

# Theory of overparametrization in quantum neural networks

Martín Larocca,<sup>1,2,\*</sup> Nathan Ju,<sup>1,\*</sup> Diego García-Martín,<sup>1,3,4</sup> Patrick J. Coles,<sup>1</sup> and M. Cerezo<sup>5,1,6</sup>

<sup>1</sup>*Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA*

<sup>2</sup>*Departamento de Física “J. J. Giambiagi” and IFIBA, FCEyN, Universidad de Buenos Aires, 1428 Buenos Aires, Argentina*

<sup>3</sup>*Barcelona Supercomputing Center, Barcelona 08034, Spain*

<sup>4</sup>*Instituto de Física Teórica, UAM-CSIC, Madrid 28049, Spain*

<sup>5</sup>*Information Sciences, Los Alamos National Laboratory, Los Alamos, NM 87545, USA*

<sup>6</sup>*Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA*

The prospect of achieving quantum advantage with Quantum Neural Networks (QNNs) is exciting. Understanding how QNN properties (e.g., the number of parameters  $M$ ) affect the loss landscape is crucial to the design of scalable QNN architectures. Here, we rigorously analyze the overparametrization phenomenon in QNNs with periodic structure. We define overparametrization as the regime where the QNN has more than a critical number of parameters  $M_c$  that allows it to explore all relevant directions in state space. Our main results show that the dimension of the Lie algebra obtained from the generators of the QNN is an upper bound for  $M_c$ , and for the maximal rank that the quantum Fisher information and Hessian matrices can reach. Underparametrized QNNs have spurious local minima in the loss landscape that start disappearing when  $M \geq M_c$ . Thus, the overparametrization onset corresponds to a computational phase transition where the QNN trainability is greatly improved by a more favorable landscape. We then connect the notion of overparametrization to the QNN capacity, so that when a QNN is overparametrized, its capacity achieves its maximum possible value. We run numerical simulations for eigensolver, compilation, and autoencoding applications to showcase the overparametrization computational phase transition. We note that our results also apply to variational quantum algorithms and quantum optimal control.

## I. INTRODUCTION

The development of Neural Networks (NNs) and Machine Learning (ML) is one of the greatest scientific revolutions of the twentieth century. Traditionally, computers were explicitly programmed to solve a task, so that a user-created code would take an input and produce a desired output. In ML, however, one follows a fundamentally different approach. Here, a computer is trained to learn from data, with the goal that it can accurately solve the problem when presented with new and previously unseen cases [1]. Currently, ML is used in virtually all areas of science, with applications such as drug discovery [2], new materials exploration [3], and self-driving cars [4].

Despite their tremendous success, training NNs is a difficult task that has even been shown to be NP-hard [5–7]. Thus, finding ways to improve the NNs trainability and generalization capacity has always been a coveted goal. Towards this end, one of the most surprising phenomena in ML is that of overparametrization. Here, one trains a NN with a capacity larger than that which is necessary to represent the distribution of the training data [8]. Usually, this implies having a number of parameters in the NN

that is much larger than the number of training points [9]. Naively, one could expect that a model with a large capacity would have training difficulties and also have overfitting (poor generalization). However, overparametrizing a NN can improve its performance and reduce its training and generalization errors [9–13], and even lead to provable convergence results [14, 15].

The advent of quantum computers [16, 17] has brought a tremendous interest in using these devices for data science. Here, researchers have embedded ML into the framework of quantum mechanics, with the new, generalized theory being called Quantum Machine Learning (QML) [18–20]. With QML, the end goal is not formal generalization but rather to exploit entanglement and superposition to achieve a quantum advantage [21–24], that is, to solve the problem more efficiently than any classical algorithm run on a classical supercomputer.

Naturally, as a generalized theory, QML has the potential to exhibit many of the issues and phenomena exhibited by (classical) ML. For instance, like the classical case, it has been shown that training QML models is NP-hard [25]. Since a QML model may consist of a data embedding followed by a parametrized quantum circuit that is often called a Quantum Neural Network (QNN), its training requires optimizing the QNN’s parameters [20, 26–29]. Recently, much effort has gone towards developing so-called Quantum Landscape Theory [30], which studies the prop-

---

\* The first two authors contributed equally to this work.

erties of QML loss function landscapes. Indeed, there are results analyzing the presence of sub-optimal local minima [31, 32], the existence of barren plateaus [33–44], and how quantum noise affects the loss landscape [45–49].

Similar to classical NNs, some examples of QNNs that exhibit overparametrization have been constructed [32, 50–55]. Some of these works have heuristically shown that increasing the number of parameters in the QNN can improve its trainability and lead to faster convergence. However, there is still need for a detailed theoretical analysis of this overparametrization phenomenon. Understanding overparametrization is crucial for Quantum Landscape Theory and for engineering QNNs to enhance their trainability.

In this work we provide a theoretical framework for the overparametrization of QNNs. Our main results indicate that, for a general type of periodic-structured QNNs, one can reach an overparametrized regime by increasing the number of parameters past some threshold critical value  $M_c$  (see Fig. 1(a)). Moreover, we prove that  $M_c$  is related to the dimension of the Dynamical Lie Algebra (DLA) [56, 57] associated with the generators of the QNN.

We here define overparametrization as the QNN having enough parameters so that the quantum Fisher information matrix saturates its achievable rank. In this case, one can explore all relevant directions in the state space by varying the QNN parameters. We then relate this notion of overparametrization to different measures of the model’s capacity [24, 58], so that a model is overparametrized when its capacity is saturated. Then, as shown in Fig. 1(b), our results have direct implications in understanding why overparametrization can improve the model’s trainability, as the overparametrization onset corresponds to a computational phase transition [52]. We verify our theoretical results by performing numerical simulations. In all cases, we find the predicted computational phase transition, where the success probability of solving the optimization problem is greatly increased after a critical number of parameters.

These results provide theoretical grounds for recent observations of the overparametrization phenomenon in QML [50, 52, 59]. Moreover, our theorems have direct consequences for the field of quantum optimal control [60–63].

## II. RESULTS

### A. Quantum Neural Networks

Quantum Neural Networks (QNNs) [18–20] employ parametrized quantum circuits to allow for task-oriented programming of quantum computers. Here, one encodes the problem of interest in a loss function  $\mathcal{L}(\theta)$ , whose min-

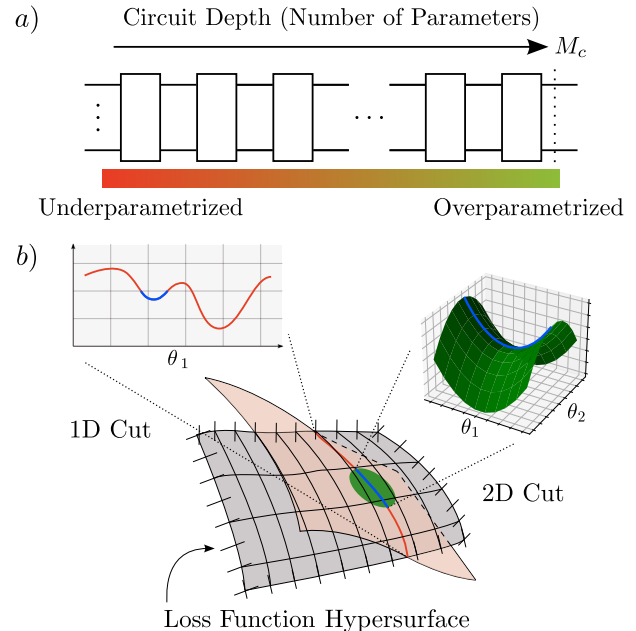


Figure 1. **Overparametrization in quantum neural networks (QNNs).** a) Quantum circuit description of the QNN. By having a low (high) number of parameters one is not able (is able) to explore all relevant directions in the Hilbert space, and thus the QNN is underparametrized (overparametrized). b) The gray surface corresponds to the unconstrained loss function landscape. An underparametrized QNN explores a low dimensional cut of the loss function (1D cut over the red lines). Here, the optimizer can get trapped in spurious local minima (blue segment) that negatively impact the parameter optimization. By increasing the number of parameters past some threshold  $M_c$ , one can explore a higher dimensional cut of the landscape (2D cut over the green region). As shown, some previous spurious local minima correspond to saddle points (blue segment), and the optimizer can escape the false trap.

ima correspond to the task’s solution. Using data from a training dataset  $\mathcal{S}$  composed of quantum states  $|\psi_\mu\rangle \in \mathcal{S}$ , one optimizes the QNN parameters to solve the problem

$$\theta_* = \arg \min_{\theta} \mathcal{L}(\theta). \quad (1)$$

Measurements on a quantum computer assist in estimating the loss function (or its gradients), while a classical optimizer is used to update the parameters and solve Eq. (1). This hybrid scheme allows the QML model to access the exponentially large dimension of the Hilbert space, with the hope that if the whole process is hard to classically simulate, then a quantum advantage could be achieved [22, 64, 65].

We consider the case when the QNN is a parametrized

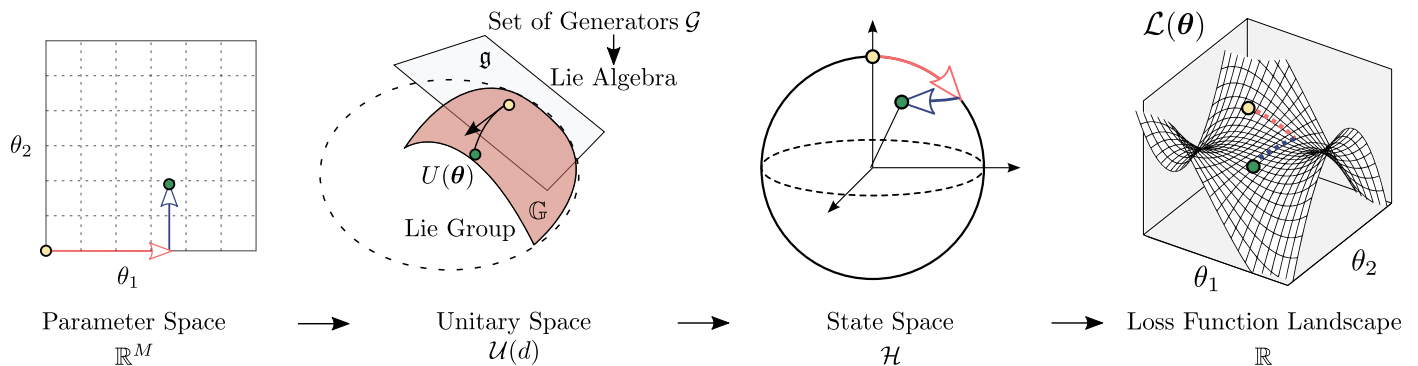


Figure 2. **Relevant mathematical spaces for QNNs.** QNNs employ a set of  $M$  trainable parameters  $\theta \in \mathbb{R}^M$ , which live in parameter space. The QNN itself is represented by a  $d$ -dimensional unitary  $U(\theta)$ , which lives in unitary space. An exemplary form of  $U(\theta)$  is that of Eq. (2), where the set of unitaries depends on the dynamical Lie algebra  $\mathfrak{g}$ , which in turn is obtained from the set of generators  $\mathcal{G}$  in Definition 1. In most QML applications, the QNN acts on an input state  $|\psi_\mu\rangle$  from a training set. Thus, the set of reachable unitaries of the ansatz translates into a set of reachable states in the Hilbert space  $\mathcal{H}$ . Finally, by performing measurements on a quantum computer one estimates the loss function or its gradient. This information is then used to navigate through the loss function landscape  $\mathcal{L}(\theta)$ . Understanding the connections between these mathematical spaces is fundamental for the theory of quantum landscapes.

quantum circuit  $U(\theta)$  that acts on the quantum states in the training set as  $U(\theta)|\psi_\mu\rangle$ . Here,  $U(\theta)$  has an  $L$ -layered periodic structure of the form

$$U(\theta) = \prod_{l=1}^L U_l(\theta_l), \quad U_l(\theta_l) = \prod_{k=1}^K e^{-i\theta_{lk} H_k}, \quad (2)$$

where the index  $l$  indicates the layer, and the index  $k$  spans the traceless Hermitian operators  $H_k$  that generate the unitaries in the ansatz. Moreover,  $\theta_l = (\theta_{l1}, \dots, \theta_{lK})$  are the parameters in a single layer, and  $\theta = \{\theta_1, \dots, \theta_L\}$  denotes the set of  $M = K \cdot L$  trainable parameters in the QNN.

As discussed in [37], Eq. (2) contains as special cases the hardware-efficient ansatz [66], quantum alternating operator ansatz (QAOA) [67, 68], Adaptive QAOA [69], Hamiltonian Variational Ansatz (HVA) [70], and Quantum Optimal Control Ansatz [71], among others [54]. As we discuss in the Methods section, due to the close connection between training a parametrized quantum circuit and the control pulses used to evolve a quantum state in a quantum optimal control protocol, all the results derived hereon can be directly applied to the field of quantum optimal control.

## B. Quantum Landscape Theory

The usefulness of a QNN for a given task hinges on several factors. First and foremost, it is crucial that a solution (or a good approximation to it) actually exists within the ansatz. Then, even if that solution exists, one must be

able to find the associated optimal parameters. The goals of Quantum Landscape Theory are to study properties of the QML loss landscape, how they emerge, and how they affect the optimization process. Here we recall the basic theoretical framework of Quantum Landscape Theory.

First, we note that there are several aspects of the problem that play a key role in how the loss function landscape arises. Specifically, as shown in Fig. 2,  $\theta$  is a vector in  $\mathbb{R}^M$ , and each set of parameters corresponds to a unitary  $U(\theta)$  in the unitary group  $U(d)$  of degree  $d$ . Then, one applies the unitary  $U(\theta)$  to an  $n$ -qubit input state  $|\psi_\mu\rangle$  (from the dataset  $\mathcal{S}$ ) in a Hilbert space  $\mathcal{H}$  of dimension  $d = 2^n$ . Finally, the loss function value  $\mathcal{L}(\theta) \in \mathbb{R}$  is determined by performing measurements over the states  $U(\theta)|\psi_\mu\rangle$ . In this sense, the action of the QML model arises from the composition of the following three maps:

$$\mathbb{R}^M \rightarrow U(d) \rightarrow \mathcal{H} \rightarrow \mathbb{R}. \quad (3)$$

Since the landscape is essentially the collection of values obtained at the end of the maps in Eq. (3), understanding each step of this process is crucial to understanding the properties of the landscape.

Let us consider the first map in Eq. (3), i.e., the map between the space of parameters and the unitary group. It has been shown that the unitaries generated by the ansatz in Eq. (2) are characterized via the so-called Dynamical Lie Algebra (DLA) [56, 57]. Specifically, consider the following definition.

**Definition 1** (Set of generators  $\mathcal{G}$ ). *Consider a parametrized quantum circuit of the form (2). The set*

of generators  $\mathcal{G} = \{H_k\}_{k=1}^K$  is defined as the set (of size  $|\mathcal{G}| = K$ ) of the Hermitian operators that generate the unitaries in a single layer of  $U(\boldsymbol{\theta})$ .

Then, the DLA is defined as follows.

**Definition 2** (Dynamical Lie Algebra (DLA)). *Consider a set of generators  $\mathcal{G}$  according to Definition 1. The DLA  $\mathfrak{g}$  is generated by repeated nested commutators of the operators in  $\mathcal{G}$ . That is,*

$$\mathfrak{g} = \text{span} \langle iH_1, \dots, iH_K \rangle_{Lie}, \quad (4)$$

where  $\langle S \rangle_{Lie}$  denotes the Lie closure, i.e., the set obtained by repeatedly taking the commutator of the elements in  $S$ .

Recall that the set of reachable unitaries  $\{U(\boldsymbol{\theta})\}_{\boldsymbol{\theta} \subseteq \mathbb{G}} \subseteq \mathbb{G} \subseteq SU(d)$  obtained from arbitrary choices of  $\boldsymbol{\theta}$  forms itself a Lie group, known as the dynamical Lie group  $\mathbb{G}$ . Then, we note that  $\mathbb{G}$  is fully obtained from the DLA as  $\mathbb{G} = e^{\mathfrak{g}}$  [37, 72]. We refer the reader to the Methods section for some intuitive understanding on the role of the DLA.

Here, we should remark that the optimal choice of ansatz (or equivalently, the best choice of generators) for a given task is still an open question. While a natural choice would be to use a QNN that is as expressible as possible [73], it has been shown that such choice can lead to trainability issues such as barren plateaus [33, 37, 38].

We can now analyze the second map in Eq. (3), i.e., the map leading to quantum states in a Hilbert space. Given the fact that  $\mathfrak{g}$  determines the set of reachable unitaries, and recalling that the QNN acts on the states  $|\psi_\mu\rangle$  in the training set  $\mathcal{S}$  as  $U(\boldsymbol{\theta})|\psi_\mu\rangle$ , then the set of reachable states (i.e., the orbit) is, in turn, also directly determined by the DLA. We note that in many cases the set of generators can have symmetries, in which case the DLA is of the form  $\mathfrak{g} = \bigoplus_\nu \mathfrak{g}_\nu$ . Here,  $\nu$  is an index over the invariant subspaces. The states in the training set need not respect some, or any, of the symmetries of the QNN. In this work, we consider the case where the states in the training set respect some of the symmetries, and we denote as  $\mathfrak{g}_\mathcal{S}$  the DLA associated with the symmetries preserved by the states in  $\mathcal{S}$ . The limiting case when the states in  $\mathcal{S}$  break all symmetries in the ansatz (or when the ansatz has no symmetries) corresponds to  $\mathfrak{g}_\mathcal{S} = \mathfrak{g}$ .

Here, one can study the set of reachable states through the action of  $U(\boldsymbol{\theta})$  on  $|\psi_\mu\rangle$  as follows. Given a set of parameters  $\boldsymbol{\theta}$  and an infinitesimal perturbation  $\boldsymbol{\delta}$  (possibly obtained from some update rule), it is useful to quantify the distance  $\mathcal{D}$  between the quantum states  $|\psi_\mu(\boldsymbol{\theta})\rangle = U(\boldsymbol{\theta})|\psi_\mu\rangle$  and  $|\psi_\mu(\boldsymbol{\theta} + \boldsymbol{\delta})\rangle = U(\boldsymbol{\theta} + \boldsymbol{\delta})|\psi_\mu\rangle$ . The second-order Taylor expansion of  $\mathcal{D}$  is given by the Fubini-Study metric [74, 75] as

$$\mathcal{D}(|\psi_\mu(\boldsymbol{\theta})\rangle, |\psi_\mu(\boldsymbol{\theta} + \boldsymbol{\delta})\rangle) = \frac{1}{2} \boldsymbol{\delta}^T \cdot F_\mu(\boldsymbol{\theta}) \cdot \boldsymbol{\delta}. \quad (5)$$

Here,  $F_\mu(\boldsymbol{\theta})$  is the Quantum Fisher Information Matrix (QFIM) for the state  $|\psi_\mu\rangle$ . The QFIM is an  $M \times M$  matrix whose elements are [76]

$$[F_\mu(\boldsymbol{\theta})]_{ij} = 4\text{Re}[\langle \partial_i \psi_\mu(\boldsymbol{\theta}) | \partial_j \psi_\mu(\boldsymbol{\theta}) \rangle - \langle \partial_i \psi_\mu(\boldsymbol{\theta}) | \psi_\mu(\boldsymbol{\theta}) \rangle \langle \psi_\mu(\boldsymbol{\theta}) | \partial_j \psi_\mu(\boldsymbol{\theta}) \rangle], \quad (6)$$

where  $|\partial_i \psi_\mu(\boldsymbol{\theta})\rangle = \partial |\psi_\mu(\boldsymbol{\theta})\rangle / \partial \theta_i = \partial_i |\psi_\mu(\boldsymbol{\theta})\rangle$  for  $\theta_i \in \boldsymbol{\theta}$ . The QFIM plays a crucial role in imaginary time evolution algorithms [77], and in quantum-aware optimizers such as the quantum natural gradient descent [78–81]. Moreover, we recall that the rank of the QFIM quantifies the number of independent directions in state space that can be explored by making an infinitesimal change in  $\boldsymbol{\theta}$ .

Finally, consider the third map in Eq. (3), i.e., the map leading to the loss function value. Similar to how the QFIM is related to the changes in state space arising by a change in the parameters, one can also quantify how much the loss function value changes by a small parameter update. In this case, one can study the curvature of the loss landscape via the Hessian matrix  $\nabla^2 \mathcal{L}(\boldsymbol{\theta})$ , an  $M \times M$  matrix whose elements are defined as

$$[\nabla^2 \mathcal{L}(\boldsymbol{\theta})]_{ij} = \partial_i \partial_j \mathcal{L}(\boldsymbol{\theta}). \quad (7)$$

Evaluating the gradient and the Hessian at a given point allows one to construct a quadratic model of the loss function, with the Hessian eigenvectors associated with positive (negative) eigenvalues determining directions of positive (negative) curvature. Thus, the rank of  $\nabla^2 \mathcal{L}(\boldsymbol{\theta})$  is related to the number of directions that lead to (second order) changes in the loss, as a zero-valued eigenvalue indicates a zero-curvature flat direction. We finally note that the Hessian has been used to characterize the loss landscapes of variational quantum algorithms [42, 82–84].

### C. Theoretical Results

Here we present our main results, where we rigorously analyze the overparametrization phenomenon in QNNs. Our results prove that: 1) there exists a critical number of parameters  $M_c$  needed to overparametrize a QNN, and 2) that  $M_c$ , and the onset of overparametrization, can be related to the dimension of the associated DLA. The proofs of our main results are sketched in the Methods section and formally derived in the Supplementary Information. For our main results in Theorem 1 and Theorem 2, we make no assumption on the loss function other than the QNN acting on the states  $|\psi_\mu\rangle$  in the training set  $\mathcal{S}$  as  $U(\boldsymbol{\theta})|\psi_\mu\rangle$  and that the loss is estimated via measurements on these evolved states. Then, for Theorem 3 we consider special cases of such loss functions.

First, consider the following definition.

**Definition 3** (Overparametrization). *A QNN is said to be overparametrized if the number of parameters  $M$  is such that the QFI matrices, for all the states in the training set, simultaneously saturate their achievable rank  $R_\mu$  at least in one point of the loss landscape. That is, if increasing the number of parameters past some minimal (critical) value  $M_c$  does not further increase the rank of any QFIM:*

$$\max_{M \geq M_c, \boldsymbol{\theta}} \text{rank}[F_\mu(\boldsymbol{\theta})] = R_\mu. \quad (8)$$

In the Methods section we give additional motivation for this definition, as well as present an equivalent definition that further highlights the geometrical nature of the overparametrization phenomenon.

According to Definition 3, when the QNN is overparametrized, one can explore all relevant and independent directions in the state space by changing the parameters of the ansatz. Evidently, since the rank of the QFIM is at most equal to  $M$ , then Definition 3 implies that  $M_c$  must be such that  $M_c \geq \max_\mu R_\mu$ . We also remark that the overparametrization is here defined for the QFIM ranks to be equal to  $R_\mu$  on a single point in the landscape. In principle, the QFIM could achieve its maximum rank in a given point, and not in others. However, as we numerically verify (see Supplementary Information), at the overparametrization onset the QFIM saturates its rank almost everywhere in the landscape simultaneously. Here, increasing the number of parameters will not further increase the number of accessible directions in state space. However, it can still be beneficial to add more parameters as this will lead to global minima with higher degeneracy [46, 62, 63, 85].

In light of Definition 3, overparametrization has implications for the trainability of the QNN parameters. If the QNN is underparametrized, the loss landscape can exhibit spurious, or false, local minima [86–89]. However, by increasing the number of parameters and overparametrizing the QNN, one can explore more directions in state space, and hence the optimizer is able to escape these false minima. As such, crossing the overparametrization threshold can be considered as a computational phase transition [52] where a more favorable landscape ameliorates the optimization.

Then, our first main result to understand how overparametrization can improve the trainability is as follows.

**Theorem 1.** *For each state  $|\psi_\mu\rangle$  in the training set  $\mathcal{S}$ , the maximum rank  $R_\mu$  of its associated QFIM (defined in Eq. (6)) is upper bounded as*

$$R_\mu \leq \dim(\mathfrak{g}_\mathcal{S}). \quad (9)$$

We remark that  $\dim(\mathfrak{g}_\mathcal{S}) \leq \dim(\mathfrak{g})$ , and hence  $\dim(\mathfrak{g})$  also upper bounds  $R_\mu$ . Theorem 1 shows that, at most, the QNN can explore  $\dim(\mathfrak{g}_\mathcal{S})$  relevant and independent directions in the state space. And thus, a sufficient condition for overparametrization is that

$$M \geq \dim(\mathfrak{g}_\mathcal{S}). \quad (10)$$

Note here that the number of parameters for overparametrization depends on the data in  $\mathcal{S}$  and on the set of generators  $\mathcal{G}$ . The latter implies that: 1) Different ansatzes for the QNN can be overparametrized for different depths even when using the same dataset, 2) The same QNN ansatz can reach overparametrization for different depths when used for two different datasets.

Then, as shown in the numerical results below, in many cases the QNN is found to be overparametrized when

$$M \sim \dim(\mathfrak{g}_\mathcal{S}). \quad (11)$$

Evidently,  $M$  in Eq. (11) can be intractable for ansatzes where  $\dim(\mathfrak{g}_\mathcal{S}) \in \mathcal{O}(b^n)$  with  $b > 1$  (e.g. controllable systems [37]). More promising, however, are QNNs where  $\dim(\mathfrak{g}_\mathcal{S}) \in \mathcal{O}(\text{poly}(n))$ , as here the QNN can be overparametrized for a number of parameters  $M \in \mathcal{O}(\text{poly}(n))$ . Below we show examples of ansatzes that can achieve overparametrization with polynomially deep circuits.

Here we note that Definition 3 allows us to connect the notion of overparametrization to that of the QNN’s capacity. We recall that the capacity (or power) of a QNN quantifies the breadth of functions that it can capture [90]. While there is no unique definition of capacity, we here consider two definitions for the so-called effective quantum dimension, which measures the power of the QNN. First, following [58], we can define the average effective quantum dimension of a QNN:

$$D_1(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{i=1}^M \mathcal{I}(\lambda_\mu^i(\boldsymbol{\theta})) \right], \quad (12)$$

where  $\lambda_\mu^i(\boldsymbol{\theta})$  are the eigenvalues of the QFIM for the state  $|\psi_\mu\rangle$ , and where  $\mathcal{I}(x) = 0$  for  $x = 0$ , and  $\mathcal{I}(x) = 1$  for  $x \neq 0$ . Here the expectation value is taken over the probability distribution that samples input states from the dataset.

The second definition follows from [24]. In the  $n \rightarrow \infty$  limit, the effective quantum dimension of [24] converges to

$$D_2 = \max_{\boldsymbol{\theta}} \left( \text{rank} \left[ \tilde{F}(\boldsymbol{\theta}) \right] \right), \quad (13)$$

where  $\tilde{F}(\boldsymbol{\theta})$  is the classical Fisher Information matrix obtained as

$$\tilde{F}(\boldsymbol{\theta}) = \mathbb{E} \left[ \frac{\partial \log(p(|\psi_\mu\rangle, y_\mu; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \frac{\partial \log(p(|\psi_\mu\rangle, y_\mu; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}^T \right]. \quad (14)$$

Here,  $p(|\psi\rangle, y; \boldsymbol{\theta})$ , describes the joint relationship between an input  $|\psi\rangle$  and an output  $y$  of the QNN. In addition, the expectation value is taken over the probability distribution that samples input states from the dataset.

Then, the following theorem holds.

**Theorem 2.** *The model capacity, as quantified by the effective dimensions of Eqs. (12) or (13), is upper bounded as*

$$D_1(\boldsymbol{\theta}) \leq \dim(\mathfrak{g}_S), \quad D_2 \leq \dim(\mathfrak{g}_S). \quad (15)$$

Moreover, when the QNN is overparametrized according to Definition 3,  $D_1(\boldsymbol{\theta})$  achieves its maximum value on at least one point of the landscape.

Theorem 2 provides an operational meaning to the overparametrization definition in terms of the model's capacity. Specifically, the onset of the overparametrization arises when the model's capacity in Eq. (12) can get saturated. Moreover, we here see that increasing the number of parameters can never increase the model capacity beyond  $\dim(\mathfrak{g}_S)$ .

Note that Definition 3 relates the overparametrization phenomenon with the rank of the QFIM and the possibility of exploring all relevant directions in the state space. One can also relate the notion of overparametrization with the rank of the Hessian and the relevant directions in the loss function landscape. Consider the case when the loss function is of the form

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{|\psi_\mu\rangle \in \mathcal{S}} c_\mu \text{Tr}[U(\boldsymbol{\theta})|\psi_\mu\rangle\langle\psi_\mu|U^\dagger(\boldsymbol{\theta})O], \quad (16)$$

where  $c_\mu$  are real coefficients associated with each state  $|\psi_\mu\rangle$  in  $\mathcal{S}$ , and where  $O$  is a Hermitian operator. Such loss functions arise for supervised quantum machine learning [43, 44, 91], autoencoding [92], principal component analysis [93–95], dynamical simulation [96–99], and, more generally, for variational quantum algorithms [20, 100]. Then, the following theorem holds.

**Theorem 3.** *Let  $\nabla^2\mathcal{L}(\boldsymbol{\theta}_*)$  be the Hessian for a loss function of the form of Eq. (16) evaluated at the optimal set of parameters  $\boldsymbol{\theta}_*$ . Then, its rank is upper bounded as*

$$\text{rank}[\nabla^2\mathcal{L}(\boldsymbol{\theta}_*)] \leq \min\{\dim(\mathfrak{g}_S), 2dr - r^2 - r\}, \quad (17)$$

where  $r = \min\{\text{rank}[\sum_\mu c_\mu |\psi_\mu\rangle\langle\psi_\mu|], \text{rank}[O]\}$ , and  $d$  is the Hilbert space dimension.

Theorem 3 shows that the maximum number of relevant directions around the global minima of the optimization problem is always smaller than  $\dim(\mathfrak{g}_S)$ . Here, we again

numerically find that in the overparametrization regime adding more parameters only adds zero-valued eigenvalues to the Hessian. We finally remark that Theorem 3 imposes a maximal rank on the Hessian when evaluated at the solution, but in general the Hessian can have a rank larger than  $\dim(\mathfrak{g}_S)$  at other points in the landscape.

Note that in principle one can define overparametrization as the rank of the Hessian being saturated at the solution. However, as discussed in the Methods, this definition could have potential issues.

## D. Numerical Results

Here we numerically illustrate the overparametrization phenomenon and the associated computational phase transition. We consider three different optimization tasks: the Variational Quantum Eigensolver (VQE), unitary compilation, and quantum autoencoding. We note that the overparametrization phenomenon has been empirically observed for the first two tasks respectively in [50, 59] and [52]. The simulations were performed with the open-source library Qibo [101, 102], and the details can be found in the Supplemental Information.

### 1. Variational Quantum Eigensolver

First, we use the VQE algorithm [103–105] to minimize the loss function

$$E(\boldsymbol{\theta}) = \langle\psi(\boldsymbol{\theta})|H_{\text{TFIM}}|\psi(\boldsymbol{\theta})\rangle, \quad (18)$$

and find the ground state of the Hamiltonian of the transverse field Ising model  $H_{\text{TFIM}}$ . Here,  $|\psi(\boldsymbol{\theta})\rangle = U(\boldsymbol{\theta})|+\rangle^{\otimes n}$  and  $H_{\text{TFIM}} = -\sum_{i=1}^{n_f} \sigma_i^z \sigma_{i+1}^z - h \sum_{i=1}^n \sigma_i^x$ , where  $\sigma_i^\mu$  denotes the  $\mu$ -Pauli matrix (with  $\mu = x, z$ ) acting on qubit  $i$ , and  $h$  is the strength of the transverse field. We set  $h = 1$  and consider both open ( $n_f = n - 1$ ) and closed ( $n_f = n$ ) boundary conditions. In the latter,  $\sigma_{n+1}^\mu = \sigma_1^\mu$ . We employ a Hamiltonian variational ansatz for the QNN [50, 70]. This ansatz has two parameters per layer and is precisely of the form in (2) (see Methods for a detailed description of the ansatz).

As shown in [37], the dimension of the DLA associated with the ansatz is given by [106]

$$\dim(\mathfrak{g}_S^{\text{closed}}) = \frac{3}{2}n, \quad \dim(\mathfrak{g}_S^{\text{open}}) = n^2, \quad (19)$$

where the superscripts indicate closed and open boundary conditions in the ansatz and in  $H_{\text{TFIM}}$ . Hence, from our theoretical results, we expect that both of these ansatzes

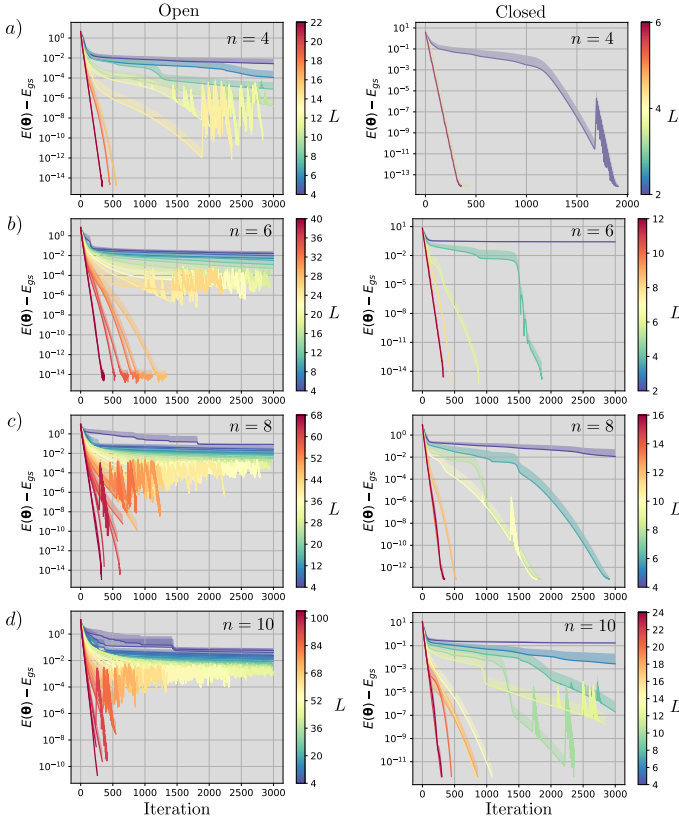


Figure 3. **Training curves for VQE implementation.** The loss function value minus the exact ground-state energy ( $E_{gs}$ ) is plotted versus iteration. We used a Hamiltonian variational ansatz with open (left) and closed (right) boundary conditions to solve the VQE task in Eq. (18) for a)  $n = 4$ , b)  $n = 6$ , c)  $n = 8$ , and d)  $n = 10$  qubits. Solid lines represent the average over 50 random initialized runs while the shaded regions correspond to the standard deviation.

can be overparametrized with only a polynomial number of parameters.

Figure 3 shows the results of minimizing the loss in Eq. (18), for problem sizes of  $n = 4, 6, 8, 10$  qubits and for ansatzes with different depths  $L$  (i.e.,  $2L$  parameters), with both open and closed boundary conditions. In all cases, we averaged over 50 random parameter initializations. First, we note that one can always observe the onset of overparametrization through a computational phase transition whereby the convergence of the optimization dramatically increases when increasing the number of parameters past some threshold. That is, for a small number of layers, the algorithm is unable to accurately find the ground state, while for a large number of layers the algorithm always rapidly converges to the solution. In fact, we observe that

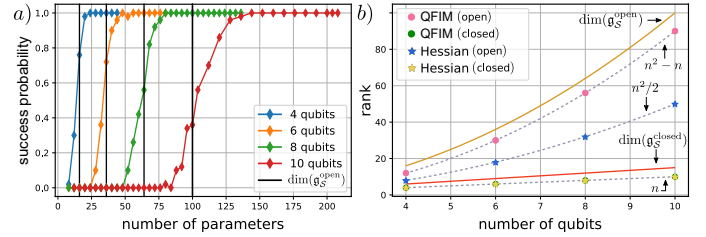


Figure 4. **Overparametrization threshold for VQE implementation.** a) The success probability (i.e., fraction of instances that converged to the global optimum within an error of  $10^{-7}$ ) is plotted versus number of parameters. Results are obtained from 50 randomly-initialized instances, for  $n = 4 - 10$  qubits and for the Hamiltonian variational ansatz with open boundary conditions. The vertical black lines indicate the dimension of the DLA. b) The rank of the QFIM and Hessian is plotted versus number of qubits. The rank of the QFIM was evaluated at the optima and at random points in the landscape. The rank of the Hessian was evaluated at the optima. Dashed lines indicate the functional dependence of the computed ranks. Here, the ansatz had a number of parameters for which overparametrization had been fully achieved.

the loss function decreases exponentially with each optimization step when the number of layers is large enough.

To analyze the number of parameters for which the overparametrization occurs, Figure 4(a) shows the success probability, i.e., the fraction of randomly-initialized instances that converged within  $10^{-7}$  of the true solution. Here, one can see the phase transition at the onset of overparametrization. Indeed, at  $M \sim \dim(\mathfrak{g}_S)$ , the success probability rapidly goes to one. This is due to the fact that the optimization hypersurface becomes more favorable by the removal of false local minima, and thus one can obtain higher-quality solutions with less iterations. Figure 4(a) also shows that further increasing the number of parameters past  $\dim(\mathfrak{g}_S)$  can in fact lead to the QNN having a higher probability of converging to the solution. There exists a point, however, for which the overparametrization saturates and there is no visible improvement in convergence speed or quality of the solution found. We found that the saturation number of parameters grows linearly for closed boundary conditions and quadratically for open boundary conditions, and thus these saturation numbers have the same scaling as their corresponding  $\dim(\mathfrak{g}_S)$ .

Finally, Fig. 4(b) shows the computations of the ranks of the QFIM and Hessian at the overparametrization threshold for Hamiltonians and ansatzes with open and closed boundary conditions. The rank of the QFIM was computed at the global optima, and also at random points in the landscape, and it was found to be the same in all cases. The rank of the Hessian was computed at the global

optima. First, let us note that these results show that Theorem 1 and Theorem 3 hold, as the dimension of the associated DLA is always an upper bound for the ranks. As shown in the figure by the dashed lines, we can find the explicit dependence for the ranks as a function of the system size. In the Supplementary Information we present additional plots for the ranks of the QFIM and Hessian.

## 2. Unitary compilation

Let us now consider a unitary compilation task. Unitary compilation refers to decomposing a target unitary into a sequence of control pulses or quantum gates that can be directly implemented on quantum hardware [107–111].

In variational unitary compiling [110, 111], one trains a parametrized quantum circuit  $U(\theta)$  so that its action matches that of a target unitary  $V$  (up to a global phase). Thus, one minimizes the loss function

$$\mathcal{L}(\theta) = 1 - |\text{Tr}[V^\dagger U(\theta)]|^2 / d^2. \quad (20)$$

Here,  $\mathcal{L}(\theta)$  can be efficiently evaluated on a quantum computer with the Hilbert-Schmidt test [110]. While  $\mathcal{L}(\theta)$  is not exactly of the form in (16), we also prove in the Methods section a theorem showing that the rank of its Hessian is also upper bounded by  $\dim(\mathfrak{g}_S)$  at the global optima.

We employ a hardware efficient ansatz [66] for  $U(\theta)$  composed of alternating layers of single qubit rotations and entangling gates. The number of parameters is therefore  $M = 2n + L(4n - 4)$  (see Methods for a detailed description of the ansatz). We sample the target unitary  $V$  from the Haar measure in the unitary group of degree  $d$ . As shown in [37], the dimension of the DLA associated with this ansatz is

$$\dim(\mathfrak{g}_S) = 4^n, \quad (21)$$

and thus grows exponentially with the number of qubits.

Figure 5(a) shows the results of minimizing the loss function in Eq. (20), for problem sizes of  $n = 2, 3, 4, 5$  qubits, and for ansatzes with different depths  $L$ . In all cases we averaged over 50 random parameter initializations. Here we can again observe that as the depth of the circuit increases, the convergence towards the global optimum improves dramatically until reaching a saturation point. Figure 5(b) plots the success probability for randomly-initialized instances. Similar to the VQE implementation, one finds that around  $\dim(\mathfrak{g}_S)$  parameters are required to consistently find high-quality solutions, and that the probability of convergence to the global optimum undergoes a drastic phase transition when the number of parameters is around  $\dim(\mathfrak{g}_S)$ . This result again implies a simplification of the

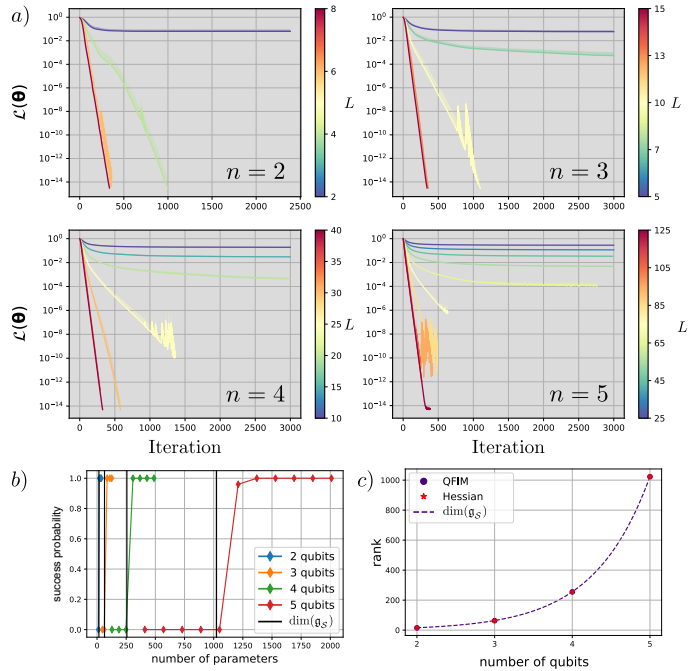


Figure 5. **Unitary compilation implementation.** a) The loss function is plotted versus iteration. The results are obtained for problem sizes of  $n = 2, 3, 4, 5$  qubits and for an  $L$ -layered hardware efficient ansatz. For each point, we averaged the results over 50 randomly-initialized problem instances. b) The success probability (i.e., fraction of instances that converged to the global optimum within an error of  $10^{-7}$ ) is plotted versus number of parameters. The vertical lines indicate the dimension of the DLA. c) The rank of the QFIM and Hessian is plotted versus number of qubits. The rank of the QFIM was evaluated at the optimum and at random points in the landscape. The rank of the Hessian was evaluated at the global optima. Here, the ansatz had a number of parameters for which overparametrization had been fully achieved.

optimization landscape, where local traps disappear. We also numerically verify Theorem 1 and Theorem 3. Namely, Fig. 5(c) plots the rank of the QFIM and the Hessian. The QFIM was evaluated at the global optima and at random points in the landscape, while the Hessian was evaluated at the global optima. For all cases we found that the ranks are equal to  $\dim(\mathfrak{g}) - 1$ .

## 3. Quantum autoencoding

Finally, we present results for the archetypal QML task of quantum autoencoding [92, 112]. A quantum autoencoder is a special type of QNN that can be used to com-



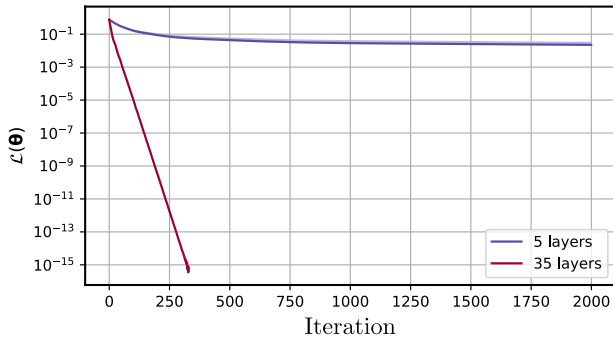


Figure 6. **Quantum autoencoder implementation.** The loss function value is plotted versus the number of iterations. Results are obtained for a  $n = 4$  qubit problem when using a layered hardware efficient ansatz with  $L = 5$  and  $L = 35$  layers.

press quantum information. Analogously to classical autoencoders, the idea is to reduce the dimensionality of the states in a dataset through the action of an encoder  $U(\theta)$ . Once compressed, the states belong to a smaller dimensional Hilbert space known as the latent space. The states compressed into this space can subsequently be recovered with high fidelity at a later time by a decoder  $U^\dagger(\theta)$ .

Consider a bipartite quantum system  $AB$  of  $n_A$  and  $n_B$  qubits, respectively, and let  $|\psi_\mu\rangle$  be states from the training set  $\mathcal{S}$ . The goal of the quantum autoencoder is to train an encoding parametrized quantum circuit  $U(\theta)$  to compress the states in  $\mathcal{S}$  onto subsystem  $A$ , so that one can discard the qubits in subsystem  $B$  without losing information. A possible loss function here is given by [92]

$$\mathcal{L}(\theta) = \sum_{|\psi_\mu\rangle \in \mathcal{S}} \left( 1 - \frac{1}{|\mathcal{S}|} \text{Tr} \left[ U(\theta) |\psi_\mu\rangle \langle \psi_\mu| U^\dagger(\theta) O \right] \right), \quad (22)$$

with  $O = |0\rangle\langle 0|^{\otimes n_B} \otimes \mathbb{1}_A$ , and where  $\mathbb{1}_A$  denotes the identity on subsystem  $A$ . We note that an alternative local version of this loss function was proposed in [34] to avoid barren plateaus issues. However, since we here consider small problem sizes, we use the loss in (22).

Our results were obtained for a system of  $n = 4$  qubits (with  $n_B=2$ ) and for the same hardware efficient ansatz used for unitary compilation. The dataset consisted of four states drawn from the NTangled dataset [113], a quantum dataset composed of states with different amounts and types of multipartite entanglement. As shown in Fig. 6, we can again see that an overparametrized QNN is able to accurately reach the global optima in few iterations. When computing the rank of the QFIM at random points of the landscape, we found that the rank is always 30. This is in contrast to the dimension of the DLA, which is  $\dim(\mathfrak{g}_{\mathcal{S}}) = 256$ , and thus the latter leads to some hope

that overparametrization can be achieved with a number of parameters that is much smaller than  $\dim(\mathfrak{g}_{\mathcal{S}})$ .

Let us here note a crucial difference between the results obtained for the Hamiltonian variational ansatz and for the hardware efficient ansatz. Namely, for the Hamiltonian variational ansatz the dimension of the DLA scaled polynomially with  $n$ , whereas for the hardware efficient ansatz the dimension of the DLA grows exponentially with the system size. Thus, in the first case the system becomes overparametrized at a polynomial number of parameters, while the latter case one can require an exponentially large one. This makes it so that overparametrization can be unachievable in practice for large problem sizes when using an ansatz with an exponentially large DLA. In addition, it has been shown that the ansatzes with exponentially large DLAs can exhibit barren plateaus [37], thus further preventing their practical use.

### III. DISCUSSION

Quantum Machine Learning (QML) is an emerging field that aims to analyze (either classical or quantum) data with significant speedup over classical Machine Learning (ML). However, like classical ML, QML also has trainability issues associated with non-convex landscapes, local minima, and the overall NP-hardness of the optimization.

Classical ML has benefited from the discovery of the overparameterization phenomenon, whereby increasing the number of parameters beyond some threshold causes many local minima to disappear (e.g., as in Fig. 1). Similarly, preliminary evidence of overparameterization in QML has been discovered for specific constructions of Quantum Neural Networks (QNNs). However, prior to our work, no general theory existed for the precise properties of QNNs that lead to overparameterization.

In this work, we provide the first general analysis of overparameterization for a broad class of QNNs (i.e., those with periodic structure). We find that the Dynamical Lie Algebra (DLA) obtained from the set of generators of the QNN plays a crucial role in determining properties of the QNN and the ensuing landscape. To our knowledge, our work is the first *algebraic* theory of overparameterization. This represents an important contribution to Quantum Landscape Theory, i.e., the understanding of QML loss function landscapes and how to engineer them.

We defined overparameterization as the QNN having more than a critical number of parameters that allow it to explore all independent and relevant directions in the state space. This translates to the Quantum Fisher Information Matrices (QFIMs) having reached their maximum achievable rank. This definition has direct implications for

the loss functions of under- and over-parametrized QNNs. Underparametrized QNNs can exhibit spurious, or false, local minima that disappear when one increases the number of parameters and reaches the overparametrization regime. Since the existence of false local minima negatively affect the QNN’s trainability, the overparametrization onset corresponds to a computational phase transition where the QNN parameter optimization improves due to a more favorable landscape.

We found that the critical number of parameters needed to overparametrize the QNN is directly linked to the dimension of the associated DLA  $\mathfrak{g}_S$ . Our theorems showed that the rank of the QFIM (across the whole landscape) and the rank of the Hessian (evaluated at the optima) are upper bounded by  $\dim(\mathfrak{g}_S)$ . Thus, one can potentially reach overparametrization if the QNN has  $\dim(\mathfrak{g}_S)$  parameters. This result is particularly interesting for QNN constructions where  $\dim(\mathfrak{g}_S) \in \mathcal{O}(\text{poly}(n))$ . Thus, our results show that there can exist QNNs that are overparametrized for a polynomial number of parameters.

We verified our theoretical results by performing numerical simulations of problems where the overparametrization had been heuristically observed [50, 52]. Here, our theoretical framework allowed us to shed new light and explain some of the observations in these prior works.

We note that most ansatzes used for QNNs in the literature are ultimately hardware efficient ansatzes. These are known to exhibit barren plateaus, and in view of our recent results, they may require an exponential number of parameters to be overparametrized (e.g., see Fig. 5). These results indicate that the search for scalable and trainable ansatzes should be a priority for the field.

In this sense, our results provide additional guidance to develop QNN architectures with extremely favorable landscapes: overparametrization and absence of barren plateaus. In this context, good candidates are architectures with polynomially large DLAs.

#### IV. METHODS

In this section, we provide additional details and intuition for the results in the main text, as well as a sketch of the proofs for our main theorems. More detailed proofs of our theorems are given in the Supplementary Information.

##### Intuition behind the dynamical Lie algebra

According to Definition 2, the DLA  $\mathfrak{g}$  is obtained from the nested commutators of the elements in the set of generators. To understand why this is the case, let us

consider a single-layered unitary  $U(\boldsymbol{\theta})$  generated by two Hermitian operators, so that  $\mathcal{G} = \{H_1, H_2\}$ . From the Baker–Campbell–Hausdorff formula, we have that

$$U(\boldsymbol{\theta}) = e^{i\theta_1 H_1} e^{i\theta_2 H_2} = e^{K_1(\boldsymbol{\theta})}, \quad (23)$$

where

$$K_1(\boldsymbol{\theta}) = i(\theta_1 H_1 + \theta_2 H_2 + \frac{i\theta_1 \theta_2}{2} [H_1, H_2] - \frac{\theta_1^2 \theta_2}{12} [H_1, [H_1, H_2]] + \dots). \quad (24)$$

In Eq. (24) we can see that by combining  $e^{i\theta_1 H_1}$  and  $e^{i\theta_2 H_2}$  into a single term, the new evolution is generated by an operator  $K_1(\boldsymbol{\theta})$  that depends on both  $\theta_1$  and  $\theta_2$ , and which contains the nested commutators between  $H_1$  and  $H_2$ . Here, it is also worth noting that the set formed by the operators  $\{iH_1, iH_2, i[H_1, [H_1, H_2]], \dots\}$  will eventually be closed under the commutation operation in the sense that not all elements will be linearly independent, but rather there will be a finite basis. This is precisely what the DLA is. It is the space spanned by the  $\dim(\mathfrak{g})$  operators that form a basis of the nested commutators.

When the QNN has multiple layers, that is, when  $U(\boldsymbol{\theta}) = \prod_{l=1}^L e^{i\theta_{l1} H_1} e^{i\theta_{l2} H_2}$ , one can recursively apply the Baker–Campbell–Hausdorff formula to express the action of the QNN as being generated by a single parametrized operator  $K_L(\boldsymbol{\theta})$ . That is, to have  $U(\boldsymbol{\theta}) = e^{K_L(\boldsymbol{\theta})}$ . Evidently, both  $K_1$  and  $K_L$  are obtained from the nested commutators of  $H_1$  and  $H_2$ , and thus both operators are elements of  $\mathfrak{g}$ . However, while  $K_1$  depends on only two parameters  $\theta_1$  and  $\theta_2$ ,  $K_L$  is parametrized by all  $2L$  elements in the vector  $\boldsymbol{\theta} = \{\theta_{l1}, \theta_{l2}\}_{l=1}^L$ . Having these additional parameters allows for a more fine-tuned control of the action of  $U(\boldsymbol{\theta})$ . Intuitively, to hope for a locally surjective map between parameter space and  $\mathfrak{g}$ , we need to place at least  $\dim(\mathfrak{g})$  parameters. Here, there will come a point where further adding parameters does not further increase one’s control of the action of  $U(\boldsymbol{\theta})$ .

We finally note that the analysis for a QNN with more than two unitaries in  $\mathcal{G}$  follows readily.

##### Motivation for the definition of overparametrization

Let us here motivate our definition of overparametrization. First, we recall that we are considering the case where the QNN  $U(\boldsymbol{\theta})$  acts on the states of the training set as  $U(\boldsymbol{\theta}) |\psi_\mu\rangle$ , and that the loss function is estimated via measurement outcomes on such evolved states.

In Definition 3, we defined overparametrization as a property of the QNN (independently of how the loss function is defined). More specifically, we consider a QNN to

be overparametrized if the QNN can explore all relevant directions in the state space. This definition is justified from the fact that, irrespective of how the loss function is estimated via measurements on  $U(\boldsymbol{\theta})|\psi_\mu\rangle$ , the accessible space in the Hilbert space is ultimately defined by the action of the QNN in the states of the training set.

Here, one could also potentially define overparametrization in terms of exploring all relevant directions in the loss landscape. However, this could have some issues. For instance, consider a QML model where one measures the evolved states  $U(\boldsymbol{\theta})|\psi_\mu\rangle$  in the computational basis and evaluates the loss function as  $\mathcal{L}(\boldsymbol{\theta}) = \sum_{\mu, \mathbf{z}} p(\mathbf{z}|\psi_\mu)/|\mathcal{S}|$ , where  $p(\mathbf{z}|\psi_\mu)$  is the probability of measuring the bitstring  $\mathbf{z}$  at the output of the QNN when sending the state  $|\psi_\mu\rangle$  as input. Evidently, here  $\mathcal{L}(\boldsymbol{\theta}) = 1$  for all  $\boldsymbol{\theta}$ , and independently of how the QNN is defined. Thus, the loss landscape is always flat, and the Hessian is trivially given by the zero matrix.

The previous example shows that while a QNN can be considered as overparametrized in the state space, this might not be relevant in the loss landscape space. In view of this issue, we have opted to define overparametrization in the state space, as the map leading to states in the Hilbert space (third map in Eq. (3) and Fig. 2) is more fundamental than the map leading to the loss landscape (fourth map). Evidently, we also expect that arguments can be made in favor of defining overparametrization in terms of the loss landscape, or even the unitary space. However, for the setting presently analyzed, Definition 3 can be considered as a first step toward better understanding the overparametrization phenomenon.

### Sketch of the Proof of Theorem 1

Let us here consider for simplicity the case when the states in the training set do not respect any symmetries in the QNN (i.e.,  $\mathfrak{g}_\mathcal{S} = \mathfrak{g}$ ). From Eq. (2) we have that

$$|\partial_j \psi_\mu(\boldsymbol{\theta})\rangle = \partial_j (U(\boldsymbol{\theta})|\psi_\mu\rangle) = -iU(\boldsymbol{\theta})\tilde{H}_j|\psi\rangle \quad (25)$$

where we defined  $\tilde{H}_j = U_1^\dagger \cdots U_j^\dagger H_j U_j \cdots U_1$ . Note that here the explicit dependence of  $\tilde{H}_j$  in the parameters  $\boldsymbol{\theta}$  is omitted. Replacing (25) in Eq. (6) we find that the elements of the QFIM can be written as

$$[F_\mu(\boldsymbol{\theta})]_{ij} = 4 \sum_{m \neq \psi_\mu} \text{Re}[\langle \psi_\mu | \tilde{H}_i | m \rangle \langle m | \tilde{H}_j | \psi_\mu \rangle]. \quad (26)$$

From here, one finds that the QFIM can be expressed as

$$F_\mu(\boldsymbol{\theta}) = -2 \sum_{m \neq \psi} (\mathbf{R}_{m\psi_\mu} \cdot \mathbf{R}_{m\psi_\mu}^\top + \mathbf{I}_{m\psi_\mu} \cdot \mathbf{I}_{m\psi_\mu}^\top), \quad (27)$$

where we have introduced the vectors  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  with components

$$R_{mn}(i) = \text{Re}[\langle m | \tilde{H}_i | n \rangle], \quad I_{mn}(i) = \text{Im}[\langle m | \tilde{H}_i | n \rangle]. \quad (28)$$

Eq. (27) allows us to write the QFIM as a sum of  $2d - 2$  rank-one matrices.

Here we recall that, by definition,  $H_j$  are elements in the DLA  $\mathfrak{g}$ . Then, since the unitaries  $U$  are elements of the dynamical Lie group  $\mathbb{G}$  generated by  $\mathfrak{g}$ , conjugating  $H_j$  by any unitary  $U$  results in another element in  $\mathfrak{g}$ . That is:  $\forall U \in \mathbb{G}$ , and  $\forall H_i \in \mathfrak{g}$  we have  $UH_jU^\dagger \in \mathfrak{g}$ . Then, by repeating this argument  $j$  times, we find that  $\tilde{H}_j \in \mathfrak{g}$ .

Letting  $\{S_\nu\}_{\nu=1}^{\dim(\mathfrak{g})}$  be a basis of  $\mathfrak{g}$ , we can express

$$\tilde{H}_j = \sum_{\nu=1}^{\dim(\mathfrak{g})} a_\nu(j) S_\nu, \quad (29)$$

where  $a_\nu(j)$  are real coefficients. From Eq. (29) we can find

$$\mathbf{R}_{mn} = \sum_{\nu=1}^{\dim(\mathfrak{g}_\mathcal{S})} \text{Re}[\langle m | S_\nu | n \rangle] \mathbf{a}_\nu, \quad (30)$$

$$\mathbf{I}_{mn} = \sum_{\nu=1}^{\dim(\mathfrak{g}_\mathcal{S})} \text{Im}[\langle m | S_\nu | n \rangle] \mathbf{a}_\nu. \quad (31)$$

Equations (30), and (31) show that the vectors  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  can be expressed as a linear combination of  $\dim(\mathfrak{g}_\mathcal{S})$  other vectors  $\{\mathbf{a}_\nu\}$ . Then, while the  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  generate the  $2d - 2$  rank-one matrices in the QFIM, we have that  $F_\mu(\boldsymbol{\theta})$  has a support on a subspace with a basis that has, at most,  $\dim(\mathfrak{g})$  elements. Thus, we find  $\text{rank}[F_\mu(\boldsymbol{\theta})] \leq \dim(\mathfrak{g})$ . The latter hence proves Theorem 1.

Here we note that Eqs. (30), and (31) do not take into consideration what the state  $|\psi_\mu\rangle$  is. However, from (27), the QFIM is actually expressed in terms of  $\mathbf{R}_{m\psi_\mu}$  and  $\mathbf{I}_{m\psi_\mu}$ . Then, from the definitions in Eq. (28), one can see that the state plays a role in the terms  $\tilde{H}_i|\psi_\mu\rangle$ . From a closer inspection, one can see that  $\tilde{H}_i|\psi_\mu\rangle$  is the action of some elements of the Lie algebra over the state  $|\psi_\mu\rangle$ . Thus, since  $\tilde{H}_i$  are directions in the Lie group, we have that  $\tilde{H}_i|\psi_\mu\rangle$  are directions in the state space.

In fact, the expressible states obtained by acting with the QNN on  $|\psi_\mu\rangle$  form a submanifold of the Hilbert space known as the state space orbit, which is defined by  $\mathbb{G}|\psi_\mu\rangle = \{U|\psi_\mu\rangle, \forall U \in \mathbb{G}\}$  [56]. The latter has a very important implications. Since the rank of the QFIM quantifies the number of independent directions in the state space that are accessible via arbitrary infinitesimal variations of the parameter vector  $\boldsymbol{\theta}$ , it cannot be larger than the dimension of the state space orbit. Thus, one can tighten

the bound in Theorem 1 as

$$\text{rank}[F_\mu(\boldsymbol{\theta})] \leq \dim(\mathbb{G}|\psi_\mu\rangle), \quad (32)$$

where we recall that  $\dim(\mathbb{G}|\psi_\mu\rangle)$  is upper bounded by  $\dim(\mathfrak{g})$ .

Thus, one can define the overparametrization as

**Definition 4** (Overparametrization). *A QNN is said to be overparametrized if the number of parameters  $M$  is such that the QFIM has rank equal to the dimension of the orbits given by the action of  $\mathbb{G}$  on the states in the training set.*

Evidently, a sufficient condition for the QNN to be overparametrized is that  $M \geq \max_\mu \dim(\mathbb{G}|\psi_\mu\rangle)$ . For example, we have numerically verified that this occurs in the VQE implementation, where the overparametrization onset occurs when  $M = \max_\mu \dim(\mathbb{G}|\psi_\mu\rangle)$ .

### Sketch of the Proof of Theorem 3

The proof of Theorem 3 follows similarly to that of Theorem 1. Specifically, one can show that the Hessian evaluated at  $\boldsymbol{\theta}_*$  can be expressed as

$$\nabla^2 \mathcal{L}(\boldsymbol{\theta}_*) = 2 \sum_{m,n=1}^d \kappa_{mn} (\mathbf{R}'_{mn} \cdot (\mathbf{R}'_{mn})^\top + \mathbf{I}'_{mn} \cdot (\mathbf{I}'_{mn})^\top), \quad (33)$$

where  $\kappa_{mn}$  are real coefficients, and where now the vectors  $\mathbf{R}'_{mn}$  and  $\mathbf{I}'_{mn}$  have components

$$\begin{aligned} R'_{mn}(i) &= \text{Re}[\langle m| Q \tilde{H}_j Q^\dagger |n\rangle], \\ I'_{mn}(i) &= \text{Im}[\langle m| Q \tilde{H}_j Q^\dagger |n\rangle]. \end{aligned}$$

Where  $Q$  is the matrix that diagonalizes the operator  $\sigma = \sum_\mu c_\mu |\psi_\mu\rangle\langle\psi_\mu|$ .

Then following a similar argument as the one previously used for Theorem 1, we can again show that the rank of the Hessian is such that  $\text{rank}[\nabla^2 \mathcal{L}(\boldsymbol{\theta}_*)] \leq \min\{\dim(\mathfrak{g}_S)\}$ , we leave for the Supplementary Information the rest of the proof, where the quantity  $r = \min\{\text{rank}[\sigma, \text{rank}[O]]\}$  comes into play.

In the Supplemental Information we also provide a proof for the following theorem

**Theorem 4.** *Consider the loss functions for a unitary compilation task*

$$\mathcal{L}_1(\boldsymbol{\theta}) = 2d - 2\text{Re}[T(\boldsymbol{\theta})], \quad \text{and} \quad \mathcal{L}_2(\boldsymbol{\theta}) = 1 - \frac{1}{d^2} |T(\boldsymbol{\theta})|^2,$$

where  $T(\boldsymbol{\theta}) = \text{Tr}[V^\dagger U(\boldsymbol{\theta})]$  for a target unitary  $V$ . Then, let  $H_1(\boldsymbol{\theta}_*)$  and  $H_2(\boldsymbol{\theta}_*)$  be the Hessian for the loss functions  $\mathcal{L}_1(\boldsymbol{\theta})$  and  $\mathcal{L}_2(\boldsymbol{\theta})$ , respectively evaluated at their solutions  $U(\boldsymbol{\theta}_*) = V$  and  $U(\boldsymbol{\theta}_*) = e^{i\phi}V$ . Then, the maximal rank of  $\nabla^2 \mathcal{L}_1(\boldsymbol{\theta}_*)$  and  $\nabla^2 \mathcal{L}_2(\boldsymbol{\theta}_*)$  is such that  $\text{rank}[\nabla^2 \mathcal{L}_1(\boldsymbol{\theta}_*)], \text{rank}[\nabla^2 \mathcal{L}_2(\boldsymbol{\theta}_*)] \leq \dim(\mathfrak{g}_S)$

### Sketch of the Proof of Theorem 2

Let us here consider for simplicity the case when the dataset contains a single state  $|\psi\rangle$  which does not respect the symmetries of the ansatz ( $\mathfrak{g}_S = \mathfrak{g}$ ). The more general proof is presented in the Supplementary Information.

Now, the capacities of Eq. (12) and (13) are

$$D_1(\boldsymbol{\theta}) = \sum_{i=1}^M \mathcal{I}(\lambda^i(\boldsymbol{\theta})), \quad (34)$$

where  $\lambda^i(\boldsymbol{\theta})$  are the eigenvalues of the QFIM for the state  $|\psi\rangle$ , and

$$D_2 = \max_{\boldsymbol{\theta}} \left( \text{rank} \left[ \tilde{F}(\boldsymbol{\theta}) \right] \right), \quad (35)$$

for  $\tilde{F}(\boldsymbol{\theta})$  the classical Fisher information for the input state  $|\psi\rangle$ .

First, we note that, by definition,  $D_1(\boldsymbol{\theta}) = \text{rank}[F(\boldsymbol{\theta})]$ , so that the inequality  $D_1(\boldsymbol{\theta}) \leq \dim(\mathfrak{g})$  follows readily from Theorem 1. Moreover, by the definition of overparametrization in Definition 3, the capacity  $D_1(\boldsymbol{\theta})$  is saturated on at least one point of the landscape.

Now we need to show that  $D_2(\boldsymbol{\theta}) \leq \dim(\mathfrak{g})$ . Here we recall that the quantum and classical Fisher information matrices are such that [75]

$$\tilde{F}(\boldsymbol{\theta}) \leq F(\boldsymbol{\theta}) \quad (36)$$

for all  $\boldsymbol{\theta}$ . Then, using that fact that if  $A$  and  $B$  are two Hermitian matrices such that  $A \leq B$ , then  $A^q \leq B^q$  for all  $q \in [0, 1]$  [114]. Thus, we have that  $\tilde{F}^q(\boldsymbol{\theta}) \leq F^q(\boldsymbol{\theta})$ , and choosing  $q = 0$  leads to

$$\text{supp}[\tilde{F}(\boldsymbol{\theta})] \leq \text{supp}[F(\boldsymbol{\theta})], \quad (37)$$

where here  $\text{supp}(\cdot)$  denotes the support of a matrix. Taking the trace on both sides allows us to obtain

$$\text{rank}[\tilde{F}(\boldsymbol{\theta})] \leq \text{rank}[F(\boldsymbol{\theta})]. \quad (38)$$

Finally, combining Theorem 1 with the definition of overparametrization in Definition 3 and the definition of the capacity  $D_2(\boldsymbol{\theta})$  in Eq. (35), it follows that  $D_2(\boldsymbol{\theta}) \leq \dim(\mathfrak{g})$ .

### Implications for Quantum Optimal Control

In Quantum Optimal Control (QOC) [115–126] one is typically interested in controlling the dynamics of a quantum state  $|\psi\rangle$  evolving through a functional time-dependent Hamiltonian

$$H(t, \{\theta_k(t)\}) = H_0 + \sum_{k=1}^K \theta_k(t) H_k \quad (39)$$

that defines the continuous-in-time equation of motion

$$\frac{dU(t)}{dt} = -iH(t, \{\theta_k(t)\})U(t), \quad \text{with } U(0) = \mathbb{I}. \quad (40)$$

Here, the idea is that the functions  $\{\theta_k(t)\}$ , known as control fields, can be trained to pursue some desired evolution.

Interestingly, it has been shown some that QOC and the field of variational quantum algorithms can be unified into a single framework where the evolution of a quantum system is controlled at the pulse level (QOC), or at the gate level (QNN) [37, 127]. Most importantly, irregardless of the choice of controls, the unitaries that are expressible by a QOC ansatz  $U(t)$  are, like in the QNN case, contained in the group generated by the DLA  $\mathfrak{g}$  (see Definition 2) that is determined by the set of generators  $\mathcal{G} = \{H_k\}_{k=0}^K$ . Since all of the results presented in this manuscript are stated in terms of the DLA of a given QNN, they can be straightforwardly adapted to the QOC setting. For example, the maximum rank achievable by some QFIM associated with an ansatz of the form in Eq. (40) will be upper bounded by the dimension of  $\mathfrak{g}$  (equivalent to Theorem 1). That is, a QOC and a QNN ansatz that share the same set of generators  $\mathcal{G}$  can be expected to have the same saturation value for their respective QFIM matrices.

Similarly, in analogy with the results in Theorems 3 and 4, the Hessian under a QOC ansatz can be expected to be upper bounded by  $\dim(\mathfrak{g})$  when evaluated at a solution. While the existence of bounds on the rank of the Hessian at solutions is well known in the control literature [129, 130], these results analyze the case when the ansatz is controllable (i.e., when  $\mathfrak{g} = \mathfrak{su}(d)$ ) and thus the bounds found are exponentially large. For example, the rank of the Hessian for unitary compilation tasks (see Eq.(20)), has been shown to be upper bounded by  $\dim(\mathfrak{g}) = \dim(\mathfrak{su}(d)) = d^2 - 1$ . Hence, the results in this work generalize these previous studies to the case of general  $\mathfrak{g}$  (i.e. uncontrollable systems). Let us note that the existence of a fundamental bound on the rank of the Hessian at the global minima is directly connected to another interesting phenomenon: the arisal of continuous submanifolds of degenerate solutions [62, 63, 85].

Although historically the quantum control community has mainly focused on controllable systems, the importance of studying uncontrollable ones, in particular those with  $\dim(\mathfrak{g}) = \mathcal{O}(\text{poly}(n))$ , has been evidenced in [37]. Here, it has been ascertained that control systems with exponentially large DLAs may encounter scalability issues, like the presence of barren plateaus in their optimization landscapes. Conversely, systems with polynomially large DLAs can avoid barren plateau issues and be scalable. Thus, the results in the present manuscript should also be considered as an additional motivation for QOC systems with polynomially sized algebras, as these will achieve overparametrization with  $\mathcal{O}(\text{poly}(n))$  parameters.

Finally, we remark that our results also provide a new insight into the existence of false traps in the control landscape [83, 86–89]. In QOC, false traps are usually analyzed through the rank of the Jacobian matrix of the map  $\{\theta_k(t)\} \rightarrow U(t)$ . Here, false traps are critical points in the landscape that are not related to local minima of the loss function itself, but to points where this map is not locally surjective. In our context, this is precisely what a rank-deficient QFIM means: points in parameter space where all possible variations of parameters do not translate into all possible directions in the state space orbit.

### Ansatzes for the numerical simulations

In this section we present the details of the two QNN ansatzes used in our numeric simulations. Let us remark that both ansatzes are of the form in Eq. (18), i.e. a periodic structured parametrized circuit defined by a given set of generators  $\mathcal{G}$ .

Let us first describe the so-called Hamiltonian variational ansatz (HVA) [50, 70]. Consider a VQE task where one wants to minimize a Hamiltonian of the form

$$H = \sum_{k=1}^N a_k A_k, \quad (41)$$

where  $A_k$  are Hermitian operators and  $a_k$  real numbers. The basic idea in the HVA ansatz is to use, as generators, the individual terms in the Hamiltonian that is being minimized, i.e.  $\mathcal{G} = \{A_k\}_{k=1}^N$ . For instance, we show in Fig. 7(a) the ansatz used to find the ground-state of the transverse field Ising model with open boundary conditions. Here, the generators are  $A_0 = \frac{1}{2} \sum_i \sigma_i^x$  and  $A_1 = \frac{1}{2} \sum_i \sigma_i^z \sigma_{i+1}^z$ .

As a second choice of ansatz, let us introduce the hardware efficient ansatz [66] used in the unitary compilation and autoencoding tasks. As shown in Fig. 7(b), this ansatz is composed of single qubit rotations followed by

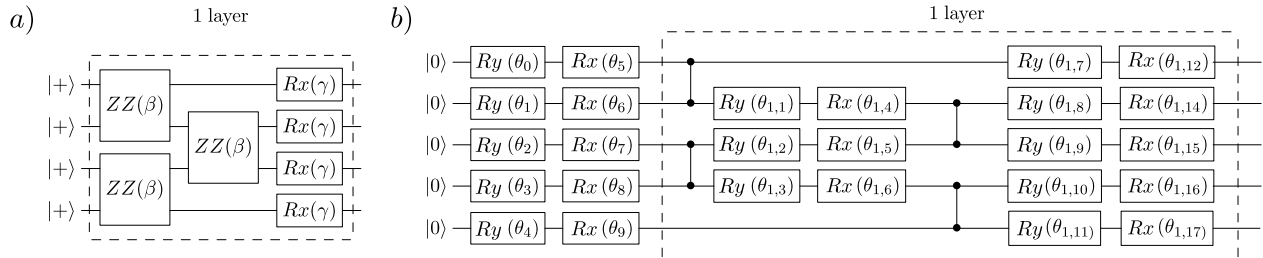


Figure 7. **QNN ansatzes for the numerical simulations.** a) Hamiltonian variational ansatz for the VQE task. Here we show a single layer of the ansatz for  $n = 4$  qubits. For closed boundary conditions, there is an extra  $ZZ(\beta)$  gate acting on the uppermost and lowermost qubits. A  $ZZ(\beta)$  gate on qubits  $i, j$  corresponds to the operator  $e^{-i\frac{\beta}{2}\sigma_i^z\sigma_j^z}$ , and it may be decomposed into two CNOTs and one  $Rz$  rotation [128]. The input state was  $|+\rangle^{\otimes n}$ . b) Hardware efficient ansatz for the unitary compilation and autoencoding tasks. Here we show a single layer of the ansatz for  $n = 5$  qubits. Notice that there is an extra  $Ry$  and  $Rx$  rotation on each qubit at the beginning of the circuit.

CZ gates acting on alternating pairs of qubits. Here we can see that the number of parameters in the ansatz is  $M = 2n + L(4n - 4)$ .

## ACKNOWLEDGEMENTS

We thank Patrick deNiverville, Julia Nakhleh, Stavros Efthymiou, Louis Schatzki and Marco Farinati for useful conversations. NJ and DGM were supported by the U.S. DOE through a quantum computing program sponsored by the Los Alamos National Laboratory (LANL) Information Science & Technology Institute. DGM acknowledges partial financial support from project QuantumCAT (ref. 001- P-001644), co-funded by the Generalitat de Catalunya and the European Union Regional Development Fund within the ERDF Operational Program of Catalunya, and from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 951911 (AI4Media). PJC and MC were initially supported by Laboratory Directed Research and Development (LDRD) program of LANL under project number 20190065DR. PJC also acknowledges support from the LANL ASC Beyond Moore’s Law project. MC also acknowledges support from

the Center for Nonlinear Studies at LANL. This work was supported by the U.S. DOE, Office of Science, Office of Advanced Scientific Computing Research, under the Accelerated Research in Quantum Computing (ARQC) program.

## AUTHOR CONTRIBUTIONS

The project was conceived by ML, PJC and MC. The manuscript was written by NJ, ML, DGM, PJC, and MC. Theoretical results were proved by NJ, ML, PJC, and MC. Numerical implementations were performed by DGM.

## DATA AVAILABILITY

Data generated and analyzed during current study are available from the corresponding author upon reasonable request.

## COMPETING INTERESTS

The authors declare no competing interests.

- 
- [1] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning* (MIT Press, 2018).
- [2] J. Vamathevan, D. Clark, P. Czodrowski, I. Dunham, E. Ferran, G. Lee, B. Li, A. Madabhushi, P. Shah, M. Spitzer, *et al.*, Applications of machine learning in drug discovery and development, *Nature Reviews Drug Discovery* **18**, 463 (2019).
- [3] J. Schmidt, M. R. Marques, S. Botti, and M. A. Marques, Recent advances and applications of machine learning in solid-state materials science, *npj Computational Materials* **5**, 1 (2019).
- [4] S. Grigorescu, B. Trasnea, T. Cocias, and G. Macesanu, A survey of deep learning techniques for autonomous driving, *Journal of Field Robotics* **37**, 362 (2020).

- [5] A. L. Blum and R. L. Rivest, Training a 3-node neural network is np-complete, *Neural Networks* **5**, 117 (1992).
- [6] A. Daniely, Complexity theoretic limitations on learning halfspaces, in *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing* (2016) pp. 105–117.
- [7] D. Boob, S. S. Dey, and G. Lan, Complexity of training relu neural network, *Discrete Optimization*, 100620 (2020).
- [8] B. Neyshabur, Z. Li, S. Bhojanapalli, Y. LeCun, and N. Srebro, The role of over-parametrization in generalization of neural networks, in *International Conference on Learning Representations* (2018).
- [9] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, Understanding deep learning (still) requires rethinking generalization, *Communications of the ACM* **64**, 107 (2021).
- [10] Z. Allen-Zhu, Y. Li, and Z. Song, A convergence theory for deep learning via over-parameterization, in *International Conference on Machine Learning* (PMLR, 2019) pp. 242–252.
- [11] Z. Allen-Zhu, Y. Li, and Y. Liang, Learning and generalization in overparameterized neural networks, going beyond two layers, *Advances in neural information processing systems* (2019).
- [12] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai, Gradient descent finds global minima of deep neural networks, in *International Conference on Machine Learning* (PMLR, 2019) pp. 1675–1685.
- [13] R.-D. Buhai, Y. Halpern, Y. Kim, A. Risteski, and D. Sontag, Empirical study of the benefits of overparameterization in learning latent variable models, in *International Conference on Machine Learning* (PMLR, 2020) pp. 1211–1219.
- [14] S. S. Du, X. Zhai, B. Póczos, and A. Singh, Gradient descent provably optimizes over-parameterized neural networks, in *International Conference on Learning Representations* (2019).
- [15] A. Brutzkus, A. Globerson, E. Malach, and S. Shalev-Shwartz, SGD learns over-parameterized networks that provably generalize on linearly separable data, in *International Conference on Learning Representations* (2018).
- [16] M. A. Nielsen and I. L. Chuang, *Quantum Computation and Quantum Information* (Cambridge University Press, 2000).
- [17] J. Preskill, Quantum computing in the nisq era and beyond, *Quantum* **2**, 79 (2018).
- [18] M. Schuld, I. Sinayskiy, and F. Petruccione, An introduction to quantum machine learning, *Contemporary Physics* **56**, 172 (2015).
- [19] J. Biamonte, P. Wittek, N. Pancotti, P. Rebentrost, N. Wiebe, and S. Lloyd, Quantum machine learning, *Nature* **549**, 195 (2017).
- [20] M. Cerezo, A. Arrasmith, R. Babbush, S. C. Benjamin, S. Endo, K. Fujii, J. R. McClean, K. Mitarai, X. Yuan, L. Cincio, and P. J. Coles, Variational quantum algorithms, *Nature Reviews Physics* **1**, 19 (2021).
- [21] H.-Y. Huang, R. Kueng, and J. Preskill, Information-theoretic bounds on quantum advantage in machine learning, *Phys. Rev. Lett.* **126**, 190505 (2021).
- [22] H.-Y. Huang, M. Broughton, M. Mohseni, R. Babbush, S. Boixo, H. Neven, and J. R. McClean, Power of data in quantum machine learning, *Nature Communications* **12**, 1 (2021).
- [23] J. M. Kübler, S. Buchholz, and B. Schölkopf, The inductive bias of quantum kernels, *arXiv preprint arXiv:2106.03747* (2021).
- [24] A. Abbas, D. Sutter, C. Zoufal, A. Lucchi, A. Figalli, and S. Woerner, The power of quantum neural networks, *Nature Computational Science* **1**, 403 (2021).
- [25] L. Bittel and M. Kliesch, Training variational quantum algorithms is np-hard, *Phys. Rev. Lett.* **127**, 120502 (2021).
- [26] M. Benedetti, E. Lloyd, S. Sack, and M. Fiorentini, Parameterized quantum circuits as machine learning models, *Quantum Science and Technology* **4**, 043001 (2019).
- [27] K. Beer, D. Bondarenko, T. Farrelly, T. J. Osborne, R. Salzmann, D. Scheiermann, and R. Wolf, Training deep quantum neural networks, *Nature Communications* **11**, 808 (2020).
- [28] I. Cong, S. Choi, and M. D. Lukin, Quantum convolutional neural networks, *Nature Physics* **15**, 1273 (2019).
- [29] E. Farhi and H. Neven, Classification with quantum neural networks on near term processors, *arXiv preprint arXiv:1802.06002* (2018).
- [30] A. Arrasmith, Z. Holmes, M. Cerezo, and P. J. Coles, Equivalence of quantum barren plateaus to cost concentration and narrow gorges, *arXiv preprint arXiv:2104.05868* (2021).
- [31] J. Rivera-Dean, P. Huembeli, A. Acín, and J. Bowles, Avoiding local minima in variational quantum algorithms with neural networks, *arXiv preprint arXiv:2104.02955* (2021).
- [32] D. Wierichs, C. Gogolin, and M. Kastoryano, Avoiding local minima in variational quantum eigensolvers with the natural gradient optimizer, *Physical Review Research* **2**, 043246 (2020).
- [33] J. R. McClean, S. Boixo, V. N. Smelyanskiy, R. Babbush, and H. Neven, Barren plateaus in quantum neural network training landscapes, *Nature communications* **9**, 1 (2018).
- [34] M. Cerezo, A. Sone, T. Volkoff, L. Cincio, and P. J. Coles, Cost function dependent barren plateaus in shallow parametrized quantum circuits, *Nature Communications* **12**, 1791 (2021).
- [35] C. O. Marrero, M. Kieferová, and N. Wiebe, Entanglement induced barren plateaus, *arXiv preprint arXiv:2010.15968* (2020).
- [36] T. L. Patti, K. Najafi, X. Gao, and S. F. Yelin, Entanglement devised barren plateau mitigation, *Physical Review Research* **3**, 033090 (2021).
- [37] M. Larocca, P. Czarnik, K. Sharma, G. Muraleedharan, P. J. Coles, and M. Cerezo, Diagnosing barren plateaus with tools from quantum optimal control, *arXiv preprint arXiv:2105.14377* (2021).

- [38] Z. Holmes, K. Sharma, M. Cerezo, and P. J. Coles, Connecting ansatz expressibility to gradient magnitudes and barren plateaus, [arXiv preprint arXiv:2101.02138 \(2021\)](#).
- [39] M. Cerezo and P. J. Coles, Higher order derivatives of quantum neural networks with barren plateaus, [Quantum Science and Technology 6, 035006 \(2021\)](#).
- [40] Z. Holmes, A. Arrasmith, B. Yan, P. J. Coles, A. Albrecht, and A. T. Sornborger, Barren plateaus preclude learning scramblers, [Physical Review Letters 126, 190501 \(2021\)](#).
- [41] A. Arrasmith, M. Cerezo, P. Czarnik, L. Cincio, and P. J. Coles, Effect of barren plateaus on gradient-free optimization, [arXiv preprint arXiv:2011.12245 \(2020\)](#).
- [42] P. Huembeli and A. Dauphin, Characterizing the loss landscape of variational quantum circuits, [Quantum Science and Technology 6, 025011 \(2021\)](#).
- [43] A. Pesah, M. Cerezo, S. Wang, T. Volkoff, A. T. Sornborger, and P. J. Coles, Absence of barren plateaus in quantum convolutional neural networks, [arXiv preprint arXiv:2011.02966 \(2020\)](#).
- [44] K. Sharma, M. Cerezo, L. Cincio, and P. J. Coles, Trainability of dissipative perceptron-based quantum neural networks, [arXiv preprint arXiv:2005.12458 \(2020\)](#).
- [45] S. Wang, E. Fontana, M. Cerezo, K. Sharma, A. Sone, L. Cincio, and P. J. Coles, Noise-induced barren plateaus in variational quantum algorithms, [arXiv preprint arXiv:2007.14384 \(2020\)](#).
- [46] E. Fontana, M. Cerezo, A. Arrasmith, I. Rungger, and P. J. Coles, Optimizing parametrized quantum circuits via noise-induced breaking of symmetries, [arXiv preprint arXiv:2011.08763 \(2020\)](#).
- [47] S. Wang, P. Czarnik, A. Arrasmith, M. Cerezo, L. Cincio, and P. J. Coles, Can error mitigation improve trainability of noisy variational quantum algorithms?, [arXiv preprint arXiv:2109.01051 \(2021\)](#).
- [48] E. Campos, D. Rabinovich, V. Akshay, and J. Biamonte, Training saturation in layerwise quantum approximate optimization, [Phys. Rev. A 104, L030401 \(2021\)](#).
- [49] D. S. Franca and R. Garcia-Patron, Limitations of optimization algorithms on noisy quantum devices, [arXiv preprint arXiv:2009.05532 \(2020\)](#).
- [50] R. Wiersema, C. Zhou, Y. de Sereville, J. F. Carrasquilla, Y. B. Kim, and H. Yuen, Exploring entanglement and optimization within the hamiltonian variational ansatz, [PRX Quantum 1, 020319 \(2020\)](#).
- [51] S. Zhang and W. Cui, Overparametrization in qaoa, [Written Report \(2020\)](#).
- [52] B. T. Kiani, S. Lloyd, and R. Maity, Learning unitaries by gradient descent, [arXiv preprint arXiv:2001.11897 \(2020\)](#).
- [53] L. Funcke, T. Hartung, K. Jansen, S. Kühn, M. Schneider, and P. Stornati, Best-approximation error for parametric quantum circuits, [arXiv preprint arXiv:2107.07378 \(2021\)](#).
- [54] J. Lee, A. B. Magann, H. A. Rabitz, and C. Arenz, Progress toward favorable landscapes in quantum combinatorial optimization, [Physical Review A 104, 032401 \(2021\)](#).
- [55] E. R. Anschuetz, Critical points in hamiltonian agnostic variational quantum algorithms, [arXiv preprint arXiv:2109.06957 \(2021\)](#).
- [56] D. D'Alessandro, *Introduction to Quantum Control and Dynamics*, Chapman & Hall/CRC Applied Mathematics & Nonlinear Science (Taylor & Francis, 2007).
- [57] R. Zeier and T. Schulte-Herbrüggen, Symmetry principles in quantum systems theory, [Journal of mathematical physics 52, 113510 \(2011\)](#).
- [58] T. Haug, K. Bharti, and M. Kim, Capacity and quantum geometry of parametrized quantum circuits, [arXiv preprint arXiv:2102.01659 \(2021\)](#).
- [59] J. Kim, J. Kim, and D. Rosa, Universal effectiveness of high-depth circuits in variational eigenproblems, [Physical Review Research 3, 023203 \(2021\)](#).
- [60] D. d'Alessandro, *Introduction to quantum control and dynamics* (CRC press, 2007).
- [61] R. Chakrabarti and H. Rabitz, Quantum control landscapes, [International Reviews in Physical Chemistry 26, 671 \(2007\)](#).
- [62] M. Larocca, E. Calzetta, and D. A. Wisniacki, Exploiting landscape geometry to enhance quantum optimal control, [Phys. Rev. A 101, 023410 \(2020\)](#).
- [63] M. Larocca, E. Calzetta, and D. Wisniacki, Fourier compression: A customization method for quantum control protocols, [Phys. Rev. A 102, 033108 \(2020\)](#).
- [64] S. Boixo, S. V. Isakov, V. N. Smelyanskiy, R. Babbush, N. Ding, Z. Jiang, M. J. Bremner, J. M. Martinis, and H. Neven, Characterizing quantum supremacy in near-term devices, [Nature Physics 14, 595 \(2018\)](#).
- [65] F. Arute, K. Arya, R. Babbush, D. Bacon, *et al.*, Quantum supremacy using a programmable superconducting processor, [Nature 574, 505 \(2019\)](#).
- [66] A. Kandala, A. Mezzacapo, K. Temme, M. Takita, M. Brink, J. M. Chow, and J. M. Gambetta, Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets, [Nature 549, 242 \(2017\)](#).
- [67] E. Farhi, J. Goldstone, and S. Gutmann, A quantum approximate optimization algorithm, [arXiv preprint arXiv:1411.4028 \(2014\)](#).
- [68] S. Hadfield, Z. Wang, B. O'Gorman, E. G. Rieffel, D. Venturelli, and R. Biswas, From the quantum approximate optimization algorithm to a quantum alternating operator ansatz, [Algorithms 12, 34 \(2019\)](#).
- [69] L. Zhu, H. L. Tang, G. S. Barron, N. J. Mayhall, E. Barnes, and S. E. Economou, An adaptive quantum approximate optimization algorithm for solving combinatorial problems on a quantum computer, [arXiv preprint arXiv:2005.10258 \(2020\)](#).
- [70] D. Wecker, M. B. Hastings, and M. Troyer, Progress towards practical quantum variational algorithms, [Phys. Rev. A 92, 042303 \(2015\)](#).
- [71] A. Choquette, A. Di Paolo, P. K. Barkoutsos, D. Sénéchal, I. Tavernelli, and A. Blais, Quantum-optimal-control-inspired ansatz for variational quantum algorithms, [Physical Review Research 3, 023092 \(2021\)](#).



- [72] M. E. Morales, J. Biamonte, and Z. Zimborás, On the universality of the quantum approximate optimization algorithm, *Quantum Information Processing* **19**, 1 (2020).
- [73] S. Sim, P. D. Johnson, and A. Aspuru-Guzik, Expressibility and entangling capability of parameterized quantum circuits for hybrid quantum-classical algorithms, *Advanced Quantum Technologies* **2**, 1900070 (2019).
- [74] R. Cheng, Quantum geometric tensor (fubini-study metric) in simple quantum system: A pedagogical introduction, *arXiv preprint arXiv:1012.1337* (2010).
- [75] J. J. Meyer, Fisher Information in Noisy Intermediate-Scale Quantum Applications, *Quantum* **5**, 539 (2021).
- [76] J. Liu, H. Yuan, X.-M. Lu, and X. Wang, Quantum fisher information matrix and multiparameter estimation, *Journal of Physics A: Mathematical and Theoretical* **53**, 023001 (2019).
- [77] S. McArdle, T. Jones, S. Endo, Y. Li, S. C. Benjamin, and X. Yuan, Variational ansatz-based quantum simulation of imaginary time evolution, *npj Quantum Information* **5**, 1 (2019).
- [78] J. Stokes, J. Izaac, N. Killoran, and G. Carleo, Quantum natural gradient, *Quantum* **4**, 269 (2020).
- [79] B. Koczor and S. C. Benjamin, Quantum natural gradient generalised to non-unitary circuits, *arXiv preprint arXiv:1912.08660* (2019).
- [80] J. Gacon, C. Zoufal, G. Carleo, and S. Woerner, Simultaneous perturbation stochastic approximation of the quantum fisher information, *arXiv preprint arXiv:2103.09232* (2021).
- [81] T. Haug and M. Kim, Natural parameterized quantum circuit, *arXiv preprint arXiv:2107.14063* (2021).
- [82] J. Kim and Y. Oz, Quantum energy landscape and vqa optimization, *arXiv preprint arXiv:2107.10166* (2021).
- [83] M. Dalgaard, J. Sherson, and F. Motzoi, Predicting quantum dynamical cost landscapes with deep learning, *arXiv preprint arXiv:2107.00008* (2021).
- [84] M. Dalgaard, F. Motzoi, J. H. M. Jensen, and J. Sherson, Hessian-based optimization of constrained quantum control, *Physical Review A* **102**, 042612 (2020).
- [85] K. W. Moore and H. Rabitz, Exploring constrained quantum control landscapes, *The Journal of chemical physics* **137**, 134113 (2012).
- [86] R.-B. Wu, R. Long, J. Dominy, T.-S. Ho, and H. Rabitz, Singularities of quantum control landscapes, *Physical Review A* **86**, 013405 (2012).
- [87] G. Riviello, C. Brif, R. Long, R.-B. Wu, K. M. Tibbetts, T.-S. Ho, and H. Rabitz, Searching for quantum optimal control fields in the presence of singular critical points, *Physical Review A* **90**, 013404 (2014).
- [88] N. Rach, M. M. Müller, T. Calarco, and S. Montangero, Dressing the chopped-random-basis optimization: A bandwidth-limited access to the trap-free landscape, *Physical Review A* **92**, 062343 (2015).
- [89] M. Larocca, P. M. Poggi, and D. A. Wisniacki, Quantum control landscape for a two-level system near the quantum speed limit, *Journal of Physics A: Mathematical and Theoretical* **51**, 385305 (2018).
- [90] P. J. Coles, Seeking quantum advantage for neural networks, *Nature Computational Science* **1**, 389 (2021).
- [91] V. Havlíček, A. D. Córcoles, K. Temme, A. W. Harrow, A. Kandala, J. M. Chow, and J. M. Gambetta, Supervised learning with quantum-enhanced feature spaces, *Nature* **567**, 209 (2019).
- [92] J. Romero, J. P. Olson, and A. Aspuru-Guzik, Quantum autoencoders for efficient compression of quantum data, *Quantum Science and Technology* **2**, 045001 (2017).
- [93] R. LaRose, A. Tikku, É. O’Neel-Judy, L. Cincio, and P. J. Coles, Variational quantum state diagonalization, *npj Quantum Information* **5**, 1 (2019).
- [94] C. Bravo-Prieto, D. García-Martín, and J. I. Latorre, Quantum singular value decomposer, *Phys. Rev. A* **101**, 062310 (2020).
- [95] M. Cerezo, K. Sharma, A. Arrasmith, and P. J. Coles, Variational quantum state eigensolver, *arXiv preprint arXiv:2004.01372* (2020).
- [96] X. Yuan, S. Endo, Q. Zhao, Y. Li, and S. C. Benjamin, Theory of variational quantum simulation, *Quantum* **3**, 191 (2019).
- [97] C. Cirstoiu, Z. Holmes, J. Iosue, L. Cincio, P. J. Coles, and A. Sornborger, Variational fast forwarding for quantum simulation beyond the coherence time, *npj Quantum Information* **6**, 1 (2020).
- [98] B. Commeau, M. Cerezo, Z. Holmes, L. Cincio, P. J. Coles, and A. Sornborger, Variational hamiltonian diagonalization for dynamical quantum simulation, *arXiv preprint arXiv:2009.02559* (2020).
- [99] J. Gibbs, K. Gili, Z. Holmes, B. Commeau, A. Arrasmith, L. Cincio, P. J. Coles, and A. Sornborger, Long-time simulations with high fidelity on quantum hardware, *arXiv preprint arXiv:2102.04313* (2021).
- [100] K. Bharti, A. Cervera-Lierta, T. H. Kyaw, T. Haug, S. Alperin-Lea, A. Anand, M. Degroote, H. Heimonen, J. S. Kottmann, T. Menke, W.-K. Mok, S. Sim, L.-C. Kwek, and A. Aspuru-Guzik, Noisy intermediate-scale quantum (nisq) algorithms, *arXiv preprint arXiv:2101.08448* (2021).
- [101] S. Efthymiou, S. Ramos-Calderer, C. Bravo-Prieto, A. Pérez-Salinas, D. García-Martín, A. Garcia-Saez, J. I. Latorre, and S. Carrazza, Qibo: a framework for quantum simulation with hardware acceleration, *arXiv preprint arXiv:2009.01845* (2020).
- [102] S. Efthymiou, S. Carrazza, S. Ramos, bpcarlos, AdrianPerezSalinas, D. García-Martín, Paul, J. Serrano, and atomicprinter, *qiboteam/qibo: Qibo 0.1.6-rc1* (2021).
- [103] A. Peruzzo, J. McClean, P. Shadbolt, M.-H. Yung, X.-Q. Zhou, P. J. Love, A. Aspuru-Guzik, and J. L. O’Brien, A variational eigenvalue solver on a photonic quantum processor, *Nature communications* **5**, 1 (2014).
- [104] C. Bravo-Prieto, J. Lumbrecas-Zarapico, L. Tagliacozzo, and J. I. Latorre, Scaling of variational quantum circuit depth for condensed matter systems, *Quantum* **4**, 272 (2020).
- [105] M. Consiglio, W. J. Chetcuti, C. Bravo-Prieto, S. Ramos-Calderer, A. Minguzzi, J. I. Latorre, L. Amico, and

- T. J. G. Apollaro, Variational quantum eigensolver for  $su(n)$  fermions, [arXiv preprint arXiv:2106.15552 \(2021\)](#).
- [106] Equations (19) are actually upper bounds for  $\dim(g_S)$ , however, the bounds get quickly saturated as  $n$  increases.
- [107] F. T. Chong, D. Franklin, and M. Martonosi, Programming languages and compiler design for realistic quantum hardware, *Nature* **549**, 180 (2017).
- [108] T. Häner, D. S. Steiger, K. Svore, and M. Troyer, A software methodology for compiling quantum programs, *Quantum Science and Technology* **3**, 020501 (2018).
- [109] D. Venturelli, M. Do, E. Rieffel, and J. Frank, Compiling quantum circuits to realistic hardware architectures using temporal planners, *Quantum Science and Technology* **3**, 025004 (2018).
- [110] S. Khatri, R. LaRose, A. Poremba, L. Cincio, A. T. Sornborger, and P. J. Coles, Quantum-assisted quantum compiling, *Quantum* **3**, 140 (2019).
- [111] K. Sharma, S. Khatri, M. Cerezo, and P. J. Coles, Noise resilience of variational quantum compiling, *New Journal of Physics* **22**, 043006 (2020).
- [112] H. Ma, C.-J. Huang, C. Chen, D. Dong, Y. Wang, R.-B. Wu, and G.-Y. Xiang, On compression rate of quantum autoencoders: Control design, numerical and experimental realization, [arXiv preprint arXiv:2005.11149 \(2020\)](#).
- [113] L. Schatzki, A. Arrasmith, P. J. Coles, and M. Cerezo, Entangled datasets for quantum machine learning, [arXiv preprint arXiv:2109.03400 \(2021\)](#).
- [114] N. Chan and M. K. Kwong, Hermitian matrix inequalities and a conjecture, *The American Mathematical Monthly* **92** (1985).
- [115] S. J. Glaser, U. Boscain, T. Calarco, C. P. Koch, W. Köckenberger, R. Kosloff, I. Kuprov, B. Luy, S. Schirmer, T. Schulte-Herbrüggen, *et al.*, Training schrödinger's cat: quantum optimal control, *The European Physical Journal D* **69**, 1 (2015).
- [116] A. Acín, I. Bloch, H. Buhrman, T. Calarco, C. Eichler, J. Eisert, D. Esteve, N. Gisin, S. J. Glaser, F. Jelezko, *et al.*, The quantum technologies roadmap: a european community view, *New Journal of Physics* **20**, 080201 (2018).
- [117] Z.-C. Yang, A. Rahmani, A. Shabani, H. Neven, and C. Chamon, Optimizing variational quantum algorithms using pontryagin's minimum principle, *Physical Review X* **7**, 021027 (2017).
- [118] D. Lu, K. Li, J. Li, H. Katiyar, A. J. Park, G. Feng, T. Xin, H. Li, G. Long, A. Brodutch, *et al.*, Enhancing quantum control by bootstrapping a quantum processor of 12 qubits, *npj Quantum Information* **3**, 1 (2017).
- [119] P. Rembold, N. Oshnik, M. M. Müller, S. Montangero, T. Calarco, and E. Neu, Introduction to quantum optimal control for quantum sensing with nitrogen-vacancy centers in diamond, *AVS Quantum Science* **2**, 024701 (2020).
- [120] J. P. Peterson, H. Katiyar, and R. Laflamme, Fast simulation of magnetic field gradients for optimization of pulse sequences, [arXiv preprint arXiv:2006.10133 \(2020\)](#).
- [121] D. Bluvstein, A. Omran, H. Levine, A. Keesling, G. Semeghini, S. Ebadi, T. T. Wang, A. A. Michailidis, N. Maskara, W. W. Ho, *et al.*, Controlling quantum many-body dynamics in driven rydberg atom arrays, *Science* **371**, 1355 (2021).
- [122] S. Ebadi, T. T. Wang, H. Levine, A. Keesling, G. Semeghini, A. Omran, D. Bluvstein, R. Samajdar, H. Pichler, W. W. Ho, *et al.*, Quantum phases of matter on a 256-atom programmable quantum simulator, *Nature* **595**, 227 (2021).
- [123] A. B. Magann, K. M. Rudinger, M. D. Grace, and M. Sarovar, Feedback-based quantum optimization, [arXiv preprint arXiv:2103.08619 \(2021\)](#).
- [124] M. Larocca and D. Wisniacki, Krylov-subspace approach for the efficient control of quantum many-body dynamics, *Physical Review A* **103**, 023107 (2021).
- [125] L. T. Brady, C. L. Baldwin, A. Bapat, Y. Kharkov, and A. V. Gorshkov, Optimal protocols in quantum annealing and quantum approximate optimization algorithm problems, *Physical Review Letters* **126**, 070505 (2021).
- [126] N. Wittler, F. Roy, K. Pack, M. Werninghaus, A. S. Roy, D. J. Egger, S. Filipp, F. K. Wilhelm, and S. Machnes, Integrated tool set for control, calibration, and characterization of quantum devices applied to superconducting qubits, *Physical Review Applied* **15**, 034080 (2021).
- [127] A. B. Magann, C. Arenz, M. D. Grace, T.-S. Ho, R. L. Kosut, J. R. McClean, H. A. Rabitz, and M. Sarovar, From pulses to circuits and back again: A quantum optimal control perspective on variational quantum algorithms, *PRX Quantum* **2**, 010101 (2021).
- [128] A. Smith, M. Kim, F. Pollmann, and J. Knolle, Simulating quantum many-body dynamics on a current digital quantum computer, *npj Quantum Information* **5**, 1 (2019).
- [129] M. Hsieh, R. Wu, and H. Rabitz, Topology of the quantum control landscape for observables, *The Journal of chemical physics* **130**, 104109 (2009).
- [130] T.-S. Ho, J. Dominy, and H. Rabitz, Landscape of unitary transformations in controlled quantum dynamics, *Physical Review A* **79**, 013422 (2009).
- [131] J. J. Meyer, Fisher Information in Noisy Intermediate-Scale Quantum Applications, *Quantum* **5**, 539 (2021).
- [132] D. P. Kingma and J. Ba, Adam: A method for stochastic optimization, in *Proceedings of the 3rd International Conference on Learning Representations (ICLR)* (2015).
- [133] K. Mitarai, M. Negoro, M. Kitagawa, and K. Fujii, Quantum circuit learning, *Physical Review A* **98**, 032309 (2018).
- [134] M. Schuld, V. Bergholm, C. Gogolin, J. Izaac, and N. Killoran, Evaluating analytic gradients on quantum hardware, *Physical Review A* **99**, 032331 (2019).
- [135] A. Mari, T. R. Bromley, and N. Killoran, Estimating the gradient and higher-order derivatives on quantum hardware, *Phys. Rev. A* **103**, 012405 (2021).

**SUPPLEMENTARY INFORMATION FOR “THEORY OF OVERPARAMETRIZATION IN QUANTUM NEURAL NETWORKS”**

In this Supplementary Information, we present detailed proofs of the theorems, and corollaries presented in the manuscript “*Theory of overparametrization in quantum neural networks*”. In addition, here we provide additional details and results for the numerical simulations.

**I. PRELIMINARIES**

Let us start by recalling that we consider the case when the QNN  $U(\boldsymbol{\theta})$  is a parametrized quantum circuit with an  $L$ -layered periodic structure of the form

$$U(\boldsymbol{\theta}) = \prod_{l=1}^L U_l(\boldsymbol{\theta}_l), \quad U_l(\boldsymbol{\theta}_l) = \prod_{k=1}^K e^{-i\theta_{lk}H_k}, \quad (42)$$

where the index  $l$  indicates the layer, and the index  $k$  spans the traceless Hermitian operators  $H_k$  that generate the unitaries in the ansatz. Here,  $\boldsymbol{\theta}_l = (\theta_{l1}, \dots, \theta_{lK})$  are the parameters in a single layer, and  $\boldsymbol{\theta} = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_L\}$  denotes the set of  $M = K \cdot L$  trainable parameters in the QNN. In this Supplementary Information, we make use of the following relabelling of the parameters  $\theta_{lk}$  and operators  $H_k$ :

$$U(\boldsymbol{\theta}) \equiv \prod_{j=1}^{LK} e^{-i\theta_j H_j}. \quad (43)$$

For convenience we also recall the following definitions:

**Definition 1** (Set of generators  $\mathcal{G}$ ). *Consider a parametrized quantum circuit of the form (42). The set of generators  $\mathcal{G} = \{H_k\}_{k=1}^K$  is defined as the set (of size  $|\mathcal{G}| = K$ ) of the Hermitian operators that generate the unitaries in a single layer of  $U(\boldsymbol{\theta})$ .*

And, the definition for the dynamical Lie Algebra:

**Definition 2** (Dynamical Lie Algebra (DLA)). *Consider a set of generators  $\mathcal{G}$  according to Definition 1. The DLA  $\mathfrak{g}$  is generated by repeated nested commutators of the operators in  $\mathcal{G}$ . That is,*

$$\mathfrak{g} = \text{span} \langle iH_1, \dots, iH_K \rangle_{Lie}, \quad (44)$$

where  $\langle S \rangle_{Lie}$  denotes the Lie closure, i.e., the set obtained by repeatedly taking the commutator of the elements in  $S$ .

**Invariant subspaces.** Consider now the case when the elements in the DLA share a symmetry (for simplicity we assume only one symmetry, although generalization to multiple symmetries is straightforward). That is, there exists a Hermitian operator  $\Sigma$  such that  $[\Sigma, g] = 0$  for all  $g \in \mathfrak{g}$ . If  $\Sigma$  has  $N$  distinct eigenvalues, then the DLA has the form  $\mathfrak{g} = \bigoplus_{m=1}^N \mathfrak{g}_m$ . This imposes a partition of Hilbert space  $\mathcal{H} = \bigoplus_{m=1}^N \mathcal{H}_m$  where each subspace  $\mathcal{H}_m$  of dimension  $d_m$  is invariant under  $\mathfrak{g}$ .

Let us introduce some notation. Consider the  $d \times d_m$  matrix that results from horizontally stacking the eigenvectors of  $\Sigma$  associated with the  $m$ -th eigenvalue (of degeneracy  $g_m$ )

$$Q_m^\dagger = \begin{bmatrix} \vdots & \vdots & \vdots \\ |v_1\rangle, & |v_2\rangle, & \dots, |v_{g_m}\rangle \\ \vdots & \vdots & \vdots \end{bmatrix}, \quad (45)$$

such that  $Q_m$  maps vectors from  $\tilde{H}$  to  $\tilde{H}_m$ . These satisfy

$$Q_m Q_n^\dagger = \mathbb{1}_{d_m} \delta_{mn}, \quad Q_m^\dagger Q_m = \mathbb{P}_m, \quad (46)$$

where  $\mathbb{P}_m$  are projectors onto the  $m$ -th eigenspace, such that  $\sum_{m=1}^N \mathbb{P}_m = \mathbb{1}$ . Let us now use the notation

$$|\psi\rangle^{(m)} = Q_m |\psi\rangle, \quad A^{(m)} = Q_m A Q_m^\dagger, \quad (47)$$

to denote the  $d_m$ -dimensional reduced states and operators, respectively. Recall that, since any unitary  $U \in \mathbb{G}$  produced by such a system is block diagonal, we can write  $U = \sum_m \mathbb{P}_m U \mathbb{P}_m$ . Also, let us note that if  $A$  is Hermitian, then  $A^{(k)}$  is also Hermitian.

## II. PROOF OF THEOREM 1

In the following we provide a proof for Theorem 1. Let us first recall the definition of the Quantum Fisher Information Matrix (QFIM). The QFIM entries are given by

$$[F(\boldsymbol{\theta})]_{jk} = 4\text{Re}[\langle \partial_j \psi(\boldsymbol{\theta}) | \partial_k \psi(\boldsymbol{\theta}) \rangle - \langle \partial_j \psi(\boldsymbol{\theta}) | \psi(\boldsymbol{\theta}) \rangle \langle \psi(\boldsymbol{\theta}) | \partial_k \psi(\boldsymbol{\theta}) \rangle], \quad (48)$$

where  $|\psi(\boldsymbol{\theta})\rangle = U(\boldsymbol{\theta}) |\psi\rangle$ . Here we also denote where  $|\partial_i \psi(\boldsymbol{\theta})\rangle = \partial |\psi(\boldsymbol{\theta})\rangle / \partial \theta_i = \partial_i |\psi(\boldsymbol{\theta})\rangle$  for  $\theta_i \in \boldsymbol{\theta}$ .

We now restate Theorem 1 for convenience.

**Theorem 1.** *For each state  $|\psi_\mu\rangle$  in the training set  $\mathcal{S}$ , the maximum rank  $R_\mu$  of its associated QFIM (defined in Eq. (48)) is upper bounded as*

$$R_\mu \leq \dim(\mathfrak{g}_\mathcal{S}). \quad (49)$$

*Proof.* Let us first note that the partial derivatives of the parametrized state are

$$|\partial_j \psi(\boldsymbol{\theta})\rangle = \partial_j (U(\boldsymbol{\theta}) |\psi\rangle) = -iU(\boldsymbol{\theta}) \tilde{H}_j |\psi\rangle, \quad 1 \leq j \leq M, \quad (50)$$

where

$$\tilde{H}_j = U_{1 \rightarrow j}^\dagger H_j U_{1 \rightarrow j}, \quad \text{and where } U_{1 \rightarrow j} = U_j \cdots U_1. \quad (51)$$

Thus,  $U_{1 \rightarrow j}$  is the propagator up to the  $j$ -th layer in the circuit and we are labeling  $H_j$  modulo  $|\mathcal{G}|$ , e.g.  $H_{|\mathcal{G}|} = H_1$ , i.e. the first generator.

Next, let us consider the case then the DLA has a symmetry (see the Preliminaries section above) and that all states in the training set belong to the  $m$ -th invariant subspace of the symmetry. We denote by  $\mathfrak{g}_\mathcal{S}$  the DLA associated with said symmetry respected by the training set, by  $\mathcal{H}_\mathcal{S}$  the corresponding Hilbert space, and by  $|\psi\rangle^{(m)} = Q_m |\psi\rangle$  the projected state according to Eq. (47). Then, for any pair of states  $|\phi\rangle, |\chi\rangle \in \mathcal{S}$ , their overlap can be described in terms of the overlap between the corresponding  $d_\mathcal{S}$ -dimensional reduced states

$$\langle \chi | \phi \rangle = \langle \chi^{(\mathcal{S})} | \phi^{(\mathcal{S})} \rangle. \quad (52)$$

Then, it is straightforward to see that  $|\psi\rangle, |\psi(\boldsymbol{\theta})\rangle$  and  $|\partial_j \psi(\boldsymbol{\theta})\rangle$  also belong in  $\mathcal{H}_\mathcal{S}$ . Hence, the overlaps in Eq. (48) can be computed in terms of their reduced counterparts  $|\psi(\boldsymbol{\theta})^{(\mathcal{S})}\rangle$  and  $|\partial_j \psi(\boldsymbol{\theta})^{(\mathcal{S})}\rangle$ . In the following, we will work with everything reduced to such subspace, but to simplify the notation, we will omit the  $\mathcal{S}$  superscript everywhere. For example, whenever we write operator  $O$ , we actually mean  $O^{(\mathcal{S})} \in \mathbb{C}^{d_\mathcal{S} \times d_\mathcal{S}}$ .

Using the explicit expression for the partial derivatives, we find that the first term in Eq. (48) is

$$\text{Re}[\langle \partial_j \psi(\boldsymbol{\theta}) | \partial_k \psi(\boldsymbol{\theta}) \rangle] = \text{Re}[i(-i) \langle \psi | \tilde{H}_j U(\boldsymbol{\theta})^\dagger U(\boldsymbol{\theta}) \tilde{H}_k | \psi \rangle] = \text{Re}[\langle \psi | \tilde{H}_j \tilde{H}_k | \psi \rangle]. \quad (53)$$

Choosing an orthonormal basis containing  $|\psi\rangle$  we can rewrite this term as

$$\text{Re}[\langle \partial_j \psi(\boldsymbol{\theta}) | \partial_k \psi(\boldsymbol{\theta}) \rangle] = \sum_m \text{Re}[\langle \psi | \tilde{H}_j | m \rangle \langle m | \tilde{H}_k | \psi \rangle] = \text{Re}[\langle \psi | \tilde{H}_j | \psi \rangle \langle \psi | \tilde{H}_k | \psi \rangle] + \sum_{m \neq \psi} \text{Re}[\langle \psi | \tilde{H}_j | m \rangle \langle m | \tilde{H}_k | \psi \rangle]. \quad (54)$$

Proceeding similarly, we find for the second term in Eq. (48)

$$\begin{aligned} \text{Re}[\langle \partial_j \psi(\boldsymbol{\theta}) | \psi(\boldsymbol{\theta}) \rangle \langle \psi(\boldsymbol{\theta}) | \partial_k \psi(\boldsymbol{\theta}) \rangle] &= \text{Re}[-i] \langle \psi | \tilde{H}_j U(\boldsymbol{\theta})^\dagger | \psi(\boldsymbol{\theta}) \rangle \langle \psi(\boldsymbol{\theta}) | U(\boldsymbol{\theta}) \tilde{H}_k | \psi \rangle \\ &= \text{Re}[\langle \psi | \tilde{H}_j | \psi \rangle \langle \psi | \tilde{H}_k | \psi \rangle]. \end{aligned} \quad (55)$$

Combining these results we can express the matrix elements of the QFIM as

$$[F(\boldsymbol{\theta})]_{jk} = 4\text{Re}[\langle \psi | \tilde{H}_j \tilde{H}_k | \psi \rangle] - 4\text{Re}[\langle \psi | \tilde{H}_j | \psi \rangle \langle \psi | \tilde{H}_k | \psi \rangle] = 4 \sum_{m \neq \psi} \text{Re}[\langle \psi | \tilde{H}_j | m \rangle \langle m | \tilde{H}_k | \psi \rangle]. \quad (56)$$

Note here that this equations also allows us to express the QFIM elements as  $[F]_{jk} = 4\text{Re}[\text{Cov}_{|\psi\rangle}(\tilde{H}_j, \tilde{H}_k)]$ . Then, defining the vectors  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  with components

$$R_{mn}(j) = \text{Re}[\langle m | \tilde{H}_j | n \rangle], \quad I_{mn}(j) = \text{Im}[\langle m | \tilde{H}_j | n \rangle], \quad (57)$$

we can express (56) as

$$[F(\boldsymbol{\theta})]_{jk} = 4 \sum_{m \neq \psi} R_{\psi m}(j) R_{m \psi}(k) - I_{\psi m}(j) I_{m \psi}(k) = 4 \sum_{m \neq \psi} R_{\psi m}(j) R_{\psi m}(k) + I_{\psi m}(j) I_{\psi m}(k), \quad (58)$$

where the second equality follows from the fact that  $R_{mn}(j) = R_{nm}(j)$ , while  $I_{mn}(j) = -I_{nm}(j)$ . Thus, one can finally express the QFIM as a sum of  $2d - 2$  rank-one matrices

$$F(\boldsymbol{\theta}) = -2 \sum_{m=1, m \neq \psi}^d (\mathbf{R}_{m\psi} \cdot \mathbf{R}_{m\psi}^\top + \mathbf{I}_{m\psi} \cdot \mathbf{I}_{m\psi}^\top). \quad (59)$$

Here we recall that, by definition,  $H_j$  are elements in the DLA  $\mathfrak{g}_S$ . Then, since the unitaries  $U$  are elements of the dynamical Lie group  $\mathbb{G}_S$  generated by  $\mathfrak{g}_S$ , conjugating  $H_j$  by any unitary  $U$  results in another element in  $\mathfrak{g}_S$ . That is:  $\forall U \in \mathbb{G}_S$ , and  $\forall H_i \in \mathfrak{g}_S$  we have  $UH_j U^\dagger \in \mathfrak{g}_S$ . Then, by repeating this argument  $j$  times, we find that  $\tilde{H}_j \in \mathfrak{g}_S$ , where  $\tilde{H}_j$  was defined in Eq. (51).

Letting  $\{S_\nu\}_{\nu=1}^{\dim(\mathfrak{g})}$  be a basis of  $\mathfrak{g}$ , we can express

$$\tilde{H}_j = \sum_{\nu=1}^{\dim(\mathfrak{g}_S)} a_\nu(j) S_\nu, \quad (60)$$

with  $a_\nu$  real coefficients. Using this fact, we can expand  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  in the following ways:

$$\begin{aligned} R_{mn}(j) &= \text{Re} \left[ \sum_{\nu}^{\dim(\mathfrak{g}_S)} \langle m | a_\nu(j) S_\nu | n \rangle \right] = \sum_{\nu}^{\dim(\mathfrak{g}_S)} \text{Re}[\langle m | S_\nu | n \rangle] a_\nu(j), \\ I_{mn}(j) &= \text{Im} \left[ \sum_{\nu}^{\dim(\mathfrak{g}_S)} \langle m | a_\nu(j) S_\nu | n \rangle \right] = \sum_{\nu}^{\dim(\mathfrak{g}_S)} \text{Im}[\langle m | S_\nu | n \rangle] a_\nu(j). \end{aligned} \quad (61)$$

More succinctly, we find

$$\mathbf{R}_{mn} = \sum_{\nu=1}^{\dim(\mathfrak{g}_S)} \text{Re} \langle m | S_\nu | n \rangle \mathbf{a}_\nu, \quad \text{and} \quad \mathbf{I}_{mn} = \sum_{\nu=1}^{\dim(\mathfrak{g}_S)} \text{Im} \langle m | S_\nu | n \rangle \mathbf{a}_\nu. \quad (62)$$

These equations show that the vectors  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  can be expressed as a linear combination of  $\dim(\mathfrak{g}_S)$  other vectors  $\{\mathbf{a}_\nu\}$ . Then, since the  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  generate the  $2d - 2$  rank-one matrices in the QFIM, we have that  $F(\boldsymbol{\theta})$  has a support on a subspace with a basis that has, at most,  $\dim(\mathfrak{g}_S)$  elements. Thus, we find

$$\text{rank}[F_\mu(\boldsymbol{\theta})] \leq \dim(\mathfrak{g}_S), \quad (63)$$

where we have recovered the  $\mu$  dependence of the QFIM.  $\square$

Let us here note that the proof of Theorem 1 holds for all states in the training set  $|\psi_\mu\rangle \in \mathcal{S}$ . Thus, for all  $|\psi_\mu\rangle$  we know that the associated QFIM  $F_\mu(\boldsymbol{\theta})$  has a column space contained within some fixed  $\dim(\mathfrak{g}_S)$  dimensional space. More precisely, from the previous proof, we have that the following Proposition holds.

**Proposition 1.** *There is some vector space spanned by  $\dim(\mathfrak{g}_S)$  vectors  $\{\mathbf{a}_\nu\}_{\nu=1}^{\dim(\mathfrak{g}_S)}$ , such that for any state in the training set  $|\psi_\mu\rangle \in \mathcal{S}$ , the associated QFIM  $F_\mu(\boldsymbol{\theta})$  has a column space contained within this vector space.*

We will make use of this proposition in the following section.

### III. PROOF OF THEOREM 2

In the following, we provide a proof for Theorem 2. For convenience, we here recall the definition of overparametrization as well as the statement for Theorem 2.

**Definition 3** (Overparametrization). *A QNN is said to be overparametrized if the number of parameters  $M$  is such that the QFI matrices, for all the states in the training set, simultaneously saturate their achievable rank  $R_\mu$  at least in one point of the loss landscape. That is, if increasing the number of parameters past some minimal (critical) value  $M_c$  does not further increase the rank of any QFIM:*

$$\max_{M \geq M_c, \boldsymbol{\theta}} \text{rank}[F_\mu(\boldsymbol{\theta})] = R_\mu. \quad (64)$$

Let us also recall two definitions of a QNN's effective dimension. First, following [58], we can define the average effective quantum dimension of a QNN:

$$D_1(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{i=1}^M \mathcal{I}(\lambda_\mu^i(\boldsymbol{\theta})) \right], \quad (65)$$

where  $\lambda_\mu^i(\boldsymbol{\theta})$  are the eigenvalues of  $F_\mu(\boldsymbol{\theta})$ , and where  $\mathcal{I}(x) = 0$  for  $x = 0$ , and  $\mathcal{I}(x) = 1$  for  $x \neq 0$ . Here the expectation value is taken over the probability distribution that samples input states from the dataset.

The second definition follows from [24]. In the  $n \rightarrow \infty$  limit, the effective quantum dimension of [24] converges to

$$D_2 = \max_{\boldsymbol{\theta}} \left( \text{rank} \left[ \tilde{F}(\boldsymbol{\theta}) \right] \right), \quad (66)$$

where  $\tilde{F}(\boldsymbol{\theta})$  is the classical Fisher Information matrix obtained as

$$\tilde{F}(\boldsymbol{\theta}) = \mathbb{E} \left[ \frac{\partial \log(p(|\psi_\mu\rangle, y_\mu; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \frac{\partial \log(p(|\psi_\mu\rangle, y_\mu; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}^T \right]. \quad (67)$$

Here,  $p(|\psi\rangle, y; \boldsymbol{\theta})$ , describes the joint relationship between an input  $|\psi\rangle$  and an output  $y$  of the QNN. In addition, the expectation value is taken over the probability distribution that samples input states from the dataset.

Then, consider the following theorem in the main text.

**Theorem 2.** *The model capacity, as quantified by the effective dimensions of Eqs. (65) or (66), is upper bounded as*

$$D_1(\boldsymbol{\theta}) \leq \dim(\mathfrak{g}_S), \quad D_2 \leq \dim(\mathfrak{g}_S). \quad (68)$$

Moreover, when the QNN is overparametrized according to Definition 3,  $D_1(\boldsymbol{\theta})$  achieves its maximum value on at least one point of the landscape.

*Proof.* When the QNN is overparametrized according to Definition 3, there exists some  $\theta$  such that the ranks of the QFIMs are maximized. That is,  $\text{rank}[F_\mu(\theta)] = R_\mu$ .

Note that the effective dimension  $D_1(\theta)$  can be expressed as  $D_1(\theta) = \mathbb{E}[\text{rank}[F_\mu(\theta)]]$ . Then, since the ranks are maximal at the overparametrization, so is  $D_1(\theta)$ . More precisely, we have  $D_1(\theta) = \mathbb{E}[R_\mu]$ . Additionally, as shown in Theorem 1,  $\text{rank}[F_\mu(\theta)] \leq \dim(\mathfrak{g}_S)$  for all  $\theta$  and  $|\psi_\mu\rangle \in \mathcal{S}$ , so  $D_1(\theta) \leq \dim(\mathfrak{g}_S)$ .

Next we consider the effective dimension  $D_2$  (Eq. (66)).  $D_2$  specifically quantifies the maximal rank of the expectation value of the classical Fisher information matrices  $\tilde{F}_\mu(\theta)$ . Because the operator  $F_\mu(\theta) - \tilde{F}_\mu(\theta)$  is positive semidefinite ([131], Section 5), the following holds:

$$\tilde{F}(\theta) = \mathbb{E}_\mu[\tilde{F}_\mu(\theta)] \leq \mathbb{E}_\mu[F_\mu(\theta)]. \quad (69)$$

In addition for any two symmetric matrices  $A$  and  $B$ , having  $A \leq B$  implies that  $A^0 \leq B^0$  ([114], Theorem 3). Thus,

$$(\tilde{F}(\theta))^0 \leq (\mathbb{E}_\mu[F_\mu(\theta)])^0, \quad (70)$$

implying that

$$\text{rank}[\tilde{F}(\theta)] = \text{Tr}[(\tilde{F}(\theta))^0] \leq \text{Tr}[(\mathbb{E}_\mu[F_\mu(\theta)])^0] = \text{rank}[\mathbb{E}_\mu[F_\mu(\theta)]]. \quad (71)$$

By applying Proposition 1 to  $\mathbb{E}_\mu[F_\mu(\theta)]$ , we arrive at the desired result:

$$\text{rank}[\tilde{F}(\theta)] \leq \text{rank}[\mathbb{E}_\mu[F_\mu(\theta)]] \leq \dim(\mathfrak{g}_S). \quad (72)$$

□

#### IV. PROOF OF THEOREM 3

In the following we prove Theorem 3, which bounds the rank of the Hessian for an observable minimization loss function of the form

$$\mathcal{L}(\theta) = \sum_{|\psi_\mu\rangle \in \mathcal{S}} c_\mu \text{Tr}[U(\theta) |\psi_\mu\rangle \langle \psi_\mu| U^\dagger(\theta) O], \quad (73)$$

at its optimum. Here, the terms  $c_\mu$  are real coefficients associated with each state  $|\psi_\mu\rangle$  in  $\mathcal{S}$ , and where the operator  $O$  is Hermitian. Let us restate the theorem for convenience.

**Theorem 3.** *Let  $\nabla^2 \mathcal{L}(\theta_*)$  be the Hessian for a loss function of the form of Eq. (73) evaluated at the optimum set of parameters  $\theta_*$ . Then, its rank is upper bounded as*

$$\text{rank}[\nabla^2 \mathcal{L}(\theta_*)] \leq \min\{\dim(\mathfrak{g}_S), 2dr - r^2 - r\}, \quad (74)$$

where  $r = \min\{\text{rank}[\sum_\mu c_\mu |\psi_\mu\rangle \langle \psi_\mu|], \text{rank}[O]\}$ , and  $d$  is the Hilbert space dimension.

*Proof.* Let us define  $\rho = \sum_\mu c_\mu |\psi_\mu\rangle \langle \psi_\mu|$ . First, the gradient has the following form:

$$\partial_j \mathcal{L}(\theta) = -\text{Tr}[\partial_j U(\theta) \rho U^\dagger(\theta) O - U(\theta) \rho \partial_j U^\dagger(\theta) O] = i \text{Tr}[\tilde{H}_j, \rho] O_f, \quad (75)$$

where  $\tilde{H}_j$  is defined in Eq. (51), and where we defined  $O_f = U(\theta)^\dagger O U(\theta)$  (we henceforth drop the explicit dependence on  $\theta$ ). Going forward, we similarly drop the explicit dependence on  $\theta$  on terms which are not being differentiated.

If we assume that  $i \leq j$ , and note that  $\partial_i \tilde{H}_j = i[\tilde{H}_i, \tilde{H}_j]$  in this case, then we can express the matrix elements of the Hessian as

$$\begin{aligned}
\partial_i \partial_j \mathcal{L}(\boldsymbol{\theta}) &= i \partial_i \text{Tr} \left[ U(\boldsymbol{\theta}) \tilde{H}_j(\boldsymbol{\theta}) \rho U^\dagger(\boldsymbol{\theta}) O - U(\boldsymbol{\theta}) \rho \tilde{H}_j(\boldsymbol{\theta}) U^\dagger(\boldsymbol{\theta}) O \right] \\
&= \text{Tr} \left[ \tilde{H}_i \tilde{H}_j \rho O_f \right] - \text{Tr} \left[ [\tilde{H}_i, \tilde{H}_j] \rho O_f \right] - \text{Tr} \left[ \tilde{H}_j \rho \tilde{H}_i O_f \right] \\
&\quad - \text{Tr} \left[ \tilde{H}_i \rho \tilde{H}_j O_f \right] + \text{Tr} \left[ \rho [\tilde{H}_i, \tilde{H}_j] O_f \right] + \text{Tr} \left[ \rho \tilde{H}_j \tilde{H}_i O_f \right] \\
&= 2 \text{Re} \left[ \text{Tr} \left[ \rho \tilde{H}_i \tilde{H}_j O_f \right] \right] - 2 \text{Re} \left[ \text{Tr} \left[ \tilde{H}_i \rho \tilde{H}_j O_f \right] \right].
\end{aligned} \tag{76}$$

We now evaluate the Hessian at the optimum  $\boldsymbol{\theta}_*$ . Here, the propagator has the form [61]  $U(\boldsymbol{\theta}_*) = R^\dagger Q$  for unitaries  $R$  and  $Q$  that respectively diagonalize  $\rho$  and  $O$ , i.e.  $\rho = Q^\dagger e Q$  and  $O = R^\dagger o R$ , such that  $e$  ( $o$ ) is a diagonal matrix containing the eigenvalues of  $\rho$  ( $O$ ) in decreasing (increasing) order. Therefore, we can rewrite Eq. (76) at the optimum as

$$\begin{aligned}
\partial_i \partial_j \mathcal{L}(\boldsymbol{\theta}_*) &= 2 \text{Re} \left[ \text{Tr} \left[ \tilde{H}_i \tilde{H}_j O_f \rho \right] \right] - 2 \text{Re} \left[ \text{Tr} \left[ \tilde{H}_i \rho \tilde{H}_j O_f \right] \right] \\
&= 2 \text{Re} \left[ \text{Tr} \left[ \tilde{H}_i \tilde{H}_j Q^\dagger o e Q \right] \right] - 2 \text{Re} \left[ \text{Tr} \left[ \tilde{H}_i Q^\dagger e Q \tilde{H}_j Q^\dagger o Q \right] \right] \\
&= 2 \sum_{m,n=1}^d (o_m e_m - o_m e_n) \text{Re} \left[ \langle m | Q \tilde{H}_i Q^\dagger | n \rangle \langle n | Q \tilde{H}_j Q^\dagger | m \rangle \right] \\
&= 2 \sum_{m,n=1}^d (o_m e_m - o_m e_n) (R'_{mn}(i) R'_{mn}(j) + I'_{mn}(i) I'_{mn}(j)),
\end{aligned} \tag{77}$$

where we have used  $O_f = Q^\dagger o Q$  at the optimum and defined  $R'_{mn}(j) = \text{Re}[\langle m | Q \tilde{H}_j Q^\dagger | n \rangle]$  and  $I'_{mn}(j) = \text{Im}[\langle m | Q \tilde{H}_j Q^\dagger | n \rangle]$ . Because Eq. (77) is symmetric in indices  $i$  and  $j$ , we can remove the assumption that  $i \leq j$ . By following a proof similar to that in Theorem 1, we have an upper bound of  $\dim(\mathfrak{g}_S)$  on the rank of the Hessian  $\nabla^2 \mathcal{L}(\boldsymbol{\theta}_*)$  because  $R'_{mn}(\cdot)$  and  $I'_{mn}(\cdot)$  reside in a  $\dim(\mathfrak{g}_S)$  dimensional space; see Eq. (62).

We will now establish the additional  $2dr - r^2 - r$  upper bound stated in the theorem. We will use the short hand  $r(o)$  and  $r(e)$  for ranks of  $o$  and  $e$ , respectively. Assume that  $r(e) \leq r(o)$  (the case of  $r(o) \leq r(e)$  proceeds similarly). We would like to split Eq. (77) into disjoint summations over  $m$  and  $n$ . Toward that goal, let us define  $t_{mn}(i, j) = R'_{mn}(i) R'_{mn}(j) + I'_{mn}(i) I'_{mn}(j)$  to rewrite Eq. (77):

$$\begin{aligned}
\frac{1}{2} \partial_i \partial_j \mathcal{L}(\boldsymbol{\theta}_*) &= \sum_m o_m e_m t_{mm}(i, j) + \sum_{m>n} o_m e_m t_{mn}(i, j) + \sum_{m<n} o_m e_m t_{mn}(i, j) \\
&\quad - \sum_m o_m e_n t_{mn}(i, j) - \sum_{m>n} o_m e_n t_{mn}(i, j) - \sum_{m<n} o_m e_n t_{mn}(i, j) \\
&= \sum_{m>n, m \leq r(e)} o_m e_m t_{mn}(i, j) + \sum_{m<n, m \leq r(e)} o_m e_m t_{mn}(i, j) - \sum_{m>n, m \leq r(o), n \leq r(e)} o_m e_n t_{mn}(i, j) \\
&\quad - \sum_{m<n, n \leq r(e)} o_m e_n t_{mn}(i, j),
\end{aligned} \tag{78}$$

where in the first equality we have simply split the sums among  $m = n$ ,  $m > n$ , and  $m < n$ . Then, in the second sum we have attached more specific subscripts to the summation and used the fact that  $r(e) \leq r(o)$ . We now combine the sums over  $m > n$  (and the same for  $m < n$ , separately) to arrive at



$$\begin{aligned}
\frac{1}{2}\partial_i\partial_j\mathcal{L}(\boldsymbol{\theta}_*) &= \sum_{m>n;m,n\leq r(e)} (o_m e_m - o_m e_n)t_{mn}(i,j) + \sum_{m<n;m,n\leq r(e)} (o_m e_m - o_m e_n)t_{mn}(i,j) \\
&- \sum_{m>n,r(e)<m\leq r(o),n\leq r(e)} o_m e_n t_{mn}(i,j) + \sum_{m<n,m\leq r(e)} o_m e_m t_{mn}(i,j),
\end{aligned} \tag{79}$$

where the second summation contains the leftover terms when combining over  $m > n$  and the fourth summations contains the leftover terms when combining over  $m < n$ . Note that  $R'_{mn} = R'_{nm}$  and  $I'_{mn} = -I'_{nm}$ . This means that we can also combine more terms between the first and second summations, and also combine terms between the third and fourth summations. By rewriting terms so that  $m > n$  and combining, we arrive at

$$\begin{aligned}
\frac{1}{2}\partial_i\partial_j\mathcal{L}(\boldsymbol{\theta}_*) &= \sum_{m>n;m,n\leq r(e)} (o_m e_m - o_m e_n + o_n e_n - o_n e_m)t_{mn}(i,j) \\
&- \sum_{m>n,r(e)<m\leq r(o),n\leq r(e)} (-o_m e_n + o_n e_n)t_{mn}(i,j) \\
&+ \sum_{m>n,m>r(o),n\leq r(e)} o_n e_n t_{mn}(i,j).
\end{aligned} \tag{80}$$

As a result, the Hessian can be expressed as

$$\begin{aligned}
\frac{1}{2}\nabla^2\mathcal{L}(\boldsymbol{\theta}_*) &= \sum_{m>n;m,n\leq r(e)} (o_m e_m - o_m e_n + o_n e_n - o_n e_m)(\mathbf{R}_{mn} \cdot \mathbf{R}_{mn}^\top + \mathbf{I}_{mn} \cdot \mathbf{I}_{mn}^\top) \\
&- \sum_{m>n,r(e)<m\leq r(o),n\leq r(e)} (-o_m e_n + o_n e_n)(\mathbf{R}_{mn} \cdot \mathbf{R}_{mn}^\top + \mathbf{I}_{mn} \cdot \mathbf{I}_{mn}^\top) \\
&+ \sum_{m>n,m>r(o),n\leq r(e)} o_n e_n (\mathbf{R}_{mn} \cdot \mathbf{R}_{mn}^\top + \mathbf{I}_{mn} \cdot \mathbf{I}_{mn}^\top),
\end{aligned} \tag{81}$$

where the  $j$ 'th entry of  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  are  $R'_{mn}(j)$  and  $I'_{mn}(j)$ , respectively. Now each summation is completely disjoint over  $(m, n)$  pairs, so the remaining projectors,  $\mathbf{R}_{mn} \cdot \mathbf{R}_{mn}^\top$  and  $\mathbf{I}_{mn} \cdot \mathbf{I}_{mn}^\top$ , are those such that  $m > n$  and  $n \leq r(e)$ . This gives an upper bound on the rank of the Hessian when  $r(e) \leq r(o)$  as

$$\text{rank}(\nabla^2\mathcal{L}(\boldsymbol{\theta}_*)) \leq 2\binom{r(e)}{2} + 2r(e)(d - r(e)) = r(e)^2 - r(e) + 2r(e)(d - r(e)). \tag{82}$$

A similar analysis for the case of  $r(o) \leq r(e)$  reveals

$$\text{rank}(\nabla^2\mathcal{L}(\boldsymbol{\theta}_*)) \leq r(o)^2 - r(o) + 2r(o)(d - r(o)). \tag{83}$$

Thus, defining  $r = \min\{r(e), r(o)\}$ , we have an upper bound on the rank of the Hessian as

$$\text{rank}(\nabla^2\mathcal{L}(\boldsymbol{\theta}_*)) \leq 2dr - r^2 - r. \tag{84}$$

□

## V. PROOF OF THEOREM 4

Here we present a proof for Theorem 4, which upper bounds the rank of the Hessian (evaluated at the solution) for a unitary compilation task. Here the goal is to train a QNN so that its action matches that of a target unitary  $V$ . We

consider two possible loss functions for this task

$$\mathcal{L}_1(\boldsymbol{\theta}) = 2d - 2\text{Re}[\text{Tr}[V^\dagger U(\boldsymbol{\theta})]], \quad \text{and} \quad \mathcal{L}_2(\boldsymbol{\theta}) = 1 - \frac{1}{d^2} |\text{Tr}[V^\dagger U(\boldsymbol{\theta})]|^2. \quad (85)$$

Where  $\mathcal{L}_1(\boldsymbol{\theta})$  is minimized if  $U(\boldsymbol{\theta}) = V$ , while  $\mathcal{L}_2(\boldsymbol{\theta})$  is minimized if  $U(\boldsymbol{\theta}) = e^{i\phi}V$ , for some any phase  $\phi$ .

We recall now the statement of Theorem 4,:

**Theorem 4.** *Consider the loss functions for a unitary compilation task*

$$\mathcal{L}_1(\boldsymbol{\theta}) = 2d - 2\text{Re}[T(\boldsymbol{\theta})], \quad \text{and} \quad \mathcal{L}_2(\boldsymbol{\theta}) = 1 - \frac{1}{d^2} |T(\boldsymbol{\theta})|^2,$$

where  $T(\boldsymbol{\theta}) = \text{Tr}[V^\dagger U(\boldsymbol{\theta})]$  for a target unitary  $V$ . Then, let  $\nabla^2 \mathcal{L}_1(\boldsymbol{\theta}_*)$  and  $\nabla^2 \mathcal{L}_2(\boldsymbol{\theta}_*)$  be the Hessians for the loss functions  $\mathcal{L}_1(\boldsymbol{\theta})$  and  $\mathcal{L}_2(\boldsymbol{\theta})$ , respectively evaluated at their solutions  $U(\boldsymbol{\theta}_*) = V$  and  $U(\boldsymbol{\theta}_*) = e^{i\phi}V$ . Then, the maximal ranks of  $\nabla^2 \mathcal{L}_1(\boldsymbol{\theta}_*)$  and  $\nabla^2 \mathcal{L}_2(\boldsymbol{\theta}_*)$  are such that  $\text{rank}[\nabla^2 \mathcal{L}_1(\boldsymbol{\theta}_*)], \text{rank}[\nabla^2 \mathcal{L}_2(\boldsymbol{\theta}_*)] \leq \dim(\mathfrak{g})$ .

*Proof.* Let us begin with  $\mathcal{L}_1(\boldsymbol{\theta})$ . The gradient of this loss function is

$$\begin{aligned} \nabla \mathcal{L}_1(\boldsymbol{\theta}) &= -2\text{Re}[\nabla T(\boldsymbol{\theta})] \\ &= -2\text{Re}[\text{Tr}[V^\dagger \nabla U(\boldsymbol{\theta})]] \\ &= 2\text{Im}[\text{Tr}[V^\dagger U(\boldsymbol{\theta}) \tilde{\mathbf{H}}(\boldsymbol{\theta})]], \end{aligned} \quad (86)$$

where  $\tilde{\mathbf{H}}(\boldsymbol{\theta}) = (\tilde{H}_1(\boldsymbol{\theta}), \dots, \tilde{H}_M(\boldsymbol{\theta}))^\top$ , and where  $\tilde{H}_j(\boldsymbol{\theta})$  was defined in Eq. (51). Similarly, assuming  $i \leq j$ , we find for the matrix elements of the Hessian

$$\begin{aligned} \partial_i \partial_j \mathcal{L}_1(\boldsymbol{\theta}) &= 2\text{Im}[\text{Tr}[V^\dagger \partial_i U(\boldsymbol{\theta}) \tilde{H}_j + V^\dagger U(\boldsymbol{\theta}) \partial_i \tilde{H}_j]] \\ &= -2\text{Re}[\text{Tr}[V^\dagger U(\boldsymbol{\theta}) \tilde{H}_i \tilde{H}_j - V^\dagger U(\boldsymbol{\theta}) [\tilde{H}_i, \tilde{H}_j]]] \\ &= -2\text{Re}[\text{Tr}[V^\dagger U(\boldsymbol{\theta}) \tilde{H}_i \tilde{H}_j]]. \end{aligned} \quad (87)$$

Evaluating at any optimal set of parameters, that is, such that  $U(\boldsymbol{\theta}_*) = V$ , we find that Eq. (87) is symmetric in indices  $i$  and  $j$ . Thus, we can remove the assumption that  $i \leq j$  and express the Hessian more succinctly:

$$\nabla^2 \mathcal{L}_1(\boldsymbol{\theta}_*) = -2\text{Re}[\text{Tr}[\tilde{\mathbf{H}} \cdot \tilde{\mathbf{H}}^\top]] = -2 \sum_{m,n=1}^d \mathbf{R}_{mn} \cdot \mathbf{R}_{mn}^\top + \mathbf{I}_{mn} \cdot \mathbf{I}_{mn}^\top \quad (88)$$

where  $\mathbf{R}_{mn} = \text{Re}[\langle m | \tilde{\mathbf{H}} | n \rangle]$  and  $\mathbf{I}_{mn} = \text{Im}[\langle m | \tilde{\mathbf{H}} | n \rangle]$ . Hence, we again find that the Hessian is a sum of  $d^2$  rank-one matrices. We note that from here onward, we drop the  $\boldsymbol{\theta}$  dependence of  $\tilde{\mathbf{H}}$ .

Then, following a proof similar to the one used in proving Theorem 1, we know that each of the vectors generating the matrices  $\mathbf{R}_{mn}$  and  $\mathbf{I}_{mn}$  can be written as a linear combination of  $\dim(\mathfrak{g})$  other vectors  $\{\mathbf{a}_\nu\}_{\nu=1}^{\dim(\mathfrak{g})}$ . Thus, the rank of the Hessian of  $\nabla^2 \mathcal{L}_1(\boldsymbol{\theta})$  at the optimum is upper bounded by larger than  $\dim(\mathfrak{g})$ .

Now, let us derive the result for  $\mathcal{L}_2(\boldsymbol{\theta})$ . The gradient of the loss function is

$$\nabla \mathcal{L}_2(\boldsymbol{\theta}) = -\frac{2}{d^2} \text{Re}[\nabla T(\boldsymbol{\theta}) T^*(\boldsymbol{\theta})] = -2\text{Re}[\text{Tr}[V^\dagger \nabla U(\boldsymbol{\theta}) T^*(\boldsymbol{\theta})]] = 2\text{Im}[\text{Tr}[V^\dagger U(\boldsymbol{\theta}) \tilde{\mathbf{H}}] T^*(\boldsymbol{\theta})] \quad (89)$$

and the Hessian

$$\begin{aligned} \nabla^2 \mathcal{L}_2(\boldsymbol{\theta}) &= \frac{2}{d^2} \text{Im}[\nabla \cdot (\text{Tr}[V^\dagger U(\boldsymbol{\theta}) \tilde{\mathbf{H}}^\top]) T^*(\boldsymbol{\theta}) + \nabla T^*(\boldsymbol{\theta}) \cdot \text{Tr}[V^\dagger U(\boldsymbol{\theta}) \tilde{\mathbf{H}}^\top]] \\ &= -\frac{2}{d^2} \text{Re}[\text{Tr}[V^\dagger U(\boldsymbol{\theta}) \tilde{\mathbf{H}} \cdot \tilde{\mathbf{H}}^\top] T^*(\boldsymbol{\theta}) + \text{Tr}[\tilde{\mathbf{H}} U^\dagger(\boldsymbol{\theta}) V] \cdot \text{Tr}[V^\dagger U(\boldsymbol{\theta}) \tilde{\mathbf{H}}^\top]]. \end{aligned} \quad (90)$$

Evaluating at a solution  $U(\boldsymbol{\theta}) = e^{i\phi}V$

$$\begin{aligned}\nabla^2\mathcal{L}_2(\boldsymbol{\theta}_*) &= -\frac{2}{d^2}\text{Re}\left[\text{Tr}\left[\tilde{\mathbf{H}}\cdot\tilde{\mathbf{H}}^\top\right]d+\text{Tr}\left[\tilde{\mathbf{H}}\right]\cdot\text{Tr}\left[\tilde{\mathbf{H}}^\top\right]\right] \\ &= -\frac{2}{d}\sum_{m,n=1}^d\mathbf{R}_{mn}\cdot\mathbf{R}_{mn}^\top+\mathbf{I}_{mn}\cdot\mathbf{I}_{mn}^\top+\frac{1}{d}\left(\mathbf{R}_{mm}\cdot\mathbf{R}_{nn}^\top-\mathbf{I}_{mm}\cdot\mathbf{I}_{nn}^\top\right).\end{aligned}\quad (91)$$

Again, this is a sum of rank-one matrices that live in the span of  $\{\mathbf{a}_\nu\}$ , and following a proof similar to that used in deriving Theorem 1, the rank of  $\nabla^2\mathcal{L}_2(\boldsymbol{\theta}_*)$  is at upper bounded by  $\dim(\mathfrak{g})$ .  $\square$

## VI. DETAILS OF THE NUMERICAL SIMULATIONS

The simulations in the main text were carried out in double precision using the open-source library `Qibo` [101, 102] (version 0.1.6). All circuits have been run on CPU because the overhead of transferring the state vector between the host and the device makes the usage of GPUs not suitable for circuits with less than 15-20 qubits. This is specially true for the case of VQAs, where back and forth communication between host and device results in a deteriorated performance. Simulations have been performed using single-thread multiprocessing to parallelize the execution of different instances of circuits with different number of qubits and depths in multiple cores. In particular, IntelCore i7-9750H, IntelCore i7-10750H and IntelCore i9-9900K cores have been employed.

The optimization method chosen in all cases has been the Adaptive Moment Estimation (Adam) algorithm [132], which is a variant of Stochastic Gradient Descent (SGD) widely used in classical machine learning, that adaptively adjusts the learning rate for each optimization parameter based on information coming from first and second moments of the gradients. This choice has been motivated by the fact that the works that have reported overparametrization in VQAs used this algorithm [50, 52], and also because we consider that a gradient-based optimizer is an appropriate choice to probe relevant features of the optimization landscape, like the disappearance of suboptimal local minima. In order to leverage automatic differentiation for the computation of gradients, the simulation backend in `Qibo` has been set to `tensorflow`. This backend, although slower than the `qibojit` and `qibotf` custom backends, allows to seamlessly deploy `Tensorflow`'s implementation of the Adam optimizer. The hyper-parameter values employed in all cases are: initial learning rate =  $10^{-2}$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\hat{\epsilon} = 10^{-7}$ . The optimization was stopped whenever we reached machine precision.

The minimizations have been carried out in all cases without considering sampling noise, *i.e.* using the full state vector in the simulation to compute expectation values of observables. The main reason for this is that we are here interested in the optimization landscape itself, and not in the stochasticity introduced by finite sampling.

Finally, we mention that parameter-shift rules [133, 134] have been employed in all cases for the computation of the quantum Fisher information and Hessian matrices, and that the simulation backend was switched to the faster `qibojit` for that.

## VII. FORMULAS FOR COMPUTING THE QFIM AND HESSIAN

We present here the explicit formulas employed in the computation of the Quantum Fisher Information Matrix,  $F(\boldsymbol{\theta})$ , and the Hessian,  $\nabla^2\mathcal{L}(\boldsymbol{\theta})$ , in each of the examples in our numerical simulations. For convenience, we recall the definitions of the elements of these two matrices,

$$[F_\mu(\boldsymbol{\theta})]_{ij} = 4\text{Re}[\langle\partial_i\psi_\mu(\boldsymbol{\theta})|\partial_j\psi_\mu(\boldsymbol{\theta})\rangle - \langle\partial_i\psi_\mu(\boldsymbol{\theta})|\psi_\mu(\boldsymbol{\theta})\rangle\langle\psi_\mu(\boldsymbol{\theta})|\partial_j\psi_\mu(\boldsymbol{\theta})\rangle] \quad , \quad [\nabla^2\mathcal{L}(\boldsymbol{\theta})]_{ij} = \partial_i\partial_j\mathcal{L}(\boldsymbol{\theta}), \quad (92)$$

where we use the notation  $\partial_i = \frac{\partial}{\partial\theta_i}$  and where the subscript  $\mu$  indicates the quantum state  $|\psi_\mu\rangle$  the QNN acts on. The QFIM can be interpreted, at each point  $\boldsymbol{\theta}$  in the landscape (and up to a constant factor), as the Hessian matrix of a pure state transfer problem where the target state is  $|\psi(\boldsymbol{\theta})\rangle$  itself. This opens up the possibility of employing quantum circuits to evaluate the QFIM on quantum hardware using parameter shift rules [133, 134], which may be useful for instance when

computing the natural gradient during the optimization of a VQA [78], or when doing variational quantum simulation of imaginary-time evolution [77]. The parameter shift rules are simple recipes to analytically compute partial derivatives of a given loss function with respect to parametrized quantum gates. In the case of the QFIM, its elements are given by [135]

$$[F_\mu(\boldsymbol{\theta})]_{ij} = \frac{1}{2} \partial_i \partial_j \mathcal{L}(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}' = \boldsymbol{\theta}} \quad (93)$$

where

$$\mathcal{L}(\boldsymbol{\theta}) = 1 - |\langle \psi(\boldsymbol{\theta}') | \psi(\boldsymbol{\theta}) \rangle|^2. \quad (94)$$

We start by computing the QFIM and the Hessian for the Hamiltonian Variational Ansatz (HVA) employed in the Variational Quantum Eigensolver (VQE) implementation. We recall that a HVA is an ansatz of the form in Eq. (42) where the generators  $\mathcal{G}$ , for a given Hamiltonian  $H = \sum_{k=1}^N a_k A_k$  (with  $A_k$  Hermitian operators and  $a_k$  real numbers), are simply  $\mathcal{G} = \{A_k\}_{k=1}^N$ . We employed this type of ansatz in the main text to minimize the loss function

$$E(\boldsymbol{\theta}) = \langle \psi(\boldsymbol{\theta}) | H_{\text{TFIM}} | \psi(\boldsymbol{\theta}) \rangle, \quad (95)$$

where  $H_{\text{TFIM}}$  is the Hamiltonian of the Transverse Field Ising Model (TFIM). Making use of the fact that for a HVA applied to the TFIM Hamiltonian, all the terms commute within a given  $e^{-i\theta_{lk} H_k}$  operator (where  $H_k = \frac{1}{2} \sum_i \sigma_i^z \sigma_{i+1}^z$  or  $H_k = \frac{1}{2} \sum_i \sigma_i^x$ ), we find that

$$\partial_{lk} e^{-i\theta_{lk} H_k} = \partial_{lk} \prod_i e^{-i\theta_{lk} H_{ki}} = \sum_i \prod_{j \neq i} e^{-i\theta_{lk} H_{kj}} \partial_{lk} e^{-i\theta_{lk} H_{ki}} = \sum_i \prod_{j \neq i} e^{-i\theta_{lkj} H_{kj}} \partial_{lk} e^{-i\theta_{lki} H_{ki}}, \quad (96)$$

where  $H_{ki} = \sigma_i^z \sigma_{i+1}^z$  or  $H_{ki} = \sigma_i^x$ . For convenience, we have additionally introduced the notation  $\theta_{lki}$  to denote the parameter in the  $l$ -th layer, that parametrizes the  $i$ -th term of the  $k$ -th generator. Note that in a periodic-structured ansatz,  $\theta_{lki} = \theta_{lk}$  for all  $i$ . Now, the partial derivative  $\partial_{lk} e^{-i\theta_{lk} H_{ki}}$  can be obtained by applying the parameter shift rule (since  $H_{ki}$  has only two distinct non-zero eigenvalues), and hence the partial derivative of a loss function with respect to  $\theta_{lk}$  is

$$\partial_{lk} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \sum_i \left( \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lki}}, \theta_{lki}^{\frac{\pi}{2}} \right) - \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lki}}, \theta_{lki}^{-\frac{\pi}{2}} \right) \right), \quad (97)$$

where  $\overline{lki}$  denotes all the indices distinct from  $l, k, i$ , and  $\theta_{lki}^s = \theta_{lki} + s$ . Therefore, applying the parameter shift rule twice, the matrix elements of the QFIM are given by

$$[F(\boldsymbol{\theta})]_{lk, l'k'} = \frac{1}{8} \sum_{i,j} \left( \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lki, l'k'j}}, \theta_{lki}^{\frac{\pi}{2}}, \theta_{l'k'j}^{\frac{\pi}{2}} \right) + \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lki, l'k'j}}, \theta_{lki}^{-\frac{\pi}{2}}, \theta_{l'k'j}^{-\frac{\pi}{2}} \right) \right. \\ \left. - \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lki, l'k'j}}, \theta_{lki}^{\frac{\pi}{2}}, \theta_{l'k'j}^{-\frac{\pi}{2}} \right) - \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lki, l'k'j}}, \theta_{lki}^{-\frac{\pi}{2}}, \theta_{l'k'j}^{\frac{\pi}{2}} \right) \right). \quad (98)$$

To analytically compute the Hessian matrix of the loss function  $E(\boldsymbol{\theta})$  for the HVA, we can also apply twice the parameter shift rule. The matrix elements  $\nabla^2 E(\boldsymbol{\theta})_{lk, l'k'} = \partial_{lk} \partial_{l'k'} E(\boldsymbol{\theta})$  of the Hessian are thus given by (see *e.g.* [39, 42])

$$[\nabla^2 E(\boldsymbol{\theta})]_{lk, l'k'} = \frac{1}{4} \sum_{i,j} \left( E \left( \boldsymbol{\theta}_{\overline{lki, l'k'j}}, \theta_{lki}^{\frac{\pi}{2}}, \theta_{l'k'j}^{\frac{\pi}{2}} \right) + E \left( \boldsymbol{\theta}_{\overline{lki, l'k'j}}, \theta_{lki}^{-\frac{\pi}{2}}, \theta_{l'k'j}^{-\frac{\pi}{2}} \right) \right. \\ \left. - E \left( \boldsymbol{\theta}_{\overline{lki, l'k'j}}, \theta_{lki}^{\frac{\pi}{2}}, \theta_{l'k'j}^{-\frac{\pi}{2}} \right) - E \left( \boldsymbol{\theta}_{\overline{lki, l'k'j}}, \theta_{lki}^{-\frac{\pi}{2}}, \theta_{l'k'j}^{\frac{\pi}{2}} \right) \right). \quad (99)$$

We now turn our attention to the Hardware Efficient Ansatz (HEA) that we employed for unitary compilation and quantum autoencoding in the main text. In this case, every  $Rx$  or  $Ry$  gate in the ansatz is a generator of the form  $e^{-i\frac{\theta}{2}\sigma^k}$  (where  $\sigma^k \in \{\sigma^x, \sigma^y\}$  has eigenvalues  $\pm 1$ ), with an independent angle. Hence, the QFIM elements are simply

$$[F(\boldsymbol{\theta})]_{lk,l'k'} = \frac{1}{8} \left( \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{\frac{\pi}{2}}, \theta_{l'k'}^{\frac{\pi}{2}} \right) + \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{-\frac{\pi}{2}}, \theta_{l'k'}^{-\frac{\pi}{2}} \right) - \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{\frac{\pi}{2}}, \theta_{l'k'}^{-\frac{\pi}{2}} \right) - \mathcal{L} \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{-\frac{\pi}{2}}, \theta_{l'k'}^{\frac{\pi}{2}} \right) \right). \quad (100)$$

We recall that the QFIM is independent of the loss function, and hence the above formula is valid for both the unitary compilation task and the quantum autoencoder. The difference between these two is that the initial states are different, and this has an impact on the QFIM. The Hessian matrices, on the contrary, depend on the loss function and thus are different in each case, but in our simulations we only computed the Hessian for the unitary compilation case. The loss here is given by  $\mathcal{L}(\boldsymbol{\theta}) = 1 - \frac{1}{d^2} |\text{Tr}(W^\dagger U(\boldsymbol{\theta}))|^2$ , where  $W$  is the unitary being compiled and  $d = 2^n$ . In this case, the term  $L(\boldsymbol{\theta}) = \frac{1}{d^2} |\text{Tr}(W^\dagger U(\boldsymbol{\theta}))|^2$  can be directly evaluated on a quantum computer, and so its second partial derivative  $\partial_{lk} \partial_{l'k'} L(\boldsymbol{\theta})$  is

$$\partial_{lk} \partial_{l'k'} L(\boldsymbol{\theta}) = \frac{1}{4} \left( L \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{\frac{\pi}{2}}, \theta_{l'k'}^{\frac{\pi}{2}} \right) + L \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{-\frac{\pi}{2}}, \theta_{l'k'}^{-\frac{\pi}{2}} \right) - L \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{\frac{\pi}{2}}, \theta_{l'k'}^{-\frac{\pi}{2}} \right) - L \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{-\frac{\pi}{2}}, \theta_{l'k'}^{\frac{\pi}{2}} \right) \right). \quad (101)$$

Then, applying the chain rule twice on  $\mathcal{L}(\boldsymbol{\theta}) = 1 - L(\boldsymbol{\theta})$  as

$$\partial_{lk} \partial_{l'k'} \mathcal{L}(\boldsymbol{\theta}) = \mathcal{L}'(L(\boldsymbol{\theta})) \partial_{lk} \partial_{l'k'} L(\boldsymbol{\theta}) + \mathcal{L}''(L(\boldsymbol{\theta})) \partial_{lk} L(\boldsymbol{\theta}) \partial_{l'k'} L(\boldsymbol{\theta}), \quad (102)$$

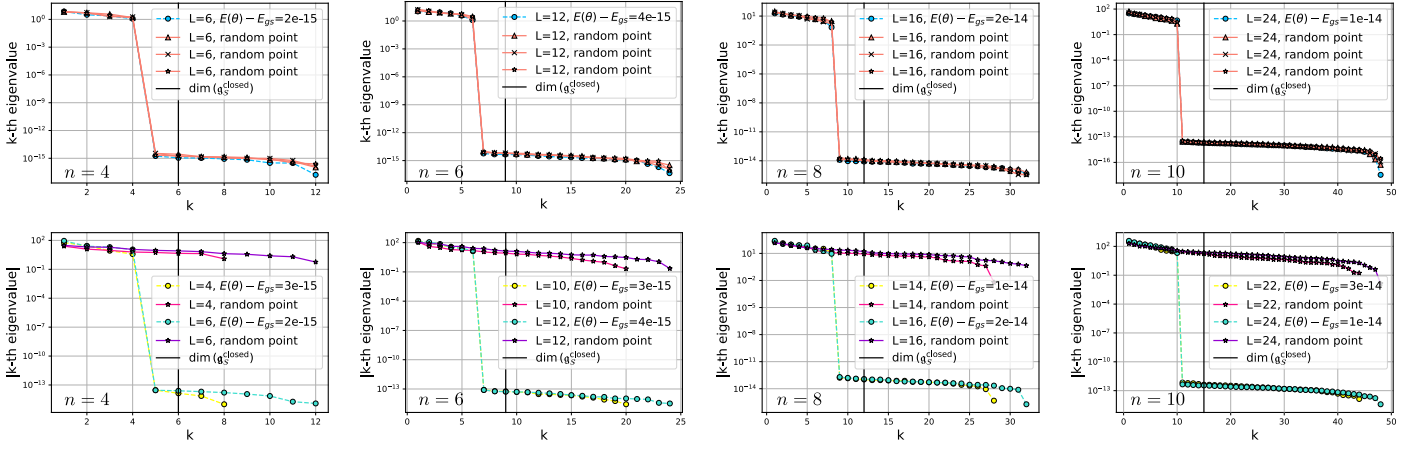
where  $\mathcal{L}'$  ( $\mathcal{L}''$ ) is the first (second) derivative of  $\mathcal{L}$  with respect to  $L$ , the expression for the Hessian matrix is found to be

$$[\nabla^2 \mathcal{L}(\boldsymbol{\theta})]_{lk,l'k'} = -\frac{1}{4} \left( L \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{\frac{\pi}{2}}, \theta_{l'k'}^{\frac{\pi}{2}} \right) + L \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{-\frac{\pi}{2}}, \theta_{l'k'}^{-\frac{\pi}{2}} \right) - L \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{\frac{\pi}{2}}, \theta_{l'k'}^{-\frac{\pi}{2}} \right) - L \left( \boldsymbol{\theta}_{\overline{lk,l'k'}}, \theta_{lk}^{-\frac{\pi}{2}}, \theta_{l'k'}^{\frac{\pi}{2}} \right) \right). \quad (103)$$

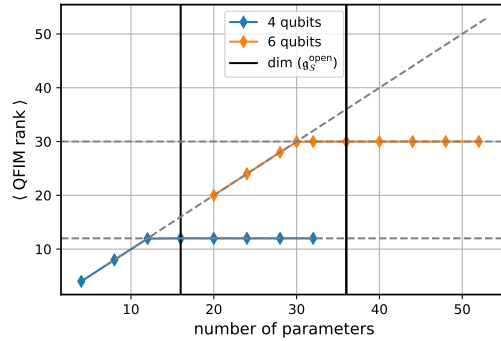
## VIII. ADDITIONAL NUMERICAL RESULTS

We present here some additional numerical results that we obtained in simulations. In Sup. Fig. 1, we show the eigenvalues of the QFIM and Hessian computed at the global optimum for the loss function  $E(\boldsymbol{\theta}) = \langle \psi(\boldsymbol{\theta}) | H_{\text{TFIM}} | \psi(\boldsymbol{\theta}) \rangle$  and the Hamiltonian variational ansatz with closed boundary conditions, for  $n = 10$  qubits. The interest in showing these plots is that therein one can better appreciate that the spectrum does not form a continuum, but rather, that there is a large gap between the non-zero and the zero eigenvalues, so that there is no ambiguity when defining the rank, stemming from numerical precision issues. Moreover, we computed the eigenvalues at random points in the landscape, where the rank of the QFIM is also bounded by  $\dim(\mathfrak{g}_S)$ , unlike the rank of the Hessian. The spectra of the Hessian at random points in the landscape further informed us that the landscape is highly non-convex, as it contains both positive and negative eigenvalues. It is also interesting to note that in order to obtain a rank of the Hessian that is bounded by  $\dim(\mathfrak{g}_S)$ , one needs to compute it at the global minimum; otherwise, we encountered fairly-good local minima that did not fulfill this result. We remark that all these features were found in all cases where we computed and diagonalized the QFIM and the Hessian.

Furthermore, in Sup. Fig. 2 we computed the rank of the QFIM at 30 random points in the landscape for the HVA with open boundaries and different depths (*i.e.* number of parameters), for  $n = 4, 6$  qubits. This figure shows that the rank quickly saturates at all points once a critical number of parameters  $M_c$  is reached, suggesting that, at least in this case, overparametrization largely arises simultaneously across the entire landscape. We note as well that before having  $M_c$  parameters, the average rank is equal to the number of parameters, *i.e.* there is a perfectly linear relation between the two quantities. Adding more parameters beyond  $M_c$  however seems to have no effect on the rank, as it only adds null eigenvalues.



SUP FIG. 1. **Spectra of the QFIM and the Hessian for the VQE implementations.** Top row: QFIM spectra for the Hamiltonian variational ansatz with closed boundary conditions, for  $n = 4, 6, 8, 10$  qubits and  $L = 6, 12, 16, 24$  layers, both at the global optima and at three random points in the landscape. Bottom row: Hessian spectra for the Hamiltonian variational ansatz with closed boundary conditions, for  $n = 4, 6, 8, 10$  qubits and  $L = 4, 6, 10, 12, 14, 16, 22, 24$  layers, both at the global optima and at a random point in the landscape.



SUP FIG. 2. **Average QFIM rank versus number of parameters for the VQE implementations.** Average rank of the QFIM across 30 random points in the landscape, for the Hamiltonian variational ansatz with open boundary conditions and  $n = 4, 6$  qubits. The horizontal dashed lines mark the maximal ranks that the QFIMs achieve, and the tilted dashed line is the line  $\langle \text{QFIM rank} \rangle = \text{number of parameters}$ . The vertical black lines correspond to the respective  $\dim(\mathfrak{g}_S)$  (leftmost for  $n = 4$ , and rightmost for  $n = 6$ ). We remark that the standard deviation is exactly 0 for all the points in the plot, except for the case  $n = 4$  ( $6$ ) and  $M = 12$  ( $30$ ) parameters, where  $\text{rank}[F(\theta)] = 12$  ( $30$ ) for 29 out of 30 points and  $\text{rank}[F(\theta)] = 11$  ( $29$ ) for the remaining one.