# Fragmentation is Efficiently Learnable by Quantum Neural Networks

MIKHAIL MINTS, California Institute of Technology, USA

ERIC R. ANSCHUETZ, California Institute of Technology, USA

*Hilbert space fragmentation* is a phenomenon in which the Hilbert space of a quantum system is dynamically decoupled into exponentially many Krylov subspaces. We can define the *Schur transform* as a unitary operation mapping some set of preferred bases of these Krylov subspaces to computational basis states labeling them. We prove that this transformation can be efficiently learned via gradient descent from a set of training data using quantum neural networks, provided that the fragmentation is sufficiently strong such that the summed dimension of the unique Krylov subspaces is polynomial in the system size. To demonstrate this, we analyze the loss landscapes of random quantum neural networks constructed out of Hilbert space fragmented systems. We prove that in this setting, it is possible to eliminate barren plateaus and poor local minima, suggesting efficient trainability when using gradient descent. Furthermore, as the algebra defining the fragmentation is not known a priori and not guaranteed to have sparse algebra elements, to the best of our knowledge there are no existing efficient classical algorithms generally capable of simulating expectation values in these networks. Our setting thus provides a rare example of a physically motivated quantum learning task with no known dequantization.

## 1 Introduction

Quantum algorithms are known to outperform the best known classical algorithms on many computational problems exhibiting some sort of algebraic structure [10, 17, 18, 33]. A significant research direction today is focused on trying to understand whether the unique advantages offered by quantum computers can help in machine learning tasks and whether *quantum neural networks* (QNNs) can be superior to classical methods. Classical neural networks are known to behave like Gaussian processes in the asymptotic limit [27], which provides theoretical guarantees on their efficient trainability in a wide range of applications. On the other hand, QNNs are known to have poorly-behaved loss landscapes. They are often dominated by *barren plateaus* [12, 23, 30], which generally prevent gradients from being estimated efficiently, and *poor local minima*, which generally prevent gradient descent from reaching the global optimum [1, 6, 38]. While there do exist settings where QNNs are efficient to train via gradient descent—such as when they have many symmetries [4, 5, 15, 24, 28, 32, 37]—in these settings there often also exist efficient classical simulation algorithms that render the use of a quantum computer for the task unnecessary outside of potentially an initial data acquisition phase [3, 11, 16].

Recent work [2] formulated the Jordan Algebraic Wishart System (JAWS) framework, which is the first full theoretical characterization of QNN loss landscapes that provides the conditions under which they can become efficiently trainable. In this work, we use these tools to formally demonstrate a specific setting where QNNs *can* in fact be efficiently trained to solve a physically-motivated problem for which no efficient classical simulation algorithm is known. Our approach is to showcase a setting where the QNN has a high-dimensional symmetry group but, crucially, the end-user does not know what that group is and—to the best of our knowledge—has no efficient strategy for learning it.

### 1.1 Fragmentation and the Schur Basis

The problem that we study in this work is that of classifying states in systems which exhibit *Hilbert space fragmentation* [21, 29, 31], a well-studied barrier to ergodicity in quantum many-body

systems conceptually similar to quantum many-body scars [9, 34]. The mechanism underlying Hilbert space fragmentation has been recently mathematically characterized [25], and we review this mathematical characterization here.

Suppose for each system size $L$ we can define a set of Hamiltonians $\mathcal{S}$ acting on a Hilbert space $\mathcal{H}$ such that $\mathcal{S}$ has the following property: there exist a set of Hermitian local operators $\boldsymbol{h}_1, \ldots, \boldsymbol{h}_m \in \mathcal{L}(\mathcal{H})$ with $m = \text{poly}(L)$ where for any $\boldsymbol{H} \in \mathcal{S}$,

$$\boldsymbol{H} = \sum_{i=1}^{m} c_i \boldsymbol{h}_i \tag{1}$$

for some $\{c_i \in \mathbb{R}\}$. Let $\mathcal{A} = \langle \boldsymbol{h}_1, \ldots, \boldsymbol{h}_m \rangle$ be the associative algebra generated by these operators under addition and matrix multiplication. Let $C$ be the commutant algebra, consisting of all operators in $\mathcal{L}(\mathcal{H})$ that commute with all elements of $\mathcal{A}$. By the Von Neumann bicommutant theorem, we know that $\mathcal{A}$ and $C$ are each other's centralizers, and we can decompose the Hilbert space as

$$\mathcal{H} = \bigoplus_{\lambda=1}^{\Lambda} \mathcal{H}_\lambda^{(\mathcal{A})} \otimes \mathcal{H}_\lambda^{(C)}, \tag{2}$$

where each $\mathcal{H}_\lambda^{(\mathcal{A})}$ is an irreducible representation of $\mathcal{A}$ and each $\mathcal{H}_\lambda^{(C)}$ is an irreducible representation of $C$. Let $N_\lambda = \dim\left(\mathcal{H}_\lambda^{(\mathcal{A})}\right), N_\lambda' = \dim\left(\mathcal{H}_\lambda^{(C)}\right)$. In this decomposition, each $\mathcal{H}_\lambda^{(\mathcal{A})}$ is a *Krylov subspace* of the system, preserved by the action of Hamiltonians constructed from the generators $\boldsymbol{h}_1, \ldots, \boldsymbol{h}_m$. Each of these subspaces has $N_\lambda'$ degenerate copies. The representation in $\mathcal{H}$ of any operator $A \in \mathcal{A}$ can be written (in a slight abuse of notation identifying the operator with its representation) in the form

$$\boldsymbol{A} = \bigoplus_{\lambda=1}^{\Lambda} \boldsymbol{A}^\lambda \otimes \text{Id}_{N_\lambda'}. \tag{3}$$

That is, any such operator acts as the identity on the "multiplicity labels" in $\mathcal{H}_\lambda^{(C)}$, meaning that it acts identically on each of the $N_\lambda'$ degenerate Krylov subspaces isomorphic to $\mathcal{H}_\lambda^{(\mathcal{A})}$.

The phenomenon of *Hilbert space fragmentation* is said to occur when the total number of Krylov subspaces, $\sum_\lambda N_\lambda'$, is exponential in the system size $L$ [26]. We are interested in the stronger condition where the dimension of $\mathcal{A}$ is polynomial in the system size. If a system does not have this property, we may be able to restrict to a subspace that does: we pick $\Lambda'$ as a function of $L$ to be such that

$$N = \sum_{\lambda=1}^{\Lambda'} N_\lambda = \text{poly}(L). \tag{4}$$

Note that this $\Lambda'$ may be a constant or it may be equal to $\Lambda$ depending on the system in question and the degree of degeneracy in the Krylov subspaces. Now, we can define the space

$$\mathcal{H}' = \bigoplus_{\lambda=1}^{\Lambda'} \mathcal{H}_\lambda^{(\mathcal{A})} \otimes \mathcal{H}_\lambda^{(C)}. \tag{5}$$

The representation in $\mathcal{H}'$ of any operator $A \in \mathcal{A}$ can then be written as

$$\boldsymbol{A} = \bigoplus_{\lambda=1}^{\Lambda'} \boldsymbol{A}^\lambda \otimes \text{Id}_{N_\lambda'}. \tag{6}$$

In analogy with the terminology for permutation-invariant quantum systems [7], we can construct the *Schur basis* for $\mathcal{H}'$ as the set of states $|\lambda, q_\lambda, p_\lambda\rangle = |\lambda, q_\lambda\rangle \otimes |\lambda, p_\lambda\rangle$ such that the states
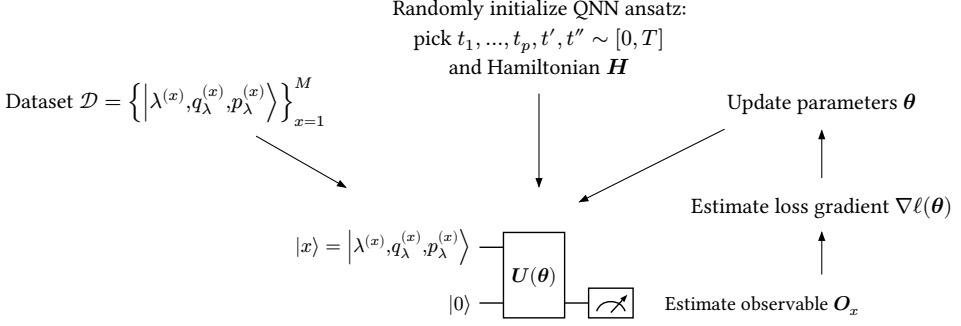
Fig. 1. Diagram of the QNN training process. We are given a dataset of Schur basis states and randomly sample a QNN ansatz architecture and an initial vector of parameters. We then compile the QNN into a quantum circuit and repeatedly run on the input states to estimate the gradient of the loss function, which is then used to adjust the parameters.

$|\lambda, q_\lambda\rangle$ form an orthonormal basis for $\mathcal{H}_\lambda^{(\mathcal{A})}$ and the states $|\lambda, p_\lambda\rangle$ form an orthonormal basis for $\mathcal{H}_\lambda^{(C)}$.

## 1.2 The Classification Task

Given a quantum state that is a basis state promised to be in the Schur basis, we would like to perform some kind of quantum measurement to determine which state exactly it is, ignoring the degeneracy. In some special settings—such as when the system is permutation invariant, and one is interested in the associated Schur basis of that Krylov decomposition—this can be done efficiently using a quantum algorithm known as the Schur transform [7]. However, generally, we know only the generators $\boldsymbol{h}_1, \ldots, \boldsymbol{h}_m$. Instead, we hope in effect to *learn* the Schur transform associated with $\mathcal{S}$, assuming we are also given copies of Schur basis states $|\lambda, q_\lambda, p_\lambda\rangle$. Specifically, we want to learn a transformation $\boldsymbol{U}$ such that

$$\mathrm{Tr}_{\mathcal{H}'}\left[\boldsymbol{U}(|\lambda, q_\lambda, p_\lambda\rangle\langle\lambda, q_\lambda, p_\lambda| \otimes |\mathbf{0}\rangle\langle\mathbf{0}|)\boldsymbol{U}^\dagger\right] \approx |\boldsymbol{\lambda}, \boldsymbol{q_\lambda}\rangle\langle\boldsymbol{\lambda}, \boldsymbol{q_\lambda}|. \tag{7}$$

Here, we have added an ancillary register of $n_a$ qubits with Hilbert space of dimension $N_a = 2^{n_a}$ and initialized it in the computational basis state $|\mathbf{0}\rangle$. The unitary $U$ will be implemented as a quantum circuit that should write a bitstring label $\boldsymbol{\lambda}, \boldsymbol{q_\lambda}$ to this ancillary register in the computational basis. Note that here we are ambivalent to the multiplicity label $p_\lambda$; there are exponentially many such labels, and we cannot hope to learn them efficiently. Instead, we also require that our network *generalizes* its classification of any $|\lambda, q_\lambda, p_\lambda\rangle$ given training examples of $\left|\lambda, q_\lambda, p'_\lambda\right\rangle$ for $p_\lambda$ not necessarily equal to $p'_\lambda$.

We implement this unitary as a QNN ansatz built from matrix exponentials of Hamiltonians constructed from the given generators $\boldsymbol{h}_1, \ldots, \boldsymbol{h}_m$, parameterized by $\boldsymbol{\theta} \in \mathbb{R}^p$. The training process is shown in Figure 1. We use gradient descent to find optimal values of the parameters to minimize the loss function:

$$\ell(\boldsymbol{\theta}) = \frac{1}{M} \sum_{x=1}^{M} \langle x| \boldsymbol{U}(\boldsymbol{\theta})\boldsymbol{O}_x\boldsymbol{U}(\boldsymbol{\theta})^\dagger |x\rangle, \tag{8}$$

where each $|x\rangle$ is of the form $|\lambda, q_\lambda, p_\lambda\rangle \otimes |\mathbf{0}\rangle$ and the corresponding objective observable is

$$\boldsymbol{O}_x = -\mathrm{Id}_{\mathcal{H}'} \otimes |\boldsymbol{\lambda}, \boldsymbol{q_\lambda}\rangle\langle\boldsymbol{\lambda}, \boldsymbol{q_\lambda}|. \tag{9}$$

### 1.3 Classical Hardness

While there exist classical techniques for performing Lie-algebraic simulations in situations where the dimension of a system's dynamical Lie algebra grows polynomially with the system size [16], these techniques are not applicable to the problem posed above. This is because we are not given a full description of the algebraic structure of the system. Furthermore, there is no known method to efficiently learn this structure in general, even with a quantum computer. For instance, one might hope to estimate expectation values of algebra elements with respect to basis elements of the algebra. However, even if the $h_i$ have a sparse representation—say, they are given as sparse Pauli decompositions—it is not necessarily the case that products of poly($L$) of them are sparse, and therefore in general there is no way to efficiently measure these coefficients on a quantum computer. In other words, we are not promised that there exists a basis of the algebra that is sparse in the quantum representation. Classical algorithms for simulating symmetric quantum systems, such as in Anschuetz et al. [3], also cannot be applied to this problem, since the symmetry group corresponding to permutations of the degenerate Krylov subspaces is not initially known and difficult to compute from just the set of generators. As we demonstrate that this problem can be efficiently solved using quantum neural networks, this establishes one of the few known settings (outside of sampling-based tasks [19]) where a quantum machine learning task has no known dequantization.

### 1.4 Trainability Conditions

To formally prove that QNNs can efficiently solve this problem, we need to demonstrate the two main conditions needed for a QNN to be efficiently trainable. The following are standard criteria of QNN trainability in the literature [2]:

(1) **The absence of barren plateaus**. In the asymptotic limit of large system size, the distribution of the loss landscape over the random initialization of the QNN architecture converges to one where the derivatives only decay polynomially with the system size, thus allowing the gradient to be efficiently estimated on a quantum computer in polynomial time. We formally prove this in Theorem 2.

(2) **The absence of poor local minima**. When the number of parameters $p$ is sufficiently high, the loss landscape enters an *overparameterized regime* where all spurious local minima disappear, and the only remaining minima are *degenerate*, existing on some submanifold of the parameter space. We follow the approach of Larocca et al. [22] and show the Hessian is not full-rank in order to demonstrate a lack of poor local minima. The value of $p$ needed for this to occur is only polynomially large in the system size, making the QNN in the overparameterized phase efficiently implementable on a quantum computer. We formally prove this in Theorem 3.

## 2 The QNN Ansatz Construction

Suppose that we have a dataset

$$\mathcal{D} = \{(|x\rangle, O_x)\}_{x \in [M]}, \tag{10}$$

of size $M \leq N$ where each $|x\rangle$ is of the form $|\lambda, q_\lambda, p_\lambda\rangle \otimes |0\rangle$, and the corresponding $O_x$ is of the form

$$O_x = -\operatorname{Id}_{\mathcal{H}'} \otimes |\lambda, q_\lambda\rangle\langle\lambda, q_\lambda|, \tag{11}$$

where $\lambda, q_\lambda$ are bitstring representations of $\lambda, q_\lambda$ on the ancillary register of $n_a = \lceil \log_2 M \rceil$ qubits. We assume we are given at least one example $x$ with a given $\lambda, q_\lambda$ pair, though we will later see (in Theorem 1) that we do not require an example for each $p_\lambda$. Let $N_a = 2^{n_a}$ be the dimension of the

Hilbert space of the ancillary register. Let

$$\mathcal{H}^* = \mathcal{H}' \otimes \mathbb{C}^{N_a} = \bigoplus_{\lambda=1}^{\Lambda'} \mathcal{H}_\lambda^* \otimes \mathcal{H}_\lambda^{(C)}, \tag{12}$$

where we define $\mathcal{H}_\lambda^* = \mathcal{H}_\lambda^{(\mathcal{A})} \otimes \mathbb{C}^{N_a}$, of dimension $N_\lambda^* = N_\lambda N_a$.

Now, we construct a randomized QNN ansatz parameterized by $\boldsymbol{\theta} = (\theta_1, \ldots, \theta_p) \in \mathbb{R}^p$ as follows:

$$U(\boldsymbol{\theta}) = e^{\mathrm{i}Ht''}\left(\prod_{i=1}^{p} e^{-\mathrm{i}Ht_i} e^{\mathrm{i}A\theta_i} e^{\mathrm{i}Ht_i}\right) e^{-\mathrm{i}Ht'}, \tag{13}$$

where we select $t_1, \ldots, t_p, t', t''$ i.i.d. uniformly from $[0, T]$ for some fixed $T$, we let $H$ be a Hamiltonian on the system and the ancillary register:

$$H = \left(\sum_{i=1}^{m} c_i \boldsymbol{h}_i\right) \otimes \left(\sum_{j=1}^{n_a} c'_{j,x} \sigma_j^x + c'_{j,y} \sigma_j^y + c'_{j,z} \sigma_j^z\right), \tag{14}$$

and we let $A$ be a local operator which we assume without loss of generality is positive semidefinite to simplify our analysis—if it is not, we just add a sufficiently large multiple of the identity to $\boldsymbol{h}_1 \otimes \mathrm{Id}_{N_a}$.

We assume that coefficients in Equation (14) are chosen such that the Hamiltonian obeys the *full Eigenstate Thermalization Hypothesis* (ETH) [14] in each Krylov subspace, which we assume is possible. The full ETH is an ansatz for higher-order time-averaged correlation functions of local observables, and has been shown to imply a quantitative scrambling behavior [13]. More formally, we assume the following [13, Eq. (5)]:

*Hypothesis 1 (Full ETH, free cumulant definition [13, Eq. (5)]).* We say a sequence of Hamiltonians $H$ describing a size-$L$ system and acting on an $N$-dimensional Hilbert space obeys the full ETH if the following is true. For any $k$, there exists some $t_k = O(\mathrm{poly}(L, k))$ such that choosing $t$ uniformly from $[0, t_k]$ and defining the free cumulant $\kappa_{2k}$ with respect to this distribution, for any local operators $A_i(t)$ time-evolved under $H$ and any local operators $B_i$:

$$\kappa_{2k}\left[A_1(t), B_1, A_2(t), B_2, \ldots, A_k(t), B_k\right] = O(1/N). \tag{15}$$

Observe that we can block-diagonalize all operators in Equation (13), writing them as a sum over the irreducible representations (algebraic sectors) of $\mathcal{A}$. We identify the operators with their representations in $\mathcal{H}^*$:

$$U(\boldsymbol{\theta}) = \bigoplus_{\lambda=1}^{\Lambda'} U^\lambda(\boldsymbol{\theta}), \tag{16}$$

where

$$U^\lambda(\boldsymbol{\theta}) = e^{\mathrm{i}H^\lambda t''}\left(\prod_{i=1}^{p} e^{-\mathrm{i}H^\lambda t_i} e^{\mathrm{i}A^\lambda \theta_i} e^{\mathrm{i}H^\lambda t_i}\right) e^{-\mathrm{i}H^\lambda t'}. \tag{17}$$

Define the loss function

$$\ell(\boldsymbol{\theta}) = \frac{1}{M} \sum_{x=1}^{M} \ell_x(\boldsymbol{\theta}) = \frac{1}{M} \sum_{x=1}^{M} \langle x | U(\boldsymbol{\theta}) O_x U(\boldsymbol{\theta})^\dagger | x \rangle. \tag{18}$$

We can also decompose this by sector. Note that since each input state is restricted to one of the sectors, we can write

$$\ell(\boldsymbol{\theta}) = \frac{1}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \ell_x^\lambda(\boldsymbol{\theta}) = \frac{1}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \langle x | U^\lambda(\boldsymbol{\theta}) O_x^\lambda U^\lambda(\boldsymbol{\theta})^\dagger | x \rangle. \tag{19}$$

Note that the indexing of $x$ is separate for each $\lambda$, so we have slightly abused notation here and are letting each $x$ in the above equation corresponds to some $q_\lambda$, where $M_\lambda$ is the number of dataset entries in the $\lambda$ block.

We initially randomly choose the QNN ansatz by sampling $t_1, \ldots, t_p, t', t''$ and choosing the Hamiltonian $H$. Then, these become fixed, and only the parameterized parts of the form $e^{iA\theta_i}$ will be changed during training as we adjust the parameters $\theta_i$. To analyze the gradient of the loss function, we will consider its distribution at any particular value of $\boldsymbol{\theta}$ over the random choice of the QNN ansatz. We perform a sequence of reductions, first reducing the distribution to the case where we only have to consider $\boldsymbol{\theta} = 0$ and then writing the loss derivatives in terms of Gaussian-distributed random variables in the asymptotic limit, allowing us to prove the result about the variance of the gradient.

## 3   Formal Proofs of the Main Results

First, we establish the prior claim that indeed the training of our QNN ansatz does not require a data set of Schur basis states over all multiplicity labels.

THEOREM 1. *If our dataset $\mathcal{D}$ contains an instance $|\lambda, q_\lambda, p_\lambda\rangle$ for each choice of $\lambda$ and $q_\lambda$, then if we take the (exponential-size) dataset $\mathcal{D}' = \{|\lambda, q_\lambda, p_\lambda\rangle\}_{\lambda, q_\lambda, p_\lambda}$ consisting of every Schur basis state, the loss $\ell_{\mathcal{D}'}(\boldsymbol{\theta})$ with respect to that dataset is upper-bounded as $\ell_{\mathcal{D}'}(\boldsymbol{\theta}) \leq \max_{\lambda, x} \ell_x^\lambda(\boldsymbol{\theta})$.*

PROOF. This follows from the fact that every layer of the QNN is constructed from matrix exponentials of elements of $\mathcal{A}$. Thus, conjugating by $U(\boldsymbol{\theta})$ will act in the same way on each degenerate Krylov subspace, which means that the loss contribution from the new datapoints for a given $\lambda, q_\lambda$ will be exactly the same as that for the point in the original dataset. So, we can write

$$\ell_{\mathcal{D}'}(\boldsymbol{\theta}) = \frac{1}{|\mathcal{D}'|} \sum_{\lambda=1}^{\Lambda'} N_\lambda' \sum_{x=1}^{M_\lambda} \ell_x^\lambda(\boldsymbol{\theta}) \leq \frac{\max_{\lambda, x} \ell_x^\lambda(\boldsymbol{\theta})}{|\mathcal{D}'|} \sum_{\lambda=1}^{\Lambda'} N_\lambda' M_\lambda \leq \max_{\lambda, x} \ell_x^\lambda(\boldsymbol{\theta}). \tag{20}$$

$\square$

Now, we want to analyze the derivatives of this loss function.

LEMMA 1. *At $\boldsymbol{\theta} = 0$, we have that*

$$\partial_i \ell(\boldsymbol{\theta})\Big|_{\boldsymbol{\theta}=0} = \frac{\mathbf{i}}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \langle x| e^{iH^\lambda t''} \left[ e^{-iH^\lambda t_i} A^\lambda e^{iH^\lambda t_i}, \ e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} \right] e^{-iH^\lambda t''} |x\rangle. \tag{21}$$

PROOF.

$$\frac{\partial}{\partial \theta_i} \ell(\boldsymbol{\theta}) = \frac{\mathbf{i}}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \frac{\partial}{\partial \theta_i} \langle x| U^\lambda(\boldsymbol{\theta}) O_x^\lambda U^\lambda(\boldsymbol{\theta})^\dagger |x\rangle =$$

$$= \frac{\mathbf{i}}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \frac{\partial}{\partial \theta_i} \langle x| e^{iH^\lambda t''} \left( \prod_{j=1}^{p} e^{-iH^\lambda t_j} e^{\theta_j A^\lambda} e^{iH^\lambda t_j} \right) e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} \left( \prod_{j=p}^{1} e^{-iH^\lambda t_j} e^{-\theta_j A^\lambda} e^{iH^\lambda t_j} \right) e^{-iH^\lambda t''} |x\rangle =$$

$$= \frac{\mathbf{i}}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \langle x| \left( e^{iH^\lambda t''} \left( \prod_{j=1}^{p} e^{-iH^\lambda t_j} e^{\theta_j A^\lambda} \left( A^\lambda \right)^{\delta_{ij}} e^{iH^\lambda t_j} \right) e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} \left( \prod_{j=p}^{1} e^{-iH^\lambda t_j} e^{-\theta_j A^\lambda} e^{iH^\lambda t_j} \right) e^{-iH^\lambda t''} - \right.$$

$$\left. - e^{iH^\lambda t''} \left( \prod_{j=1}^{p} e^{-iH^\lambda t_j} e^{\theta_j A^\lambda} e^{iH^\lambda t_j} \right) e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} \left( \prod_{j=p}^{1} e^{-iH^\lambda t_j} e^{-\theta_j A^\lambda} \left( A^\lambda \right)^{\delta_{ij}} e^{iH^\lambda t_j} \right) e^{-iH^\lambda t''} \right) |x\rangle. \tag{22}$$

At $\boldsymbol{\theta} = 0$, this simplifies to

$$\partial_i \ell(\boldsymbol{\theta})\Big|_{\boldsymbol{\theta}=0} = \frac{\mathbf{i}}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \langle x | e^{\mathbf{i}H^\lambda t''} \left[ e^{-\mathbf{i}H^\lambda t_i} A^\lambda e^{\mathbf{i}H^\lambda t_i}, \ e^{-\mathbf{i}H^\lambda t'} O_x^\lambda e^{\mathbf{i}H^\lambda t'} \right] e^{-\mathbf{i}H^\lambda t''} |x\rangle . \quad (23)$$

$\square$

Now, we want to analyze the distribution of the gradient of the loss function. To do this, we consider joint distributions of individual terms that contribute to the expression for the gradient.

LEMMA 2. *Consider a set of nonnegative integers* $\{k_{\lambda,x,i,1}, k_{\lambda,x,i,2} : \lambda \in [\Lambda'], x \in M_\lambda, i \in [p]\}$. *Let*

$$k = \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \sum_{i=1}^{p} \left( k_{\lambda,x,i,1} + k_{\lambda,x,i,2} \right). \quad (24)$$

*Let* $\boldsymbol{H}$ *be a Hamiltonian satisfying the full Eigenstate Thermalization Hypothesis (Hypothesis 1). Then, there exists some* $T = \mathrm{poly}(L, k)$ *such that*

$$\mathbb{E}_{\substack{t_1,\ldots,t_p, \\ t',t'' \sim \mathcal{U}(0,T)}} \prod_{\lambda=1}^{\Lambda'} \prod_{x=1}^{M_\lambda} \prod_{i=1}^{p} \mathrm{Tr} \left( e^{-\mathbf{i}H^\lambda t''} |x\rangle\langle x| e^{\mathbf{i}H^\lambda t''} e^{-\mathbf{i}H^\lambda t_i} A^\lambda e^{\mathbf{i}H^\lambda t_i} e^{-\mathbf{i}H^\lambda t'} O_x^\lambda e^{\mathbf{i}H^\lambda t'} \right)^{k_{\lambda,x,i,1}}$$

$$\mathrm{Tr} \left( e^{-\mathbf{i}H^\lambda t''} |x\rangle\langle x| e^{\mathbf{i}H^\lambda t''} e^{-\mathbf{i}H^\lambda t'} O_x^\lambda e^{\mathbf{i}H^\lambda t'} e^{-\mathbf{i}H^\lambda t_i} A^\lambda e^{\mathbf{i}H^\lambda t_i} \right)^{k_{\lambda,x,i,2}} = \quad (25)$$

$$= O(1/N_{\min}^*) + \mathbb{E}_{\substack{g_1,\ldots,g_p, \\ g',g'' \sim \mathrm{Haar}}} \prod_{\lambda=1}^{\Lambda'} \prod_{x=1}^{M_\lambda} \prod_{i=1}^{p} \mathrm{Tr} \left( g''^\lambda |x\rangle\langle x| g''^{\lambda\dagger} g_i^\lambda A^\lambda g_i^{\lambda\dagger} g'^\lambda O_x^\lambda g'^{\lambda\dagger} \right)^{k_{\lambda,x,i,1}}$$

$$\mathrm{Tr} \left( g''^\lambda |x\rangle\langle x| g''^{\lambda\dagger} g'^\lambda O_x^\lambda g'^{\lambda\dagger} g_i^\lambda A^\lambda g_i^{\lambda\dagger} \right)^{k_{\lambda,x,i,2}} ,$$

*where*

$$N_{\min}^* = \min_\lambda N_\lambda^*. \quad (26)$$

*Here, the Haar-random distribution is over the unitaries with the block-diagonal structure of the Krylov subspaces: equivalently, each* $g_i^\lambda, g'^\lambda, g''^\lambda$ *is chosen to be Haar-random in each of the sectors.*

PROOF. We can justify the above claim using generally accepted physical assumptions as follows. Arguing similarly to Appendix C of Fava et al. [13], if we assume Hypothesis 1, we can say that for

a sufficiently large time $T$,

$$\underset{\substack{t_1,\ldots,t_p,\\ t',t''\sim\mathcal{U}(0,T)}}{\mathbb{E}} \prod_{\lambda=1}^{\Lambda'}\prod_{x=1}^{M_\lambda}\prod_{i=1}^{p} \text{Tr}\left(e^{-iH^\lambda t''}|x\rangle\langle x|e^{iH^\lambda t''}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}\right)^{k_{\lambda,x,i,1}}$$

$$\text{Tr}\left(e^{-iH^\lambda t''}|x\rangle\langle x|e^{iH^\lambda t''}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}\right)^{k_{\lambda,x,i,2}} =$$

$$= \prod_{\lambda=1}^{\Lambda'}\left(O(1/N_\lambda^*) + \prod_{x=1}^{M_\lambda}\prod_{i=1}^{p}\underset{\substack{t_1,\ldots,t_p,\\ t',t''\sim\mathcal{U}(0,T)}}{\mathbb{E}}\left[\text{Tr}\left(e^{-iH^\lambda t''}|x\rangle\langle x|e^{iH^\lambda t''}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}\right)\right]^{k_{\lambda,x,i,1}}\right.$$

$$\left.\underset{\substack{t_1,\ldots,t_p,\\ t',t''\sim\mathcal{U}(0,T)}}{\mathbb{E}}\left[\text{Tr}\left(e^{-iH^\lambda t''}|x\rangle\langle x|e^{iH^\lambda t''}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}\right)\right]^{k_{\lambda,x,i,2}}\right) =$$

$$= O(1/N_{\min}^*) + \prod_{\lambda=1}^{\Lambda'}\prod_{x=1}^{M_\lambda}\prod_{i=1}^{p}\underset{\substack{t_1,\ldots,t_p,\\ t',t''\sim\mathcal{U}(0,T)}}{\mathbb{E}}\left[\text{Tr}\left(e^{-iH^\lambda t''}|x\rangle\langle x|e^{iH^\lambda t''}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}\right)\right]^{k_{\lambda,x,i,1}}$$

$$\underset{\substack{t_1,\ldots,t_p,\\ t',t''\sim\mathcal{U}(0,T)}}{\mathbb{E}}\left[\text{Tr}\left(e^{-iH^\lambda t''}|x\rangle\langle x|e^{iH^\lambda t''}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}\right)\right]^{k_{\lambda,x,i,2}}.$$

(27)

Now, using the results of Fava et al. [13], we can write

$$\underset{\substack{t_1,\ldots,t_p,\\ t',t''\sim\mathcal{U}(0,T)}}{\mathbb{E}} \text{Tr}\left(e^{-iH^\lambda t''}|x\rangle\langle x|e^{iH^\lambda t''}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}\right) = \frac{1}{(N_\lambda^*)^2}\text{Tr}(A^\lambda)\text{Tr}(O_x^\lambda).$$

(28)

Now, when we replace the time evolution with the Haar-random unitaries, by free independence in the asymptotic limit [36] we will get the same result as above up to an $O(1/N_{\min}^*)$ error.  □

LEMMA 3. *We can approximate, at any value of $\boldsymbol{\theta}$,*

$$\partial_i\ell_x^\lambda(\boldsymbol{\theta}) \rightsquigarrow \langle x|\, g''^{\lambda\dagger}\left[g_i^\lambda A^\lambda g_i^{\lambda\dagger},\; g'^\lambda O_x^\lambda g'^{\lambda\dagger}\right]g''^\lambda|x\rangle,$$

(29)

*where $g_1,\ldots,g_p, g', g'' \sim$ Haar up to an $o(1)$ error the Lévy-Prokhorov metric (Definition 1), if we choose a sufficiently large $T$. Furthermore, this approximation is valid for the joint distribution of any two gradient contributions $\partial_i\ell_x^\lambda(\boldsymbol{\theta})$ and $\partial_j\ell_y^{\lambda'}(\boldsymbol{\theta})$: their joint distribution can be approximated up to an $o(1)$ Lévy-Prokhorov error by replacing the Hamiltonian time evolution with Haar-random matrices.*

PROOF. Consider the joint distribution $\mathfrak{p}$ of the terms in the expression for the $\partial_i\ell_x^\lambda$ for $\boldsymbol{\theta} = 0$, as in Lemma 1, considering the commutator expression as two separate terms. The terms are of the form

$$\langle x|\, e^{iH^\lambda t''}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}e^{-iH^\lambda t''}|x\rangle$$

(30)

and

$$\langle x|\, e^{iH^\lambda t''}e^{-iH^\lambda t'}O_x^\lambda e^{iH^\lambda t'}e^{-iH^\lambda t_i}A^\lambda e^{iH^\lambda t_i}e^{-iH^\lambda t''}|x\rangle.$$

(31)

and there are $d = 2$ of these terms, so the distribution is in $\mathbb{R}^2$. We can also consider the joint distribution of the terms contributing to both $\partial_i\ell_x^\lambda(\boldsymbol{\theta})$ and $\partial_l\ell_y^{\lambda'}(\boldsymbol{\theta})$, in which case the distribution will be in $\mathbb{R}^4$. Similarly, let $\mathfrak{q}$ be the joint distribution of the same terms, but with the time-evolved Hamiltonian replaced with the Haar-random matrices $g_i, g_j, g', g''$. Now, from Lemma 2, we know that we can pick $T = \text{poly}(L, k)$ such that mixed moments of order $k$ of the distributions $\mathfrak{p}$ and

$\mathfrak{q}$ differ by an additive error of $\epsilon = 1/N^*_{\min}$ for any fixed. The random variables corresponding to the terms are bounded, and thus they are subgaussian [35] and thus their moments of order $k$ are upper bounded by $\left(C\sqrt{k}\right)^k$ for some constant $C$ independent of $k$. Now, applying Lemma 9 gives us that the distribution $\mathfrak{p}$ converges to $\mathfrak{q}$ up to an error of

$$O\left(\frac{\log \log N_{\min}}{\log N_{\min}} + \frac{\log k}{\sqrt{k}}\right) = o(1) \tag{32}$$

in the Lévy-Prokhorov metric (note that the $\mu$ term is subleading). Now, consider taking a general value of $\boldsymbol{\theta}$. Then, each of the terms becomes conjugated by some unitary matrix (see Lemma 1). By the unitary invariance property of the Haar ensemble, the mixed moments of the distribution $\mathfrak{q}$ remain the same, which means that it approximates the joint distribution of the derivative terms at any value of $\boldsymbol{\theta}$. since the Lévy-Prokhorov distance is preserved under the linear transformation of summing the terms to form the commutator expressions and is only reduced when dividing by the factor of $M$, the proof is complete. □

LEMMA 4. *We can approximate, at any value of $\boldsymbol{\theta}$,*

$$\partial_i \ell_x^\lambda(\boldsymbol{\theta}) \rightsquigarrow \frac{\mathbf{i}}{MN_\lambda^{*2}} \langle x| \, \boldsymbol{g}''^{\lambda\dagger} \left[\tilde{\boldsymbol{g}}_i^\lambda \boldsymbol{A}^\lambda \tilde{\boldsymbol{g}}_i^{\lambda\dagger}, \; \tilde{\boldsymbol{g}}'^\lambda \boldsymbol{O}_x^\lambda \tilde{\boldsymbol{g}}'^{\lambda\dagger}\right] \boldsymbol{g}''^\lambda |x\rangle, \tag{33}$$

*where $\tilde{\boldsymbol{g}}_1^\lambda, \ldots, \tilde{\boldsymbol{g}}_p^\lambda, \tilde{\boldsymbol{g}}'^\lambda$ have i.i.d. standard Gaussian entries, up to an $o(1)$ error in the Lévy-Prokhorov metric, and the same approximation is valid for the joint distribution of any two loss contributions $\partial_i \ell_x^\lambda(\boldsymbol{\theta})$ and $\partial_j \ell_y^{\lambda'}(\boldsymbol{\theta})$. From these results, we incur an additional error of*

$$O\left(\sqrt{\frac{\log N^*}{N^*}}\right) \tag{34}$$

PROOF. We first perform the reduction to Haar-random matrices as in Lemma 3. Now, we can proceed similarly to Lemma 27 of Anschuetz [2], based on the results of Jiang [20]. □

We know by construction that $\boldsymbol{A}$ is positive semidefinite. We can diagonalize it, writing each block as

$$\boldsymbol{A}^\lambda = \sum_{\mu=1}^{N_\lambda^*} a_\mu^\lambda |\mu\rangle\langle\mu| \tag{35}$$

Now, we want to perform an additional reduction to replace $\boldsymbol{A}$ with a *semi-isotropic* $\tilde{\boldsymbol{A}}$, which we define as:

$$\tilde{\boldsymbol{A}}^\lambda = \frac{\text{Tr}(\boldsymbol{A}^\lambda)}{r_A^\lambda} \sum_{\mu=1}^{r_A^\lambda} |\mu\rangle\langle\mu|, \tag{36}$$

where

$$r_A^\lambda = \frac{\text{Tr}(\boldsymbol{A}^\lambda)^2}{\text{Tr}((\boldsymbol{A}^\lambda)^2)}. \tag{37}$$

That is, we want to replace $\boldsymbol{A}^\lambda$ with a matrix all of whose nonzero eigenvalues are the same, preserving the average eigenvalue of $\boldsymbol{A}^\lambda$.

Observe that

$$\text{Tr}\left(\tilde{\boldsymbol{A}}^\lambda\right) = \text{Tr}\left(\boldsymbol{A}^\lambda\right) \tag{38}$$

and

$$\operatorname{Tr}\left((\tilde{A}^\lambda)^2\right) = \frac{\operatorname{Tr}(A^\lambda)^2}{r_A^\lambda} = \operatorname{Tr}((A^\lambda)^2). \tag{39}$$

LEMMA 5. *Replacing $A$ with $\tilde{A}$, we can write*

$$\partial_i \ell_x^\lambda(\theta) \rightsquigarrow \hat{\ell}_{x;i}^\lambda := \frac{\mathbf{i}}{MN_\lambda^{*2}} \langle x| g''^{\lambda\dagger} \left[\tilde{g}_i^\lambda \tilde{A}^\lambda \tilde{g}_i^{\lambda\dagger}, \ \tilde{g}'^\lambda O_x^\lambda \tilde{g}'^{\lambda\dagger}\right] g''^\lambda |x\rangle, \tag{40}$$

*up to an $o(1)$ error in the Lévy-Prokhorov metric. The same is true for joint distributions of two derivative contributions as before.*

PROOF. Consider each of the terms contributing to the gradient as previously. Replace each $A^\lambda$ term with $A^\lambda - \tilde{A}^\lambda$, so we have terms like

$$\frac{1}{MN_\lambda^{*2}} \langle x| \tilde{g}_i^\lambda \left(A^\lambda - \tilde{A}^\lambda\right) \tilde{g}_i^{\lambda\dagger} \tilde{g}'^\lambda O_x^\lambda \tilde{g}'^{\lambda\dagger} |x\rangle \tag{41}$$

and

$$\frac{1}{MN_\lambda^{*2}} \langle x| \tilde{g}'^\lambda O_x^\lambda \tilde{g}'^{\lambda\dagger} \tilde{g}_i^\lambda \left(A^\lambda - \tilde{A}^\lambda\right) \tilde{g}_i^{\lambda\dagger} |x\rangle \tag{42}$$

All these terms are of the form

$$L_k^\lambda = \frac{1}{MN_\lambda^{*2}} \operatorname{Tr}\left(M_k^\lambda X_k^\lambda \left(A^\lambda - \tilde{A}^\lambda\right) X_k^{\lambda\dagger}\right), \tag{43}$$

where each $X_k^\lambda$ has i.i.d. Gaussian entries but the $X_k^\lambda$ are not necessarily independent from each other. We can write

$$\frac{1}{MN_\lambda^{*2}} X_k^\lambda \left(A^\lambda - \tilde{A}^\lambda\right) X_k^{\lambda\dagger} = \frac{1}{MN_\lambda^{*2}} \sum_{i=1}^{N_\lambda^*} (a_i^\lambda - \tilde{a}_i^\lambda)|x_{k,i}^\lambda\rangle\langle x_{k,i}^\lambda|, \tag{44}$$

where the $a_i^\lambda, \tilde{a}_i^\lambda$ are the eigenvalues of $A^\lambda, \tilde{A}^\lambda$ in non-increasing order and $\left|x_{k,i}^\lambda\right\rangle$ is the $i$th column of $X_k^\lambda$ in the basis in which $A^\lambda$ is diagonal. Now, for any normalized vector $|\psi\rangle$, observe that $\left\langle x_{k,i}^\lambda \middle| \psi\right\rangle$ is gamma-distributed since it is a sum of squared Gaussian random variables. Then, by the version of Bernstein's inequality found in Vershynin [35], we have that

$$\Pr\left[\left|\frac{1}{MN_\lambda^{*2}} \sum_{i=1}^{N_\lambda^*} (a_i^\lambda - \tilde{a}_i^\lambda) \left|\left\langle x_{k,i}^\lambda \middle| \psi\right\rangle\right|^2\right| \geq \epsilon\right] \leq$$

$$\leq 2 \exp\left[-cMN_\lambda^{*2} \min\left(\frac{\epsilon}{\max_i \left|a_i^\lambda - \tilde{a}_i^\lambda\right|}, \frac{\epsilon^2}{\frac{1}{MN_\lambda^{*2}} \sum_i \left(a_i^\lambda - \tilde{a}_i^\lambda\right)^2}\right)\right] \tag{45}$$

for some constant $c > 0$. Now, observe that

$$\frac{1}{MN_\lambda^{*2}} \sum_i \left(a_i^\lambda - \tilde{a}_i^\lambda\right)^2 = \frac{1}{MN_\lambda^{*2}} \operatorname{Tr}\left[\left(A^\lambda - \tilde{A}^\lambda\right)^2\right] =$$

$$= \frac{1}{MN_\lambda^{*2}} \operatorname{Tr}\left[\left(A^\lambda\right)^2 + \left(\tilde{A}^\lambda\right)^2 - A^\lambda \tilde{A}^\lambda - \tilde{A}^\lambda A^\lambda\right] = \frac{2}{MN_\lambda^{*2}} \left(\operatorname{Tr}((A^\lambda)^2) - \operatorname{Tr}\left(A^\lambda \tilde{A}^\lambda\right)\right). \tag{46}$$

Now, since

$$\mathrm{Tr}\left(A^\lambda \tilde{A}^\lambda\right) = \frac{\mathrm{Tr}(A^\lambda)}{r_A^\lambda} \sum_{i=1}^{r_A^\lambda} a_i^\lambda \geq \frac{\mathrm{Tr}((A^\lambda)^2)}{\mathrm{Tr}(A^\lambda)} \left(\mathrm{Tr}(A^\lambda) - N_\lambda^* a_{r_A^\lambda+1}^\lambda\right), \tag{47}$$

we must have that

$$\frac{1}{MN_\lambda^{*2}} \sum_i \left(a_i^\lambda - \tilde{a}_i^\lambda\right)^2 \leq \frac{2}{MN_\lambda^{*2}} \frac{\mathrm{Tr}((A^\lambda)^2)}{\mathrm{Tr}(A^\lambda)} N_\lambda^* a_{r_A^\lambda+1}^\lambda = o(1), \tag{48}$$

which holds since $A$ is a local operator and its eigenvalues do not depend on the system size. We also have that

$$\max_i \left|a_i^\lambda - \tilde{a}_i^\lambda\right| = O\left(a_1^\lambda - \frac{\mathrm{Tr}((A^\lambda)^2)}{\mathrm{Tr}(A^\lambda)} a_{r_A^\lambda+1}\right) = O(1). \tag{49}$$

Now, we can just apply the union bound on all terms (of which there at most 4 since we are considering the joint distribution of at most 2 contributions) to say that

$$\Pr\left[\exists \lambda, k : \ L_k^\lambda \geq \epsilon\right] \leq 8 \exp\left[-cMN_{\min}^{*2} \min(\epsilon, \epsilon^2)\right], \tag{50}$$

Now, to solve for the Ky Fan metric (which upper-bounds the Lévy-Prokhorov metric), we need to find an upper bound on

$$\inf\{\epsilon > 0 : \Pr\left[\exists \lambda, k : \ L_k^\lambda > \epsilon\right] \leq \epsilon\}. \tag{51}$$

For the case where $\epsilon < \epsilon^2$, equality occurs when

$$8 \exp\left(-cMN_{\min}^{*2}\epsilon\right) = \epsilon \tag{52}$$

$$\epsilon \exp\left(cMN_{\min}^{*2}\epsilon\right) = 8 \tag{53}$$

$$cMN_{\min}^{*2}\epsilon \exp\left(cMN_{\min}^{*2}\epsilon\right) = 8cMN_{\min}^{*2} \tag{54}$$

$$cMN_{\min}^{*2}\epsilon = W(8cMN_{\min}^{*2}), \tag{55}$$

where $W$ is the Lambert $W$ function. Since it is upper-bounded by the logarithm, an upper bound for the Ky Fan metric is

$$cMN_{\min}^{*2}\epsilon \leq \log(8cMN_{\min}^{*2})$$
$$\epsilon \leq \frac{\log(8cMN_{\min}^{*2})}{cMN_{\min}^{*2}} = O\left(\frac{\log(MN_{\min})}{MN_{\min}^2}\right). \tag{56}$$

When $\epsilon^2 < \epsilon$ (when $\epsilon < 1$), equality occurs when

$$8 \exp\left(-cMN_{\min}^{*2}\epsilon^2\right) = \epsilon \tag{57}$$

$$cMN_{\min}^{*2}\epsilon^2 = W(8cMN_{\min}^{*2}\epsilon), \tag{58}$$

so the upper bound for the Ky Fan metric (and thus the Lévy-Prokhorov metric) is

$$cMN_{\min}^{*2}\epsilon^2 \leq \log(8cMN_{\min}^{*2}\epsilon) \leq \log(8cMN_{\min}^{*2}) \tag{59}$$

$$\epsilon = O\left(\frac{\sqrt{\log(MN_{\min})}}{MN_{\min}^2}\right) = o(1). \tag{60}$$

Since $\epsilon < 1$ for large enough system sizes, this is the correct asymptotic bound. □

Now that we have performed these reductions, we can, up to a vanishing error in Lévy-Prokhorov metric, approximate the distribution of the gradient of the loss function as some $\widehat{\nabla \ell}$.

THEOREM 2.

$$\mathbb{E}\left[\left\|\widehat{\nabla\ell}\right\|^2\right] = \sum_{\lambda=1}^{\Lambda'} \frac{2M_\lambda(N_a-1)\,\mathrm{Tr}((A^\lambda)^2)p}{M^2 N_a^4 N_\lambda^2} = \Omega\left(\frac{1}{\mathrm{poly}(L)}\right). \tag{61}$$

PROOF. We know by the reductions above that the contribution to the derivative with respect to $\theta_i$ from entry $x$ in sector $\lambda$ can be approximated as

$$\hat{\ell}^\lambda_{x;i} = \frac{\mathbf{i}}{MN_\lambda^{*2}} \langle x|\, \tilde{g}''^{\lambda\dagger} \left[\tilde{g}_i^\lambda \tilde{A}^\lambda \tilde{g}_i^{\lambda\dagger},\ \tilde{g}'^\lambda O_x^\lambda \tilde{g}'^{\lambda\dagger}\right] \tilde{g}''^\lambda\, |x\rangle. \tag{62}$$

By unitary invariance, we can replace $\tilde{g}_i^\lambda$ with $\tilde{g}'^\lambda \tilde{g}_i^\lambda$ while maintaining the same probability distribution. But now, since

$$\left[\tilde{g}'^\lambda \tilde{g}_i^\lambda \tilde{A}^\lambda \tilde{g}_i^{\lambda\dagger} \tilde{g}'^{\lambda\dagger},\ \tilde{g}'^\lambda O_x^\lambda \tilde{g}'^{\lambda\dagger}\right] = \tilde{g}'^\lambda \tilde{g}_i^\lambda \tilde{A}^\lambda \tilde{g}_i^{\lambda\dagger} O_x^\lambda \tilde{g}'^{\lambda\dagger} - \tilde{g}'^\lambda O_x^\lambda \tilde{g}_i^\lambda \tilde{A}^\lambda \tilde{g}_i^{\lambda\dagger} \tilde{g}'^{\lambda\dagger} = \tilde{g}'^\lambda \left[\tilde{g}_i^\lambda \tilde{A}^\lambda \tilde{g}_i^{\lambda\dagger},\ O_x^\lambda\right] \tilde{g}'^{\lambda\dagger}, \tag{63}$$

we can absorb $g''^\lambda$ into $g'^\lambda$ and write

$$\hat{\ell}^\lambda_{x;i} = \frac{\mathbf{i}}{MN_\lambda^{*2}} \langle x|\, \tilde{g}'^\lambda \left[\tilde{g}_i^\lambda \tilde{A}^\lambda \tilde{g}_i^{\lambda\dagger},\ O_x^\lambda\right] \tilde{g}'^{\lambda\dagger}\, |x\rangle =$$

$$= \frac{\mathbf{i}}{MN_\lambda^{*2}} \frac{\mathrm{Tr}((A^\lambda)^2)}{\mathrm{Tr}(A^\lambda)} \sum_{\mu=1}^{r_A^\lambda} \langle x|\, \tilde{g}'^\lambda \left[\tilde{g}_i^\lambda|\mu\rangle\langle\mu|\tilde{g}_i^{\lambda\dagger},\ O_x^\lambda\right] \tilde{g}'^{\lambda\dagger}\, |x\rangle \tag{64}$$

Now, we can take a matrix $X^\lambda$ of dimension $N_\lambda^* \times (M_\lambda + pr_A^\lambda)$ with i.i.d. Gaussian entries, such that

$$X^\lambda\, |x\rangle = \tilde{g}'^{\lambda\dagger}\, |x\rangle \tag{65}$$

and

$$X^\lambda\, |i,\mu\rangle = \tilde{g}_i^\lambda\, |\mu\rangle. \tag{66}$$

Note that we are using $\{x\}$ and $\{(i,\mu)\}$ to label the $M_\lambda + pr_A^\lambda$ columns of $X^\lambda$.

Now, we can write

$$
\begin{aligned}
\hat{\ell}^\lambda_{x;i} &= \frac{i \operatorname{Tr}((A^\lambda)^2)}{M N^{*2}_\lambda \operatorname{Tr}(A^\lambda)} \sum_{\mu=1}^{r^\lambda_A} \langle x| X^{\lambda\dagger} \left[ X^\lambda |i,\mu\rangle\langle i,\mu| X^{\lambda\dagger}, O^\lambda_x \right] X^\lambda |x\rangle = \\
&= \frac{i \operatorname{Tr}((A^\lambda)^2)}{M N^{*2}_\lambda \operatorname{Tr}(A^\lambda)} \sum_{\mu=1}^{r^\lambda_A} \left( \langle x|^\lambda X^{\lambda\dagger} X^\lambda |i,\mu\rangle \langle i,\mu| X^{\lambda\dagger} O^\lambda_x X^\lambda |x\rangle - \langle x| X^{\lambda\dagger} O^\lambda_x X^\lambda |i,\mu\rangle \langle i,\mu| X^{\lambda\dagger} X^\lambda |x\rangle \right) = \\
&= \frac{i \operatorname{Tr}((A^\lambda)^2)}{M N^{*2}_\lambda \operatorname{Tr}(A^\lambda)} \sum_{\mu=1}^{r^\lambda_A} \left( \langle x| W^\lambda |i,\mu\rangle \langle i,\mu| W^{\lambda,x} |x\rangle - \langle x| W^{\lambda,x} |i,\mu\rangle \langle i,\mu| W^\lambda |x\rangle \right) = \\
&= \frac{i \operatorname{Tr}((A^\lambda)^2)}{M N^{*2}_\lambda \operatorname{Tr}(A^\lambda)} \sum_{\mu=1}^{r^\lambda_A} \sum_{y=1}^{M_\lambda+1} \left( \langle x| W^{\lambda,y} |i,\mu\rangle \langle i,\mu| W^{\lambda,x} |x\rangle - \langle x| W^{\lambda,x} |i,\mu\rangle \langle i,\mu| W^{\lambda,y} |x\rangle \right) = \\
&= \frac{i \operatorname{Tr}((A^\lambda)^2)}{M N^{*2}_\lambda \operatorname{Tr}(A^\lambda)} \sum_{\mu=1}^{r^\lambda_A} \sum_{y=1}^{M_\lambda+1} (-2i) \operatorname{Im} \left( \langle x| W^{\lambda,x} |i,\mu\rangle \langle i,\mu| W^{\lambda,y} |x\rangle \right) = \\
&= \frac{2 \operatorname{Tr}((A^\lambda)^2)}{M N^{*2}_\lambda \operatorname{Tr}(A^\lambda)} \sum_{\mu=1}^{r^\lambda_A} \sum_{y=1}^{M_\lambda+1} \operatorname{Im} \left( \langle x| W^{\lambda,x} |i,\mu\rangle \langle i,\mu| W^{\lambda,y} |x\rangle \right).
\end{aligned}
$$

(67)

Here, we define $W^{\lambda,x} = X^{\lambda\dagger} O^\lambda_x X^\lambda$ for $1 \le x \le M_\lambda$, and we define

$$
W^{\lambda,M_\lambda+1} = W^\lambda - \sum_{x=1}^{M_\lambda} W^{\lambda,x}
$$

(68)

and $W^\lambda = X^{\lambda\dagger} X^\lambda$.

Now, we can write this in terms of the entries of $X^{\lambda,x}$ which we take to be an $r_x \times (M_\lambda + p r^\lambda_A)$ submatrix of $X$. Here for $1 \le x \le M_\lambda$ we have $r_x = \operatorname{rank}(O^\lambda_x) = N_\lambda$, and $r_{M_\lambda+1} = N^*_\lambda - M_\lambda N_\lambda$.

$$
\hat{\ell}^\lambda_{x;i} = \frac{2 \operatorname{Tr}((A^\lambda)^2)}{M N^{*2}_\lambda \operatorname{Tr}(A^\lambda)} \sum_{\mu=1}^{r^\lambda_A} \sum_{y \ne x} \operatorname{Im} \left( \sum_{j=1}^{r_x} \left( X^{\lambda,x}_{j,x} \right)^* X^{\lambda,x}_{j,(i,\mu)} \sum_{k=1}^{r_y} \left( X^{\lambda,y}_{k,(i,\mu)} \right)^* X^{\lambda,y}_{k,x} \right).
$$

(69)

Now, observe that for $\lambda \ne \lambda'$ or $x \ne x'$, we must have by a symmetry argument that

$$
\mathbb{E} \left[ (\hat{\ell}^\lambda_{x;i})(\hat{\ell}^{\lambda'}_{x';i}) \right] = 0
$$

(70)

since each term in the expansion of this will have a factor that is a Gaussian random variable independent from the rest. Also recall that this is a valid approximation for the corresponding term in the actual loss since we ensured in our reductions that the joint distribution of two contributions such as this converges in distribution to our approximation.

Now,

$$
\mathbb{E}\left[(\hat{\ell}_{x;i})^2\right] = \left(\frac{2\operatorname{Tr}((A^\lambda)^2)}{MN_\lambda^{*2}\operatorname{Tr}(A^\lambda)}\right)^2 \sum_{\mu=1}^{r_A^\lambda} \sum_{\substack{y=1\\y\neq x}}^{N_\lambda+1} \mathbb{E}\left[\operatorname{Im}\left(\sum_{j=1}^{r_x}\left(X_{j,x}^{\lambda,x}\right)^* X_{j,(i,\mu)}^{\lambda,x} \sum_{k=1}^{r_y}\left(X_{k,(i,\mu)}^{\lambda,y}\right)^* X_{k,x}^{\lambda,y}\right)^2\right] =
$$

$$
= \frac{1}{2}\left(\frac{2\operatorname{Tr}((A^\lambda)^2)}{MN_\lambda^{*2}\operatorname{Tr}(A^\lambda)}\right)^2 \sum_{\mu=1}^{r_A^\lambda} \sum_{\substack{y=1\\y\neq x}}^{N_\lambda+1} \mathbb{E}\left[\left|\sum_{j=1}^{r_x}\left(X_{j,x}^{\lambda,x}\right)^* X_{j,(i,\mu)}^{\lambda,x} \sum_{k=1}^{r_y}\left(X_{k,(i,\mu)}^{\lambda,y}\right)^* X_{k,x}^{\lambda,y}\right|^2\right].
$$

(71)

Note that all cross terms in the above calculation get eliminated due to symmetry in the Gaussian distributions of the unmatched matrix elements.

Now, observe that

$$
\mathbb{E}\left[\left|\sum_{j=1}^{r_x}\left(X_{j,x}^{\lambda,x}\right)^* X_{j,(i,\mu)}^{\lambda,x} \sum_{k=1}^{r_y}\left(X_{k,(i,\mu)}^{\lambda,y}\right)^* X_{k,x}^{\lambda,y}\right|^2\right] =
$$

$$
= \mathbb{E}\left[\left|\sum_{j=1}^{r_x}\left(X_{j,x}^{\lambda,x}\right)^* X_{j,(i,\mu)}^{\lambda,x}\right|^2\right]\mathbb{E}\left[\left|\sum_{k=1}^{r_y}\left(X_{k,(i,\mu)}^{\lambda,y}\right)^* X_{k,x}^{\lambda,y}\right|^2\right] =
$$

(72)

$$
= \mathbb{E}\left[\sum_{j=1}^{r_x}\left|X_{j,x}^{x,\lambda}\right|^2\left|X_{j,(i,\mu)}^{\lambda,x}\right|^2\right]\mathbb{E}\left[\sum_{k=1}^{r_y}\left|X_{k,(i,\mu)}^{\lambda,y}\right|^2\left|X_{k,x}^{\lambda,y}\right|^2\right] = r_x r_y = N_\lambda r_y.
$$

Then, since

$$
\sum_{\substack{y=1\\y\neq x}}^{N_\lambda+1} r_y = N_\lambda^* - N_\lambda = N_\lambda(N_a - 1),
$$

(73)

we have that

$$
\mathbb{E}\left[(\hat{\ell}_{x;i})^2\right] = \frac{1}{2}\left(\frac{2\operatorname{Tr}((A^\lambda)^2)}{MN_\lambda^{*2}\operatorname{Tr}(A^\lambda)}\right)^2 r_A^\lambda(N_a-1)N_\lambda^2 =
$$

(74)

$$
= \frac{2\operatorname{Tr}((A^\lambda)^2)^2}{M^2N_\lambda^{*4}\operatorname{Tr}(A^\lambda)^2}\frac{\operatorname{Tr}(A^\lambda)^2}{\operatorname{Tr}((A^\lambda)^2)}(N_a-1)N_\lambda^2 = \frac{2(N_a-1)\operatorname{Tr}((A^\lambda)^2)}{M^2N_a^4N_\lambda^2}
$$

Thus,

$$
\mathbb{E}\left[\left\|\widehat{\nabla\ell}\right\|^2\right] = \sum_{i=1}^p \mathbb{E}\left[(\hat{\ell}_i)^2\right] = \sum_{i=1}^p \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \mathbb{E}\left[(\hat{\ell}_{x;i})^2\right] = \sum_{\lambda=1}^{\Lambda'}\frac{2M_\lambda(N_a-1)\operatorname{Tr}((A^\lambda)^2)p}{M^2N_a^4N_\lambda^2} = \Omega\left(\frac{1}{\operatorname{poly}(L)}\right).
$$

(75)

$\square$

By proving Theorem 2, we have formally established one of the two main claims: the absence of barren plateaus. Since the gradient of the loss function only scales inversely polynomially with $L$ and the quantum circuit for our QNN ansatz is polynomial-size, it is efficient to estimate gradients on a quantum computer by performing repeated runs and measurements with varied parameters, and thus gradient descent can be efficiently performed.

We now argue that the model enters the overparameterized regime, as in Larocca et al. [22], by showing that the Hessian of the loss landscape stops being full-rank when the number of parameters $p$ is only polynomial in the system size.

THEOREM 3. *For some $p = O(\text{poly}(L))$, the Hessian matrix $\hat{H}$ at any value of $\boldsymbol{\theta}$ (where $\langle i| \hat{H} |j\rangle = \partial_i \partial_j \ell(\boldsymbol{\theta})$) does not have full rank.*

PROOF. We can write the second derivative of the loss function as:

$$\partial_i \partial_j \ell(\boldsymbol{\theta}) = \frac{1}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \langle x| U^\lambda(\boldsymbol{\theta}) O_x^\lambda U^\lambda(\boldsymbol{\theta})^\dagger |x\rangle =$$

$$= -\frac{1}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \langle x| e^{iH^\lambda t''} \left( \left( \prod_{k=1}^{p} e^{-iH^\lambda t_k} e^{\theta_k A^\lambda} \left( A^\lambda \right)^{\delta_{ik}+\delta_{jk}} e^{iH^\lambda t_k} \right) e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} \left( \prod_{k=p}^{1} e^{-iH^\lambda t_k} e^{-\theta_k A^\lambda} e^{iH^\lambda t_k} \right) - \right.$$

$$- \left( \prod_{k=1}^{p} e^{-iH^\lambda t_k} e^{\theta_k A^\lambda} \left( A^\lambda \right)^{\delta_{ik}} e^{iH^\lambda t_k} \right) e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} \left( \prod_{k=p}^{1} e^{-iH^\lambda t_k} e^{-\theta_k A^\lambda} \left( A^\lambda \right)^{\delta_{jk}} e^{iH^\lambda t_k} \right) -$$

$$- \left( \prod_{k=1}^{p} e^{-iH^\lambda t_k} e^{\theta_k A^\lambda} \left( A^\lambda \right)^{\delta_{jk}} e^{iH^\lambda t_k} \right) e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} \left( \prod_{k=p}^{1} e^{-iH^\lambda t_k} e^{-\theta_k A^\lambda} \left( A^\lambda \right)^{\delta_{ik}} e^{iH^\lambda t_k} \right) +$$

$$+ \left. \left( \prod_{k=1}^{p} e^{-iH^\lambda t_k} e^{\theta_k A^\lambda} e^{iH^\lambda t_k} \right) e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} \left( \prod_{k=p}^{1} e^{-iH^\lambda t_k} e^{-\theta_k A^\lambda} \left( A^\lambda \right)^{\delta_{ik}+\delta_{jk}} e^{iH^\lambda t_k} \right) \right) e^{-iH^\lambda t''} |x\rangle . \tag{76}$$

Now, for every $i$, let

$$V_i = \prod_{k=1}^{i} e^{-iH^\lambda t_k} e^{\theta_k A^\lambda} e^{iH^\lambda t_k}. \tag{77}$$

Then, assuming $i \geq j$ (without loss of generality, since $\partial_i \partial_j \ell(\boldsymbol{\theta}) = \partial_j \partial_i \ell(\boldsymbol{\theta})$), we can rewrite the above as follows:

$$\partial_i \partial_j \ell(\boldsymbol{\theta}) = -\frac{1}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \langle x| e^{iH^\lambda t''} \left( V_j e^{-iH^\lambda t_j} A^\lambda e^{iH^\lambda t_j} V_j^\dagger V_i e^{-iH^\lambda t_i} A^\lambda e^{iH^\lambda t_i} V_i^\dagger V_p e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} V_p^\dagger - \right.$$

$$- V_i e^{-iH^\lambda t_i} A^\lambda e^{iH^\lambda t_i} V_i^\dagger V_p e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} V_p^\dagger V_j e^{-iH^\lambda t_j} A^\lambda e^{iH^\lambda t_j} V_j^\dagger -$$

$$- V_j e^{-iH^\lambda t_j} A^\lambda e^{iH^\lambda t_j} V_j^\dagger V_p e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} V_p^\dagger V_i e^{-iH^\lambda t_i} A^\lambda e^{iH^\lambda t_i} V_i^\dagger +$$

$$+ V_p e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} V_p^\dagger V_i e^{-iH^\lambda t_i} A^\lambda e^{iH^\lambda t_i} V_i^\dagger V_j e^{-iH^\lambda t_j} A^\lambda e^{iH^\lambda t_j} V_j^\dagger \left. \right) e^{-iH^\lambda t''} |x\rangle =$$

$$= -\frac{1}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \langle x| e^{iH^\lambda t''} \left[ V_j e^{-iH^\lambda t_j} A^\lambda e^{iH^\lambda t_j} V_j^\dagger, \left[ V_i e^{-iH^\lambda t_i} A^\lambda e^{iH^\lambda t_i} V_i^\dagger, V_p e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} V_p^\dagger \right] \right] e^{-iH^\lambda t''} |x\rangle . \tag{78}$$

Now, analogously to Theorem 2, we let $X^\lambda$ be a matrix of dimension $N_\lambda^* \times (M_\lambda + pN_\lambda^*)$ (note that we are no longer doing the semi-isotropic reduction on $A$) such that:

$$X^\lambda |x\rangle = e^{-iH^\lambda t''} |x\rangle \tag{79}$$

and

$$X^\lambda |i, \mu\rangle = V_i e^{-iH^\lambda t_i} |\mu\rangle . \tag{80}$$

Let $\tilde{O}_x^\lambda = V_p e^{-iH^\lambda t'} O_x^\lambda e^{iH^\lambda t'} V_p^\dagger$, $W^\lambda = X^{\lambda\dagger} X^\lambda$, and $W^{\lambda, x} = X^{\lambda\dagger} \tilde{O}_x^\lambda X^\lambda$. Then, we can write:

$$\partial_i \partial_j \ell(\theta) = -\frac{1}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \sum_{\mu=1}^{N_\lambda^*} \sum_{\nu=1}^{N_\lambda^*} a_\mu^\lambda a_\nu^\lambda \langle x| X^{\lambda\dagger} \left[ X^\lambda |j,\nu\rangle \langle j,\nu| X^{\lambda\dagger}, \left[ X^\lambda |i,\mu\rangle \langle i,\mu| X^{\lambda\dagger}, \tilde{O}_x^\lambda \right] \right] X^\lambda |x\rangle =$$

$$= -\frac{1}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \sum_{\mu=1}^{N_\lambda^*} \sum_{\nu=1}^{N_\lambda^*} a_\mu^\lambda a_\nu^\lambda \Big( \langle x| W^\lambda |j,\nu\rangle \langle j,\nu| W^\lambda |i,\mu\rangle \langle i,\mu| W^{\lambda,x} |x\rangle -$$

$$- \langle x| W^\lambda |j,\nu\rangle \langle j,\nu| W^{\lambda,x} |i,\mu\rangle \langle i,\mu| W^\lambda |x\rangle -$$

$$- \langle x| W^\lambda |i,\mu\rangle \langle i,\mu| W^{\lambda,x} |j,\nu\rangle \langle j,\nu| W^\lambda |x\rangle +$$

$$+ \langle x| W^{\lambda,x} |i,\mu\rangle \langle i,\mu| W^\lambda |j,\nu\rangle \langle j,\nu| W^\lambda |x\rangle \Big) =$$

$$= \frac{2}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \sum_{\mu=1}^{N_\lambda^*} \sum_{\nu=1}^{N_\lambda^*} a_\mu^\lambda a_\nu^\lambda \operatorname{Re} \Big( \langle x| W^\lambda |i,\mu\rangle \langle i,\mu| W^{\lambda,x} |j,\nu\rangle \langle j,\nu| W^\lambda |x\rangle -$$

$$- \langle x| W^{\lambda,x} |i,\mu\rangle \langle i,\mu| W^\lambda |j,\nu\rangle \langle j,\nu| W^\lambda |x\rangle \Big) =$$

$$= \frac{2}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \sum_{\mu=1}^{N_\lambda^*} \sum_{\nu=1}^{N_\lambda^*} a_\mu^\lambda a_\nu^\lambda \operatorname{Re} \Big( S_{i,j}^{\lambda,x,\mu,\nu,1} W_{i,j}^{\lambda,x,\mu,\nu} - S_{i,j}^{\lambda,x,\mu,\nu,2} W_{i,j}^{\lambda,\mu,\nu} \Big),$$

$$(81)$$

where

$$W_{i,j}^{\lambda,x,\mu,\nu} = \langle i,\mu| W^{\lambda,x} |j,\nu\rangle \tag{82}$$

$$W_{i,j}^{\lambda,\mu,\nu} = \langle i,\mu| W^\lambda |j,\nu\rangle \tag{83}$$

$$S_{i,j}^{\lambda,x,\mu,\nu,1} = \langle x| W^\lambda |i,\mu\rangle \langle j,\nu| W^\lambda |x\rangle \tag{84}$$

$$S_{i,j}^{\lambda,x,\mu,\nu,2} = \langle x| W^{\lambda,x} |i,\mu\rangle \langle j,\nu| W^\lambda |x\rangle. \tag{85}$$

So, we can write the Hessian as

$$\hat{H} = \frac{2}{M} \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \sum_{\mu=1}^{N_\lambda^*} \sum_{\nu=1}^{N_\lambda^*} a_\mu^\lambda a_\nu^\lambda \operatorname{Re} \Big( S^{\lambda,x,\mu,\nu,1} \odot W^{\lambda,x,\mu,\nu} - S^{\lambda,x,\mu,\nu,2} \odot W^{\lambda,\mu,\nu} \Big). \tag{86}$$

Now, we know that

$$\operatorname{rank}(W^{\lambda,x,\mu,\nu}) \le \operatorname{rank}(W^{\lambda,\mu,\nu}) \le \operatorname{rank}(W^\lambda) \le N_\lambda^*. \tag{87}$$

Additionally, it is clear that $S^{\lambda,x,\mu,\nu,1}$ and $S^{\lambda,x,\mu,\nu,2}$ both have rank 1: their entries are products of a term only depending on $i$ and a term only depending on $j$, so they can be written as outer products of two vectors. Now, it is known that the rank of a Hadamard product of two matrices is upper-bounded by the product of their ranks, which means that

$$\operatorname{rank} \hat{H} \le \sum_{\lambda=1}^{\Lambda'} \sum_{x=1}^{M_\lambda} \sum_{\mu=1}^{r_A^\lambda} \sum_{\nu=1}^{r_A^\lambda} 2N_\lambda^* = \sum_{\lambda=1}^{\Lambda'} 2M_\lambda (N_\lambda^*)^3 \le 2MN^{*3} = 2MN_a^3 N^3. \tag{88}$$

Thus, if we choose $p \ge \operatorname{rank}(\hat{H})$, it must be the case that $\hat{H}$, which is a $p \times p$ matrix, does not have full rank. $\qquad\square$

## 4 Numerical Results

We numerically verify our predicted loss landscape properties for a QNN constructed using the Temperley–Lieb model, a canonical example of Hilbert space fragmentation where the algebra is known and thus the dynamics are exactly solvable [26]. The local terms generating the family of Hamiltonians are projectors of the form

$$|00\rangle\langle00| + |00\rangle\langle11| + |11\rangle\langle00| + |11\rangle\langle11| \tag{89}$$

on adjacent qubits. The simulations were ran on a 4-qubit system using the dataset

$$\left\{|0\rangle \otimes \left|\Phi^+\right\rangle \otimes |1\rangle, \qquad \left|\Phi^+\right\rangle \otimes \left|\Phi^+\right\rangle\right\} \tag{90}$$

and on an 8-qubit system using the dataset

$$\left\{\left|\Phi^+\right\rangle \otimes |010101\rangle\right\}, \tag{91}$$

where

$$\left|\Phi^+\right\rangle = \frac{1}{\sqrt{2}}(|00\rangle + |11\rangle). \tag{92}$$

The simulations were performed using the `qutip` library in Python. Several optimizations, such as a precomputing a table of matrix exponentials, were introduced to make the simulation faster. However, since our simulations were performed using exact diagonalization, our methods become impractically long for larger system sizes. As this is a special case where the algebra is known a priori, we hope to strengthen our simulations of the performance of this neural network in the future through the use of existing classical simulation algorithms for symmetric quantum systems [3, 11, 16]. That said, even at small system sizes our results—reported in Figure 2 and Figure 3—support the correctness of the QNN construction and illustrate the phenomenon of the overparameterized regime: there exists a number of parameters for which the QNN will almost always reach the global optimum loss during training.
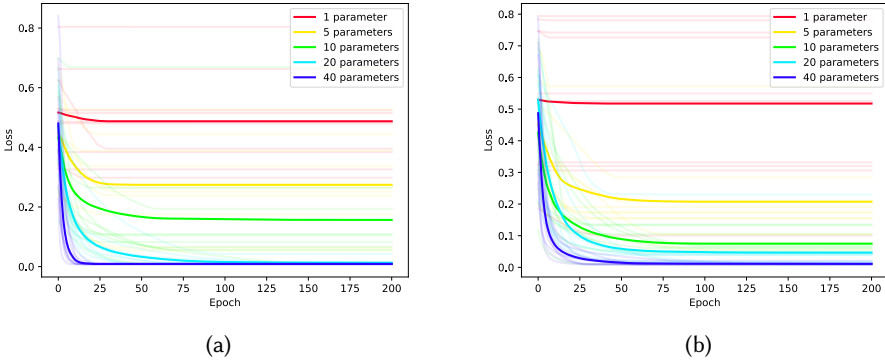


Fig. 2. Training curves for QNN models with 1, 5, 10, 20, and 40 parameters on the 4-qubit Temperley-Lieb dataset (a) and the 8-qubit Temperley-Lieb dataset (b). For each number of parameters, 10 QNNs were randomly initialized and trained for 200 epochs with a learning rate of 0.1. Training was stopped if the loss was decaying slower than 5% every 5 epochs or if it reached less than 0.01. The faded lines show the individual training curves, and the solid lines show their averages. For these diagrams, the loss value is adjusted to scale from 0 to 1 instead of from −1 to 0.
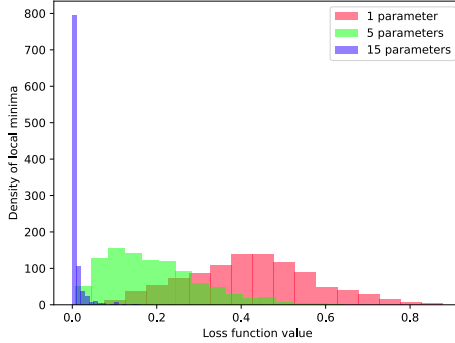
Fig. 3. Distribution of local minima reached during training for the 4-qubit Temperley-Lieb dataset. For this plot, 1000 QNNs were randomly initialized for each number of parameters and trained for 200 epochs with a learning rate of 0.1. The data plotted on the $x$-axis is the final loss achieved during training. Since we can assume that during gradient descent the parameters converge to the "closest" local minimum, this is a proxy for the distribution of the local minima of the loss landscape. We can see that with 15 parameters, the QNN enters the overparameterized regime since the distribution of minima becomes peaked at zero.

## 5   Conclusion

In this work, we have demonstrated a physically-motivated problem setting where 1) quantum neural networks efficiently solve the problem and 2) there are no known existing classical algorithms for simulating these networks, despite the substantial recent progress in classical simulation methods for this task [3, 11, 16]. At a high level, this followed by finding a setting where the quantum network was able to take advantage of algebraic structure present in the problem, but for which this structure is not given a convenient classical description for a classical simulator to take advantage of. We invite others to test our results by attempting to develop classical simulation methods in this more restricted setting, as empirical tests of the tractability of this task has implications for not only the potential utility of the algorithm presented here, but also more generally for determining unknown algebraic properties of symmetric quantum systems.

The problem we presented as a showcase for our QNN construction is physically motivated, and we hope may be of practical interest: as physicists study new quantum systems exhibiting fragmentation, they might construct a physical source of some known states to act as the training set and analyze these states using a quantum computer on which the QNN training can be performed. Then, when the optimal QNN parameters are known, the resulting quantum circuit can be used to classify new quantum states and identify them as copies of ones previously observed in a different degenerate subspace of the Hilbert space. This can help physicists better understand the properties and structure of fragmented physical systems.

While we examine one specific setting with this property, we believe this strategy for finding physically-motivated problems for quantum learning algorithms can be more generally leveraged. We hope this work motivates future studies of settings where quantum neural networks can take advantage of symmetries that are intractable to compute classically.

## Acknowledgments

## References

[1] Eric Ricardo Anschuetz. 2022. Critical Points in Quantum Generative Models. In *International Conference on Learning Representations*. https://openreview.net/forum?id=2f1z55GVQN

[2] Eric R. Anschuetz. 2025. A Unified Theory of Quantum Neural Network Loss Landscapes. In *International Conference on Learning Representations*. OpenReview. https://openreview.net/forum?id=fv8TTt9srF

[3] Eric R. Anschuetz, Andreas Bauer, Bobak T. Kiani, and Seth Lloyd. 2023. Efficient classical algorithms for simulating symmetric quantum systems. *Quantum* 7 (Nov. 2023), 1189. doi:10.22331/q-2023-11-28-1189

[4] Eric R. Anschuetz and Xun Gao. 2024. Arbitrary Polynomial Separations in Trainable Quantum Machine Learning. arXiv:2402.08606 [quant-ph]

[5] Eric R. Anschuetz, Hong-Ye Hu, Jin-Long Huang, and Xun Gao. 2023. Interpretable Quantum Advantage in Neural Sequence Learning. *PRX Quantum* 4 (6 2023), 020338. Issue 2. doi:10.1103/PRXQuantum.4.020338

[6] Eric R. Anschuetz and Bobak T. Kiani. 2022. Quantum variational algorithms are swamped with traps. *Nat. Commun.* 13 (2022), 7760. Issue 1. doi:10.1038/s41467-022-35364-5

[7] Dave Bacon, Isaac L. Chuang, and Aram W. Harrow. 2006. Efficient Quantum Circuits for Schur and Clebsch-Gordan Transforms. *Phys. Rev. Lett.* 97 (Oct. 2006), 170502. Issue 17. doi:10.1103/PhysRevLett.97.170502

[8] Istvan Berkes and Walter Philipp. 1979. Approximation Thorems for Independent and Weakly Dependent Random Vectors. *The Annals of Probability* 7, 1 (1979), 29 – 54. doi:10.1214/aop/1176995146

[9] Hannes Bernien, Sylvain Schwartz, Alexander Keesling, Harry Levine, Ahmed Omran, Hannes Pichler, Soonwon Choi, Alexander S Zibrov, Manuel Endres, Markus Greiner, et al. 2017. Probing many-body dynamics on a 51-atom quantum simulator. *Nature* 551, 7682 (2017), 579–584. doi:10.1038/nature24622

[10] Dominic W Berry, Yuan Su, Casper Gyurik, Robbie King, Joao Basso, Alexander Del Toro Barba, Abhishek Rajput, Nathan Wiebe, Vedran Dunjko, and Ryan Babbush. 2024. Analyzing prospects for quantum advantage in topological data analysis. *PRX Quantum* 5, 1 (2024), 010319.

[11] M. Cerezo, Martin Larocca, Diego García-Martín, N. L. Diaz, Paolo Braccia, Enrico Fontana, Manuel S. Rudolph, Pablo Bermejo, Aroosa Ijaz, Supanut Thanasilp, Eric R. Anschuetz, and Zoë Holmes. 2025. Does provable absence of barren plateaus imply classical simulability? *Nat. Commun.* 16, 1 (2025), 7907. doi:10.1038/s41467-025-63099-6

[12] Marco Cerezo, Akira Sone, Tyler Volkoff, Lukasz Cincio, and Patrick J Coles. 2021. Cost function dependent barren plateaus in shallow parametrized quantum circuits. *Nat. Commun.* 12, 1 (2021), 1791–1802. doi:10.1038/s41467-021-21728-w

[13] Michele Fava, Jorge Kurchan, and Silvia Pappalardi. 2025. Designs via Free Probability. *Phys. Rev. X* 15 (Feb 2025), 011031. Issue 1. doi:10.1103/PhysRevX.15.011031

[14] Laura Foini and Jorge Kurchan. 2019. Eigenstate thermalization hypothesis and out of time order correlators. *Phys. Rev. E* 99 (April 2019), 042139. Issue 4. doi:10.1103/PhysRevE.99.042139

[15] Joe Gibbs, Zoë Holmes, Matthias C. Caro, Nicholas Ezzell, Hsin-Yuan Huang, Lukasz Cincio, Andrew T. Sornborger, and Patrick J. Coles. 2024. Dynamical simulation via quantum machine learning with provable generalization. *Phys. Rev. Res.* 6 (March 2024), 013241. Issue 1.

[16] Matthew L. Goh, Martin Larocca, Lukasz Cincio, M. Cerezo, and Frédéric Sauvage. 2025. Lie-algebraic classical simulations for quantum computing. *Physical Review Research* 7, 3 (Sept. 2025). doi:10.1103/3y65-f5w6

[17] Casper Gyurik, Alexander Schmidhuber, Robbie King, Vedran Dunjko, and Ryu Hayakawa. 2024. Quantum computing and persistence in topological data analysis. *arXiv preprint arXiv:2410.21258* (2024).

[18] Aram W Harrow, Avinatan Hassidim, and Seth Lloyd. 2009. Quantum algorithm for linear systems of equations. *Physical review letters* 103, 15 (2009), 150502.

[19] Hsin-Yuan Huang, Yunchao Liu, Michael Broughton, Isaac Kim, Anurag Anshu, Zeph Landau, and Jarrod R. McClean. 2024. Learning Shallow Quantum Circuits. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing* (Vancouver, BC, Canada) *(STOC 2024)*. Association for Computing Machinery, New York, NY, USA, 1343–1351. doi:10.1145/3618260.3649722

[20] Tiefeng Jiang. 2010. The Entries of Haar-Invariant Matrices from the Classical Compact Groups. *Journal of Theoretical Probability* 23, 4 (01 Dec 2010), 1227–1243. doi:10.1007/s10959-009-0241-7

[21] Vedika Khemani, Michael Hermele, and Rahul Nandkishore. 2020. Localization from Hilbert space shattering: From theory to physical realizations. *Phys. Rev. B* 101 (May 2020), 174204. Issue 17. doi:10.1103/PhysRevB.101.174204

[22] Martín Larocca, Nathan Ju, Diego García-Martín, Patrick J. Coles, and Marco Cerezo. 2023. Theory of overparametrization in quantum neural networks. *Nature Computational Science* 3, 6 (01 Jun 2023), 542–551. doi:10.1038/s43588-023-00467-6

[23] Jarrod R McClean, Sergio Boixo, Vadim N Smelyanskiy, Ryan Babbush, and Hartmut Neven. 2018. Barren plateaus in quantum neural network training landscapes. *Nat. Commun.* 9, 1 (2018), 4812. doi:10.1038/s41467-018-07090-4

[24] Johannes Jakob Meyer, Marian Mularski, Elies Gil-Fuster, Antonio Anna Mele, Francesco Arzani, Alissa Wilms, and Jens Eisert. 2023. Exploiting Symmetry in Variational Quantum Machine Learning. *PRX Quantum* 4 (3 2023), 010328. Issue 1. doi:10.1103/PRXQuantum.4.010328

[25] Sanjay Moudgalya and Olexei I. Motrunich. 2022. Hilbert Space Fragmentation and Commutant Algebras. *Phys. Rev. X* 12 (March 2022), 011050. Issue 1. doi:10.1103/PhysRevX.12.011050

[26] Sanjay Moudgalya and Olexei I. Motrunich. 2024. Exhaustive Characterization of Quantum Many-Body Scars Using Commutant Algebras. *Physical Review X* 14, 4 (Dec. 2024). doi:10.1103/physrevx.14.041069

[27] Radford M. Neal. 1996. *Priors for Infinite Networks.* Springer New York, New York, NY, 29–53. doi:10.1007/978-1-4612-0745-0_2

[28] Quynh T. Nguyen, Louis Schatzki, Paolo Braccia, Michael Ragone, Patrick J. Coles, Frédéric Sauvage, Martín Larocca, and M. Cerezo. 2024. Theory for Equivariant Quantum Neural Networks. *PRX Quantum* 5 (May 2024), 020328. Issue 2. doi:10.1103/PRXQuantum.5.020328

[29] Shriya Pai, Michael Pretko, and Rahul M. Nandkishore. 2019. Localization in Fractonic Random Circuits. *Phys. Rev. X* 9 (April 2019), 021003. Issue 2. doi:10.1103/PhysRevX.9.021003

[30] Michael Ragone, Bojko N Bakalov, Frédéric Sauvage, Alexander F Kemper, Carlos Ortiz Marrero, Martín Larocca, and M Cerezo. 2024. A Lie algebraic theory of barren plateaus for deep parameterized quantum circuits. *Nat. Commun.* 15, 1 (2024), 7172. doi:10.1038/s41467-024-49909-3

[31] Pablo Sala, Tibor Rakovszky, Ruben Verresen, Michael Knap, and Frank Pollmann. 2020. Ergodicity Breaking Arising from Hilbert Space Fragmentation in Dipole-Conserving Hamiltonians. *Phys. Rev. X* 10 (Feb. 2020), 011047. Issue 1. doi:10.1103/PhysRevX.10.011047

[32] Louis Schatzki, Martin Larocca, Quynh T Nguyen, Frederic Sauvage, and Marco Cerezo. 2024. Theoretical guarantees for permutation-equivariant quantum neural networks. *npj Quantum Inf.* 10, 1 (2024), 12. doi:10.1038/s41534-024-00804-1

[33] Peter W Shor. 1999. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. *SIAM review* 41, 2 (1999), 303–332.

[34] Christopher J Turner, Alexios A Michailidis, Dmitry A Abanin, Maksym Serbyn, and Zlatko Papić. 2018. Weak ergodicity breaking from quantum many-body scars. *Nat. Phys.* 14, 7 (2018), 745–749. doi:10.1038/s41567-018-0137-5

[35] Roman Vershynin. 2018. *High-Dimensional Probability: An Introduction with Applications in Data Science.* Cambridge University Press.

[36] Dan Voiculescu. 1991. Limit laws for Random matrices and free products. *Inventiones mathematicae* 104, 1 (01 Dec 1991), 201–220. doi:10.1007/BF01245072

[37] Maxwell T. West, Jamie Heredge, Martin Sevior, and Muhammad Usman. 2024. Provably Trainable Rotationally Equivariant Quantum Machine Learning. *PRX Quantum* 5 (July 2024), 030320. Issue 3. doi:10.1103/PRXQuantum.5.030320

[38] Xuchen You, Shouvanik Chakrabarti, and Xiaodi Wu. 2022. A Convergence Theory for Over-parameterized Variational Quantum Eigensolvers. arXiv:2205.12481 [quant-ph]

## Appendices

## A    Details on the Lévy-Prokhorov Metric and Error Bounds

*Definition 1 (Lévy-Prokhorov Metric).*  Consider two probability measures $F, G$ over $\mathbb{R}^d$, we define the Lévy-Prokhorov distance between $F$ and $G$ relative to a norm $\|\cdot\|$ on $\mathbb{R}^d$ as

$$\pi(F, G) = \inf\{\epsilon > 0 : F(A) < G(A^\epsilon) + \epsilon \ \forall A\}, \tag{93}$$

where $A$ ranges over the Borel sets in $\mathbb{R}^d$ and $A^\epsilon$ is an $\epsilon$-neighborhood of $A$ in the given norm.

LEMMA 6.  $\pi(F, G) = \pi(G, F)$.

PROOF.  Suppose that for all Borel sets $A$, we have $F(A) < G(A^\epsilon) + \epsilon$. Then,

$$F(A^c) < G((A^c)^\epsilon) + \epsilon. \tag{94}$$

For $\epsilon > 0$, define $A^{-\epsilon}$ to be the set of all points $x$ such that the open $\epsilon$-ball centered at $x$ is contained in $A$. Then, observe that

$$\overline{(A^c)^\epsilon} = \overline{(A^{-\epsilon})^c}, \tag{95}$$

which means that we can write

$$1 - F(A) < 1 - G(A^{-\epsilon}) + \epsilon \tag{96}$$

$$G(A^{-\epsilon}) < F(A) + \epsilon \tag{97}$$

Now, since for all $A$, we have $A \subseteq (A^\epsilon)^{-\epsilon}$, so we can substitute

$$G(A) \leq G((A^\epsilon)^{-\epsilon}) < F(A^\epsilon) + \epsilon. \tag{98}$$

Thus, $\pi(F, G) = \pi(G, F)$.  □

*Definition 2 (Convolution).*  If $\varphi, \gamma$ are the probability densities of the distributions $F, G$, then the probability density of the convolution $F * G$ is defined to be

$$\rho(t) = \int_{\mathbb{R}^d} \varphi(x)\gamma(t - x) \, dx. \tag{99}$$

*Definition 3 (Characteristic Function).*  The characteristic function of a distribution $F$ with density $\varphi$ is defined as

$$f(t) = \mathbb{E}[e^{itX}] = \int_{\mathbb{R}} e^{itx}\varphi(x) \, dx, \tag{100}$$

where in the above $X \sim F$ is a random variable.

LEMMA 7.  *Let $F, G$, and $H$ be probability distributions on $\mathbb{R}^d$, and let $r > 0$. Then, the following inequality holds:*

$$\pi(F, G) \leq \pi(F * H, G * H) + 2 \max\{r, H(\|x\| \geq r)\}, \tag{101}$$

PROOF.  Let $X \sim F, Y \sim H$ be independent random variables. Let $A, K \subseteq \mathbb{R}^d$. We know that

$$\begin{aligned}
(F * H)(A) = \Pr[X + Y \in A] &\geq \\
&\geq \Pr[X + k \in A \ \forall k \in K] \Pr[Y \in K] \geq \\
&\geq \Pr[X + k \in A \ \forall k \in K] + \Pr[Y \in K] - 1 = \\
&= F(\{x : x + K \subseteq A\}) + H(K) - 1.
\end{aligned} \tag{102}$$

Thus, letting $K$ be the open ball of radius $\epsilon$, we have that

$$(F * H)(A^\epsilon) \geq F(A) + H(\|x\| < \epsilon) - 1 = F(A) - H(\|x\| \geq \epsilon). \tag{103}$$

Note that we use the shorthand notation $H(\|x\| < \epsilon)$ to mean $H(\{x : \|x\| < \epsilon\})$.

Now, we can also argue that

$$
\begin{aligned}
(F * H)(A) &= \Pr[X + Y \in A] \leq \\
&\leq \Pr[X \in A - K] \Pr[Y \in K] + \Pr[X \notin A - K] \Pr[Y \in K^c] \leq \\
&\leq \Pr[X \in A - K] + \Pr[Y \in K^c] = F(A - K) + H(K^c).
\end{aligned}
\tag{104}
$$

This gives us the inequality

$$
(F * H)(A) \leq F(A^\epsilon) + H(\|x\| \geq \epsilon).
\tag{105}
$$

The same holds when replacing $F$ with $G$. Now, if we let $L = \pi(F * H, G * H)$, then we must have by the above and by the definition of the Lévy-Prokhorov metric that

$$
\begin{aligned}
F(A) &\leq (F * H)(A^r) + H(\|x\| \geq r) \leq \\
&\leq (G * H)(A^{r+L}) + H(\|x\| \geq r) + L \leq \\
&\leq G(A^{2r+L}) + 2H(\|x\| \geq r) + L \leq \\
&\leq G(A^{L+2\max\{r, H(\|x\| \geq r)\}}) + 2H(\|x\| \geq r) + L + 2\max\{r, H(\|x\| \geq r)\},
\end{aligned}
\tag{106}
$$

which proves that

$$
\pi(F, G) \leq \pi(F * H, G * H) + 2\max\{r, H(\|x\| \geq r)\}.
\tag{107}
$$

<div align="right">□</div>

LEMMA 8. *(This lemma is adapted from Lemma 2.2 of Berkes and Philipp [8], where we have filled in several details). Let $F$ and $G$ be probability distributions on $\mathbb{R}^d$ with characteristic functions $f$ and $g$, respectively. Then, for some sufficiently large $T = O(d)$, we have that the Lévy-Prokhorov metric relative to the Euclidean norm is bounded as*

$$
\pi_2(F, G) < \left(\frac{T}{\pi}\right)^d \int_{\|u\|_2 \leq T} |f(u) - g(u)| \, \mathrm{d}u + F(\{x : \|x\|_2 \geq T/2\}) + 16d \frac{\log T}{T}.
\tag{108}
$$

PROOF. Let $H$ be a distribution on $\mathbb{R}^d$ with density $v(x)$ and characteristic function $h \in L^1$. Let $F_1 = F * H$ and $G_1 = G * H$. Let $\varphi$ and $\gamma$ be the probability densities of $F_1$ and $G_1$, respectively. Their characteristic functions are $f_1 = fh$ and $g_1 = gh$. Now, using the definition of convolution and applying the inverse Fourier transform, we have

$$
\begin{aligned}
|\varphi(x) - \gamma(x)| &= (2\pi)^{-d} \left| \int_{\mathbb{R}^d} e^{-\mathrm{i}\langle u, x\rangle} (f_1(u) - g_1(u)) \, \mathrm{d}u \right| \leq \\
&\leq (2\pi)^{-d} \int_{\mathbb{R}^d} |(f(u) - g(u))| \, |h(u)| \, \mathrm{d}u \leq \\
&\leq (2\pi)^{-d} \left( \int_{\|u\|_2 \leq T} |(f(u) - g(u))| \, \mathrm{d}u + 2 \int_{\|u\|_2 \geq T} |h(u)| \, \mathrm{d}u \right),
\end{aligned}
\tag{109}
$$

where $T$ is any real number. Then, for any Borel set $B \in \mathbb{R}^d$,

$$
\begin{aligned}
F_1(B) - G_1(B) &\leq F_1(B \cap \{\|x\|_2 \leq T\}) - G_1(B \cap \{\|x\|_2 \leq T\}) + F_1(\|x\|_2 \geq T) \leq \\
&\leq \int_{\|x\|_2 \leq T} |\varphi(x) - \gamma(x)| \, \mathrm{d}x + F(\|x\|_2 \geq T/2) + H(\|x\|_2 \geq T/2) \leq \\
&\leq \left(\frac{T}{\pi}\right)^d \left( \int_{\|u\|_2 \leq T} |(f(u) - g(u))| \, \mathrm{d}u + 2 \int_{\|u\|_2 \geq T} |h(u)| \, \mathrm{d}u \right) + F(\|x\|_2 \geq T/2) + H(\|x\|_2 \geq T/2).
\end{aligned}
\tag{110}
$$

Since for all $B$

$$F_1(B) = G_1(B) + (F_1(B) - G_1(B)) \leq G_1(B^{(F_1(B) - G_1(B))}) + (F_1(B) - G_1(B)), \tag{111}$$

we have that

$$\pi(F_1, G_1) \leq |F_1(B) - G_1(B)|. \tag{112}$$

We also know from Lemma 7 that for any $r > 0$,

$$\pi(F, G) \leq \pi(F_1, G_1) + 2 \max\{r, H(\|x\|_2 \geq r)\}. \tag{113}$$

Now, if we choose $r \leq T/2$, we can put this together to get

$$\pi(F, G) \leq \left(\frac{T}{\pi}\right)^d \left(\int_{\|u\|_2 \leq T} |(f(u) - g(u))| \, du + 2 \int_{\|u\|_2 \geq T} |h(u)| \, du\right) + \\ + F(\|x\|_2 \geq T/2) + 3 \max\{r, H(\|x\|_2 \geq r)\}. \tag{114}$$

Now, let $\sigma = 3d^{1/2}T^{-1} \log^{1/2} T$ and $r = 5dT^{-1} \log T$. Let $H$ be a distribution with probability density

$$v(x) = (2\pi\sigma^2)^{-d/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{j=1}^d x_j^2\right). \tag{115}$$

Then,

$$h(u) = \int_{\mathbb{R}^d} e^{\mathbf{i}\langle u, x\rangle} v(x) \, dx = \\
= (2\pi\sigma^2)^{-d/2} \prod_{j=1}^d \int_{\mathbb{R}} e^{\mathbf{i}u_j x - \frac{1}{2\sigma^2} x^2} \, dx = \\
= (2\pi\sigma^2)^{-d/2} \prod_{j=1}^d \int_{\mathbb{R}} e^{-\frac{1}{2\sigma^2}(x^2 - 2\mathbf{i}\sigma^2 u_j x - \sigma^4 u_j^2 + \sigma^4 u_j^2)} \, dx = \\
= (2\pi\sigma^2)^{-d/2} \prod_{j=1}^d \int_{\mathbb{R}} e^{-\frac{1}{2\sigma^2}(x - \mathbf{i}\sigma^2 u_j)^2 - \frac{1}{2}\sigma^2 u_j^2} \, dx = \\
= (2\pi\sigma^2)^{-d/2} \exp\left(-\frac{1}{2}\sigma^2 \sum_{j=1}^d u_j^2\right) \left(\int_{\mathbb{R}} e^{-\frac{1}{2\sigma^2} x^2} \, dx\right)^d = \\
= \exp\left(-\frac{1}{2}\sigma^2 \sum_{j=1}^d u_j^2\right). \tag{116}$$

Note that we are able to remove the $\mathbf{i}\sigma^2 u_j$ term in the exponent in the integrand, because the integrand is an entire function and so its integral around a rectangle with one side on the real axis and one side shifted by $-\mathbf{i}\sigma^2 u_j$ must be zero. Since the the contribution of the other two sides of the rectangle goes to zero, the value of the integral remains the same when moved to the real axis.

So now, we have that

$$\pi(F, G) \leq \left(\frac{T}{\pi}\right)^d \left(\int_{\|u\|_2 \leq T} |(f(u) - g(u))| \, du + 2 \int_{\|u\|_2 \geq T} \exp\left(-\frac{1}{2}\sigma^2 \sum_{j=1}^d u_j^2\right) du\right) + F(\|x\|_2 \geq T/2) + \\
+ 3 \max\left\{r, (2\pi\sigma^2)^{-d/2} \int_{\|u\|_2 \geq r} \exp\left(-\frac{1}{2\sigma^2} \sum_{j=1}^d u_j^2\right) du\right\}. \tag{117}$$

Now, observe that for any $A$,

$$\int_{\|u\|_2 \geq A} \exp\left(-\frac{1}{2}\sum_{j=1}^{d} u_j^2\right) \mathrm{d}u = (2\pi)^{d/2} \Pr[\chi_d^2 \geq A^2] = (2\pi)^{d/2} \Pr\left[e^{\frac{3}{8}\chi_d^2} \geq e^{\frac{3}{8}A^2}\right] \leq$$
$$\leq (2\pi)^{d/2} e^{-3A^2/8} \mathbb{E}\left[e^{\frac{3}{8}\chi_d^2}\right] = (2\pi)^{d/2} e^{-3A^2/8} 2^d. \tag{118}$$

Going back to our previous expression, we can substitute $v_j = \sigma u_j$ for the second integral and $v_j = u_j/\sigma$ for the third, which gives us

$$\pi(F,G) \leq \left(\frac{T}{\pi}\right)^d \left(\int_{\|u\|_2 \leq T} |(f(u)-g(u))| \,\mathrm{d}u + \frac{2}{\sigma^d}\int_{|v| \geq \sigma T} \exp\left(-\frac{1}{2}\sum_{j=1}^{d} v_j^2\right)\mathrm{d}v\right) + F(\|x\|_2 \geq T/2) +$$

$$+ 3\max\left\{r, (2\pi)^{-d/2}\int_{|v| \geq r/\sigma} \exp\left(-\frac{1}{2}\sum_{j=1}^{d} v_j^2\right)\mathrm{d}v\right\} =$$

$$= \left(\frac{T}{\pi}\right)^d \left(\int_{\|u\|_2 \leq T} |(f(u)-g(u))| \,\mathrm{d}u + \frac{2}{\sigma^d}(2\pi)^{d/2}e^{-3T^2\sigma^2/8}2^d\right) + F(\|x\|_2 \geq T/2) + 3\max\left\{r, e^{-\frac{3r^2}{8\sigma^2}}2^d\right\}. \tag{119}$$

Now, since

$$\frac{3r^2}{8\sigma^2} = \frac{3}{8}\frac{25d^2T^{-2}\log^2 T}{9dT^{-2}\log T} = \frac{25d\log T}{24}, \tag{120}$$

so the second term in the maximum becomes

$$T^{-\frac{25d}{24}}2^d \tag{121}$$

which becomes subleading to

$$r = \frac{5d\log T}{T} \tag{122}$$

for large $d$. Also, we have that

$$\frac{2}{\sigma^d}(2\pi)^{d/2}e^{-3T^2\sigma^2/8}2^d = \frac{2^{d+1}(2\pi)^{d/2}}{\left(3d^{1/2}T^{-1}\log^{1/2}T\right)^d}T^{-\frac{27}{8}d}, \tag{123}$$

which becomes subleading if $T$ is sufficiently large. Increasing the coefficient for $r$ by one to strictly dominate the subleading terms, we arrive at the final expression:

$$\pi(F,G) \leq \left(\frac{T}{\pi}\right)^d \int_{\|u\|_2 \leq T} |(f(u)-g(u))| \,\mathrm{d}u + F(\|x\|_2 \geq T/2) + \frac{16d\log T}{T}. \tag{124}$$

$\square$

COROLLARY 1. *Let $F$ and $G$ be two distributions on $\mathbb{R}^d$ with density functions $\varphi$ and $\gamma$ and characteristic functions $f$ and $g$ respectively. Assume that there exists some $C$ such that*

$$\int_{\|x\|_\infty \geq C} f(x)\,\mathrm{d}x \leq \mu. \tag{125}$$

*Then, there exists a universal constant $K$ such that for all $T \geq \max(2C\sqrt{d}, Kd)$, we have that*

$$\pi(F,G) \leq \left(\frac{T}{\pi}\right)^d \int_{\|u\|_\infty \leq T} |(f(u)-g(u))| \,\mathrm{d}u + \frac{16d\log T}{T} + \mu. \tag{126}$$

Proof. We need to convert from the Euclidean norm used in Lemma 8 to the infinity norm. We know that $\|u\|_\infty \le \|u\|_2 \le \sqrt{d}\,\|u\|_\infty$, so after performing that change, the claim immediately follows. □

Lemma 9. (This has been adapted from Corollary 38 of Anschuetz [2]). Let $F, G$ be distributions on $\mathbb{R}^d$ with densities $\varphi, \gamma$ and characteristic functions $f, g$. Assume that each moment of order $k' \le k$ of $F$ differs from that of $G$ by an additive error of at most $\epsilon > 0$ and assume that all moments of order $k' > k$ are bounded by $(C\sqrt{k'})^{k'}$ for some constant $C$. Let $\zeta > 0$ be a sufficiently small constant. Assume that

$$\int_{\|x\|_\infty \ge \frac{1}{2\sqrt{d}} \min\left\{\frac{1-\zeta}{d}\log(\epsilon^{-1}), \frac{\sqrt{k+1}}{2eCd}\right\}} f(x)\,\mathrm{d}x \le \mu. \tag{127}$$

Also assume that

$$d^2 = o\left(\min\left\{\frac{\log(\epsilon^{-1})}{\log\log(\epsilon^{-1})}, \frac{\sqrt{k}}{\log k}\right\}\right). \tag{128}$$

Then, the Lévy-Prokhorov distance is bounded by

$$\pi(F, G) = O\left(\frac{d^2 \log\log(\epsilon^{-1})}{\log(\epsilon^{-1})} + \frac{d^2 \log k}{\sqrt{k}} + \mu\right). \tag{129}$$

Proof. Let

$$T = \min\left\{\frac{1-\zeta}{d}\log(\epsilon^{-1}), \frac{\sqrt{k+1}}{2eCd}\right\}. \tag{130}$$

This, with the stated assumptions, satisfies the conditions for applying Corollary 1, which gives us that

$$\pi(F, G) \le \tilde{O}(T^{2d}) \sup_{\|u\|_\infty \le T} |f(u) - g(u)| + O\left(\frac{d \log T}{T}\right) + O(\mu). \tag{131}$$

Since

$$dT \le \frac{\sqrt{t+1}}{2eC}, \tag{132}$$

we can apply Lemma 37 of Anschuetz [2] to say that

$$\sup_{\|u\|_\infty \le T} |f(u) - g(u)| \le \sup_{\|u\|_1 \le dT} |f(u) - g(u)| \le \epsilon e^{dT} + 2^{1-k} \le \epsilon^\zeta + 2^{1-t}. \tag{133}$$

We also have that

$$T^{2d} = \exp(2d \log T) = \exp\left(o(\min\{\log(\epsilon^{-1}), \sqrt{t}\})\right). \tag{134}$$

So,

$$\tilde{O}(T^{2d}) \sup_{\|u\|_\infty \le T} |f(u) - g(u)| + O\left(\frac{d \log T}{T}\right) \le O\left(\epsilon^{-1+\zeta} + e^{\sqrt{k}+(1-k)\log 2}\right). \tag{135}$$

Now, looking at the second term, we have that

$$O\left(\frac{d \log T}{T}\right) \le O\left(d \min\left\{\log\log(\epsilon^{-1}), \log k\right\} \max\left\{\frac{d}{\log(\epsilon^{-1})}, \frac{d}{\sqrt{k}}\right\}\right) \le O\left(\frac{d^2 \log\log(\epsilon^{-1})}{\log(\epsilon^{-1})} + \frac{d^2 \log k}{\sqrt{k}}\right). \tag{136}$$

This second term dominates the first, which means that

$$\pi(F, G) = O\left(\frac{d^2 \log\log(\epsilon^{-1})}{\log(\epsilon^{-1})} + \frac{d^2 \log k}{\sqrt{k}} + \mu\right). \tag{137}$$

□