [59] can be optimized towards the global minima. Moreover, VQAs with deep circuits have more expressive power than that of shallow circuits [60–62], which implies the potential to handle more complex tasks in quantum machine learning and related fields.

Inspired by various advantages of deep VQAs, some approaches have been proposed recently for solving the related barren plateau problem in practice. For example, the block-identity strategy [63] initializes gate blocks in pairs and sets parameters inversely, such that the initial circuit is equivalent to the identity circuit with zero depth. Since shallow circuits have no vanishing gradient problem, the corresponding VQA is trainable with guarantees at the first step. However, we remark that the block-identity condition would not hold after the first step, and the structure of the circuit needs to be designed properly. The layerwise training method [64] trains parameters in the circuit layers by layers, such that the depth of trainable part is limited. However, this method implements circuits with larger depth than that of the origin circuit, and parameters in the first few layers are not optimized. A recent work provides theoretical guarantees on the trainability of deep circuits with certain structures [65]. However, the proposed theory only suits VQAs with local observables, but many practical applications such as finding the ground state of molecules and the quantum compiling [66, 67] apply global observables.

2 Notations and quantum computing basics

We denote by $[N]$ the set $\{1, \dots, N\}$. The form $\|\cdot\|_2$ represents the ℓ_2 norm for the vector and the spectral norm for the matrix, respectively. We denote by a_j the j -th component of the vector \mathbf{a} . The tensor product operation is denoted as “ \otimes ”. The conjugate transpose of a matrix A is denoted as A^\dagger . The trace of a matrix A is denoted as $\text{Tr}[A]$. We denote $\nabla_{\boldsymbol{\theta}} f$ as the gradient of the function f with respect to the variable $\boldsymbol{\theta}$. We employ notations \mathcal{O} to describe complexity notions.

Now we introduce quantum computing knowledge and notations. The pure state of a qubit could be written as $|\phi\rangle = a|0\rangle + b|1\rangle$, where $a, b \in \mathbb{C}$ satisfy $|a|^2 + |b|^2 = 1$, and $|0\rangle = (1, 0)^T$, $|1\rangle = (0, 1)^T$. The N -qubit space is formed by the tensor product of N single-qubit spaces. For pure states, the corresponding density matrix is defined as $\rho = |\phi\rangle\langle\phi|$, in which $\langle\phi| = (|\phi\rangle)^\dagger$. We use the density matrix to represent general mixed quantum states, i.e., $\rho = \sum_k c_k |\phi_k\rangle\langle\phi_k|$, where $c_k \in \mathbb{R}$ and $\sum_k c_k = 1$. A single-qubit operation to the state behaves like the matrix-vector multiplication and can be referred to as the gate \square in the quantum circuit language. Specifically, single-qubit operations are often used as $R_X(\theta) = e^{-i\theta X}$, $R_Y(\theta) = e^{-i\theta Y}$, and $R_Z(\theta) = e^{-i\theta Z}$, where

$$X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, Y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, Z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$

Pauli matrices will be referred to as $\{I, X, Y, Z\} = \{\sigma_0, \sigma_1, \sigma_2, \sigma_3\}$ for the convenience. Moreover, two-qubit operations, such as the CZ gate and the $\sqrt{i\text{SWAP}}$ gate, are employed for generating quantum entanglement:

$$\text{CZ} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix}, \sqrt{i\text{SWAP}} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1/\sqrt{2} & i/\sqrt{2} & 0 \\ 0 & i/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

We could obtain information from the quantum system by performing measurements, for example, measuring the state $|\phi\rangle = a|0\rangle + b|1\rangle$ generates 0 and 1 with probability $p(0) = |a|^2$ and $p(1) = |b|^2$, respectively. Such a measurement operation could be mathematically referred to as calculating the average of the observable $O = \sigma_3$ under the state $|\phi\rangle$:

$$\langle\phi|O|\phi\rangle \equiv \text{Tr}[\sigma_3|\phi\rangle\langle\phi|] = |a|^2 - |b|^2 = p(0) - p(1).$$

Mathematically, quantum observables are Hermitian matrices. Specifically, the average of a unitary observable under arbitrary states is bounded by $[-1, 1]$. We remark that $\mathcal{O}(\frac{1}{\epsilon^2})$ times of measurements could provide an $\epsilon\|O\|_2$ -error estimation to the value $\text{Tr}[O\rho]$.

3 Framework of general VQAs

In this section, we introduce the framework of general VQAs and corresponding notations. A typical variational quantum algorithm can be viewed as the optimization of the function f , which is defined as the expectation of observables. The expectation varies for different initial states and different parameters $\boldsymbol{\theta}$ used in quantum circuits. Throughout this paper, we define

$$f(\boldsymbol{\theta}) = \text{Tr}[OV(\boldsymbol{\theta})\rho_{\text{in}}V(\boldsymbol{\theta})^\dagger] \quad (1)$$

as the loss function of VQAs, where $V(\theta)$ denotes the parameterized quantum circuit, the hermitian matrix O denotes the observable, and ρ_{in} denotes the density matrix of the input state. Next, we explain observables, input states, and parameterized quantum circuits in detail.

Both the observable and the density matrix could be decomposed under the Pauli basis. We define the *locality* of a quantum observable as the maximum number of non-identity Pauli matrices in the tensor product, such that the corresponding coefficient is not zero. Thus, the observable with the constant locality is said to be *local*, and the observable that acts on all qubits is said to be *global*.

The observable and the input state in VQAs could have various formulations for specific tasks. For the quantum simulation or the quantum chemistry scenario, observables are constrained to be the system Hamiltonians, while input states are usually prepared as computational basis states. For example, $(|0\rangle\langle 0|)^{\otimes N}$ is used frequently in quantum simulations [17, 18]. Hartree–Fock (HF) states [9, 10], which are prepared by the tensor product of $\{|0\rangle, |1\rangle\}$, serve as good initial states in quantum chemistry tasks [9, 11–13]. For quantum machine learning (QML) tasks, initial states encode the information of the training data, which could have a complex form. Many encoding strategies have been introduced in the literature [24, 68, 69]. In contrary with the complex initial states, observables employed in QML are quite simple. For example, $\pi_0 = |0\rangle\langle 0|$ serves as the observable in most QML tasks related with the classification [24–26] or the dimensional reduction [70].

Apart from the input states and the observable choices, parameterized quantum circuits employed in different variational quantum algorithms have various structures, which are also known as *ansatzes* [71–73]. Specifically, the ansatz in the VQA denotes the initial guess on the circuit structure. For example, alternating-layered ansatzes [71, 74] are proposed for approximating the Hamiltonian evolution. Recently, hardware efficient ansatzes [7, 75] and tensor-network based ansatzes [76, 77], which could utilize parameters efficiently on noisy quantum computers, have been developed for various tasks, including quantum simulations and quantum machine learning. For quantum chemistry tasks, unitary coupled cluster ansatzes [78, 79] are preferred since they preserve the number of electrons corresponding to circuit states.

In practice, ansatz is deployed as the sequence of single-qubit rotations $\{e^{-i\theta\sigma_k}, k \in \{1, 2, 3\}\}$ and two-qubit gates. We remark that the gradient of the VQA satisfies the parameter-shift rule [51, 80, 81]; namely, for independently deployed parameters θ_j , the corresponding partial derivative is

$$\frac{\partial f}{\partial \theta_j} = f(\theta_+) - f(\theta_-), \quad (2)$$

where θ_+ and θ_- are different from θ only at the j -th parameter: $\theta_j \rightarrow \theta_j \pm \frac{\pi}{4}$. Thus, the gradient of f could be estimated efficiently, which allows optimizing VQAs with gradient-based methods [82–84].

4 Theoretical results about Gaussian initialized VQAs

In this section, we provide our theoretical guarantees on the trainability of deep quantum circuits through proper designs for the initial parameter distribution. In short, we prove that the gradient of the L -layer N -qubit circuit is upper bounded by $1/\text{poly}(L, N)$, if initial parameters are sampled from a Gaussian distribution with $\mathcal{O}(1/L)$ variance. Our bounds significantly improve existing results of the gradients of VQAs, which have the order $2^{-\mathcal{O}(L)}$ for shallow circuits and the order $2^{-\mathcal{O}(N)}$ for deep circuits. We prove different results for the local and global observable cases in Section 4.1 and Section 4.2, respectively.

4.1 Independent parameters with local observables

First, we introduce the Gaussian initialization of parameters for the local observable case. We use the quantum circuit illustrated in Figure 1 as the ansatz in this section. The circuit in Figure 1 performs L layers of single qubit rotations and CZ gates on the input state ρ_{in} , followed by a R_X layer and a R_Y layer. We denote the single-qubit gate on the n -th qubit of the ℓ -th layer as $e^{-i\theta_{\ell,n}G_{\ell,n}}$, $\forall \ell \in \{1, \dots, L+2\}$ and $n \in \{1, \dots, N\}$, where $\theta_{\ell,n}$ is the corresponding parameter and $G_{\ell,n}$ is a Hermitian unitary. To eliminate degenerate parameters, we require that single-qubit gates in the first L layers do not commute with the CZ gate. After gates operations, we measure the observable

$$\sigma_{\mathbf{i}} = \sigma_{(i_1, i_2, \dots, i_N)} = \sigma_{i_1} \otimes \sigma_{i_2} \otimes \dots \otimes \sigma_{i_N}, \quad (3)$$

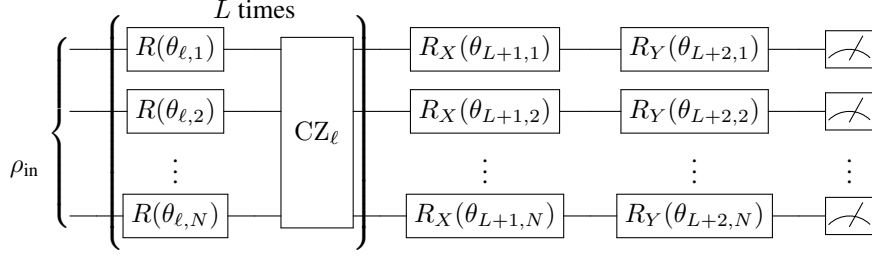


Figure 1: The quantum circuit framework for the local observable case. The circuit performs L layers of single qubit rotations and CZ layers on the input state ρ_{in} , followed by a R_X layer and a R_Y layer. In the ℓ -th single qubit layer, we employ the gate $e^{-i\theta_{\ell,n}G_{\ell,n}}$ for all qubits $n \in [N]$, where $G_{\ell,n}$ is a Hermitian unitary, which anti-commutes with σ_3 for $\ell \in [L]$. In each CZ_ℓ layer, CZ gates are employed between arbitrary qubit pairs. The measurement is performed on S qubits where the observable acts nontrivially on these qubits.

where $i_j \in \{0, 1, 2, 3\}, \forall j \in \{1, \dots, N\}$, and \mathbf{i} contains S non-zero elements. Figure 1 provides a general framework of VQAs with local observables, which covers various ansatzes proposed in the literature [65, 85, 61, 64]. The bound of the gradient norm of the Gaussian initialized variational quantum circuit is provided in Theorem 4.1 with the proof in Appendix.

Theorem 4.1. Consider the L -layer N -qubit variational quantum circuit $V(\boldsymbol{\theta})$ defined in Figure 1 and the cost function $f(\boldsymbol{\theta}) = \text{Tr}[\sigma_{\mathbf{i}}V(\boldsymbol{\theta})\rho_{\text{in}}V(\boldsymbol{\theta})^\dagger]$, where the observable $\sigma_{\mathbf{i}}$ follows the definition (3). Then,

$$\mathbb{E}_{\boldsymbol{\theta}} \|\nabla_{\boldsymbol{\theta}} f\|^2 \geq \frac{L}{S^S(L+2)^{S+1}} \text{Tr}[\sigma_{\mathbf{j}}\rho_{\text{in}}]^2, \quad (4)$$

where S is the number of non-zero elements in \mathbf{i} , and the index $\mathbf{j} = (j_1, j_2, \dots, j_N)$ such that $j_m = 0, \forall i_m = 0$ and $j_m = 3, \forall i_m \neq 0$. The expectation is taken with the Gaussian distribution $\mathcal{N}\left(0, \frac{1}{4S(L+2)}\right)$ for the parameters $\boldsymbol{\theta}$.

Compared to existing works [44, 45, 47, 48, 65], Theorem 4.1 provides a larger lower bound of the gradient norm, which improves the complexity exponentially with the depth of trainable circuits. Different from unitary 2-design distributions [44, 45, 47, 48] or the uniform distribution in the parameter space [52, 86, 65] that were employed in existing works, we analyze the expectation of the gradient norm under a depth-induced Gaussian distribution. This change follows a natural idea that the trainability is not required in the whole parameter space or the entire circuit space, but only on the parameter trajectory during the training. Moreover, large norm of gradients could only guarantee the trainability in the beginning stage, instead of the whole optimization, since a large gradient for trained parameters corresponds to non-convergence. Thus, the barren plateau problem could be crucial if initial parameters have vanishing gradients, which has been proved for deep VQAs with uniform initializations. In contrary, we could solve the barren plateau problem if parameters are initialized properly with large gradients, as provided in Theorem 4.1. Finally, Gaussian initialized circuits converge to benign values if optima appear around $\boldsymbol{\theta} = \mathbf{0}$, which holds in many cases. For example, over-parameterized quantum circuits have benign local minima [58] if the number of parameters exceeds the over-parameterization threshold. Moreover, over-parameterized circuits have exponential convergence rates [87, 88] on tasks like quantum machine learning and the quantum eigensolver. These works indicate that quantum circuits with sufficient depths could find good optimums near the initial points, which is similar to the classical wide neural network case [89].

4.2 Correlated parameters with global observables

Next, we extend the Gaussian initialization framework to general quantum circuits with correlated parameters and global observables. Quantum circuits with correlated parameters have wide applications

in quantum simulations and quantum chemistry [9, 11–13]. One example is the Givens rotation

$$R^{\text{Givens}}(\theta) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos \theta & -\sin \theta & 0 \\ 0 & \sin \theta & \cos \theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} = \begin{array}{c} \boxed{\sqrt{i}\text{SWAP}} \quad \boxed{R_Z(\frac{-\theta}{2})} \quad \boxed{\sqrt{i}\text{SWAP}} \\ \boxed{R_Z(\frac{\theta+\pi}{2})} \quad \boxed{R_Z(\frac{\pi}{2})} \end{array} \quad (5)$$

which preserves the number of electrons in parameterized quantum states [11].

To analyze VQAs with correlated parameterized gates, we consider the ansatz $V(\boldsymbol{\theta}) = \prod_{j=L}^1 V_j(\theta_j)$, which consists of parameterized gates $\{V_j(\theta_j)\}_{j=1}^L$. Denote by h_j the number of unitary gates that share the same parameter θ_j . Thus, the parameterized gate $V_j(\theta_j)$ consists of a list of fixed and parameterized unitary operations

$$V_j(\theta_j) = \prod_{k=1}^{h_j} W_{jk} e^{-i \frac{\theta_j}{a_j} G_{jk}} \quad (6)$$

with the term $a_j \in \mathbb{R}/\{0\}$, where the Hamiltonian G_{jk} and the fixed gate W_{jk} are unitary $\forall k \in [h_j]$. Moreover, we consider the objective function

$$f(\boldsymbol{\theta}) = \text{Tr} \left[O \prod_{j=L}^1 V_j(\theta_j) \rho_{\text{in}} \prod_{j=1}^L V_j(\theta_j)^\dagger \right], \quad (7)$$

where ρ_{in} and O denote the input state and the observable, respectively. In practical tasks of quantum chemistry, the molecule Hamiltonian H serves as the observable O . Minimizing the function (7) provides the ground energy and the corresponding ground state of the molecule. We provide the bound of the gradient norm of the Gaussian initialized variational quantum circuit in Theorem 4.2 with the proof in Appendix. Similar to the local observable case, we could bound the norm of the gradient of Eq. (7) if parameters are initialized with $\mathcal{O}(\frac{1}{L})$ variance. Theorem 4.2 provides nontrivial bounds when the gradient at the zero point is large. This condition holds when the mean-field theory provides a good initial guess to the corresponding problems, e.g. the ground energy task in quantum chemistry and quantum many-body problems [90].

Theorem 4.2. *Consider the N -qubit variational quantum algorithms with the objective function (7). Then the following formula holds for any $\ell \in \{1, \dots, L\}$,*

$$\mathbb{E}_{\boldsymbol{\theta}} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \geq (1 - \epsilon) \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\boldsymbol{\theta}=\mathbf{0}}, \quad (8)$$

where $\mathbf{0} \in \mathbb{R}^L$ is the zero vector. The expectation is taken with Gaussian distributions $\mathcal{N}(0, \gamma_j^2)$ for parameters in $\boldsymbol{\theta} = \{\theta_j\}_{j=1}^L$, where the variance $\gamma_j^2 \leq \frac{a_j^2 \epsilon}{16h_j^2(3h_j(h_j-1)+1)L\|O\|_2^2} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\boldsymbol{\theta}=\mathbf{0}}$.

We remark that Theorem 4.2 not only provides an initialization strategy, but also guarantees the update direction during the training. Different from the classical neural network, where the gradient could be calculated accurately, the gradient of VQAs, obtained by the parameter-shift rule (2), is perturbed by the measurement noise. A guide on the size of acceptable measurement noise could be useful for the complexity analysis of VQAs. Specifically, define $\boldsymbol{\theta}^{(t-1)}$ as the parameter at the $t-1$ -th iteration. Denote by $\boldsymbol{\theta}^{(t)}$ and $\tilde{\boldsymbol{\theta}}^{(t)}$ the parameter updated from $\boldsymbol{\theta}^{(t-1)}$ for noiseless and noisy cases, respectively. Then $\tilde{\boldsymbol{\theta}}^{(t)}$ differs from $\boldsymbol{\theta}^{(t)}$ by a Gaussian error term due to the measurement noise. We expect to derive the gradient norm bound for $\tilde{\boldsymbol{\theta}}^{(t)}$, as provided in Corollary 4.3. Thus, $\frac{1}{\gamma^2} = \mathcal{O}(\frac{L}{\epsilon})$ number of measurements is sufficient to guarantee a large gradient.

Corollary 4.3. *Consider the N -qubit variational quantum algorithms with the objective function (7). Then the following formula holds for any $\ell \in \{1, \dots, L\}$,*

$$\mathbb{E}_{\boldsymbol{\delta}} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}+\boldsymbol{\delta}} \geq (1 - \epsilon) \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}. \quad (9)$$

The expectation is taken with Gaussian distributions $\mathcal{N}(0, \gamma_j^2)$ for parameters $\boldsymbol{\delta} = \{\delta_j\}_{j=1}^L$, where the variance $\gamma_j^2 \leq \frac{a_j^2 \epsilon}{16h_j^2(3h_j(h_j-1)+1)L\|O\|_2^2} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}$, $\forall j \in [L]$.

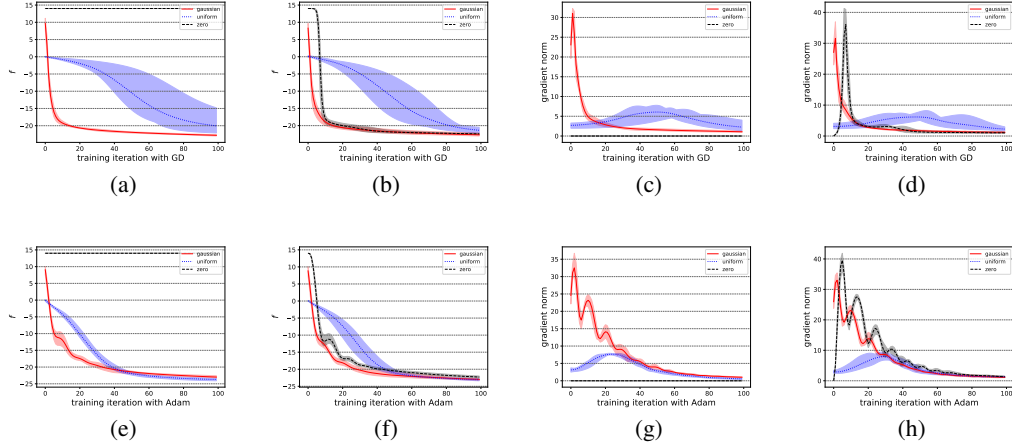


Figure 2: Numerical results of finding the ground energy of the Heisenberg model. The first row shows training results with the gradient descent optimizer, where Figures 2(a) and 2(b) illustrate the loss function corresponding to Eq.(10) during the optimization with accurate and noisy gradients, respectively. Figures 2(c) and 2(d) show the ℓ_2 norm of corresponding gradients. The second row shows training results with the Adam optimizer, where Figures 2(e) and 2(f) illustrate the loss function with accurate and noisy gradients, respectively. Figures 2(g) and 2(h) show the ℓ_2 norm of corresponding gradients. Each line denotes the average of 5 rounds of optimizations.

Corollary 4.3 is derived by analyzing the gradient of the function $g(\delta) = f(\delta + \theta^{(t)})$ via Theorem 4.2. For any number of measurements such that the corresponding Gaussian noise δ satisfies the condition in Corollary 4.3, the trainability at the updated point is guaranteed.

5 Experiments

In this section, we analyze the training behavior of two variational quantum algorithms, i.e., finding the ground energy and state of the Heisenberg model and the LiH molecule, respectively. All numerical experiments are provided using the PennyLane package [91].

5.1 Heisenberg model

In the first task, we aim to find the ground state and the ground energy of the Heisenberg model [92]. The corresponding Hamiltonian matrix is

$$H = \sum_{i=1}^{N-1} X_i X_{i+1} + Y_i Y_{i+1} + Z_i Z_{i+1}, \quad (10)$$

where N is the number of qubit, $X_i = I^{\otimes(i-1)} \otimes X \otimes I^{\otimes(N-i)}$, $Y_i = I^{\otimes(i-1)} \otimes Y \otimes I^{\otimes(N-i)}$, and $Z_i = I^{\otimes(i-1)} \otimes Z \otimes I^{\otimes(N-i)}$. We employ the loss function defined by Eq. (1) with the input state $(|0\rangle\langle 0|)^{\otimes N}$ and the observable (10). Thus, by minimizing the function (1), we can obtain the least eigenvalue of the observable (10), which is the ground energy. We adopt the ansatz with $N = 15$ qubits, which consists of $L_1 = 10$ layers of $R_Y R_X CZ$ blocks. In each block, we first employ the CZ gate to neighboring qubits pairs $\{(1, 2), \dots, (N, 1)\}$, followed by R_X and R_Y rotations for all qubits. Overall, the quantum circuit has 300 parameters. We consider three initialization methods for comparison, i.e., initializations with the Gaussian distribution $\mathcal{N}(0, \gamma^2)$ and the uniform distribution in $[0, 2\pi]$, respectively, and the zero initialization (all parameters equal to 0 at the initial point). We remark that each term in the observable (10) contains at most $S = 2$ non-identity Pauli matrices, which is consistent with the $(S, L) = (2, 18)$ case of Theorem 4.1. Thus, we expect that the Gaussian initialization with the variance $\gamma^2 = \frac{1}{4S(L+2)} = \frac{1}{160}$ could provide trainable initial parameters.

In the experiment, we train VQAs with gradient descent (GD) [93] and Adam optimizers [94], respectively. The learning rate is 0.01 and 0.01 for both GD and Adam cases. Since the estimation

of gradients on real quantum computers could be perturbed by statistical measurement noise, we compare optimizations using accurate and noisy gradients. For the latter case, we set the variance of measurement noises to be 0.01. The numerical results of the Heisenberg model are shown in the Figure 2. The loss during the training with gradient descents is shown in Figures 2(a) and 2(b) for the accurate and the noisy gradient cases, respectively. The Gaussian initialization outperforms the other two initializations with faster convergence rates. Figures 2(c) and 2(d) verify that Gaussian initialized VQAs have larger gradients in the early stage, compared to that of uniformly initialized VQAs. We notice that zero initialized VQAs cannot be trained with accurate gradient descent, since the initial gradient equals to zero. This problem is alleviated in the noisy case, as shown in Figures 2(b) and 2(d). Since the gradient is close to zero at the initial stage, the update direction mainly depends on the measurement noise, which forms the Gaussian distribution. Thus, the parameter in the noisy zero initialized VQAs is expected to accumulate enough variances, which takes around 10 iterations based on Figure 2(h). As illustrated in Figure 2(b), the loss function corresponding to the zero initialization decreases quickly after the variance accumulation stage. Results in Figures 2(e) and 2(h) show similar training behaviors using the Adam optimizer.

5.2 Quantum chemistry

In the second task, we aim to find the ground state and the ground energy of the LiH molecule. We follow settings on the ansatz in Refs. [12, 13]. For the molecule with n_e active electrons and n_o free spin orbitals, the corresponding VQA contains $N = n_o$ qubits, which employs the HF state [9, 10]

$$|\phi_{\text{HF}}\rangle = \underbrace{|1\rangle \otimes \cdots \otimes |1\rangle}_{n_e} \otimes \underbrace{|0\rangle \otimes \cdots \otimes |0\rangle}_{n_o - n_e}$$

as the input state. We construct the parameterized quantum circuit with Givens rotation gates [12], where each gate is implemented on 2 or 4 qubits with one parameter. Specifically, for the LiH molecule, the number of electrons $n_e = 2$, the number of free spin orbitals $n_o = 10$, and the number of different Givens rotations is $L = 24$ [13]. We follow the molecule Hamiltonian H_{LiH} defined in Ref. [13]. Thus, the loss function for finding the ground energy of LiH is defined as

$$f(\boldsymbol{\theta}) = \text{Tr} [H_{\text{LiH}} V_{\text{Givens}}(\boldsymbol{\theta}) |\phi_{\text{HF}}\rangle \langle \phi_{\text{HF}}| V_{\text{Givens}}(\boldsymbol{\theta})^\dagger], \quad (11)$$

where $V_{\text{Givens}}(\boldsymbol{\theta}) = \prod_{i=1}^{24} R_i^{\text{Givens}}(\theta_i)$ denotes the product of all parameterized Givens rotations of the LiH molecule. By minimizing the function (11), we can obtain the least eigenvalue of the Hamiltonian H_{LiH} , which is the ground energy of the LiH molecule.

In practice, we initialize parameters in the VQA (11) with three distributions for comparison, i.e., the Gaussian distribution $\mathcal{N}(0, \gamma^2)$, the zero distribution (all parameters equal to 0), and the uniform distribution in $[0, 2\pi]$. For 2-qubit Givens rotations, the term $(h, a) = (2, 2)$ as shown in Eq. (5). For 4-qubit Givens rotations, the term $(h, a) = (8, 8)$ [95]. Thus, we set the variance in the Gaussian distribution $\gamma^2 = \frac{8^2 \times \frac{1}{2}}{48 \times 8^4 \times 24}$, which matches the $(L, h, a, \epsilon) = (24, 8, 8, \frac{1}{2})$ case of Theorem 4.2. Similar to the task of the Heisenberg model, we consider both the accurate and the noisy gradient cases, where the variance of noises in the latter case is the constant 0.001. Moreover, we consider the noisy case with adaptive noises, where the variance of the noise on each partial derivative $\frac{\partial f}{\partial \theta_\ell} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}$ in the t -th iteration is

$$\gamma^2 = \frac{1}{96 \times 24 \times 8^2 \|H_{\text{LiH}}\|_2^2} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t-1)}}. \quad (12)$$

The variance in Eq. (12) matches the $(L, h, a, \epsilon) = (24, 8, 8, \frac{1}{2})$ case of Corollary 4.3 when the VQA is nearly converged:

$$\frac{\partial f}{\partial \theta_\ell} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \approx \frac{\partial f}{\partial \theta_\ell} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t-1)}}.$$

In the experiment, we train VQAs with gradient descent and Adam optimizers. Learning rates are set to be 0.1 and 0.01 for GD and Adam cases, respectively. The loss (11) during training iterations is shown in Figure 3. Optimization results with gradient descents are shown in Figures 3(a)-3(c) for the accurate gradient case, the adaptive noisy gradient case, and the noisy gradient case with the constant noise variance 0.001, respectively. The variance of the noise in the adaptive noisy gradient case follows Eq. (12). Figures 3(a) and 3(b) show similar performance, where the loss f with the

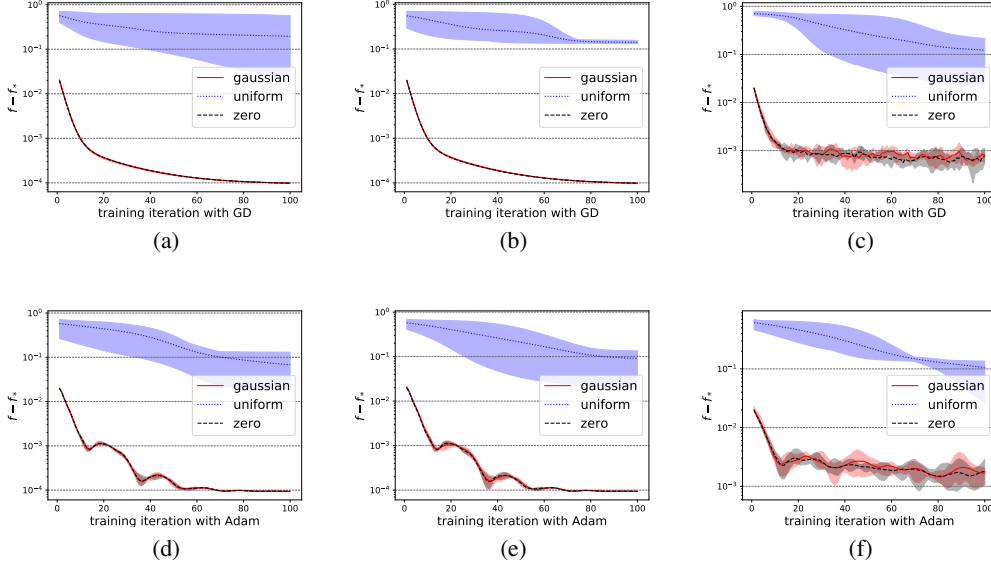


Figure 3: Numerical results of finding the ground energy of the molecule LiH. The first and second rows show training results with the gradient descent and the Adam optimizer, respectively. The left, the middle, and the right columns show results using accurate gradients, noisy gradients with adaptive-distributed noises, and noisy gradients with constant-distributed noises. The variance of noises in the middle line (Figures 3(b) and 3(e)) follows Eq. (12), while the variance of noises in the right line (Figures 3(c) and 3(f)) is 0.001. Each line denotes the average of 5 rounds of optimizations.

Gaussian initialization and the zero initialization converge to 10^{-4} over the global minimum f_* . The loss with the uniform initialization is higher than 10^{-1} over the global minimum. Figure 3(c) shows the training with constantly perturbed gradients. The Gaussian initialization and the zero initialization induce the 10^{-3} convergence, while the loss function with the uniform initialization is still higher than 10^{-1} over the global minimum. Figures 3(d)-3(f) show similar training behaviors using the Adam optimizer. Based on Figures 3(a)-3(f), the Gaussian initialization and the zero initialization outperform the uniform initialization in all cases. We notice that optimization with accurate gradients and optimization with adaptive noisy gradients have the same convergence rate and the final value of the loss function, which is better than that using constantly perturbed gradients. We remark that the number of measurements $T = \mathcal{O}(\frac{1}{\text{Var}(\text{noise})})$. Thus, finite number of measurements with the noise (12) for gradient estimation is enough to achieve the performance of accurate gradients, which verifies Theorem 4.2 and Corollary 4.3.

6 Conclusions

In this work, we provide a Gaussian initialization strategy for solving the vanishing gradient problem of deep variational quantum algorithms. We prove that the gradient norm of N -qubit quantum circuits with L layers could be lower bounded by $\text{poly}(N, L)^{-1}$, if the parameter is sampled independently from the Gaussian distribution with the variance $\mathcal{O}(\frac{1}{L})$. Our results hold for both the local and the global observable cases, and could be generalized to VQAs employing correlated parameterized gates. Compared to the local case, the bound for the global case depends on the gradient performance at the zero point. Further analysis towards the zero-case-free bound could be investigated as future directions. Moreover, we show that the necessary number of measurements, which scales $\mathcal{O}(\frac{L}{\epsilon})$, suffices for estimating the gradient during the training. We provide numerical experiments on finding the ground energy and state of the Heisenberg model and the LiH molecule, respectively. Experiments show that the proposed Gaussian initialization method outperforms the uniform initialization method with a faster convergence rate, and the training using gradients with adaptive noises shows the same convergence compared to the training using noiseless gradients.

References

- [1] Frank Arute, Kunal Arya, Ryan Babbush, Dave Bacon, Joseph C Bardin, Rami Barends, Rupak Biswas, Sergio Boixo, Fernando GSL Brandao, David A Buell, et al. Quantum supremacy using a programmable superconducting processor. *Nature*, 574(7779):505–510, 2019.
- [2] Han-Sen Zhong, Hui Wang, Yu-Hao Deng, Ming-Cheng Chen, Li-Chao Peng, Yi-Han Luo, Jian Qin, Dian Wu, Xing Ding, Yi Hu, Peng Hu, Xiao-Yan Yang, Wei-Jun Zhang, Hao Li, Yuxuan Li, Xiao Jiang, Lin Gan, Guangwen Yang, Lixing You, Zhen Wang, Li Li, Nai-Le Liu, Chao-Yang Lu, and Jian-Wei Pan. Quantum computational advantage using photons. *Science*, 370(6523):1460–1463, 2020.
- [3] John Preskill. Quantum computing in the nisq era and beyond. *Quantum*, 2:79, August 2018.
- [4] Marco Cerezo, Andrew Arrasmith, Ryan Babbush, Simon C Benjamin, Suguru Endo, Keisuke Fujii, Jarrod R McClean, Kosuke Mitarai, Xiao Yuan, Lukasz Cincio, and Patrick J. Coles. Variational quantum algorithms. *Nature Reviews Physics*, pages 1–20, 2021.
- [5] Sam McArdle, Suguru Endo, Alán Aspuru-Guzik, Simon C. Benjamin, and Xiao Yuan. Quantum computational chemistry. *Rev. Mod. Phys.*, 92:015003, Mar 2020.
- [6] Alberto Peruzzo, Jarrod McClean, Peter Shadbolt, Man-Hong Yung, Xiao-Qi Zhou, Peter J Love, Alán Aspuru-Guzik, and Jeremy L O’Brien. A variational eigenvalue solver on a photonic quantum processor. *Nat. Commun.*, 5(1):1–7, 2014.
- [7] Abhinav Kandala, Antonio Mezzacapo, Kristan Temme, Maika Takita, Markus Brink, Jerry M Chow, and Jay M Gambetta. Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets. *Nature*, 549(7671):242–246, 2017.
- [8] Oscar Higgott, Daochen Wang, and Stephen Brierley. Variational quantum computation of excited states. *Quantum*, 3:156, July 2019.
- [9] Harper R Grimsley, Sophia E Economou, Edwin Barnes, and Nicholas J Mayhall. An adaptive variational algorithm for exact molecular simulations on a quantum computer. *Nat. Commun.*, 10(1):1–9, 2019.
- [10] Cornelius Hempel, Christine Maier, Jonathan Romero, Jarrod McClean, Thomas Monz, Heng Shen, Petar Jurcevic, Ben P. Lanyon, Peter Love, Ryan Babbush, Alán Aspuru-Guzik, Rainer Blatt, and Christian F. Roos. Quantum chemistry calculations on a trapped-ion quantum simulator. *Phys. Rev. X*, 8:031022, Jul 2018.
- [11] Frank Arute, Kunal Arya, Ryan Babbush, Dave Bacon, Joseph C. Bardin, Rami Barends, Sergio Boixo, Michael Broughton, Bob B. Buckley, David A. Buell, et al. Hartree-fock on a superconducting qubit quantum computer. *Science*, 369(6507):1084–1089, 2020.
- [12] Ho Lun Tang, V.O. Shkolnikov, George S. Barron, Harper R. Grimsley, Nicholas J. Mayhall, Edwin Barnes, and Sophia E. Economou. Qubit-adapt-vqe: An adaptive algorithm for constructing hardware-efficient ansätze on a quantum processor. *PRX Quantum*, 2:020310, Apr 2021.
- [13] Alain Delgado, Juan Miguel Arrazola, Soran Jahangiri, Zeyue Niu, Josh Izaac, Chase Roberts, and Nathan Killoran. Variational quantum algorithm for molecular geometry optimization. *Phys. Rev. A*, 104:052402, Nov 2021.
- [14] I. M. Georgescu, S. Ashhab, and Franco Nori. Quantum simulation. *Rev. Mod. Phys.*, 86:153–185, Mar 2014.
- [15] Xiao Yuan, Suguru Endo, Qi Zhao, Ying Li, and Simon C. Benjamin. Theory of variational quantum simulation. *Quantum*, 3:191, October 2019.
- [16] Sam McArdle, Tyson Jones, Suguru Endo, Ying Li, Simon C Benjamin, and Xiao Yuan. Variational ansatz-based quantum simulation of imaginary time evolution. *Npj Quantum Inf.*, 5(1):1–6, 2019.
- [17] Suguru Endo, Jinzhao Sun, Ying Li, Simon C. Benjamin, and Xiao Yuan. Variational quantum simulation of general processes. *Phys. Rev. Lett.*, 125:010501, Jun 2020.

- [18] C Neill, T McCourt, X Mi, Z Jiang, MY Niu, W Mroczkiewicz, I Aleiner, F Arute, K Arya, J Atalaya, et al. Accurately computing the electronic properties of a quantum ring. *Nature*, 594(7864):508–512, 2021.
- [19] Xiao Mi, Matteo Ippoliti, Chris Quintana, Ami Greene, Zijun Chen, Jonathan Gross, Frank Arute, Kunal Arya, Juan Atalaya, Ryan Babbush, et al. Time-crystalline eigenstate order on a quantum processor. *Nature*, pages 1–1, 2021.
- [20] J. Randall, C. E. Bradley, F. V. van der Gronden, A. Galicia, M. H. Abobeih, M. Markham, D. J. Twitchen, F. Machado, N. Y. Yao, and T. H. Taminiau. Many-body localized discrete time crystal with a programmable spin-based quantum simulator. *Science*, 374(6574):1474–1478, 2021.
- [21] Amita B. Deb and Niels Kjærgaard. Observation of pauli blocking in light scattering from quantum degenerate fermions. *Science*, 374(6570):972–975, 2021.
- [22] G. Semeghini, H. Levine, A. Keesling, S. Ebadi, T. T. Wang, D. Bluvstein, R. Verresen, H. Pichler, M. Kalinowski, R. Samajdar, A. Omran, S. Sachdev, A. Vishwanath, M. Greiner, V. Vuletić, and M. D. Lukin. Probing topological spin liquids on a programmable quantum simulator. *Science*, 374(6572):1242–1247, 2021.
- [23] K. J. Satzinger, Y.-J. Liu, A. Smith, C. Knapp, M. Newman, C. Jones, Z. Chen, C. Quintana, X. Mi, A. Dunsworth, et al. Realizing topologically ordered states on a quantum processor. *Science*, 374(6572):1237–1241, 2021.
- [24] Maria Schuld, Alex Bocharov, Krysta M. Svore, and Nathan Wiebe. Circuit-centric quantum classifiers. *Phys. Rev. A*, 101:032308, Mar 2020.
- [25] Vojtěch Havlíček, Antonio D Córcoles, Kristan Temme, Aram W Harrow, Abhinav Kandala, Jerry M Chow, and Jay M Gambetta. Supervised learning with quantum-enhanced feature spaces. *Nature*, 567(7747):209–212, 2019.
- [26] Maria Schuld and Nathan Killoran. Quantum machine learning in feature hilbert spaces. *Phys. Rev. Lett.*, 122:040504, Feb 2019.
- [27] Yuxuan Du, Min-Hsiu Hsieh, Tongliang Liu, and Dacheng Tao. Expressive power of parametrized quantum circuits. *Phys. Rev. Research*, 2:033125, Jul 2020.
- [28] Samuel Yen-Chi Chen, Chao-Han Huck Yang, Jun Qi, Pin-Yu Chen, Xiaoli Ma, and Hsi-Sheng Goan. Variational quantum circuits for deep reinforcement learning. *IEEE Access*, 8:141007–141024, 2020.
- [29] Valeria Saggio, Beate E Asenbeck, Arne Hamann, Teodor Strömberg, Peter Schiansky, Vedran Dunjko, Nicolai Friis, Nicholas C Harris, Michael Hochberg, Dirk Englund, et al. Experimental quantum speed-up in reinforcement learning agents. *Nature*, 591(7849):229–233, 2021.
- [30] He-Liang Huang, Yuxuan Du, Ming Gong, Youwei Zhao, Yulin Wu, Chaoyue Wang, Shaowei Li, Futian Liang, Jin Lin, Yu Xu, Rui Yang, Tongliang Liu, Min-Hsiu Hsieh, Hui Deng, Hao Rong, Cheng-Zhi Peng, Chao-Yang Lu, Yu-Ao Chen, Dacheng Tao, Xiaobo Zhu, and Jian-Wei Pan. Experimental quantum generative adversarial networks for image generation. *Phys. Rev. Applied*, 16:024051, Aug 2021.
- [31] Yuxuan Du, Min-Hsiu Hsieh, Tongliang Liu, Shan You, and Dacheng Tao. Learnability of quantum neural networks. *PRX Quantum*, 2:040337, Nov 2021.
- [32] Michael Lubasch, Jaewoo Joo, Pierre Moinier, Martin Kiffner, and Dieter Jaksch. Variational quantum algorithms for nonlinear problems. *Phys. Rev. A*, 101:010301, Jan 2020.
- [33] Kenji Kubo, Yuya O. Nakagawa, Suguru Endo, and Shota Nagayama. Variational quantum simulations of stochastic differential equations. *Phys. Rev. A*, 103:052425, May 2021.
- [34] Yong-Xin Yao, Niladri Gomes, Feng Zhang, Cai-Zhuang Wang, Kai-Ming Ho, Thomas Iadecola, and Peter P. Orth. Adaptive variational quantum dynamics simulations. *PRX Quantum*, 2:030307, Jul 2021.

- [35] Hai-Ling Liu, Yu-Sen Wu, Lin-Chun Wan, Shi-Jie Pan, Su-Juan Qin, Fei Gao, and Qiao-Yan Wen. Variational quantum algorithm for the poisson equation. *Phys. Rev. A*, 104:022418, Aug 2021.
- [36] Oleksandr Kyriienko, Annie E. Paine, and Vincent E. Elfving. Solving nonlinear differential equations with differentiable quantum circuits. *Phys. Rev. A*, 103:052416, May 2021.
- [37] Carlos Bravo-Prieto, Ryan LaRose, Marco Cerezo, Yigit Subasi, Lukasz Cincio, and Patrick Coles. Variational quantum linear solver: a hybrid algorithm for linear systems. *Bulletin of the American Physical Society*, 65, 2020.
- [38] Xiaosi Xu, Jinzhao Sun, Suguru Endo, Ying Li, Simon C. Benjamin, and Xiao Yuan. Variational algorithms for linear algebra. *Science Bulletin*, 66(21):2181–2188, 2021.
- [39] Xin Wang, Zhixin Song, and Youle Wang. Variational quantum singular value decomposition. *Quantum*, 5:483, June 2021.
- [40] Marcello Benedetti, Erika Lloyd, Stefan Sack, and Mattia Fiorentini. Parameterized quantum circuits as machine learning models. *Quantum Science and Technology*, 4(4):043001, nov 2019.
- [41] Hugo Larochelle, Yoshua Bengio, Jérôme Louradour, and Pascal Lamblin. Exploring strategies for training deep neural networks. *Journal of machine learning research*, 10(1), 2009.
- [42] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1675–1685. PMLR, 09–15 Jun 2019.
- [43] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT’2010*, pages 177–186, Heidelberg, 2010. Physica-Verlag HD.
- [44] Jarrod R McClean, Sergio Boixo, Vadim N Smelyanskiy, Ryan Babbush, and Hartmut Neven. Barren plateaus in quantum neural network training landscapes. *Nat. Commun.*, 9(1):1–6, 2018.
- [45] Marco Cerezo, Akira Sone, Tyler Volkoff, Lukasz Cincio, and Patrick J Coles. Cost function dependent barren plateaus in shallow parametrized quantum circuits. *Nat. Commun.*, 12(1):1–12, 2021.
- [46] Aram W Harrow and Richard A Low. Random quantum circuits are approximate 2-designs. *Commun. Math. Phys.*, 291(1):257–302, 2009.
- [47] A V Uvarov and J D Biamonte. On barren plateaus and cost function locality in variational quantum algorithms. 54(24):245301, may 2021.
- [48] Arthur Pesah, M. Cerezo, Samson Wang, Tyler Volkoff, Andrew T. Sornborger, and Patrick J. Coles. Absence of barren plateaus in quantum convolutional neural networks. *Phys. Rev. X*, 11:041011, 2021.
- [49] Carlos Ortiz Marrero, Mária Kieferová, and Nathan Wiebe. Entanglement-induced barren plateaus. *PRX Quantum*, 2:040316, Oct 2021.
- [50] Zoë Holmes, Kunal Sharma, M. Cerezo, and Patrick J. Coles. Connecting ansatz expressibility to gradient magnitudes and barren plateaus. *PRX Quantum*, 3:010313, Jan 2022.
- [51] Jun Li, Xiaodong Yang, Xinhua Peng, and Chang-Pu Sun. Hybrid quantum-classical approach to quantum optimal control. *Phys. Rev. Lett.*, 118:150503, 2017.
- [52] Kaining Zhang, Min-Hsiu Hsieh, Liu Liu, and Dacheng Tao. Toward trainability of quantum neural networks. *arXiv:2011.06258*, 2020.
- [53] Samson Wang, Enrico Fontana, Marco Cerezo, Kunal Sharma, Akira Sone, Lukasz Cincio, and Patrick Coles. Noise-induced barren plateaus in variational quantum algorithms. *Nat. Commun.*, 2021.

- [54] Andrew Arrasmith, M. Cerezo, Piotr Czarnik, Lukasz Cincio, and Patrick J. Coles. Effect of barren plateaus on gradient-free optimization. *Quantum*, 5:558, 2021.
- [55] M Cerezo and Patrick J Coles. Higher order derivatives of quantum neural networks with barren plateaus. 6(3):035006, 2021.
- [56] Lennart Bittel and Martin Kliesch. Training variational quantum algorithms is np-hard. *Phys. Rev. Lett.*, 127:120502, Sep 2021.
- [57] Xuchen You and Xiaodi Wu. Exponentially many local minima in quantum neural networks. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 12144–12155. PMLR, 18–24 Jul 2021.
- [58] Eric Ricardo Anschuetz. Critical points in quantum generative models. In *International Conference on Learning Representations*, 2022.
- [59] Martín Larocca, Nathan Ju, Diego García-Martín, Patrick J. Coles, and Marco Cerezo. Theory of overparametrization in quantum neural networks. *Nature Computational Science*, 3(6):542–551, Jun 2023.
- [60] Yuxuan Du, Zhuozhuo Tu, Xiao Yuan, and Dacheng Tao. Efficient measure for the expressivity of variational quantum algorithms. *Phys. Rev. Lett.*, 128:080506, Feb 2022.
- [61] Tobias Haug, Kishor Bharti, and M.S. Kim. Capacity and quantum geometry of parametrized quantum circuits. *PRX Quantum*, 2:040309, Oct 2021.
- [62] Matthias C Caro, Hsin-Yuan Huang, M Cerezo, Kunal Sharma, Andrew Sornborger, Lukasz Cincio, and Patrick J Coles. Generalization in quantum machine learning from few training data. *Nature Communications*, 2022.
- [63] Edward Grant, Leonard Wossnig, Mateusz Ostaszewski, and Marcello Benedetti. An initialization strategy for addressing barren plateaus in parametrized quantum circuits. *Quantum*, 3:214, 2019.
- [64] Andrea Skolik, Jarrod R McClean, Masoud Mohseni, Patrick van der Smagt, and Martin Leib. Layerwise learning for quantum neural networks. *Quantum Mach. Intell.*, 3(1):1–11, 2021.
- [65] Kaining Zhang, Min-Hsiu Hsieh, Liu Liu, and Dacheng Tao. Toward trainability of deep quantum neural networks. *arXiv:2112.15002*, 2021.
- [66] Sumeet Khatri, Ryan LaRose, Alexander Poremba, Lukasz Cincio, Andrew T. Sornborger, and Patrick J. Coles. Quantum-assisted quantum compiling. *Quantum*, 3:140, may 2019.
- [67] Kunal Sharma, Sumeet Khatri, M Cerezo, and Patrick J Coles. Noise resilience of variational quantum compiling. *New Journal of Physics*, 22(4):043006, apr 2020.
- [68] Israel F Araujo, Daniel K Park, Francesco Petruccione, and Adenilton J da Silva. A divide-and-conquer algorithm for quantum state preparation. *Sci. Rep.*, 11(1):1–12, 2021.
- [69] Louis Schatzki, Andrew Arrasmith, Patrick J Coles, and M Cerezo. Entangled datasets for quantum machine learning. *arXiv:2109.03400*, 2021.
- [70] Carlos Bravo-Prieto. Quantum autoencoders with enhanced data encoding. *Machine Learning: Science and Technology*, 2(3):035028, jul 2021.
- [71] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. A quantum approximate optimization algorithm. *arXiv:1411.4028*, 2014.
- [72] Sam McArdle, Tyson Jones, Suguru Endo, Ying Li, Simon C Benjamin, and Xiao Yuan. Variational ansatz-based quantum simulation of imaginary time evolution. *Npj Quantum Inf.*, 5(1):1–6, 2019.
- [73] Stuart Hadfield, Zhihui Wang, Bryan O’Gorman, Eleanor G. Rieffel, Davide Venturelli, and Rupak Biswas. From the quantum approximate optimization algorithm to a quantum alternating operator ansatz. *Algorithms*, 12(2), 2019.

- [74] Mark Fingerhuth, Tomáš Babej, et al. A quantum alternating operator ansatz with hard and soft constraints for lattice protein folding. *arXiv:1810.13411*, 2018.
- [75] Kouhei Nakaji and Naoki Yamamoto. Expressibility of the alternating layered ansatz for quantum computation. *Quantum*, 5:434, April 2021.
- [76] Iris Cong, Soonwon Choi, and Mikhail D Lukin. Quantum convolutional neural networks. *Nat. Phys.*, 15(12):1273–1278, 2019.
- [77] Timo Felser, Simone Notarnicola, and Simone Montangero. Efficient tensor network ansatz for high-dimensional quantum many-body problems. *Phys. Rev. Lett.*, 126:170603, Apr 2021.
- [78] M-H Yung, Jorge Casanova, Antonio Mezzacapo, Jarrod McClean, Lucas Lamata, Alan Aspuru-Guzik, and Enrique Solano. From transistor to trapped-ion computers for quantum chemistry. *Sci. Rep.*, 4(1):1–7, 2014.
- [79] Yangchao Shen, Xiang Zhang, Shuaining Zhang, Jing-Ning Zhang, Man-Hong Yung, and Kihwan Kim. Quantum implementation of the unitary coupled cluster for simulating molecular electronic structure. *Phys. Rev. A*, 95:020501, Feb 2017.
- [80] Gavin E Crooks. Gradients of parameterized quantum gates using the parameter-shift rule and gate decomposition. *arXiv:1905.13311*, 2019.
- [81] David Wierichs, Josh Izaac, Cody Wang, and Cedric Yen-Yu Lin. General parameter-shift rules for quantum gradients. *Quantum*, 6:677, March 2022.
- [82] D. Zhu, N. M. Linke, M. Benedetti, K. A. Landsman, N. H. Nguyen, C. H. Alderete, A. Perdomo-Ortiz, N. Korda, A. Garfoot, C. Brecque, L. Egan, O. Perdomo, and C. Monroe. Training of quantum circuits on a hybrid quantum computer. *Sci. Adv.*, 5(10), 2019.
- [83] James Stokes, Josh Izaac, Nathan Killoran, and Giuseppe Carleo. Quantum natural gradient. *Quantum*, 4:269, May 2020.
- [84] Ryan Sweke, Frederik Wilde, Johannes Meyer, Maria Schuld, Paul K. Faehrmann, Barthélémy Meynard-Piganeau, and Jens Eisert. Stochastic gradient descent for hybrid quantum-classical optimization. *Quantum*, 4:314, 2020.
- [85] Chih-Chieh Chen, Masaya Watabe, Kodai Shiba, Masaru Sogabe, Katsuyoshi Sakamoto, and Tomah Sogabe. On the expressibility and overfitting of quantum circuit learning. *ACM Transactions on Quantum Computing*, 2(2), jul 2021.
- [86] Martin Larocca, Piotr Czarnik, Kunal Sharma, Gopikrishnan Muraleedharan, Patrick J. Coles, and M. Cerezo. Diagnosing Barren Plateaus with Tools from Quantum Optimal Control. *Quantum*, 6:824, September 2022.
- [87] Junyu Liu, Khadijeh Najafi, Kunal Sharma, Francesco Tacchino, Liang Jiang, and Antonio Mezzacapo. Analytic theory for the dynamics of wide quantum neural networks. *Phys. Rev. Lett.*, 130:150601, Apr 2023.
- [88] Xuchen You, Shouvanik Chakrabarti, and Xiaodi Wu. A convergence theory for over-parameterized variational quantum eigensolvers. *arXiv:2205.12481*, 2022.
- [89] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [90] Luigi Amico and Vittorio Penna. Dynamical mean field theory of the bose-hubbard model. *Phys. Rev. Lett.*, 80:2189–2192, Mar 1998.
- [91] Ville Bergholm, Josh Izaac, Maria Schuld, Christian Gogolin, M Sohaib Alam, Shahnawaz Ahmed, Juan Miguel Arrazola, Carsten Blank, Alain Delgado, Soran Jahangiri, et al. PennyLane: Automatic differentiation of hybrid quantum-classical computations. *arXiv:1811.04968*, 2018.

- [92] F Bonechi, E Celeghini, R Giachetti, E Sorace, and M Tarlini. Heisenberg xxz model and quantum galilei group. *Journal of Physics A: Mathematical and General*, 25(15):L939–L943, aug 1992.
- [93] Léon Bottou. Stochastic gradient descent tricks. In *Neural networks: Tricks of the trade*, pages 421–436. Springer, 2012.
- [94] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.
- [95] Juan Miguel Arrazola, Olivia Di Matteo, Nicolás Quesada, Soran Jahangiri, Alain Delgado, and Nathan Killoran. Universal quantum circuits for quantum chemistry. *Quantum*, 6:742, June 2022.
- [96] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- [97] Y Nesterov. A method of solving a convex programming problem with convergence rate $\mathcal{O}(1/k^2)$. In *Sov. Math. Dokl*, volume 27.
- [98] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(61):2121–2159, 2011.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#)
 - (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#)
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#)
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#)
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) We run each experiment independently for 5 times.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[No\]](#) Each experiment finishes quickly on CPUs with a few hours.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[N/A\]](#)
 - (b) Did you mention the license of the assets? [\[N/A\]](#)
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[N/A\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)

5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

This supplementary material contains four parts:

- Section A provides some additional experiment results.
- Section B provides some technical lemmas which are useful for proving main theorems in this work.
- Section C provides the proof of Theorem 4.1.
- Section D provides the proof of Theorem 4.2.

A Additional Experiments

A.1 Heisenberg model

In this section, we introduce additional experiment results toward finding the ground state energy of the Heisenberg model with different circuit depths and optimizers. We follow simulation details in the main text.

First, we consider the effect of different circuit depths and Gaussian initializations with different variances. The loss function has the formulation Eq. (10) with the number of qubits $N = 15$. We adopt the ansatz circuit 1 with $L_1 \in \{8, 10, 12\}$ layers of $R_Y R_X CZ$ blocks, which correspond with $L \in \{14, 18, 22\}$ case of Theorem 4.1, respectively. In the experiment, we train VQAs using gradient descent with the learning rate 0.01. Since the estimation of gradients on real quantum computers could be perturbed by statistical measurement noise, we compare optimizations using accurate and noisy gradients. For the latter case, we set the variance of measurement noises to be 0.01. We train different Gaussian initialized VQAs with variances $\{0.01\gamma, 0.1\gamma, \gamma, 10\gamma, 100\gamma\}$, where the value γ follows the formulation in Theorem 4.1.

We illustrate results in Figures 4 and 5, which correspond to the noiseless and the noisy case, respectively. As show in figures of the loss during optimizations, the Gaussian initialization with the variance γ outperforms other Gaussian initializations with faster convergence rates. Gaussian initializations with small variances $\{0.01\gamma, 0.1\gamma\}$ have similar performances with the zero initialization for the noisy training case, and Gaussian initializations with large variances $\{10\gamma, 100\gamma\}$ behave similarly with the uniform initialization presented in the main text. Moreover, circuits initialized with larger variances $\{10\gamma, 100\gamma\}$ need more iterations to converge when the depth increases, while circuits with variances $\{0.01\gamma, 0.1\gamma, \gamma\}$ show similar convergence rates for different depths.

Next, we compare different initializations with other optimizers, i.e., the gradient descent with momentum [96], the Nesterov accelerated gradient (NAG) [97], and the adaptive gradient (AdaGrad) [98]. We follow the loss function (10) with $(N, L) = (15, 18)$ and $(N, L) = (18, 38)$. The learning rate and the noise are the same as that in the experiment considering different Gaussian variances. We illustrate results in Figures 6 and 7. As shown in Figure 6 and Figure 1 in the main text, the performance of GD with momentum and the NAG is similar to that of the Adam optimizer, while the performance of the AdaGrad is similar to the GD optimizer. By comparing Figures 6 and 7, we notice that uniformly initialized circuits converge slower when the qubit number and the circuit depth increase.

A.2 Quantum chemistry

In this section, we introduce additional experiment results toward finding the ground state energy of the Heisenberg model with different circuit depths. We repeat the LiH task in the main text with the depth $L \in \{24, 48, 72\}$ by stacking the circuit V_{Givens} in Eq. (11). The noise setting follows the adaptive noise with the variance in Eq. (12). We adopt gradient descent and the Adam optimizer with learning rates 0.1 and 0.01, respectively. The result is shown in Figure 8. For the gradient descent case, the convergence rate of the loss function increases when the circuit depth grows. For the Adam case, circuits with different depths show similar convergence speeds.

B Technical Lemmas

In this section, we provide some technical lemmas.

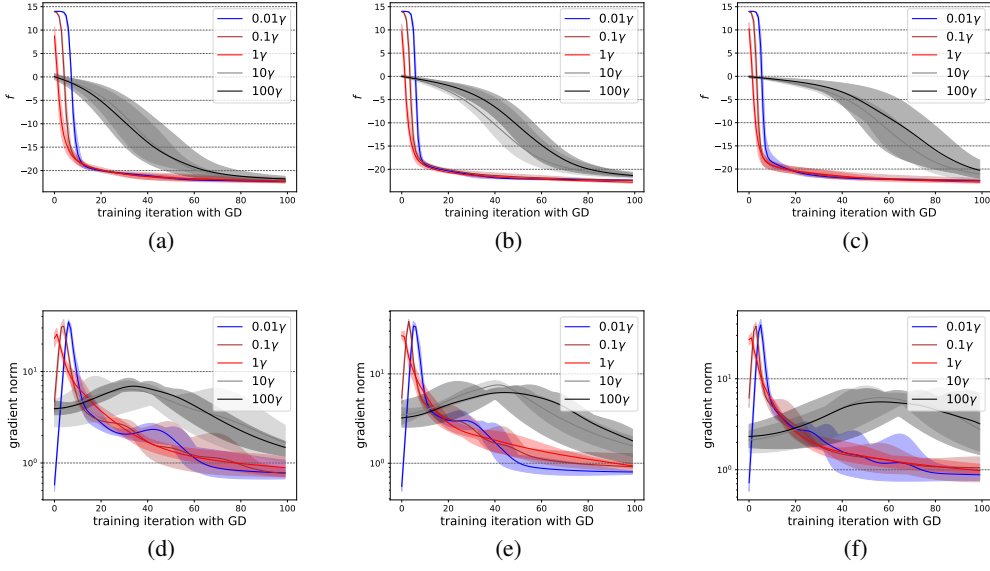


Figure 4: Numerical results of finding the ground state energy of the Heisenberg model using the noiseless gradient descent. Figures 4(a)-4(c) show the loss during optimizations for different $L \in \{14, 18, 22\}$ using the circuit 1 in the main text. For each L , we adopt Gaussian initializations with different variances $0.01\gamma, 0.1\gamma, \gamma, 10\gamma, 100\gamma$, where the value γ follows the formulation in Theorem 4.1. Figures 4(d)-4(f) show the ℓ_2 norm of corresponding gradients during the optimization. Each line illustrates the average of 5 rounds of independent experiments.

Lemma B.1. Let θ be a variable with Gaussian distribution $\mathcal{N}(0, \gamma^2)$. Let $\rho = \sum_k c_k \rho_k$ be the linear combination of density matrices $\{\rho_k\}$ with real coefficients $\{c_k\}$. Let G be a hermitian unitary and $V = e^{-i\theta G}$. Let O be an arbitrary hermitian quantum observable that anti-commutes with G . Then

$$\mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \text{Tr} [OV\rho V^\dagger]^2 \geq (1 - 4\gamma^2) \text{Tr} [O\rho]^2 + 4\gamma^2(1 - 4\gamma^2) \text{Tr} [iGO\rho]^2. \quad (13)$$

Proof. By replacing the term

$$V = e^{-i\theta G} = I \cos \theta - iG \sin \theta,$$

we have

$$\begin{aligned} \text{Tr} [OV\rho V^\dagger] &= \text{Tr} [O(I \cos \theta - iG \sin \theta)\rho(I \cos \theta + iG \sin \theta)] \\ &= \cos 2\theta \text{Tr} [O\rho] + \sin 2\theta \text{Tr} [iGO\rho], \end{aligned} \quad (14)$$

where Eq. (14) follows from the condition $OG + GO = 0$. Since O anti-commutes with G , iGO could be served as a hermitian observable. Based on Eq. (14), we have

$$\begin{aligned} \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \text{Tr} [OV\rho V^\dagger]^2 &= \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\cos 2\theta \text{Tr} [O\rho] + \sin 2\theta \text{Tr} [iGO\rho] \right)^2 \\ &= \frac{1 + e^{-8\gamma^2}}{2} \text{Tr} [O\rho]^2 + \frac{1 - e^{-8\gamma^2}}{2} \text{Tr} [iGO\rho]^2 \end{aligned} \quad (15)$$

$$\geq (1 - 4\gamma^2) \text{Tr} [O\rho]^2 + 4\gamma^2(1 - 4\gamma^2) \text{Tr} [iGO\rho]^2, \quad (16)$$

where Eq. (15) is obtained by calculating expectation terms. InEq. (16) holds since $1 - 8\gamma^2 \leq e^{-8\gamma^2} \leq 1 - 8\gamma^2 + 32\gamma^4$. Thus, we have proved Eq. (13). \square

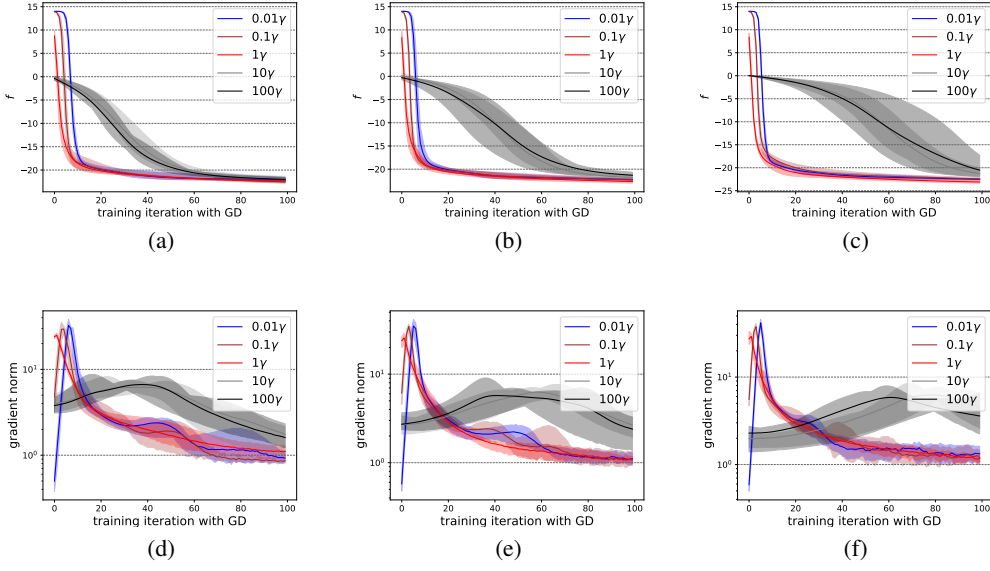


Figure 5: Numerical results of finding the ground state energy of the Heisenberg model using the noisy gradient descent. Figures 5(a)-5(c) show the loss during optimizations for different $L \in \{14, 18, 22\}$ using the circuit 1 in the main text. For each L , we adopt Gaussian initializations with different variances $0.01\gamma, 0.1\gamma, \gamma, 10\gamma, 100\gamma$, where the value γ follows the formulation in Theorem 4.1. Figures 5(d)-5(f) show the ℓ_2 norm of corresponding gradients during the optimization. Each line illustrates the average of 5 rounds of independent experiments.

Lemma B.2. Let θ be a variable with Gaussian distribution $\mathcal{N}(0, \gamma^2)$. Let ρ be the density matrix of a quantum state. Let G be a hermitian unitary and $V = e^{-i\theta G}$. Let O be an arbitrary hermitian quantum observable that anti-commutes with G . Then

$$\mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\frac{\partial}{\partial \theta} \text{Tr} [OV\rho V^\dagger] \right)^2 \geq (1 - 4\gamma^2) \left(\frac{\partial}{\partial \theta} \text{Tr} [OV\rho V^\dagger] \right)^2 \Big|_{\theta=0} + 16\gamma^2(1 - 4\gamma^2) \text{Tr} [O\rho]^2. \quad (17)$$

Proof. By calculating the gradient for both sides of Eq. (14), we obtain

$$\frac{\partial}{\partial \theta} \text{Tr} [OV\rho V^\dagger] = -2 \sin 2\theta \text{Tr} [O\rho] + 2 \cos 2\theta \text{Tr} [iGO\rho]. \quad (18)$$

Let $\theta = 0$ in Eq. (18), we obtain

$$\frac{\partial}{\partial \theta} \text{Tr} [OV\rho V^\dagger] \Big|_{\theta=0} = 2 \text{Tr} [iGO\rho]. \quad (19)$$

Now we proceed to prove Lemma B.2.

$$\text{The left part of Eq. (17)} = \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} (-2 \sin 2\theta \text{Tr} [O\rho] + 2 \cos 2\theta \text{Tr} [iGO\rho])^2$$

$$= 2(1 - e^{-8\gamma^2}) \text{Tr} [O\rho]^2 + 2(1 + e^{-8\gamma^2}) \text{Tr} [iGO\rho]^2 \quad (20)$$

$$\geq 16\gamma^2(1 - 4\gamma^2) \text{Tr} [O\rho]^2 + 4(1 - 4\gamma^2) \text{Tr} [iGO\rho]^2 \quad (21)$$

$$= (1 - 4\gamma^2) \left(\frac{\partial}{\partial \theta} \text{Tr} [OV\rho V^\dagger] \right)^2 \Big|_{\theta=0} + 16\gamma^2(1 - 4\gamma^2) \text{Tr} [O\rho]^2. \quad (22)$$

Eq. (20) is obtained by calculating expectation terms. InEq. (21) is obtained by using $1 - 8\gamma^2 \leq e^{-8\gamma^2} \leq 1 - 8\gamma^2 + 32\gamma^4$. Eq. (22) follows from Eq. (19). Thus, we have proved Eq. (17). \square

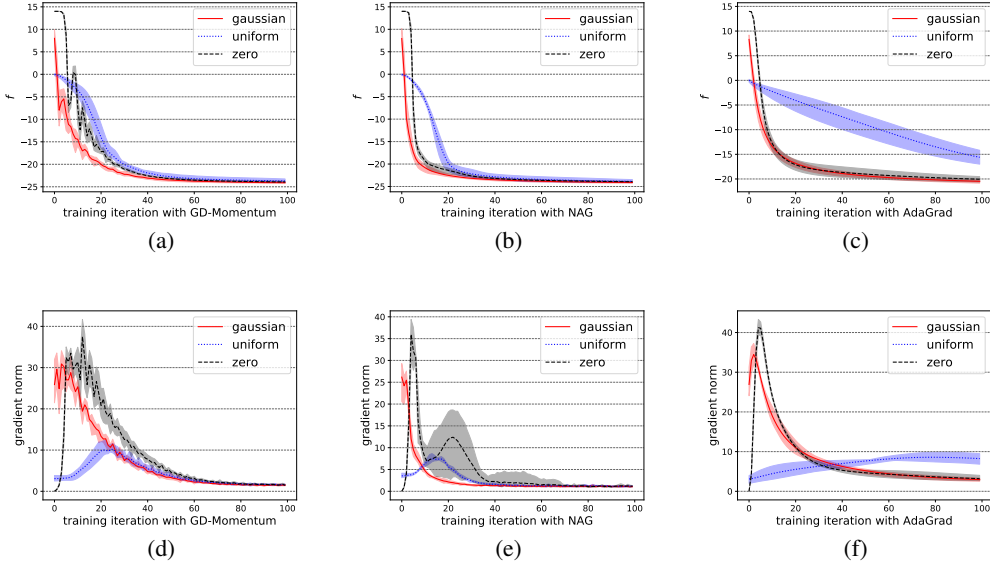


Figure 6: Numerical results of finding the ground state energy of the Heisenberg model with qubits $N = 15$ (noisy case). Figures 6(a)-6(c) show the loss during optimizations using the gradient descent with momentum, the Nesterov accelerated gradient (NAG), and the adaptive gradient (AdaGrad), respectively. Figures 6(d)-6(f) show the ℓ_2 norm of gradients during the optimization. Each line illustrates the average of 5 rounds of independent experiments.

Lemma B.3. Denote by $\rho = \sum_k c_k \rho_k$ the linear combination of density matrices $\{\rho_k\}$ with real coefficients $\{c_k\}$. Let $V_h(\theta) = W_1 e^{-i\theta G_1} W_2 \cdots W_h e^{-i\theta G_h}$, where $\{G_n\}_{n=1}^h$ is a list of hermitian unitaries and $\{W_n\}_{n=1}^h$ is a list of unitary matrices. Denote by O an arbitrary hermitian quantum observable. Then

$$\mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \text{Tr} \left[O V_h(\theta) \rho V_h(\theta)^\dagger \right]^2 \geq \text{Tr} \left[O V_h(0) \rho V_h(0)^\dagger \right]^2 - [12h(h-1) + 4] \gamma^2 \|c\|_1^2 \|O\|_2^2, \quad (23)$$

where $\|c\|_1 = \sum_k |c_k|$ denotes the ℓ_1 norm of c , $\|O\|_2$ denotes the spectral norm of O , and the variance $\gamma^2 \leq \frac{1}{12h^2}$.

Proof. Before the proof, we define several notations for convenience. We define $V_0 = I$ and

$$V_j(\theta) = W_{h+1-j} e^{-i\theta G_{h+1-j}} \cdots W_h e^{-i\theta G_h}, \forall j \in \{1, \dots, h\}. \quad (24)$$

We denote $\mathbf{0}_k$, $\mathbf{1}_k$, and $\mathbf{2}_k$ as k -dimensional vectors with components 0, 1, and 2, respectively. We define $O_{i_1, \dots, i_k}^{j_1, \dots, j_k} = O$ for the $k = 0$ case and

$$O_{i_1, \dots, i_k}^{j_1, \dots, j_k} = \begin{cases} W_k^\dagger O_{i_1, \dots, i_{k-1}}^{j_1, \dots, j_{k-1}} W_k, & \text{if } i_k = 0, j_k = 0, \\ \frac{1}{2} G_k \left\{ G_k, W_k^\dagger O_{i_1, \dots, i_{k-1}}^{j_1, \dots, j_{k-1}} W_k \right\}, & \text{if } i_k = 1, j_k = 0, \\ \frac{1}{2} G_k \left[G_k, W_k^\dagger O_{i_1, \dots, i_{k-1}}^{j_1, \dots, j_{k-1}} W_k \right], & \text{if } i_k = 2, j_k = 0, \\ i G_k O_{i_1, \dots, i_{k-1}, i_k}^{j_1, \dots, j_{k-1}, 0}, & \text{if } j_k = 1, \end{cases} \quad (25)$$

for increasing $k \in \{1, \dots, h\}$, where $i_k \in \{0, 1, 2\}$ and $j_k \in \{0, 1\}$.

For all $1 \leq k \leq \ell \leq h$, the definition (25) provides the commuting and anti-commuting parts of $O_{i_1, \dots, i_{k-1}, 0, i_{k+1}, \dots, i_\ell}^{j_1, \dots, j_{k-1}, 0, j_{k+1}, \dots, j_\ell}$ with respect to G_k , respectively, i.e.,

$$O_{i_1, \dots, i_{k-1}, 0, i_{k+1}, \dots, i_\ell}^{j_1, \dots, j_{k-1}, 0, j_{k+1}, \dots, j_\ell} = O_{i_1, \dots, i_{k-1}, 1, i_{k+1}, \dots, i_\ell}^{j_1, \dots, j_{k-1}, 0, j_{k+1}, \dots, j_\ell} + O_{i_1, \dots, i_{k-1}, 2, i_{k+1}, \dots, i_\ell}^{j_1, \dots, j_{k-1}, 0, j_{k+1}, \dots, j_\ell},$$

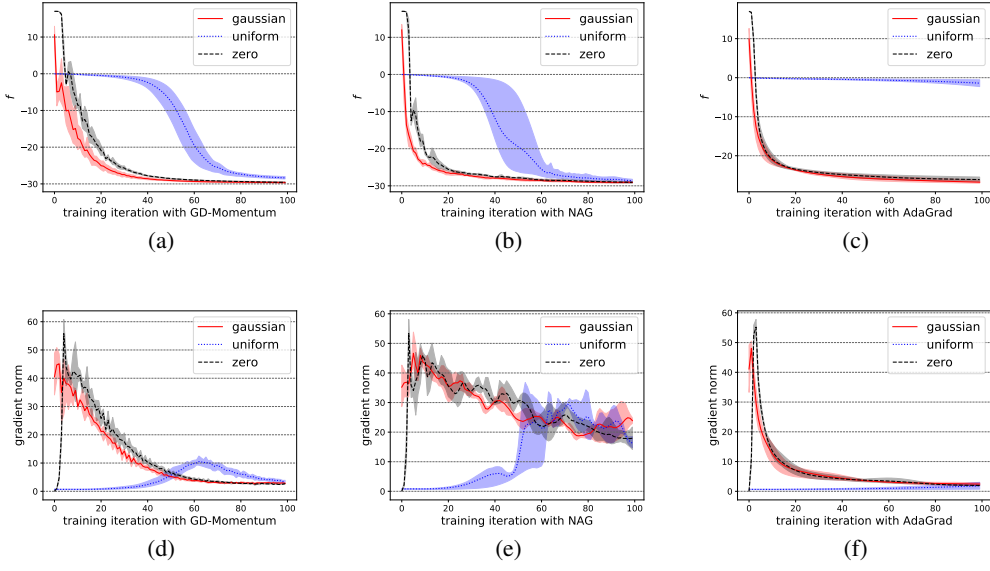


Figure 7: Numerical results of finding the ground state energy of the Heisenberg model with qubits $N = 18$ (noisy case). Figures 7(a)-7(c) show the loss during optimizations using the gradient descent with momentum, the Nesterov accelerated gradient (NAG), and the adaptive gradient (AdaGrad), respectively. Figures 7(d)-7(f) show the ℓ_2 norm of gradients during the optimization. Each line illustrates the average of 3 rounds of independent experiments.

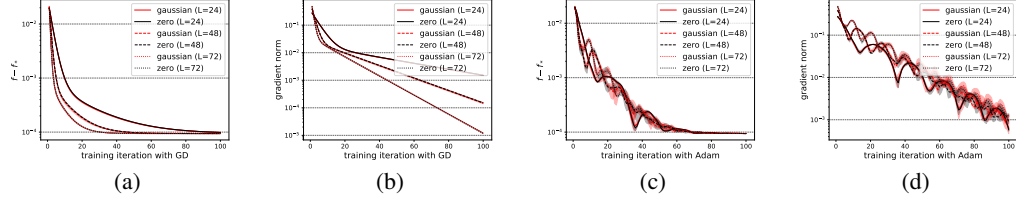


Figure 8: Numerical results of finding the ground state energy of the LiH molecule using noisy gradients. Figures 8(a) and 8(c) show the loss during optimizations for different $L \in \{24, 48, 72\}$ with the gradient descent and the Adam optimizer, respectively. Figures 8(b) and 8(d) show the ℓ_2 norm of gradients during the optimization. Each line illustrates the average of 3 rounds of independent experiments.

$$G_k O_{i_1, \dots, i_{k-1}, 1, i_{k+1}, \dots, i_\ell}^{j_1, \dots, j_{k-1}, 0, j_{k+1}, \dots, j_\ell} = O_{i_1, \dots, i_{k-1}, 1, i_{k+1}, \dots, i_\ell}^{j_1, \dots, j_{k-1}, 0, j_{k+1}, \dots, j_\ell} G_k,$$

$$G_k O_{i_1, \dots, i_{k-1}, 2, i_{k+1}, \dots, i_\ell}^{j_1, \dots, j_{k-1}, 0, j_{k+1}, \dots, j_\ell} = -O_{i_1, \dots, i_{k-1}, 2, i_{k+1}, \dots, i_\ell}^{j_1, \dots, j_{k-1}, 0, j_{k+1}, \dots, j_\ell} G_k.$$

Since for all $k \in [L]$, G_k is a unitary matrix, O_i^j is a hermitian observable for all $i \in \{0, 1, 2\}^\ell$, $j \in \{0, 1\}^\ell$, and $\ell \in [L]$. Meanwhile, the spectral norm of O_i^j is bounded,

$$\begin{aligned} \|O_{i_1, \dots, i_{h-1}, i_h}^{j_1, \dots, j_{h-1}, j_h}\|_2 &\leq \|O_{i_1, \dots, i_{h-1}, i_h}^{j_1, \dots, j_{h-1}, 0}\|_2 \leq \frac{1}{2} \|O_{i_1, \dots, i_{h-1}, 0}^{j_1, \dots, j_{h-1}, 0}\|_2 + \frac{1}{2} \|O_{i_1, \dots, i_{h-1}, 0}^{j_1, \dots, j_{h-1}, 0}\|_2 \\ &= \|O_{i_1, \dots, i_{h-1}, 0}^{j_1, \dots, j_{h-1}, 0}\|_2 = \|O_{i_1, \dots, i_{h-1}}^{j_1, \dots, j_{h-1}}\|_2 \leq \|O\|_2, \end{aligned} \quad (26)$$

where $\|A\|_2$ denotes the spectral norm of the matrix A . Moreover, for all $k, \ell \geq 0$ such that $k + \ell \leq h$, the observable $O_{i_1, \dots, i_k, 0_{h-k-\ell}, i_{h-\ell+1}, \dots, i_h}^{j_1, \dots, j_k, 0_{h-k-\ell}, j_{h-\ell+1}, \dots, j_h}$ could be recovered by

$$\sum_{n=k+1}^{h-\ell} \sum_{i_n=1}^2 O_{i_1, \dots, i_k, i_{k+1}, \dots, i_{h-\ell}, i_{h-\ell+1}, \dots, i_h}^{j_1, \dots, j_k, 0_{h-k-\ell}, j_{h-\ell+1}, \dots, j_h} = O_{i_1, \dots, i_k, 0_{h-k-\ell}, i_{h-\ell+1}, \dots, i_h}^{j_1, \dots, j_k, 0_{h-k-\ell}, j_{h-\ell+1}, \dots, j_h}. \quad (27)$$

Now we begin the proof. To analyze the expectation with respect to the parameter θ , we need the detailed formulation of $\text{Tr} \left[O V_h \rho V_h^\dagger \right]$ as the function of θ . In fact, for all $h' \in \{0, 1, \dots, h\}$ and all

$$\mathbf{i} \in \{0, 1, 2\}^{h-h'}, \mathbf{j} \in \{0, 1\}^{h-h'},$$

we have

$$\text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} V_{h'} \rho V_{h'}^\dagger \right] = \sum_{\mathbf{j}'=\mathbf{0}_{h'}}^{\mathbf{1}_{h'}} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_{h'}}^{\mathbf{2}_{h'}} (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - h'} (\sin 2\theta)^{\|\mathbf{j}'\|_1} \text{Tr} \left[O_{\mathbf{i}, \mathbf{i}'}^{\mathbf{j}, \mathbf{j}'} \rho \right], \quad (28)$$

where $\|\mathbf{i}'\|_1 \equiv \sum_{k=1}^{\dim(\mathbf{i}')} |i'_k|$ denotes the ℓ_1 norm of the vector \mathbf{i} .

Eq. (28) can be proved inductively. First, for the case $h' = 0$, Eq. (28) holds trivially. Next, we assume that Eq. (28) holds for the $h' = k$ case. Then for all

$$\mathbf{i} \in \{0, 1, 2\}^{h-k-1}, \mathbf{j} \in \{0, 1\}^{h-k-1},$$

we have

$$\begin{aligned} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} V_{k+1} \rho V_{k+1}^\dagger \right] &= \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} W_{h-k} (I \cos \theta - i G_{h-k} \sin \theta) V_k \rho V_k^\dagger (I \cos \theta + i G_{h-k} \sin \theta) W_{h-k}^\dagger \right] \\ &= \cos^2 \theta \text{Tr} \left[O_{\mathbf{i},0}^{\mathbf{j},0} V_k \rho V_k^\dagger \right] + \sin^2 \theta \text{Tr} \left[G_{h-k} O_{\mathbf{i},0}^{\mathbf{j},0} G_{h-k} V_k \rho V_k^\dagger \right] \\ &\quad + \sin \theta \cos \theta \text{Tr} \left[i G_{h-k} O_{\mathbf{i},0}^{\mathbf{j},0} V_k \rho V_k^\dagger \right] - \sin \theta \cos \theta \text{Tr} \left[O_{\mathbf{i},0}^{\mathbf{j},0} i G_{h-k} V_k \rho V_k^\dagger \right] \\ &= \text{Tr} \left[O_{\mathbf{i},1}^{\mathbf{j},0} V_k \rho V_k^\dagger \right] + \cos 2\theta \text{Tr} \left[O_{\mathbf{i},2}^{\mathbf{j},0} V_k \rho V_k^\dagger \right] + \sin 2\theta \text{Tr} \left[O_{\mathbf{i},2}^{\mathbf{j},1} V_k \rho V_k^\dagger \right], \end{aligned} \quad (29)$$

$$(30)$$

where Eqs. (29) and (30) are derived by using the definition (25). We proceed by employing the $h' = k$ case of Eq. (28), such that

$$\begin{aligned} \text{Eq. (30)} &= \sum_{\mathbf{j}'=\mathbf{0}_k}^{\mathbf{1}_k} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} (\sin 2\theta)^{\|\mathbf{j}'\|_1} \text{Tr} \left[O_{\mathbf{i},1}^{\mathbf{j},0} \rho \right] \\ &\quad + \cos 2\theta \sum_{\mathbf{j}'=\mathbf{0}_k}^{\mathbf{1}_k} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} (\sin 2\theta)^{\|\mathbf{j}'\|_1} \text{Tr} \left[O_{\mathbf{i},2}^{\mathbf{j},0} \rho \right] \\ &\quad + \sin 2\theta \sum_{\mathbf{j}'=\mathbf{0}_k}^{\mathbf{1}_k} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} (\sin 2\theta)^{\|\mathbf{j}'\|_1} \text{Tr} \left[O_{\mathbf{i},2}^{\mathbf{j},1} \rho \right] \\ &= \sum_{\mathbf{j}'=\mathbf{0}_{k+1}}^{\mathbf{1}_{k+1}} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_{k+1}}^{\mathbf{2}_{k+1}} (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k - 1} (\sin 2\theta)^{\|\mathbf{j}'\|_1} \text{Tr} \left[O_{\mathbf{i}, \mathbf{i}'}^{\mathbf{j}, \mathbf{j}'} \rho \right], \end{aligned} \quad (31)$$

which matches the formulation of the $h' = k + 1$ case of Eq. (28). Thus, Eq. (28) has been proved.

Now we begin to prove Eq. (23). Employing the $h' = h$ case of Eq. (28) could yield

$$\begin{aligned} &\mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\text{Tr} \left[O V_h \rho V_h^\dagger \right] \right)^2 \\ &= \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{\mathbf{j}=\mathbf{0}_h}^{\mathbf{1}_h} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{\mathbf{2}_h} (\cos 2\theta)^{\|\mathbf{i}\|_1 - \|\mathbf{j}\|_1 - h} (\sin 2\theta)^{\|\mathbf{j}\|_1} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right)^2 \end{aligned} \quad (32)$$

$$= \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{\mathbf{i}=\mathbf{1}_h}^{\mathbf{2}_h} (\cos 2\theta)^{\|\mathbf{i}\|_1 - h} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{0}_h} \rho \right] + \sum_{\mathbf{j} > \mathbf{0}_h}^{\mathbf{1}_h} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{\mathbf{2}_h} (\cos 2\theta)^{\|\mathbf{i}\|_1 - \|\mathbf{j}\|_1 - h} (\sin 2\theta)^{\|\mathbf{j}\|_1} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right)^2 \quad (33)$$

$$\begin{aligned}
&\geq \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{i=1_h}^{2_h} (\cos 2\theta)^{\|i\|_1 - h} \text{Tr} [O_i^{0_h} \rho] \right)^2 \\
&\quad + 2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{j > 0_h}^{1_h} \sum_{i=j+1_h}^{2_h} (\cos 2\theta)^{\|i\|_1 - \|j\|_1 - h} (\sin 2\theta)^{\|j\|_1} \text{Tr} [O_i^j \rho] \sum_{i'=1_h}^{2_h} (\cos 2\theta)^{\|i'\|_1 - h} \text{Tr} [O_{i'}^{0_h} \rho].
\end{aligned} \tag{34}$$

InEq. (34) is obtained by discarding the square of the latter term in the bracket of Eq. (33). We remark that if Eqs. (35) and (36) hold, we can prove Eq. (23) by using Eqs. (32-34).

$$\mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{i=1_h}^{2_h} (\cos 2\theta)^{\|i\|_1 - h} \text{Tr} [O_i^{0_h} \rho] \right)^2 - (\text{Tr} [OV_h(0) \rho V_h(0)^\dagger])^2 \geq -(6h - 2)\gamma^2 \|c\|_1^2 \|O\|_2^2, \tag{35}$$

$$\begin{aligned}
&\mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{j > 0_h}^{1_h} \sum_{i=j+1_h}^{2_h} (\cos 2\theta)^{\|i\|_1 - \|j\|_1 - h} (\sin 2\theta)^{\|j\|_1} \text{Tr} [O_i^j \rho] \sum_{i'=1_h}^{2_h} (\cos 2\theta)^{\|i'\|_1 - h} \text{Tr} [O_{i'}^{0_h} \rho] \\
&\geq -(6h^2 - 9h + 3) \gamma^2 \|c\|_1^2 \|O\|_2^2.
\end{aligned} \tag{36}$$

In the following proof, we would derive Eqs. (35) and (36). We focus on the Eq. (35) first. In fact, the left side of Eq. (35) is bounded by

$$\mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{i=1_h}^{2_h} [1 - (\cos 2\theta)^{\|i\|_1 - h} - 1] \text{Tr} [O_i^{0_h} \rho] \right)^2 - (\text{Tr} [O_{0_h}^{0_h} \rho])^2 \tag{37}$$

$$= \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{i=1_h}^{2_h} [1 - (\cos 2\theta)^{\|i\|_1 - h} - 1] \text{Tr} [O_i^{0_h} \rho] \right)^2 - \left(\sum_{i=1_h}^{2_h} \text{Tr} [O_i^{0_h} \rho] \right)^2 \tag{38}$$

$$\geq -2 \left| \sum_{i=1_h}^{2_h} \text{Tr} [O_i^{0_h} \rho] \right| \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left| \sum_{i=1_h}^{2_h} [1 - (\cos 2\theta)^{\|i\|_1 - h}] \text{Tr} [O_i^{0_h} \rho] \right| \tag{39}$$

$$= -2 |\text{Tr} [O_{0_h}^{0_h} \rho]| \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left| \sum_{i=1_h}^{2_h} [1 - (\cos 2\theta)^{\|i\|_1 - h}] \text{Tr} [O_i^{0_h} \rho] \right| \tag{40}$$

$$\geq -2 \|c\|_1 \|O\|_2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left| \sum_{i=1_h}^{2_h} [1 - (\cos 2\theta)^{\|i\|_1 - h}] \text{Tr} [O_i^{0_h} \rho] \right| \tag{41}$$

$$\geq -2 \|c\|_1^2 \|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} [(2 - \cos 2\theta)^h - 1]. \tag{42}$$

Eq. (37) is obtained by using the definition (25). Eq. (38) is derived by using Eq. (27). InEq. (39) is obtained by using $(a - b)^2 - b^2 \geq -2|a| \cdot |b|$. Eq. (40) yields from Eq. (27). InEq. (41) is derived by using

$$\left| \text{Tr} [O_i^j \rho] \right| = \left| \sum_k c_k \text{Tr} [O_i^j \rho_k] \right| \leq \sum_k |c_k| \left| \text{Tr} [O_i^j \rho_k] \right| \leq \sum_k |c_k| \|O_i^j\|_2 \leq \|c\|_1 \|O\|_2. \tag{43}$$

InEq. (42) is obtained by using the $h' = h$ case of InEq. (44), i.e.,

$$\left| \sum_{i'=j'+1_{h'}}^{2_{h'}} [1 - (\cos 2\theta)^{\|i'\|_1 - \|j'\|_1 - h'}] \text{Tr} [O_{i'}^{j', j'} \rho] \right| \leq [(2 - \cos 2\theta)^{h' - \|j'\|_1} - 1] \|c\|_1 \|O\|_2 \tag{44}$$

for all $h' \in \{0, 1, \dots, h\}$, $j' \in \{0, 1\}^{h'}$, $i \in \{0, 1, 2\}^{h-h'}$, and $j \in \{0, 1\}^{h-h'}$.

InEq. (44) can be proved inductively. First, for the case $h' = 0$, Eq. (44) holds trivially. Next we assume that Eq. (44) holds for the case $h' = k$. Then for all $\mathbf{i} \in \{0, 1, 2\}^{h-k-1}$ and $\mathbf{j} \in \{0, 1\}^{h-k-1}$, we have

$$\begin{aligned} & \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_{k+1}}^{\mathbf{2}_{k+1}} \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k - 1} \right] \text{Tr} \left[O_{\mathbf{i}', \mathbf{i}}^{\mathbf{j}', \mathbf{j}} \rho \right] \right| \\ &= \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \sum_{\mathbf{i}'_{k+1}=\mathbf{j}'_{k+1}+1}^2 \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 + \mathbf{i}'_{k+1} - \|\mathbf{j}'\|_1 - \mathbf{j}'_{k+1} - k - 1} \right] \text{Tr} \left[O_{\mathbf{i}', \mathbf{i}'_{k+1}, \mathbf{i}}^{\mathbf{j}', \mathbf{j}'_{k+1}, \mathbf{j}} \rho \right] \right|. \end{aligned} \quad (45)$$

For the case $\mathbf{j}'_{k+1} = 1$,

$$\begin{aligned} \text{Eq. (45)} &= \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}', 2, \mathbf{i}}^{\mathbf{j}', 1, \mathbf{j}} \rho \right] \right| \\ &\leq \left[(2 - \cos 2\theta)^{k - \|\mathbf{j}'\|_1} - 1 \right] \|\mathbf{c}\|_1 \|O\|_2 \end{aligned} \quad (46)$$

$$= \left[(2 - \cos 2\theta)^{k+1 - \|\mathbf{j}'\|_1 - \mathbf{j}'_{k+1}} - 1 \right] \|\mathbf{c}\|_1 \|O\|_2. \quad (47)$$

InEq. (46) is derived by using the $h' = k$ case of InEq. (44). Eq. (47) is derived by using $\mathbf{j}'_{k+1} = 1$. For the case $\mathbf{j}'_{k+1} = 0$,

$$\begin{aligned} \text{Eq. (45)} &= \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}', 1, \mathbf{i}}^{\mathbf{j}', 0, \mathbf{j}} \rho \right] \right. \\ &\quad \left. + \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k + 1} \right] \text{Tr} \left[O_{\mathbf{i}', 2, \mathbf{i}}^{\mathbf{j}', 0, \mathbf{j}} \rho \right] \right| \\ &= \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}', 0, \mathbf{i}}^{\mathbf{j}', 0, \mathbf{j}} \rho \right] \right. \\ &\quad \left. + (1 - \cos 2\theta) \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[(\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} - 1 + 1 \right] \text{Tr} \left[O_{\mathbf{i}', 2, \mathbf{i}}^{\mathbf{j}', 0, \mathbf{j}} \rho \right] \right| \end{aligned} \quad (48)$$

$$\begin{aligned} &\leq \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}', 0, \mathbf{i}}^{\mathbf{j}', 0, \mathbf{j}} \rho \right] \right| \\ &\quad + (1 - \cos 2\theta) \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}', 2, \mathbf{i}}^{\mathbf{j}', 0, \mathbf{j}} \rho \right] \right| \\ &\quad + (1 - \cos 2\theta) \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \text{Tr} \left[O_{\mathbf{i}', 2, \mathbf{i}}^{\mathbf{j}', 0, \mathbf{j}} \rho \right] \right| \end{aligned} \quad (49)$$

$$\begin{aligned} &\leq \left[(2 - \cos 2\theta)^{k - \|\mathbf{j}'\|_1} - 1 \right] \|\mathbf{c}\|_1 \|O\|_2 + (1 - \cos 2\theta) \left[(2 - \cos 2\theta)^{k - \|\mathbf{j}'\|_1} - 1 \right] \|\mathbf{c}\|_1 \|O\|_2 \\ &\quad + (1 - \cos 2\theta) \|\mathbf{c}\|_1 \|O\|_2 \end{aligned} \quad (50)$$

$$\leq \left[(2 - \cos 2\theta)^{k+1 - \|\mathbf{j}'\|_1 - \mathbf{j}'_{k+1}} - 1 \right] \|\mathbf{c}\|_1 \|O\|_2. \quad (51)$$

Eq. (48) is derived by using Eq. (27). InEq. (49) is obtained since $|a + b| \leq |a| + |b|$. InEq. (50) is obtained using the $h' = k$ case of Eq. (44) and Eq. (27). InEq. (51) is derived by using Eq. (26). Thus we have proved Eq. (44) since Eqs. (47) and (51) match the $h' = k + 1$ case.

Since $\cos 2\theta \geq 1 - 2\theta^2$, we could further bound Eq. (42) by

$$\text{Eq. (42)} \geq -2\|\mathbf{c}\|_1^2 \|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left[(1 + 2\theta^2)^h - 1 \right] \quad (52)$$

$$\begin{aligned}
&= -2\|\mathbf{c}\|_1^2\|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{t=1}^h \binom{h}{t} (2\theta^2)^t \\
&= -2\|\mathbf{c}\|_1^2\|O\|_2^2 \sum_{t=1}^h \binom{h}{t} (2t-1)!! (2\gamma^2)^t \tag{53}
\end{aligned}$$

$$\geq -2\|\mathbf{c}\|_1^2\|O\|_2^2 \sum_{t=1}^h h(h-1)^{t-1} 2^{t-1} \left(\frac{1}{6h^2}\right)^{t-1} (2\gamma^2) \tag{54}$$

$$\begin{aligned}
&= -2\|\mathbf{c}\|_1^2\|O\|_2^2 2h\gamma^2 \left[1 + \frac{h-1}{3h^2} \sum_{t=0}^{h-2} \left(\frac{h-1}{3h^2}\right)^t\right] \\
&\geq -(6h-2)\|\mathbf{c}\|_1^2\|O\|_2^2\gamma^2. \tag{55}
\end{aligned}$$

Eq. (53) is derived by calculating expectation terms. InEq. (54) yields from $\frac{(2t-1)!!}{t!} \leq 2^{t-1}$ and the condition $\gamma^2 \leq \frac{1}{12h^2}$. Thus, we have proved InEq. (35).

Next, we focus on the Eq. (36). The left side of Eq. (36) could be bounded by

$$\begin{aligned}
&= \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{j > \mathbf{0}_h, 2\|j\|_1}^{1_h} \sum_{i=j+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|i\|_1 - \|j\|_1 - h} (\sin 2\theta)^{\|j\|_1} \text{Tr} [O_i^j \rho] \\
&\quad \cdot \sum_{i'=\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|i'\|_1 - h} \text{Tr} [O_{i'}^{0_h} \rho] \tag{56}
\end{aligned}$$

$$\begin{aligned}
&\geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{j > \mathbf{0}_h, 2\|j\|_1}^{1_h} (\sin 2\theta)^{\|j\|_1} \left(\left| \sum_{i=j+\mathbf{1}_h}^{2_h} \left[1 - (\cos 2\theta)^{\|i\|_1 - \|j\|_1 - h}\right] \text{Tr} [O_i^j \rho] \right| \right. \\
&\quad \left. + \left| \sum_{i=j+\mathbf{1}_h}^{2_h} \text{Tr} [O_i^j \rho] \right| \right) \cdot \left(\left| \sum_{i'=\mathbf{1}_h}^{2_h} \left[1 - (\cos 2\theta)^{\|i'\|_1 - h}\right] \text{Tr} [O_{i'}^{0_h} \rho] \right| + \left| \sum_{i'=\mathbf{1}_h}^{2_h} \text{Tr} [O_{i'}^{0_h} \rho] \right| \right) \tag{57}
\end{aligned}$$

$$\begin{aligned}
&\geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{j > \mathbf{0}_h, 2\|j\|_1}^{1_h} (\sin 2\theta)^{\|j\|_1} \left(\left[(2 - \cos 2\theta)^{h - \|j\|_1} - 1 \right] \|\mathbf{c}\|_1 \|O_{0_h}^{0_h}\|_2 + \left| \text{Tr} [O_{2j}^j \rho] \right| \right) \\
&\quad \cdot \left(\left[(2 - \cos 2\theta)^h - 1 \right] \|\mathbf{c}\|_1 \|O_{0_h}^{0_h}\|_2 + \left| \text{Tr} [O_{0_h}^{0_h} \rho] \right| \right) \tag{58}
\end{aligned}$$

$$\geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{j > \mathbf{0}_h, 2\|j\|_1}^{1_h} (\sin 2\theta)^{\|j\|_1} (2 - \cos 2\theta)^{2h - \|j\|_1} \|\mathbf{c}\|_1^2 \|O\|_2^2 \tag{59}$$

$$\geq - \|\mathbf{c}\|_1^2 \|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{j > \mathbf{0}_h, 2\|j\|_1}^{1_h} (2\theta)^{\|j\|_1} (1 + 2\theta^2)^{2h - \|j\|_1}. \tag{60}$$

Eq. (56) is obtained by noticing that the expectation of $\sin^a 2\theta \cos^b 2\theta$ equals to zero, if a is odd. InEq. (57) is obtained by using $\sum_{i,j} a_i b_j \geq -(\sum_i |a_i|)(\sum_j |b_j|)$ and $|a+b| \leq |a| + |b|$. InEq. (58) is derived by using the $h' = h$ case of Eq. (44) and Eq. (27). InEq. (59) is obtained by using $\|O_{0_h}^{0_h}\| = \|O\|$ and Eq. (43). InEq. (60) is derived by using $(\sin 2\theta)^2 \leq (2\theta)^2$ and $\cos 2\theta \geq 1 - 2\theta^2$.

We proceed from InEq. (60), which could be further bounded by

$$= - \|\mathbf{c}\|_1^2 \|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{t=1}^{\lfloor h/2 \rfloor} \binom{h}{2t} (2\theta)^{2t} \sum_{m=0}^{2h-2t} \binom{2h-2t}{m} (2\theta^2)^m \tag{61}$$

$$= - \|\mathbf{c}\|_1^2 \|O\|_2^2 \sum_{t=1}^{\lfloor h/2 \rfloor} \binom{h}{2t} \sum_{m=0}^{2h-2t} \binom{2h-2t}{m} 2^{2t+m} (2t+2m-1)!! \gamma^{2t+2m} \tag{62}$$

$$\geq -\|\mathbf{c}\|_1^2 \|O\|_2^2 \sum_{t=1}^{\lfloor h/2 \rfloor} \sum_{m=0}^{2h-2t} \frac{h(h-1)^{2t-1}}{2^t t! (2t-1)!!} \frac{(2h-2)^m}{m!} 2^{2t+m} \cdot (2t-1)!! (2t+1)(2t+3) \cdots (2t+2m-1) \gamma^{2t+2m} \quad (63)$$

$$\geq -\|\mathbf{c}\|_1^2 \|O\|_2^2 \sum_{t=1}^{\lfloor h/2 \rfloor} \sum_{m=0}^{2h-2t} \frac{h(h-1)^{2t-1}}{2^t 2^{t-1}} \frac{(2h-2)^m}{m!} 2^{2t+m} (2h)^m \gamma^{2t+2m} \quad (64)$$

$$\geq -\|\mathbf{c}\|_1^2 \|O\|_2^2 \sum_{t=1}^{\lfloor h/2 \rfloor} \left(2h(h-1)^{2t-1} \gamma^{2t} + \sum_{m=1}^{2h-2t} 4h(h-1)^{2t-1} (2h-2)^m (2h)^m \gamma^{2t+2m} \right) \quad (65)$$

$$= -\|\mathbf{c}\|_1^2 \|O\|_2^2 \left(\sum_{t=1}^{\lfloor h/2 \rfloor} 2h(h-1)^{2t-1} \gamma^{2t} \right) \cdot \left(1 + 2 \sum_{m=1}^{2h-2t} [4h(h-1) \gamma^2]^m \right) \quad (66)$$

$$\geq -\|\mathbf{c}\|_1^2 \|O\|_2^2 3h(h-1) \gamma^2 (1 + 12h(h-1) \gamma^2) \quad (67)$$

$$\geq -(6h^2 - 9h + 3) \gamma^2 \|\mathbf{c}\|_1^2 \|O\|_2^2. \quad (68)$$

Here, Eq. (61) is obtained since the summation $\sum_{j>0_h}^{1_h}$ contains $\binom{h}{2t}$ terms such that $\|j\|_1 = 2t$, for all $t \in \{1, \dots, \lfloor \frac{h}{2} \rfloor\}$. Eq. (62) is derived by calculating expectation terms. InEq. (63) is obtained by using

$$\binom{h}{2t} \leq \frac{h(h-1)^{2t-1}}{(2t)!} = \frac{h(h-1)^{2t-1}}{2^t t! (2t-1)!!} \text{ and } \binom{2h-2t}{m} \leq \frac{(2h-2t)^m}{m!}.$$

InEq. (64) is derived by using $t! \geq 2^{t-1}$ and

$$(2t+2k-1)(2t+2m-2k+1) \leq (2t+m)^2 \leq (2h)^2, \forall k \in \{1, \dots, m-1\}.$$

InEq. (65) is obtained by splitting the summation \sum_m and using $m! \geq 2^{m-1}, \forall m \geq 1$. InEq. (67) is derived by calculating geometric sequences with the condition $\gamma^2 \leq \frac{1}{12h^2}$. InEq. (68) follows from the condition $\gamma^2 \leq \frac{1}{12h^2}$. Thus, we have proved Eq. (36). \square

Lemma B.4. Let ρ be the density matrix of a quantum state. Let $V_h = W_1 e^{-i\theta G_1} W_2 \cdots W_h e^{-i\theta G_h}$, where $\{G_n\}_{n=1}^h$ is a list of hermitian unitaries and $\{W_n\}_{n=1}^h$ is a list of unitary matrices. Denote by O an arbitrary hermitian quantum observable. Then

$$\mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\frac{\partial}{\partial \theta} \text{Tr} [OV_h \rho V_h^\dagger] \right)^2 \geq (1 - 4\gamma^2) \left(\frac{\partial}{\partial \theta} \text{Tr} [OV_h \rho V_h^\dagger] \Big|_{\theta=0} \right)^2 - 96h^2(h-1)\gamma^2 \|O\|_2^2 - 20h^2(h-1)(h-2)\gamma^2 \|O\|_2^2, \quad (69)$$

where $\|O\|_2$ denotes the spectral norm of O and the variance $\gamma^2 \leq \frac{1}{16h^3}$.

Proof. For convenience, we follow the notation O_i^j in Eq. (25). We can obtain the detailed formulation of $\frac{\partial}{\partial \theta} \text{Tr} [OV_h \rho V_h^\dagger]$ by using the $h' = h$ case of Eq. (28),

$$\frac{\partial}{\partial \theta} \text{Tr} [OV_h \rho V_h^\dagger] = \frac{\partial}{\partial \theta} \sum_{j=0_h}^{1_h} \sum_{i=j+1_h}^{2_h} (\cos 2\theta)^{\|i\|_1 - \|j\|_1 - h} (\sin 2\theta)^{\|j\|_1} \text{Tr} [O_i^j \rho] \quad (70)$$

$$\begin{aligned} &= 2 \sum_{j=0_h}^{1_h} \sum_{i=j+1_h}^{2_h} (h + \|j\|_1 - \|i\|_1) (\cos 2\theta)^{\|i\|_1 - \|j\|_1 - h - 1} (\sin 2\theta)^{\|j\|_1 + 1} \text{Tr} [O_i^j \rho] \\ &\quad + 2 \sum_{j=0_h}^{1_h} \sum_{i=j+1_h}^{2_h} \|j\|_1 (\cos 2\theta)^{\|i\|_1 - \|j\|_1 - h + 1} (\sin 2\theta)^{\|j\|_1 - 1} \text{Tr} [O_i^j \rho] \end{aligned} \quad (71)$$

$$= 2 \sum_{j=0_h}^{1_h} \sum_{i=j+1_h}^{2_h} (h + \|j\|_1 - \|i\|_1) (\cos 2\theta)^{\|i\|_1 - \|j\|_1 - h - 1} (\sin 2\theta)^{\|j\|_1 + 1} \text{Tr} [O_i^j \rho]$$

$$\begin{aligned}
& + 2 \sum_{\|\mathbf{j}\|_1=1} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|\mathbf{i}\|_1-h} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \\
& + 2 \sum_{\|\mathbf{j}\|_1 \geq 2}^{1_h} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} \|\mathbf{j}\|_1 (\cos 2\theta)^{\|\mathbf{i}\|_1-\|\mathbf{j}\|_1-h+1} (\sin 2\theta)^{\|\mathbf{j}\|_1-1} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right]. \tag{72}
\end{aligned}$$

Here, Eq. (70) follows from Eq. (28). Eq. (71) is derived by calculating the gradient of sine and cosine terms. By discarding the square of the sum of the first and the third term in Eq. (72), we obtain

$$\begin{aligned}
& \left(\frac{\partial}{\partial \theta} \text{Tr} \left[O V_h \rho V_h^\dagger \right] \right)^2 \geq 4 \left(\sum_{\|\mathbf{j}\|_1=1} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|\mathbf{i}\|_1-h} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right)^2 \\
& + 8 \left(\sum_{\|\mathbf{j}\|_1 \geq 2}^{1_h} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} \|\mathbf{j}\|_1 (\cos 2\theta)^{\|\mathbf{i}\|_1-\|\mathbf{j}\|_1-h+1} (\sin 2\theta)^{\|\mathbf{j}\|_1-1} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right) \\
& \cdot \left(\sum_{\|\mathbf{j}'\|_1=1} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|\mathbf{i}'\|_1-h} \text{Tr} \left[O_{\mathbf{i}'}^{\mathbf{j}'} \rho \right] \right) \\
& + 8 \left(\sum_{\mathbf{j}=\mathbf{0}_h}^{1_h} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} (h + \|\mathbf{j}\|_1 - \|\mathbf{i}\|_1) (\cos 2\theta)^{\|\mathbf{i}\|_1-\|\mathbf{j}\|_1-h-1} (\sin 2\theta)^{\|\mathbf{j}\|_1+1} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right) \\
& \cdot \left(\sum_{\|\mathbf{j}'\|_1=1} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|\mathbf{i}'\|_1-h} \text{Tr} \left[O_{\mathbf{i}'}^{\mathbf{j}'} \rho \right] \right). \tag{73}
\end{aligned}$$

Let $\theta = 0$ in Eq. (72), we obtain

$$\frac{\partial}{\partial \theta} \text{Tr} \left[O V_h \rho V_h^\dagger \right] \Big|_{\theta=0} = 2 \sum_{\|\mathbf{j}\|_1=1} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right]. \tag{74}$$

Thus, we could obtain Eq. (69) if Eqs. (75-77) hold.

$$\begin{aligned}
& \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{\|\mathbf{j}\|_1=1} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|\mathbf{i}\|_1-h} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right)^2 - (1 - 4\gamma^2) \left(\sum_{\|\mathbf{j}\|_1=1} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right)^2 \\
& \geq -\frac{13}{3} h^2 (h-1) \gamma^2 \|O\|_2^2, \tag{75}
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{\mathbf{j}=\mathbf{0}_h}^{1_h} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} (h + \|\mathbf{j}\|_1 - \|\mathbf{i}\|_1) (\cos 2\theta)^{\|\mathbf{i}\|_1-\|\mathbf{j}\|_1-h-1} (\sin 2\theta)^{\|\mathbf{j}\|_1+1} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right) \\
& \cdot \left(\sum_{\|\mathbf{j}'\|_1=1} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|\mathbf{i}'\|_1-h} \text{Tr} \left[O_{\mathbf{i}'}^{\mathbf{j}'} \rho \right] \right) \geq -\frac{59}{6} h^2 (h-1) \gamma^2 \|O\|_2^2, \tag{76}
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{\|\mathbf{j}\|_1 \geq 2}^{1_h} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} \|\mathbf{j}\|_1 (\cos 2\theta)^{\|\mathbf{i}\|_1-\|\mathbf{j}\|_1-h+1} (\sin 2\theta)^{\|\mathbf{j}\|_1-1} \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right) \\
& \cdot \left(\sum_{\|\mathbf{j}'\|_1=1} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|\mathbf{i}'\|_1-h} \text{Tr} \left[O_{\mathbf{i}'}^{\mathbf{j}'} \rho \right] \right) \geq -\frac{5}{2} h^2 (h-1)(h-2) \gamma^2 \|O\|_2^2. \tag{77}
\end{aligned}$$

We begin by proving Eq. (75). The left side of Eq. (75) can be lower bounded as

$$\geq \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{\|\mathbf{j}\|_1=1} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} \left[\cos 2\theta - (\cos 2\theta)^{\|\mathbf{i}\|_1-h} - \cos 2\theta \right] \text{Tr} \left[O_{\mathbf{i}}^{\mathbf{j}} \rho \right] \right)^2$$

$$- \left(\cos 2\theta \sum_{\|\mathbf{j}\|_1=1} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} \text{Tr} [O_{\mathbf{i}}^j \rho] \right)^2 \quad (78)$$

$$\geq -2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} (\cos 2\theta)^2 \sum_{\|\mathbf{j}\|_1=1} \left| \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} [1 - (\cos 2\theta)^{\|\mathbf{i}\|_1-1-h}] \text{Tr} [O_{\mathbf{i}}^j \rho] \right| \cdot \sum_{\|\mathbf{j}'\|_1=1} \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_h}^{2_h} \text{Tr} [O_{\mathbf{i}'}^{j'} \rho] \right| \quad (79)$$

$$\geq -2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} (\cos 2\theta)^2 \sum_{\|\mathbf{j}\|_1=1} [(2 - \cos 2\theta)^{h-1} - 1] \|O\|_2 \cdot \sum_{\|\mathbf{j}\|_1=1} \left| \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} \text{Tr} [O_{\mathbf{i}}^j \rho] \right| \quad (80)$$

$$\geq -2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} (\cos 2\theta)^2 h [(2 - \cos 2\theta)^{h-1} - 1] \|O\|_2 \cdot h \|O\|_2 \quad (81)$$

$$\geq -2h^2 \|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} [(1 + 2\theta^2)^{h-1} - 1] \quad (82)$$

$$\geq -\frac{13}{3} h^2 (h-1) \gamma^2 \|O\|_2^2. \quad (83)$$

Here, InEq. (78) follows from

$$1 - 4\gamma^2 = \mathbb{E}_\theta [1 - 4\theta^2] \leq \mathbb{E}_\theta (1 - 2\theta^2)^2 \leq \mathbb{E}_\theta (\cos 2\theta)^2.$$

InEq. (79) is obtained by using $(a-b)^2 - b^2 \geq -2|a| \cdot |b|$. InEq. (80) follows from the $h' = h$ and $\|\mathbf{j}\| = 1$ case of Eq. (44). InEq. (81) is derived by using Eq. (27). InEq. (82) is obtained by using $\cos 2\theta \geq 1 - 2\theta^2$. InEq. (83) follows from the derivation below.

$$\begin{aligned} \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} (1 + 2\theta^2)^{h-1} - 1 &= \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{t=1}^{h-1} \binom{h-1}{t} (2\theta^2)^t \\ &= \sum_{t=1}^{h-1} \binom{h-1}{t} (2t-1)!! (2\gamma^2)^t \end{aligned} \quad (84)$$

$$\begin{aligned} &\leq \sum_{t=1}^{h-1} (h-1)(h-2)^{t-1} 2^{t-1} (2\gamma^2)^t \\ &\leq 2(h-1)\gamma^2 \sum_{t=1}^{h-1} [h^3 \gamma^2]^{t-1} \end{aligned} \quad (85)$$

$$\leq \frac{13}{6} (h-1) \gamma^2, \quad (86)$$

where Eq. (84) is obtained by calculating expectation terms. InEq (85) follows from $h^3 \geq 4(h-2)$ for integer h . InEq. (86) is derived by calculating geometric sequences with the condition $\gamma^2 \leq \frac{1}{16h^3}$.

Next, we prove Eq. (76). The left side of Eq. (76) could be lower bounded by

$$\begin{aligned} &= \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{\|\mathbf{j}'\|_1=1} \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_h}^{2_h} (\cos 2\theta)^{\|\mathbf{i}'\|_1-1-h} \text{Tr} [O_{\mathbf{i}'}^{j'} \rho] \right) \\ &\quad \cdot \left(\sum_{\substack{\mathbf{j}=\mathbf{0}_h \\ 2|(\|\mathbf{j}\|_1+1)}}^{1_h} \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{2_h} (h + \|\mathbf{j}\|_1 - \|\mathbf{i}\|_1) (\cos 2\theta)^{\|\mathbf{i}\|_1-\|\mathbf{j}\|_1-h} (\sin 2\theta)^{\|\mathbf{j}\|_1+1} \text{Tr} [O_{\mathbf{i}}^j \rho] \right) \end{aligned} \quad (87)$$

$$\geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{\|\mathbf{j}'\|_1=1} \left(\left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_h}^{2_h} [1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1-1-h}] \text{Tr} [O_{\mathbf{i}'}^{j'} \rho] \right| + \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_h}^{2_h} \text{Tr} [O_{\mathbf{i}'}^{j'} \rho] \right| \right)$$

$$\cdot \sum_{\substack{\mathbf{j}=\mathbf{0}_h \\ 2|(\|\mathbf{j}\|_1+1)}}^{\mathbf{1}_h} (\sin 2\theta)^{\|\mathbf{j}\|_1+1} \left| \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{\mathbf{2}_h} (\|\mathbf{i}\|_1 - \|\mathbf{j}\|_1 - h) (\cos 2\theta)^{\|\mathbf{i}\|_1 - \|\mathbf{j}\|_1 - h} \text{Tr} [O_{\mathbf{i}}^{\mathbf{j}} \rho] \right| \quad (88)$$

$$\geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{\|\mathbf{j}'\|_1=1} (2 - \cos 2\theta)^{h-1} \|O\|_2 \sum_{\substack{\mathbf{j}=\mathbf{0}_h \\ 2|(\|\mathbf{j}\|_1+1)}}^{\mathbf{1}_h} (\sin 2\theta)^{\|\mathbf{j}\|_1+1} \left| \sum_{\mathbf{i}=\mathbf{j}+\mathbf{1}_h}^{\mathbf{2}_h} \left[(h - \|\mathbf{j}\|_1) - (\|\mathbf{i}\|_1 - \|\mathbf{j}\|_1 - h) (\cos 2\theta)^{\|\mathbf{i}\|_1 - \|\mathbf{j}\|_1 - h} - (h - \|\mathbf{j}\|_1) \right] \text{Tr} [O_{\mathbf{i}}^{\mathbf{j}} \rho] \right| \quad (89)$$

$$\geq -h \|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} (2 - \cos 2\theta)^{h-1} \cdot \sum_{\substack{\mathbf{j}=\mathbf{0}_h \\ 2|(\|\mathbf{j}\|_1+1)}}^{\mathbf{1}_h} (\sin 2\theta)^{\|\mathbf{j}\|_1+1} (h - \|\mathbf{j}\|_1) (3 - \cos 2\theta) (2 - \cos 2\theta)^{h-\|\mathbf{j}\|_1-1}. \quad (90)$$

Here, Eq. (87) is obtained by noticing that the expectation of $\sin^a 2\theta \cos^b 2\theta$ equals to zero, if a is odd. InEq. (88) is derived by using $|\sum_k a_k| \leq \sum_k |a_k|$. InEq. (89) is obtained by using the $h' = h$, $\|\mathbf{j}\|_1 = 1$ case of Eq. (44) and Eq. (27). InEq. (90) follows from Eq. (27) and the $h' = h$ case of Eq. (91), i.e.

$$\left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_{h'}}^{\mathbf{2}_{h'}} \left[(h' - \|\mathbf{j}'\|_1) - (\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - h') (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - h'} \right] \text{Tr} [O_{\mathbf{i}'}^{\mathbf{j}', \mathbf{j}} \rho] \right| \leq g(h' - \|\mathbf{j}'\|_1) \|O\|_2 \quad (91)$$

for all $h' \in \{0, 1, \dots, h\}$, $\mathbf{i} \in \{0, 1, 2\}^{h-h'}$, and $\mathbf{j} \in \{0, 1\}^{h-h'}$, where

$$g(x) = x [(3 - \cos 2\theta)(2 - \cos 2\theta)^{x-1} - 1]. \quad (92)$$

InEq. (91) can be proved inductively. First, InEq. (91) holds trivially when $h' = 0$. Next, we assume that InEq. (91) holds for the $h' = k$ case. Then, for all $\mathbf{i} \in \{0, 1, 2\}^{h-k-1}$ and $\mathbf{j} \in \{0, 1\}^{h-k-1}$, we have

$$\begin{aligned} & \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_{k+1}}^{\mathbf{2}_{k+1}} \left[(k+1 - \|\mathbf{j}'\|_1) - (\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k-1) (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k-1} \right] \text{Tr} [O_{\mathbf{i}'}^{\mathbf{j}', \mathbf{j}} \rho] \right| \\ &= \left| \sum_{\mathbf{i}'_{k+1}=j'_{k+1}+1}^2 \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[(k+1 - \|\mathbf{j}'\|_1 - j'_{k+1}) \right. \right. \\ & \quad \left. \left. - (\|\mathbf{i}'\|_1 + i'_{k+1} - \|\mathbf{j}'\|_1 - j'_{k+1} - k-1) (\cos 2\theta)^{\|\mathbf{i}'\|_1 + i'_{k+1} - \|\mathbf{j}'\|_1 - j'_{k+1} - k-1} \right] \text{Tr} [O_{\mathbf{i}', i'_{k+1}}^{\mathbf{j}', j_{k+1}, \mathbf{j}} \rho] \right|. \end{aligned} \quad (93)$$

For the case $j'_{k+1} = 1$, we have

$$\begin{aligned} \text{Eq. (93)} &= \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{\mathbf{2}_k} \left[(k - \|\mathbf{j}'\|_1) - (\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k) (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} [O_{\mathbf{i}', 2, \mathbf{i}}^{\mathbf{j}', 1, \mathbf{j}} \rho] \right| \\ &\leq g(k - \|\mathbf{j}'\|_1) \|O\|_2 \end{aligned} \quad (94)$$

$$= g(k+1 - \|\mathbf{j}'\|_1 - j'_{k+1}) \|O\|_2. \quad (95)$$

Here, Eq. (94) follows from the $h' = k$ case of InEq. (91). Eq. (95) is obtained by using $j'_{k+1} = 1$. We remark that InEq. (95) matches the $h' = k+1$ case of InEq. (91).

For the case $j'_{k+1} = 0$, the situation is more complicated. We have

$$\begin{aligned}
& \text{Eq. (93)} \\
& = \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{2_k} \left[(k+1 - \|\mathbf{j}'\|_1) - (\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k)(\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}',1,\mathbf{i}}^{\mathbf{j}',0,\mathbf{j}} \rho \right] \right. \\
& \quad + \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{2_k} \left[(k+1 - \|\mathbf{j}'\|_1) - (\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k+1)(\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k+1} \right] \text{Tr} \left[O_{\mathbf{i}',2,\mathbf{i}}^{\mathbf{j}',0,\mathbf{j}} \rho \right] \left. \right| \\
& \leq \left| \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{2_k} \text{Tr} \left[O_{\mathbf{i}',1,\mathbf{i}}^{\mathbf{j}',0,\mathbf{j}} \rho \right] \right. \\
& \quad + \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{2_k} \left[(k - \|\mathbf{j}'\|_1) - (\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k)(\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}',0,\mathbf{i}}^{\mathbf{j}',0,\mathbf{j}} \rho \right] \\
& \quad + (1 - \cos 2\theta)(k+1 - \|\mathbf{j}'\|_1) \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{2_k} \text{Tr} \left[O_{\mathbf{i}',2,\mathbf{i}}^{\mathbf{j}',0,\mathbf{j}} \rho \right] \\
& \quad + \cos 2\theta \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{2_k} \left[1 - (\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}',2,\mathbf{i}}^{\mathbf{j}',0,\mathbf{j}} \rho \right] \\
& \quad - (1 - \cos 2\theta) \sum_{\mathbf{i}'=\mathbf{j}'+\mathbf{1}_k}^{2_k} \left[(k - \|\mathbf{j}'\|_1) - (\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k)(\cos 2\theta)^{\|\mathbf{i}'\|_1 - \|\mathbf{j}'\|_1 - k} \right] \text{Tr} \left[O_{\mathbf{i}',2,\mathbf{i}}^{\mathbf{j}',0,\mathbf{j}} \rho \right] \left. \right| \\
& \leq \left\| O_{\mathbf{0}_k,1,\mathbf{i}}^{\mathbf{0}_k,0,\mathbf{j}} \right\|_2 + g(k - \|\mathbf{j}'\|_1) \|O\|_2 + (1 - \cos 2\theta)(k+1 - \|\mathbf{j}'\|_1) \left\| O_{\mathbf{0}_k,2,\mathbf{i}}^{\mathbf{0}_k,0,\mathbf{j}} \right\|_2 \\
& \quad + |\cos 2\theta| \left[(2 - \cos 2\theta)^{k - \|\mathbf{j}'\|_1} - 1 \right] \left\| O_{\mathbf{0}_k,2,\mathbf{i}}^{\mathbf{0}_k,0,\mathbf{j}} \right\|_2 + (1 - \cos 2\theta)g(k - \|\mathbf{j}'\|_1) \|O\|_2 \quad (96) \\
& \leq \left[(2 - \cos 2\theta)(k - \|\mathbf{j}'\|_1) \left[(3 - \cos 2\theta)(2 - \cos 2\theta)^{k - \|\mathbf{j}'\|_1 - 1} - 1 \right] \right. \\
& \quad \left. + (2 - \cos 2\theta)^{k - \|\mathbf{j}'\|_1} + (1 - \cos 2\theta)(k+1 - \|\mathbf{j}'\|_1) \right] \|O\|_2 \quad (97) \\
& \leq g(k+1 - \|\mathbf{j}'\| - j'_{k+1}) \|O\|_2. \quad (98)
\end{aligned}$$

Here, InEq. (96) is obtained by using the $h' = k$ case of Eq. (91). InEq. (97) is obtained by using Eqs. (26) and (92). InEq. (98) follows from Eq. (92) and the condition $j'_{k+1} = 0$. Since Eqs. (95) and (98) match the formulation of the $h' = k+1$ case of Eq. (91), we have proved Eq. (91) for general cases.

We proceed from Eq. (90), which can be lower bounded by

$$\geq -h(h-1)\|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0,\gamma^2)} \sum_{\substack{\mathbf{j}=\mathbf{0}_h \\ 2(\|\mathbf{j}\|_1+1)}}^{\mathbf{1}_h} (2\theta)^{\|\mathbf{j}\|_1+1} 2(1+2\theta^2)^{2h-1-\|\mathbf{j}\|_1} \quad (99)$$

$$\geq -2h(h-1)\|O\|_2^2 \mathbb{E}_{\theta \sim \mathcal{N}(0,\gamma^2)} \sum_{t=1}^{\lfloor \frac{h+1}{2} \rfloor} \binom{h}{2t-1} (2\theta)^{2t} \sum_{m=0}^{2h-2t} \binom{2h-2t}{m} (2\theta^2)^m \quad (100)$$

$$= -2h(h-1)\|O\|_2^2 \sum_{t=1}^{\lfloor \frac{h+1}{2} \rfloor} \binom{h}{2t-1} \sum_{m=0}^{2h-2t} \binom{2h-2t}{m} 2^{2t+m} (2t+2m-1)!! \gamma^{2t+2m} \quad (101)$$

$$\geq -\frac{59}{6} h^2 (h-1) \gamma^2 \|O\|_2^2. \quad (102)$$

Here, InEq. (99) is obtained by using $1 \geq \cos 2\theta \geq 1 - 2\theta^2$. InEq. (100) is obtained since the summation $\sum_{\mathbf{j}=\mathbf{0}_h}^{\mathbf{1}_h}$ contains $\binom{h}{2t-1}$ terms such that $\|\mathbf{j}\|_1 = 2t-1$, for all $t \in \{1, \dots, \lfloor \frac{h+1}{2} \rfloor\}$.

Eq. (101) is derived by calculating expectation terms. InEq. (102) is obtained by bounding the summation terms, i.e.

$$\begin{aligned} & \sum_{t=1}^{\lfloor \frac{h+1}{2} \rfloor} \binom{h}{2t-1} \sum_{m=0}^{2h-2t} \binom{2h-2t}{m} 2^{2t+m} (2t+2m-1)!! \gamma^{2t+2m} \\ & \leq \sum_{t=1}^{\lfloor \frac{h+1}{2} \rfloor} \frac{h(h-1)^{2t-2}}{2^{t-1}(t-1)!(2t-1)!!} \sum_{m=0}^{2h-2t} \frac{(2h-2t)^m}{m!} 2^{2t+m} \\ & \quad \cdot (2t-1)!!(2t+1)(2t+3) \cdots (2t+2m-1) \gamma^{2t+2m} \end{aligned} \quad (103)$$

$$\leq \sum_{t=1}^{\lfloor \frac{h+1}{2} \rfloor} h(h-1)^{2t-2} 2^{t+1} \gamma^{2t} \sum_{m=0}^{2h-2t} (2h-2)^m 2^m (2h)^m \gamma^{2m} \quad (104)$$

$$\leq 4h\gamma^2 \sum_{t=1}^{\lfloor \frac{h+1}{2} \rfloor} [h(h-1)^2 \gamma^2]^{t-1} \sum_{m=0}^{2h-2t} (2h^3 \gamma^2)^m \quad (105)$$

$$\leq 4h\gamma^2 \frac{16}{15} \times \frac{8}{7} \leq \frac{59}{12} h\gamma^2. \quad (106)$$

Here, InEq (103) follows from $t \geq 1$ and $(2t-1)! = 2^t t! (2t-1)!!$. InEq. (104) is obtained by using $t \geq 1$ and

$$(2t+2k-1)(2t+2m-2k+1) \leq (m+2t)^2 \leq (2h)^2, \forall k \in \{1, \dots, m\}.$$

InEq. (105) is derived by using $8h(h-1) \leq 2h^3$ for the integer h . InEq. (106) is obtained by calculating geometric sequences with the condition $\gamma^2 \leq \frac{1}{16h^3}$.

Finally, we prove Eq. (77). The left side of Eq. (77) could be lower bounded by

$$\begin{aligned} & = \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \left(\sum_{\|j'\|_1=1} \sum_{i'=j'+\mathbf{1}_h}^{2h} (\cos 2\theta)^{\|i'\|_1-h-1} \text{Tr} [O_{i'}^{j'} \rho] \right) \\ & \quad \cdot \left(\sum_{\substack{j=\mathbf{0}_h \\ \|j\|_1 \geq 2, 2|(\|j\|_1-1)}}^{1_h} \sum_{i=j+\mathbf{1}_h}^{2h} \|j\|_1 (\cos 2\theta)^{\|i\|_1-\|j\|_1-h+2} (\sin 2\theta)^{\|j\|_1-1} \text{Tr} [O_i^j \rho] \right) \end{aligned} \quad (107)$$

$$\begin{aligned} & \geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} \sum_{\|j'\|_1=1} (2 - \cos 2\theta)^{h-1} \|O\|_2 \sum_{\substack{j=\mathbf{0}_h \\ \|j\|_1 \geq 2, 2|(\|j\|_1-1)}}^{1_h} \|j\|_1 (\sin 2\theta)^{\|j\|_1-1} (\cos 2\theta)^2 \\ & \quad \left\{ \left| \sum_{i=j+\mathbf{1}_h}^{2h} [1 - (\cos 2\theta)^{\|i\|_1-\|j\|_1-h}] \text{Tr} [O_i^j \rho] \right| + \left| \sum_{i=j+\mathbf{1}_h}^{2h} \text{Tr} [O_i^j \rho] \right| \right\} \end{aligned} \quad (108)$$

$$\begin{aligned} & \geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} h(2 - \cos 2\theta)^{h-1} \|O\|_2 \\ & \quad \cdot \sum_{\substack{j=\mathbf{0}_h \\ \|j\|_1 \geq 2, 2|(\|j\|_1-1)}}^{1_h} \|j\|_1 (\sin 2\theta)^{\|j\|_1-1} (\cos 2\theta)^2 (2 - \cos 2\theta)^{h-\|j\|_1} \|O\|_2 \end{aligned} \quad (109)$$

$$\geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} h \|O\|_2^2 \sum_{\substack{j=\mathbf{0}_h \\ \|j\|_1 \geq 2, 2|(\|j\|_1-1)}}^{1_h} \|j\|_1 (2\theta)^{\|j\|_1-1} (1 + 2\theta^2)^{2h-1-\|j\|_1} \quad (110)$$

$$\geq - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} h \|O\|_2^2 \sum_{t=1}^{\lfloor \frac{h-1}{2} \rfloor} \binom{h}{2t+1} (2t+1)(2\theta)^{2t} (1 + 2\theta^2)^{2h-2-2t}. \quad (111)$$

Eq. (107) is obtained by noticing that the expectation of $\sin^a 2\theta \cos^b 2\theta$ equals to zero, if a is odd. InEq. (108) follows from the derivation (87-89). InEq. (109) is obtained by using the $h' = h$, $\|j\|_1 = 1$ case of Eq. (44). InEq. (110) follows from $1 \geq \cos 2\theta \geq 1 - 2\theta^2$ and $(\sin 2\theta)^2 \leq (2\theta)^2$. InEq. (111) is obtained since the summation $\sum_{j=0_h}^{1_h}$ contains $\binom{h}{2t+1}$ terms such that $\|j\|_1 = 2t + 1$, for all $t \in \{1, \dots, \lfloor \frac{h-1}{2} \rfloor\}$. We further bound InEq. (111) by

$$\begin{aligned} &= - \mathbb{E}_{\theta \sim \mathcal{N}(0, \gamma^2)} h \|O\|_2^2 \sum_{t=1}^{\lfloor \frac{h-1}{2} \rfloor} \binom{h}{2t+1} (2t+1) (2\theta)^{2t} \sum_{m=0}^{2h-2-2t} \binom{2h-2-2t}{m} (2\theta^2)^m \\ &\geq -h \|O\|_2^2 \sum_{t=1}^{\lfloor \frac{h-1}{2} \rfloor} \binom{h}{2t+1} (2t+1) \sum_{m=0}^{2h-2-2t} \binom{2h-2-2t}{m} 2^{2t+m} (2t+2m-1)!! \gamma^{2t+2m} \end{aligned} \quad (112)$$

$$\begin{aligned} &\geq -h \|O\|_2^2 \sum_{t=1}^{\lfloor \frac{h-1}{2} \rfloor} \frac{h(h-1)(h-2)(h-3)^{2t-2}}{(2t)!} 2^{2t} \gamma^{2t} \\ &\quad \cdot \sum_{m=0}^{2h-2-2t} (2h-2-2t)^m 2^m (2t+2m-1)!! \gamma^{2m} \end{aligned} \quad (113)$$

$$\begin{aligned} &= -h^2 (h-1)(h-2) \|O\|_2^2 \sum_{t=1}^{\lfloor \frac{h-1}{2} \rfloor} \frac{(h-3)^{2t-2}}{2^t t! (2t-1)!!} 2^{2t} \gamma^{2t} \\ &\quad \cdot \sum_{m=0}^{2h-2-2t} (2h-2-2t)^m 2^m (2t+2m-1)!! \gamma^{2m} \end{aligned} \quad (114)$$

$$= -\frac{5}{2} h^2 (h-1)(h-2) \gamma^2 \|O\|_2^2. \quad (115)$$

Here, InEq. (112) is obtained by calculating expectation terms. InEq. (113) is derived by using $t \geq 1$. Eq. (114) follows from $(2t)! = 2^t t! (2t-1)!!$. Eq. (115) is obtained by bounding the summation terms, i.e.

$$\begin{aligned} &\sum_{t=1}^{\lfloor \frac{h-1}{2} \rfloor} \frac{(h-3)^{2t-2}}{2^t t! (2t-1)!!} 2^{2t} \gamma^{2t} \sum_{m=0}^{2h-2-2t} (2h-2-2t)^m 2^m (2t+2m-1)!! \gamma^{2m} \\ &\leq \sum_{t=1}^{\lfloor \frac{h-1}{2} \rfloor} \frac{(h-3)^{2t-2}}{2^{t-1}} 2^t \gamma^{2t} \sum_{m=0}^{2h-2-2t} (2h-4)^m 2^m (2h-2)^m \gamma^{2m} \end{aligned} \quad (116)$$

$$\leq 2\gamma^2 \sum_{t=1}^{\lfloor \frac{h-1}{2} \rfloor} [(h-3)^2 \gamma^2]^{t-1} \sum_{m=0}^{2h-2-2t} (h^3 \gamma^2)^m \quad (117)$$

$$\leq 2\gamma^2 \left(\frac{16}{15}\right)^2 \leq \frac{5}{2} \gamma^2. \quad (118)$$

Here, InEq. (116) is obtained by using $t! \geq 2^{t-1}$, $\forall t \geq 1$ and

$$(2t+2k-1)(2t+2m-2k+1) \leq (m+2t)^2 \leq (2h-2)^2, \forall k \in \{1, \dots, m\}.$$

InEq. (117) follows from $h^3 \geq 8(h-1)(h-2)$ for integer h . InEq. (118) is obtained by calculating geometric sequences with the condition $\gamma^2 \leq \frac{1}{16h^3}$. Thus, we have proved Eq. (77). \square

C Proof of Theorem 4.1

Proof. Denote by $I_S := \{m | i_m \neq 0, m \in [N]\}$ the set of qubits where the observable acts non-trivially. First, we notice that the norm of the whole gradient is lower bounded by that of particle

derivatives summed over a part of parameters, i.e.

$$\mathbb{E}_{\boldsymbol{\theta}} \|\nabla_{\boldsymbol{\theta}} f\|^2 \geq \sum_{q=1}^L \sum_{n \in I_S} \mathbb{E}_{\boldsymbol{\theta}} \left(\frac{\partial f}{\partial \theta_{q,n}} \right)^2. \quad (119)$$

Thus, we could obtain the formulation in the theorem if

$$\mathbb{E}_{\boldsymbol{\theta}} \left(\frac{\partial f}{\partial \theta_{q,n}} \right)^2 \geq \frac{1}{S^{S+1}(L+2)^{S+1}} \text{Tr} [\sigma_j \rho_{\text{in}}]^2, \quad (120)$$

holds for any $q \in \{1, \dots, L\}$ and $n \in I_S$.

Now we begin to prove Eq. (120). Our main idea is to integrate the square of the partial derivative of f with respect to $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{L+2})$ by using Lemma B.1 and Lemma B.2.

We introduce several notations for convenience. Denote the variance $\gamma^2 = \frac{1}{4S(L+2)}$. Denote the ℓ -th single qubit rotations and CZ layer as $R_{\ell}(\boldsymbol{\theta}_{\ell})$ and CZ_{ℓ} , respectively, where

$$R_{\ell}(\boldsymbol{\theta}_{\ell}) = e^{-i\theta_{\ell,1}G_{\ell,1}} \otimes e^{-i\theta_{\ell,2}G_{\ell,2}} \otimes \dots \otimes e^{-i\theta_{\ell,N}G_{\ell,N}}, \quad (121)$$

and $G_{\ell,j}$ is the Hamiltonian corresponding to the parameter $\theta_{\ell,j}$. Denote by ρ_k the state after the k -th layer, $\forall k \in \{0, 1, \dots, 2L+2\}$,

$$\rho_k := \begin{cases} \left(\prod_{i=\frac{k}{2}}^1 \text{CZ}_i R_i(\boldsymbol{\theta}_i) \right) \rho_{\text{in}} \left(\prod_{i=1}^{\frac{k}{2}} R_i(\boldsymbol{\theta}_i)^{\dagger} \text{CZ}_i^{\dagger} \right) & (k = 2\ell \leq 2L), \\ R_{\frac{k+1}{2}}(\boldsymbol{\theta}_{\frac{k+1}{2}}) \rho_{k-1} R_{\frac{k+1}{2}}(\boldsymbol{\theta}_{\frac{k+1}{2}})^{\dagger} & (k = 2\ell + 1 \leq 2L + 1), \\ R_{L+2}(\boldsymbol{\theta}_{L+2}) \rho_{k-1} R_{L+2}(\boldsymbol{\theta}_{L+2})^{\dagger} & (k = 2L + 2). \end{cases} \quad (122)$$

Thus, ρ_k is parameterized by $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p\}$, where $p = \ell$ if $k = 2\ell \leq 2L$, $p = \ell + 1$ if $k = 2\ell + 1 \leq 2L + 1$, and $p = L + 2$ if $k = 2L + 2$.

Next, rewrite the formulation of Eq. (120) in detail:

$$\begin{aligned} & \mathbb{E}_{\boldsymbol{\theta}_1} \dots \mathbb{E}_{\boldsymbol{\theta}_{L+2}} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_i V(\boldsymbol{\theta}) \rho_{\text{in}} V(\boldsymbol{\theta})^{\dagger}] \right)^2 \\ &= \mathbb{E}_{\boldsymbol{\theta}_1} \dots \mathbb{E}_{\boldsymbol{\theta}_{L+2}} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_i \rho_{2L+2}] \right)^2 \end{aligned} \quad (123)$$

$$\geq [4\gamma^2(1-4\gamma^2)]^{S_1} (1-4\gamma^2)^{S_3} \mathbb{E}_{\boldsymbol{\theta}_1} \dots \mathbb{E}_{\boldsymbol{\theta}_{L+1}} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_{\mathbf{3|i;1}} \rho_{2L+1}] \right)^2 \quad (124)$$

$$\geq [4\gamma^2(1-4\gamma^2)]^{S_1+S_2} (1-4\gamma^2)^{S_1+2S_3} \mathbb{E}_{\boldsymbol{\theta}_1} \dots \mathbb{E}_{\boldsymbol{\theta}_L} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_{\mathbf{3|i}} \rho_{2L}] \right)^2 \quad (125)$$

$$\geq [4\gamma^2(1-4\gamma^2)]^S (1-4\gamma^2)^S \mathbb{E}_{\boldsymbol{\theta}_1} \dots \mathbb{E}_{\boldsymbol{\theta}_L} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_{\mathbf{3|i}} \rho_{2L}] \right)^2, \quad (126)$$

where $\mathbf{3|i}$ denotes the index by replacing non-zero elements of $\mathbf{i} = (i_1, \dots, i_N)$ with 3 and $\mathbf{3|i;1}$ denotes the index by replacing non-zero elements of $\mathbf{i} = (i_1, \dots, i_N)$ with 3 if the original value is 1. We refer to S_1 , S_2 , and S_3 as the number of 1, 2, and 3 in the index \mathbf{i} , respectively. Eq. (123) is obtained by using the notation ρ_{2L+2} defined in (122). We obtain Eqs. (124) and (125) by using Lemma B.1 for the R_Y and R_X gate case, respectively. InEq. (126) follows from $S = S_1 + S_2 + S_3$.

Then, we proceed from Eq. (126) and take the expectation for parameters in $(\boldsymbol{\theta}_L, \dots, \boldsymbol{\theta}_{q+1})$.

$$\text{Eq. (126)} = [2\gamma(1-4\gamma^2)]^{2S} \mathbb{E}_{\boldsymbol{\theta}_1} \dots \mathbb{E}_{\boldsymbol{\theta}_L} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_{\mathbf{3|i}} \text{CZ}_L R_L(\boldsymbol{\theta}_L) \rho_{2L-2} R_L(\boldsymbol{\theta}_L)^{\dagger} \text{CZ}_L^{\dagger}] \right)^2 \quad (127)$$

$$= [2\gamma(1-4\gamma^2)]^{2S} \mathbb{E}_{\boldsymbol{\theta}_1} \dots \mathbb{E}_{\boldsymbol{\theta}_L} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_{\mathbf{3|i}} R_L(\boldsymbol{\theta}_L) \rho_{2L-2} R_L(\boldsymbol{\theta}_L)^{\dagger}] \right)^2 \quad (128)$$

$$\geq [2\gamma(1-4\gamma^2)]^{2S} (1-4\gamma^2)^S \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{L-1}} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_{2L-2}] \right)^2 \quad (129)$$

$$\geq [2\gamma(1-4\gamma^2)]^{2S} (1-4\gamma^2)^{(L-q)S} \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_q} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_{2q}] \right)^2. \quad (130)$$

Eq. (127) follows from the definition of ρ_{2L} (122). Eq. (128) is obtained since

$$\text{CZ}(\sigma_j \otimes \sigma_k) \text{CZ}^\dagger = \sigma_j \otimes \sigma_k, \forall j, k \in \{0, 3\}.$$

InEq. (129) is derived by using the Lemma B.1. We repeat the derivation in Eqs. (127-129) inductively for parameters $(\theta_L, \dots, \theta_{q+1})$, which yields InEq. (130).

Next, we consider the expectation with respect to θ_q . We have

$$\begin{aligned} \text{Eq. (130)} &= [2\gamma(1-4\gamma^2)]^{2S} (1-4\gamma^2)^{(L-q)S} \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_q} \left(\frac{\partial}{\partial \theta_{q,n}} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_{2q-1}] \right)^2 \\ &\geq [2\gamma(1-4\gamma^2)]^{2S} (1-4\gamma^2)^{(L-q)S} (1-4\gamma^2)^{S-1} [4\gamma^2(1-4\gamma^2)]^4 \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{q-1}} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_{2q-2}]^2, \end{aligned} \quad (131)$$

where expectations with respect to parameters $\{\theta_{q,j}\}_{j \in I_S, j \neq n}$ are calculated via Lemma B.1 and the expectation with respect to $\theta_{q,n}$ is calculated via Lemma B.2.

Finally we proceed from Eq. (131) and take the expectation for parameters in $(\theta_{q-1}, \dots, \theta_1)$. We have

$$\mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{q-1}} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_{2q-2}]^2 = \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{q-1}} \text{Tr} [\sigma_{\mathbf{3}|i} \text{CZ}_{q-1} R_{q-1}(\theta_{q-1}) \rho_{2q-4} R_{q-1}(\theta_{q-1})^\dagger \text{CZ}_{q-1}^\dagger]^2 \quad (132)$$

$$= \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{q-1}} \text{Tr} [\sigma_{\mathbf{3}|i} R_{q-1}(\theta_{q-1}) \rho_{2q-4} R_{q-1}(\theta_{q-1})^\dagger]^2 \quad (133)$$

$$\geq (1-4\gamma^2)^S \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{q-2}} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_{2q-4}]^2 \quad (134)$$

$$\geq (1-4\gamma^2)^{(q-1)S} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_0]^2. \quad (135)$$

Eq. (132) is derived by using the definition of ρ_{2q-2} . Eq. (133) is obtained since

$$\text{CZ}(\sigma_j \otimes \sigma_k) \text{CZ}^\dagger = \sigma_j \otimes \sigma_k, \forall j, k \in \{0, 3\}.$$

InEq. (134) is derived by using Lemma B.1. We repeat the derivation in Eqs. (132-134) inductively for parameters $(\theta_{q-1}, \dots, \theta_1)$, which yields InEq. (135). Employing Eq. (135) to Eq. (131) yields

$$\text{Eq. (130)} \geq 4(4\gamma^2)^{S+1} (1-4\gamma^2)^{S(L+2)} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_0]^2 \quad (136)$$

$$= 4 \left(\frac{1}{S(L+2)} \right)^{S+1} \left(1 - \frac{1}{S(L+2)} \right)^{S(L+2)} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_0]^2 \quad (137)$$

$$\geq 4 \left(\frac{1}{S(L+2)} \right)^{S+1} \left(1 - \frac{1}{2} \right)^2 \text{Tr} [\sigma_{\mathbf{3}|i} \rho_0]^2 \quad (138)$$

$$= \frac{1}{S^{S+1}(L+2)^{S+1}} \text{Tr} [\sigma_{\mathbf{3}|i} \rho_0]^2. \quad (139)$$

Eq. (137) is derived by using the condition $\gamma^2 = \frac{1}{4S(L+2)}$. Eq. (138) is obtained by noticing that function $g(x) = (1 - \frac{1}{x})^x$ is monotonically increasing when $x \geq 2$. Thus, we have proved Eq. (120). \square

D Proof of Theorem 4.2

Proof. To begin with, we define several notations for convenience. Denote by ρ_j the state after the j -th parameterized operator, i.e.

$$\rho_j(\theta_1, \dots, \theta_j) = \left(\prod_{i=j}^1 V_i(\theta_i) \right) \rho_{\text{in}} \left(\prod_{i=1}^j V_i(\theta_i)^\dagger \right). \quad (140)$$

Denote by O_j the observable, i.e.

$$O_j = V_j(0)^\dagger \cdots V_L(0)^\dagger O V_L(0) \cdots V_j(0), \quad \forall j \in \{1, \dots, L\}. \quad (141)$$

Now we begin to prove the Theorem. First, we remark that $\forall j \in [L]$, the $a_j \neq 1$ case can be converted to the $a_j = 1$ case by using the transformation

$$\theta'_j = \frac{\theta_j}{a_j},$$

where the variance of the new and the old parameter satisfies

$$\text{Var}[\theta'_j] = \frac{1}{a_j^2} \text{Var}[\theta_j].$$

In the following proof, we assume that $a_j = 1, \forall j \in [L]$. By using the parameter-shift rule, $\frac{\partial f}{\partial \theta_\ell}$ could be written as the linear sum of $2h$ expectations on the observable O with coefficients ± 1 . Then for the case $\ell \leq L-1$, we have

$$\mathbb{E}_{\boldsymbol{\theta}} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 = \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_L} \left(\frac{\partial}{\partial \theta_\ell} \text{Tr} [O V_L(\theta_L) \rho_{L-1} V_L(\theta_L)^\dagger] \right)^2 \quad (142)$$

$$\geq \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{L-1}} \left(\frac{\partial}{\partial \theta_\ell} \text{Tr} [O V_L(0) \rho_{L-1} V_L(0)^\dagger] \right)^2 - [12h_L(h_L - 1) + 4] 4h_L^2 \gamma_L^2 \|O\|_2^2 \quad (143)$$

$$= \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{L-1}} \left(\frac{\partial}{\partial \theta_\ell} \text{Tr} [O_L \rho_{L-1}] \right)^2 - [12h_L(h_L - 1) + 4] 4h_L^2 \gamma_L^2 \|O\|_2^2 \quad (144)$$

$$= \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{L-1}} \left(\frac{\partial f}{\partial \theta_\ell}(\theta_1, \dots, \theta_{L-1}, 0) \right)^2 - [12h_L(h_L - 1) + 4] 4h_L^2 \gamma_L^2 \|O\|_2^2, \quad (145)$$

where Eq. (142) follows from the definition of ρ_j in Eq. (140). InEq. (143) is obtained by using Lemma B.3, where $\|c\|_1 = 2h$. Eq. (144) follows from the definition of O_j in Eq. (141). Eq. (145) follows from the formulation $f(\boldsymbol{\theta}) = \text{Tr}[O\rho(\boldsymbol{\theta})]$. By proceeding the derivation (142-145) for $L - \ell$ times, we have

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 &\geq \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_\ell} \left(\frac{\partial f}{\partial \theta_\ell}(\theta_1, \dots, \theta_\ell, 0, \dots, 0) \right)^2 - \sum_{j=\ell+1}^L [12h_j(h_j - 1) + 4] 4h_j^2 \gamma_j^2 \|O\|_2^2 \\ &= \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_\ell} \left(\frac{\partial}{\partial \theta_\ell} \text{Tr} [O_{\ell+1} V_\ell(\theta_\ell) \rho_{\ell-1} V_\ell(\theta_\ell)^\dagger] \right)^2 - \sum_{j=\ell+1}^L [12h_j(h_j - 1) + 4] 4h_j^2 \gamma_j^2 \|O\|_2^2 \quad (146) \end{aligned}$$

$$\begin{aligned} &\geq \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{\ell-1}} (1 - 4\gamma_\ell^2) \left(\frac{\partial}{\partial \theta_\ell} \text{Tr} [O_{\ell+1} V_\ell(\theta_\ell) \rho_{\ell-1} V_\ell(\theta_\ell)^\dagger] \right)^2 \Big|_{\theta_\ell=0} - 96h_\ell^2(h_\ell - 1)\gamma_\ell^2 \|O\|_2^2 \\ &\quad - 20h_\ell^2(h_\ell - 1)(h_\ell - 2)\gamma_\ell^2 \|O\|_2^2 - \sum_{j=\ell+1}^L [12h_j(h_j - 1) + 4] 4h_j^2 \gamma_j^2 \|O\|_2^2 \quad (147) \end{aligned}$$

$$\begin{aligned} &\geq \mathbb{E}_{\theta_1} \cdots \mathbb{E}_{\theta_{\ell-1}} \left(\frac{\partial f}{\partial \theta_\ell}(\theta_1, \dots, \theta_{\ell-1}, 0, 0, \dots, 0) \right)^2 - 4\gamma_\ell^2 (2h_\ell)^2 \|O\|_2^2 - 96h_\ell^2(h_\ell - 1)\gamma_\ell^2 \|O\|_2^2 \\ &\quad - 20h_\ell^2(h_\ell - 1)(h_\ell - 2)\gamma_\ell^2 \|O\|_2^2 - \sum_{j=\ell+1}^L [12h_j(h_j - 1) + 4] 4h_j^2 \gamma_j^2 \|O\|_2^2, \quad (148) \end{aligned}$$

where Eq. (146) follows from definitions ρ_j (140) and O_j (141). InEq. (147) is derived by using Lemma B.4. InEq. (148) follows from the parameter-shift rule. We proceed from InEq. (148) by employing the derivation (142-145) for parameters $(\theta_{\ell-1}, \dots, \theta_1)$, which yields

$$\mathbb{E}_{\boldsymbol{\theta}} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2$$

$$\begin{aligned}
&\geq \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\mathbf{o}} - \sum_{j=1}^{\ell-1} 16h_j^2 [3h_j(h_j-1)+1] \gamma_j^2 \|O\|_2^2 - \sum_{j=\ell+1}^L 16h_j^2 [3h_j(h_j-1)+1] \gamma_j^2 \|O\|_2^2 \\
&\quad - 16h_\ell^2 \gamma_\ell^2 \|O\|_2^2 - 96h_\ell^2 (h_\ell-1) \gamma_\ell^2 \|O\|_2^2 - 20h_\ell^2 (h_\ell-1)(h_\ell-2) \gamma_\ell^2 \|O\|_2^2 \tag{149}
\end{aligned}$$

$$\begin{aligned}
&\geq \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\mathbf{o}} - \sum_{j=1}^L 16h_j^2 [3h_j(h_j-1)+1] \gamma_j^2 \|O\|_2^2 \\
&\geq (1-\epsilon) \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\mathbf{o}}. \tag{150}
\end{aligned}$$

InEq. (149) is obtained by using Lemma B.3, where $\|c\|_1 = 2h$. InEq. (150) follows from the condition $\gamma_j^2 \leq \frac{a_j^2 \epsilon}{16h_j^2(3h_j(h_j-1)+1)L\|O\|_2^2} \left(\frac{\partial f}{\partial \theta_\ell} \right)^2 \Big|_{\theta=\mathbf{o}}$ and $a_j = 1, \forall j \in [L]$. Thus, we have proved the theorem. \square