

Trainability of Dissipative Perceptron-Based Quantum Neural Networks

Kunal Sharma,^{1,2,*} M. Cerezo,^{1,3,*} Lukasz Cincio,¹ and Patrick J. Coles¹

¹Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

²Hearne Institute for Theoretical Physics and Department of Physics and Astronomy,
Louisiana State University, Baton Rouge, Louisiana 70803, USA

³Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

Several architectures have been proposed for quantum neural networks (QNNs), with the goal of efficiently performing machine learning tasks on quantum data. Rigorous scaling results are urgently needed for specific QNN constructions to understand which, if any, will be trainable at a large scale. Here, we analyze the gradient scaling (and hence the trainability) for a recently proposed architecture that we call dissipative QNNs (DQNNs), where the input qubits of each layer are discarded at the layer’s output. We find that DQNNs can exhibit barren plateaus, i.e., gradients that vanish exponentially in the number of qubits. Moreover, we provide quantitative bounds on the scaling of the gradient for DQNNs under different conditions, such as different cost functions and circuit depths, and show that trainability is not always guaranteed. Our work represents the first rigorous analysis of the scalability of a perceptron-based QNN.

Introduction.—Neural networks (NN) have impacted many fields such as neuroscience, engineering, computer science, chemistry, and physics [1]. However, their historical development has seen periods of great progress interleaved with periods of stagnation, due to serious technical challenges [2]. The perceptron was introduced early on as an artificial neuron [3], but it was only realized later that a multilayer perceptron (now known as a feedforward NN) had much greater power than a single-layer one [1, 2]. Still there was the major issue of how to train multiple layers, and this was eventually addressed by the backpropagation method [4].

Motivated by the success of NNs and the advent of noisy intermediate-scale quantum devices [5], there has been tremendous effort to develop quantum neural networks (QNNs) [6]. The hope is that QNNs will harness the power of quantum computers to outperform their classical counterparts on machine learning tasks [7, 8], especially for quantum data or tasks that are inherently quantum in nature [9].

Despite several QNN proposals that have been successfully implemented [10–17], more research is needed on the advantages and limitations of specific architectures. Delving into potential scalability issues of QNNs could help to prevent a “winter” for these models, like what was seen historically for classical NNs. This has motivated recent works studying the scaling of gradients in QNNs [18, 19]. There, it was shown that variational quantum algorithms [20–30], which aim to train QNNs to accomplish specific tasks, may exhibit gradients that vanish exponentially with the system size. This so-called barren-plateau phenomenon, where the parameters cannot be efficiently trained for large implementations, was demonstrated for hardware-efficient QNNs, where quantum gates are arranged in a bricklike structure that matches the connectivity of the quantum device [18, 19].

Analyzing the existence of barren plateaus in QNNs is paramount to determining if they can lead to a quantum speedup. This is due to the fact that exponentially vanishing gradients imply that the precision needed to estimate such gradients grows exponentially. Since the standard goal of quantum algorithms is polynomial scaling as opposed to the typical exponential scaling of classical algorithms, a QNN with exponentially vanishing gradients has no hope of achieving this goal. On the other hand, a QNN with gradients that vanish polynomially means that the algorithm requires a polynomial precision, and hence that the hope of quantum speedup is preserved.

Here, we analyze the trainability and the existence of barren plateaus in a class of QNNs that we refer to as *dissipative QNNs* (DQNNs). In a DQNN each node within the network corresponds to a qubit [31], and the connections in the network are modelled by quantum perceptrons [32–37]. The term dissipative refers to the fact that ancillary qubits form the output layer, while the qubits from the input layer are discarded. This architecture has seen significant recent attention and has been proposed as a scalable approach to QNNs [37–39]. In particular, in Ref. [37], based on small scale numerical experiments, it was speculated that dissipative quantum neural networks do not suffer from the barren plateau (vanishing gradient) problem. However, contrary to Ref. [37], we here analytically prove that DQNNs are not immune to barren plateaus. For example, DQNNs with deep global perceptrons are untrainable despite the dissipative nature of the architecture.

Here we study the large-scale trainability of DQNNs. In particular, we focus on tasks where DQNNs are employed to learn a unitary matrix connecting input and output quantum states and for general supervised quantum machine learning tasks where training data consists of quantum states and corresponding classical labels. For these tasks, we show that the barren plateau phenomenon can also arise in DQNNs. We also discuss certain conditions (e.g., the structure and depth of the DQNN) under which one could avoid a barren plateau and achieve train-

* The first two authors contributed equally to this work.

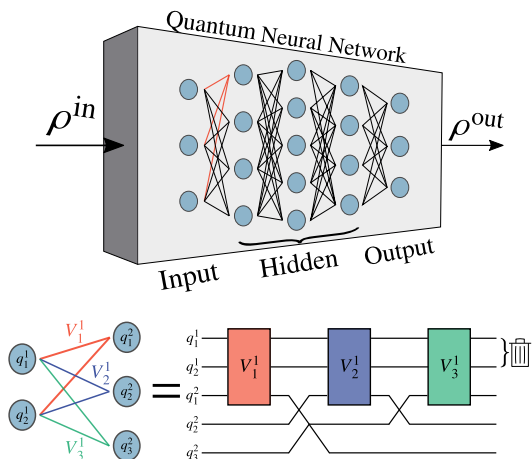


FIG. 1. Schematic diagram of a dissipative perceptron-based quantum neural network (DQNN). Top: The DQNN is composed of input, hidden, and output layers. Each node in the network corresponds to a qubit, which can be connected to qubits in adjacent layers via perceptrons (depicted as lines). The input and output of the DQNN are quantum states denoted as ρ^{in} and ρ^{out} , respectively. Bottom: Quantum circuit description of the DQNN. The j th qubit of the l th layer is denoted q_j^l . Each perceptron corresponds to a unitary operation on the qubits it connects, with V_j^l denoting the j -th perceptron in the l -th layer.

ability. In particular, our work implies that scalability is not guaranteed, and without careful thought of the structure of DQNNs, their gradients may vanish exponentially in the system size. As a by-product of our analysis of specific perceptron architectures, we also show that hardware-efficient QNNs are special cases of DQNNs. Therefore, many important results for hardware-efficient QNNs, such as the ones studied in Refs. [18, 19] also hold for DQNNs. Finally, we remark that we employ novel analytical techniques in our proofs (different from those used in Refs. [18, 19]), which were necessary to develop due to the dissipative nature of DQNNs. Our techniques may be broadly useful in the study of the scaling of other QNN architectures.

Preliminaries.— Let us first introduce the DQNN architecture. As schematically shown in Fig. 1, the DQNN is composed of a series of layers (input, hidden, and output) where the qubits at each node are connected via perceptrons. A quantum perceptron is defined as an arbitrary unitary operator with m input and k output qubits. For simplicity, we consider the case when $k = 1$, so that each perceptron acts on $m + 1$ qubits. The case of arbitrary k is presented in the Supplemental Material.

The qubits in the input layer are initialized to a state ρ^{in} , while all qubits in the hidden and output layers are initialized to a fiducial state such as $|\mathbf{0}\rangle_{\text{hid,out}} = |0 \dots 0\rangle_{\text{hid,out}}$. Henceforth we employ the notation “in”, “hid”, and “out” to indicate operators on qubits in the input, hidden, and output layers, respectively. The output state of the DQNN is a quantum state ρ^{out} (generally

mixed) which can be expressed as

$$\rho^{\text{out}} \equiv \text{Tr}_{\text{in,hid}} [V(\rho^{\text{in}} \otimes |\mathbf{0}\rangle_{\text{hid,out}} \langle \mathbf{0}|) V^\dagger], \quad (1)$$

with $V = V_{n_{\text{out}}}^{\text{out}} \dots V_{n_1}^1 \dots V_1^1$, and where V_j^l is the perceptron unitary on the l -th layer acting on the j -th output qubit. Here n_l indicates the number of qubits in the l -th layer.

Let us now make two important remarks. First, note that the order in which the perceptrons act is relevant, as in general the unitaries V_j^l will not commute. Second, we remark that for this architecture the perceptrons are applied layer by layer, meaning that once all V_j^l (for fixed l) have been applied and the information has propagated forward between layers $l - 1$ and l , one can discard the qubits in layer $l - 1$. This implies that the width of the DQNN depends on the number of qubits in two adjacent layers and not in the total number of qubits in the network.

To train the DQNN, we assume repeatable access to training data in the form of pairs $\{|\phi_x^{\text{in}}\rangle, |\phi_x^{\text{out}}\rangle\}$, with $x = 1, \dots, N$. We note that, as discussed in the Supplemental Material, our results also hold more generally for supervised quantum machine learning tasks where the training data is of the form $\{|\phi_x^{\text{in}}\rangle, y_x\}$, with y_x a label assigned to the input state $|\phi_x^{\text{in}}\rangle$ [40].

We then define a cost function (or loss function) which quantifies how well the DQNN reproduces the training data. We assume that the cost is of the form

$$C = \frac{1}{N} \sum_{x=1}^N C_x, \quad \text{with} \quad C_x = \text{Tr}[O_x \rho_x^{\text{out}}]. \quad (2)$$

As discussed below, in general there are multiple choices for the operator O_x which lead to faithful cost functions, i.e., costs that are extremized if and only if one perfectly learns the mapping on the training data. If the circuit description of output states is provided, one can employ the inverse of the corresponding unitary on the output of a DQNN [41]. Then a measurement in the computational basis estimates the cost function. Otherwise, one can employ a recently developed procedure based on classical shadows to estimate the state overlap [42].

When O_x acts non-trivially on all qubits of the output layer, we use the term *global cost function*, denoted as C^G . Here one usually compares objects (states or operators) living in exponentially large Hilbert spaces. For instance, choosing

$$O_x^G = \mathbb{1} - |\phi_x^{\text{out}}\rangle \langle \phi_x^{\text{out}}|, \quad (3)$$

leads to a global cost function that quantifies the average fidelity between each ρ_x^{out} and $|\phi_x^{\text{out}}\rangle$.

As shown in Ref. [19], local cost functions do not exhibit a barren plateau for shallow hardware-efficient QNNs. Therefore, it is important to study if local observables can also lead to trainability guarantees in DQNNs. Henceforth, we use the term *local cost function*, denoted as C^L , for the cases when the operator O_x acts non-trivially on a small number of qubits in the output layer.

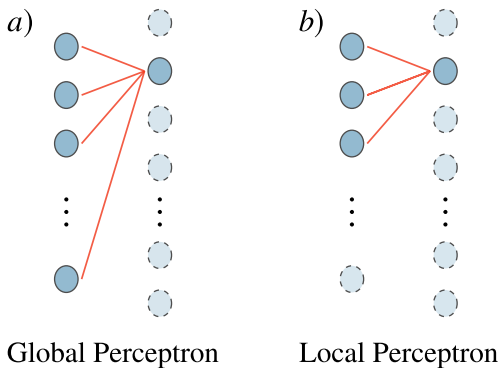


FIG. 2. Global and local perceptrons. a) The global perceptron acts non-trivially on all input qubits, i.e., $m = n$. b) The local perceptron acts non-trivially only on a small number of input qubits. For the case shown, $m = 3$.

Since the global cost in Eq. (3) is a state fidelity function, in general it will not be possible to design a corresponding faithful local cost. Therefore, we restrict ourselves to the case when $|\phi_x^{\text{out}}\rangle$ is a tensor-product state across n_{out} qubits $|\phi_x^{\text{out}}\rangle = |\psi_{x,1}^{\text{out}}\rangle \otimes \dots \otimes |\psi_{x,n_{\text{out}}}^{\text{out}}\rangle$. Then, we can define the following local observable:

$$O_x^L = \mathbb{1} - \frac{1}{n_{\text{out}}} \sum_{j=1}^{n_{\text{out}}} |\psi_{x,j}^{\text{out}}\rangle \langle \psi_{x,j}^{\text{out}}| \otimes \mathbb{1}_{\bar{j}}, \quad (4)$$

where $\mathbb{1}_{\bar{j}}$ denotes the identity over all qubits in the output layer except for qubit j . Equation (4) leads to a faithful local cost that vanishes under the same condition as the global cost defined from Eq. (3) [41, 43].

Finally let us introduce the term *global perceptron* to refer to the case when the perceptron V_j^l acts non-trivially on *all* qubits in the l -th layer, i.e., when $m = n_{l-1}$. On the other hand, a *local perceptron* is defined as a unitary V_j^l acting on a number of qubits $m \in \mathcal{O}(1)$ which is independent of n_{l-1} . Figure 2 schematically shows a global and a local perceptron.

To analyze the existence of barren plateaus and the trainability of the DQNN one needs to define an ansatz and a training method for the perceptrons. In what follows we consider two general training approaches.

Random parameterized quantum circuits.—We first consider the case where the perceptrons are parameterized quantum circuits (i.e., variational circuits) that can be expressed as a sequence of parameterized and unparameterized gates from a given gate alphabet [18, 44]. That is, the perceptrons are of the form

$$V_j^l(\theta_j^l) = \prod_{k=1}^{\eta_j^l} R_k(\theta^k) W_k, \quad (5)$$

with $R_k(\theta^k) = e^{-(i/2)\theta^k \Gamma_k}$, W_k an unparameterized unitary, and where Γ_k is a Hermitian operator with $\text{Tr}[\Gamma_k^2] \leq$

2^{n+1} . Such parameterization is widely used as it can allow for a straightforward evaluation of the cost function gradients, and since in general its quantum circuit description can be easily obtained [45–47].

A common strategy for training random parameterized quantum circuits is to randomly initialize the parameters in (5), and employ a training loop to minimize the cost function. To analyze the trainability of the DQNN we compute the variance of the partial derivative $\partial C / \partial \theta^\nu \equiv \partial_\nu C$, where θ^ν belongs to a given V_j^l

$$\text{Var}[\partial_\nu C] = \langle (\partial_\nu C)^2 \rangle - \langle \partial_\nu C \rangle^2. \quad (6)$$

Here the notation $\langle \dots \rangle$ indicates the average over all randomly initialized perceptrons. From (5), we find

$$\partial_\nu C = \frac{i}{2N} \sum_{x=1}^N \text{Tr} \left[A_j^l \tilde{\rho}_x^{\text{in}} (A_j^l)^\dagger [\mathbb{1}_{\bar{j}} \otimes \Gamma_k, (B_j^l)^\dagger \tilde{O}_x B_j^l] \right], \quad (7)$$

where we have defined

$$B_j^l = \mathbb{1}_{\bar{j}} \otimes \prod_{k=1}^{\nu-1} R_k(\theta^k) W_k, \quad A_j^l = \mathbb{1}_{\bar{j}} \otimes \prod_{k=\nu}^{\eta_j^l} R_k(\theta^k) W_k, \quad (8)$$

such that $\mathbb{1}_{\bar{j}} \otimes V_j^l = A_j^l B_j^l$, and where $\mathbb{1}_{\bar{j}}$ indicates the identity on all qubits on which V_j^l does not act. Note that the trace in (7) is over *all* qubits in the DQNN. In addition, we define

$$\tilde{\rho}_x^{\text{in}} = V_{j-1}^l \dots V_1^l (\rho_x^{\text{in}} \otimes |\mathbf{0}\rangle \langle \mathbf{0}|_{\text{hid,out}}) (V_1^l)^\dagger \dots (V_{j-1}^l)^\dagger, \\ \tilde{O}_x = (V_{j+1}^l)^\dagger \dots (V_{n_{\text{out}}}^{\text{out}})^\dagger (\mathbb{1}_{\text{in,hid}} \otimes O_x) V_{n_{\text{out}}}^{\text{out}} \dots V_{j+1}^l.$$

If the perceptron V_j^l is sufficiently random so that A_j^l , B_j^l , or both, form independent unitary 1-designs, then we find that $\langle \partial_\nu C \rangle = 0$ (see Supplemental Material). In this case, $\text{Var}[\partial_\nu C]$ quantifies (on average) how much the gradient concentrates around zero. Hence, exponentially small $\text{Var}[\partial_\nu C]$ values would imply that the slope of the cost function landscape is insufficient to provide cost-minimizing directions.

Here we recall that a t -design is a set of unitaries $\{V_y \in U(d)\}_{y \in Y}$ (of size $|Y|$) on a d -dimensional Hilbert space such that for every polynomial $P_t(V_y)$ of degree at most t in the matrix elements of V_y , and of V_y^\dagger one has [48] $\langle P_t(V) \rangle_V = \frac{1}{|Y|} \sum_{y \in Y} P_t(V_y) = \int d\mu(V) P_t(V)$, where the integral is over the unitary group $U(d)$.

Let us assume for simplicity the case when the DQNN input and output layers have the same number of qubits ($n_{\text{in}} = n_{\text{out}} = n$). As shown in the Supplemental Material, the following theorem holds.

Theorem 1. *Consider a DQNN with deep global perceptrons parameterized as in Eq. (5), such that A_j^l , B_j^l in Eq. (8) and V_j^l ($\forall j, l$) form independent 2-designs over $n+1$ qubits. Then, the variance of the partial derivative of the cost function with respect to θ^ν in V_j^l is upper bounded as*

$$\text{Var}[\partial_\nu C^G] \leq g(n), \quad \text{with } g(n) \in \mathcal{O}(1/2^{2n}), \quad (9)$$

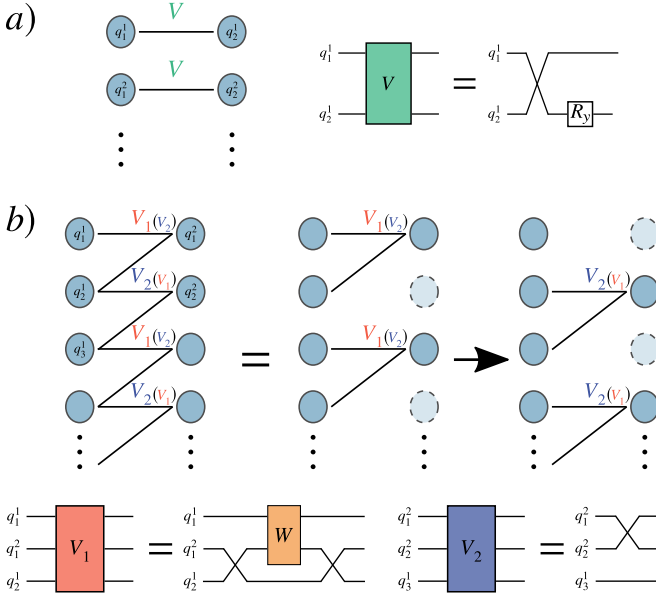


FIG. 3. Shallow local perceptrons ansatzes. a) Here $m = 1$ so that each perceptron acts on a single input and output qubit. Moreover, for all j and l we have $V_j^l = V$. The unitaries V are simply given by a SWAP operator followed by a single qubit rotation around the y axis. b) Local perceptrons V_j^l with $m = 2$. The local perceptrons are given by the unitaries V_1 , or V_2 . Specifically, for l odd on j odd (even) $V_j^l = V_1(V_2)$, while for l even and j odd (even) we have $V_j^l = V_2(V_1)$. Here we also show the order in which the perceptrons are applied so that we first implement the unitaries with j odd, followed by the unitaries with j even. The W gate in V_1 forms a local 2-design on two qubits.

if O_x is the global operator of Eq. (3), and upper bounded as

$$\text{Var}[\partial_\nu C^L] \leq h(n), \quad \text{with } h(n) \in \mathcal{O}(1/2^n), \quad (10)$$

when O_x is the local operator in (4).

Theorem 1 shows that DQNNs with deep global perceptron unitaries that form two-designs [49, 50] exhibit barren plateaus for global and local cost functions. An immediate question that follows is whether barren plateaus still arise for shallow perceptrons, which cannot form 2-designs on $n + 1$ qubits. In what follows we analyze specific cases of shallow local perceptrons for which results can be obtained.

Let us first consider the simple perceptrons of Fig. 3(a), where $m = 1$, and where R_y denotes a single qubit rotation around the y axis: $R_y(\theta^\nu) = e^{-i\theta^\nu Y/2}$ (with all angles randomly initialized). In this case one recovers the toy model example of [19], and we know that if O_x is the global operator of (3), then $\text{Var}[\partial_\nu C^G] = \frac{1}{8} \left(\frac{3}{8}\right)^{n-1}$. On the other hand, if O_x is the local operator in (4), then $\text{Var}[\partial_\nu C^L] = \frac{1}{8n^2}$.

These results suggest that DQNNs with simple shallow local perceptrons and global cost functions are untrainable when randomly initialized. On the other hand, they

also indicate that barren plateaus for DQNNs might be avoided by employing: (1) shallow (local) perceptrons, and (2) local cost functions.

Let us now consider the shallow local perceptron of Fig. 3(b), where each unitary W forms a local 2-design on two qubits. For this architecture the ensuing DQNN can be *exactly* mapped into a layered hardware-efficient ansatz as in [19], where two layers of the DQNN correspond to a single layer of the hardware-efficient ansatz [51]. Note that this mapping is not general, but rather valid for the specific architecture in Fig. 3(b). As shown in Ref. [19], when employing a global cost function, with O_x given by (3), one finds that if the number of layers is $\mathcal{O}(\text{poly}(\log(n)))$, then the DQNN cost function exhibits barren plateaus as

$$\text{Var}[\partial_\nu C^G] \leq \hat{f}(n), \quad \text{with } \hat{f}(n) \in \mathcal{O}\left((\sqrt{3}/4)^n\right). \quad (11)$$

On the other hand, for a local cost function with O_x given by (4), if the number of layers is in $\mathcal{O}(\log(n))$, then there is no barren plateau [19] as

$$\hat{g}(n) \leq \text{Var}[\partial_\nu C^L], \quad \text{with } \hat{g}(n) \in \Omega(1/\text{poly}(n)). \quad (12)$$

Here we remark that (12) was obtained following the same assumptions as those used in Corollary 2 of [19]. Note that obtaining a lower bound for the variance implies that the DQNN trainability is guaranteed.

Parameter matrix multiplication.—While in random parametrized quantum circuits one optimizes and trains a single gate angle at a time, other optimization approaches can also be considered. In what follows we analyze the trainability for a method introduced in Ref. [37] where at each time-step all perceptrons are simultaneously optimized.

In this training approach, which we call parameter matrix multiplication, the perceptrons are not explicitly decomposed into quantum circuits, but rather are treated as unitary matrices. The perceptrons $V_j^l(0)$ are randomly initialized at time-step zero, and at each step s they are updated via

$$V_j^l(s + \varepsilon) = e^{i\varepsilon H_j^l(s)} V_j^l(s). \quad (13)$$

The matrices H_j^l are such that $\text{Tr}[(H_j^l)^2] \leq 2^{n+1}$ and are parametrized as $H_j^l(s) = \sum_{\mathbf{u}\mathbf{v}} h_{j,\mathbf{u},\mathbf{v}}^l X^{\mathbf{u}} Z^{\mathbf{v}}$, with $X^{\mathbf{u}} Z^{\mathbf{v}} = X_1^{u_1} Z_1^{v_1} \otimes X_2^{u_2} Z_2^{v_2} \dots$, and where X_j and Z_j are Pauli operators on qubit j . The matrices $K_j^l(s)$ are called *parameter matrices*, and at each time-step the coefficients $h_{j,\mathbf{u},\mathbf{v}}^l$ need to be optimized. As shown in the Supplemental Material, if at least one perceptron $V_j^l(0)$ is sufficiently random so that it forms a global unitary 1-design, then we find $\langle \partial C / \partial s \rangle \equiv \langle \partial_s C \rangle = 0$.

As proved in the Supplemental Material, the following theorem holds.

Theorem 2. Consider a DQNN with deep global perceptrons, which are updated via the parameter matrix multiplication of (13). Suppose that for all j, l , the $V_j^l(0)$ perceptrons form independent 2-designs over $n + 1$ qubits. Then the variance of the partial derivative of the cost function with respect to the time-step parameter s is upper bounded as

$$\text{Var}[\partial_s C] \leq f(n), \quad \text{with } f(n) \in \mathcal{O}(1/2^n), \quad (14)$$

when O_x is the global operator of (3), or the local operator in (4).

Although the updating method in (13) simultaneously updates all perceptrons at each time-step, Theorem 2 implies that barren plateaus also arise when using the parameter matrix multiplication method.

We note that our proof techniques invoke the pure state properties of input and output states. Since the output state of a randomly initialized DQNN will be close to a maximally mixed across any bipartite cut [52], we speculate that our results can be extended to expectation values of the arbitrary Hamiltonian. We leave this question for future work.

Conclusions.—In this Letter, we analyzed the trainability of a special class of QNNs called DQNNs. We first proved that the trainability of DQNNs is not always guaranteed as they can exhibit barren plateaus in their cost function landscape. The existence of such barren plateaus was linked to the localities (i.e., the number of qubits they act non-trivially on) of the perceptrons and of the cost function. Specifically, we showed that: (i) DQNNs with deep global perceptrons are untrainable despite the dissipative nature of the architecture, and (ii) for shallow and local perceptrons, employing global cost functions leads to barren plateaus, while using local costs avoids them. We note that our results are completely general for DQNN architectures, e.g., covering arbitrary numbers of hidden layers and general perceptrons acting on any number of qubits.

In addition, we provided a specific architecture for DQNNs with local shallow perceptrons that can be exactly mapped to a layered hardware-efficient ansatz. This result not only indicates that hardware-efficient QNNs can be represented as DQNNs, but it also allows us to derive trainability guarantees for these DQNNs. In this case, since the perceptrons are local, each neuron only receives information from a small number of qubits in the

previous layer. Such architecture is reminiscent of classical convolutional neural networks, which are known to avoid some of the trainability problems of fully connected networks [53].

These results show that much work needs to be done to understand the trainability of QNNs and guarantee that they can provide a quantum speedup over classical neural networks. For instance, interesting future research directions are QNN-specific optimizers [54–57], analyzing the resilience of QNNs to noise [22, 41], and strategies to prevent barren plateaus [58–61]. Another interesting direction is to extend our results to the case when the input and output states are mixed states, particularly when the goal is to match marginals of the target output state and the output of a DQNN [62]. Furthermore, exploring architectures beyond DQNNs and hardware-efficient QNNs would be of interest, particularly if such architectures have large-scale trainability.

Acknowledgments

We thank Jarrod McClean, Tobias Osborne, and Andrew Sornborger for helpful conversations. All authors acknowledge support from LANL’s Laboratory Directed Research and Development (LDRD) program. MC was also supported by the Center for Nonlinear Studies at LANL. PJC also acknowledges support from the LANL ASC Beyond Moore’s Law project. This work was also supported by the U.S. Department of Energy (DOE), Office of Science, Office of Advanced Scientific Computing Research, under the Accelerated Research in Quantum Computing (ARQC) program.

Supplemental Material.—The Supplemental Material contains details of our proofs and References [63, 64].

Note Added.—Our work is the first to analyze barren plateaus in the context of data science applications, and also the first to consider perceptron-based quantum neural networks (QNNs). Our work has inspired more recent studies of trainability for other QNN architectures, such as quantum convolutional neural networks [65], tree-based architectures [66], and others [67–69]. We also note that our results can also be interpreted as a type of entanglement-induced barren plateau. Here, a large amount of entanglement in a parameterized quantum circuit can lead to trainability issues when qubits are discarded, and the output qubits are concentrated around the maximally mixed state. This phenomenon was further studied in [52, 70].

[1] Simon Haykin, *Neural networks: a comprehensive foundation* (Prentice Hall PTR, NJ, 1994).
 [2] Marvin Minsky and Seymour A Papert, *Perceptrons: An introduction to computational geometry* (MIT press, Cambridge, MA, 2017).

[3] Frank Rosenblatt, *The perceptron, a perceiving and recognizing automaton Project Para* (Cornell Aeronautical Laboratory, 1957).
 [4] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams, “Learning representations by back-propagating

- errors,” *Nature (London)* **323**, 533–536 (1986).
- [5] J. Preskill, “Quantum computing in the NISQ era and beyond,” *Quantum* **2**, 79 (2018).
- [6] Maria Schuld, Ilya Sinayskiy, and Francesco Petruccione, “The quest for a quantum neural network,” *Quantum Information Processing* **13**, 2567–2586 (2014).
- [7] Michael A Nielsen, *Neural networks and deep learning*, Vol. 2018 (Determination press San Francisco, CA, USA:, 2015).
- [8] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd, “Quantum machine learning,” *Nature* **549**, 195–202 (2017).
- [9] Kunal Sharma, M Cerezo, Zoë Holmes, Lukasz Cincio, Andrew Sornborger, and Patrick J Coles, “Reformulation of the No-Free-Lunch theorem for entangled data sets,” *Phys. Rev. Lett.* **128**, 070501 (2022).
- [10] J. Romero, J. P. Olson, and A. Aspuru-Guzik, “Quantum autoencoders for efficient compression of quantum data,” *Quantum Science and Technology* **2**, 045001 (2017).
- [11] Vedran Dunjko and Hans J Briegel, “Machine learning & artificial intelligence in the quantum domain: a review of recent progress,” *Reports on Progress in Physics* **81**, 074001 (2018).
- [12] Guillaume Verdon, Jason Pye, and Michael Broughton, “A universal training algorithm for quantum deep learning,” *arXiv preprint arXiv:1806.09729* (2018).
- [13] Edward Farhi and Hartmut Neven, “Classification with quantum neural networks on near term processors,” *arXiv preprint arXiv:1802.06002* (2018).
- [14] Carlo Ciliberto, Mark Herbster, Alessandro Davide Ialongo, Massimiliano Pontil, Andrea Rocchetto, Simone Severini, and Leonard Wossnig, “Quantum machine learning: a classical perspective,” *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* **474**, 20170551 (2018).
- [15] Nathan Killoran, Thomas R Bromley, Juan Miguel Arrazola, Maria Schuld, Nicolás Quesada, and Seth Lloyd, “Continuous-variable quantum neural networks,” *Physical Review Research* **1**, 033063 (2019).
- [16] Iris Cong, Soonwon Choi, and Mikhail D Lukin, “Quantum convolutional neural networks,” *Nature Physics* **15**, 1273–1278 (2019).
- [17] Zhih-Ahn Jia, Biao Yi, Rui Zhai, Yu-Chun Wu, Guang-Can Guo, and Guo-Ping Guo, “Quantum neural network states: A brief review of methods and applications,” *Advanced Quantum Technologies* **2**, 1800077 (2019).
- [18] Jarrod R McClean, Sergio Boixo, Vadim N Smelyanskiy, Ryan Babbush, and Hartmut Neven, “Barren plateaus in quantum neural network training landscapes,” *Nature communications* **9**, 4812 (2018).
- [19] Marco Cerezo, Akira Sone, Tyler Volkoff, Lukasz Cincio, and Patrick J Coles, “Cost function dependent barren plateaus in shallow parametrized quantum circuits,” *Nature communications* **12**, 1–12 (2021).
- [20] 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**, 4213 (2014).
- [21] Bela Bauer, Dave Wecker, Andrew J Millis, Matthew B Hastings, and Matthias Troyer, “Hybrid quantum-classical approach to correlated materials,” *Physical Review X* **6**, 031045 (2016).
- [22] Jarrod R McClean, Jonathan Romero, Ryan Babbush, and Alán Aspuru-Guzik, “The theory of variational hybrid quantum-classical algorithms,” *New Journal of Physics* **18**, 023023 (2016).
- [23] A. Arrasmith, L. Cincio, A. T. Sornborger, W. H. Zurek, and P. J. Coles, “Variational consistent histories as a hybrid algorithm for quantum foundations,” *Nature communications* **10**, 3438 (2019).
- [24] Tyson Jones, Suguru Endo, Sam McArdle, Xiao Yuan, and Simon C Benjamin, “Variational quantum algorithms for discovering hamiltonian spectra,” *Physical Review A* **99**, 062304 (2019).
- [25] X. Xu, J. Sun, S. Endo, Y. Li, S. C. Benjamin, and X. Yuan, “Variational algorithms for linear algebra,” *arXiv:1909.03898 [quant-ph]*.
- [26] Carlos Bravo-Prieto, Ryan LaRose, M. Cerezo, Yigit Subasi, Lukasz Cincio, and Patrick J. Coles, “Variational quantum linear solver: A hybrid algorithm for linear systems,” *arXiv:1909.05820* (2019).
- [27] Xiao Yuan, Suguru Endo, Qi Zhao, Ying Li, and Simon C Benjamin, “Theory of variational quantum simulation,” *Quantum* **3**, 191 (2019).
- [28] Cristina Cirstoiu, Zoe Holmes, Joseph Iosue, Lukasz Cincio, Patrick J Coles, and Andrew Sornborger, “Variational fast forwarding for quantum simulation beyond the coherence time,” *arXiv preprint arXiv:1910.04292* (2019).
- [29] Marco Cerezo, Alexander Poremba, Lukasz Cincio, and Patrick J Coles, “Variational quantum fidelity estimation,” *Quantum* **4**, 248 (2020).
- [30] M Cerezo, Kunal Sharma, Andrew Arrasmith, and Patrick J Coles, “Variational quantum state eigensolver,” *arXiv preprint arXiv:2004.01372* (2020).
- [31] Noriaki Kouda, Nobuyuki Matsui, Haruhiko Nishimura, and Ferdinand Peper, “Qubit neural network and its learning efficiency,” *Neural Computing & Applications* **14**, 114–121 (2005).
- [32] MV Altaisky, “Quantum neural network,” *arXiv preprint quant-ph/0107012* (2001).
- [33] Alaa Sagheer and Mohammed Zidan, “Autonomous quantum perceptron neural network,” *arXiv preprint arXiv:1312.4149* (2013).
- [34] Michael Siomau, “A quantum model for autonomous learning automata,” *Quantum information processing* **13**, 1211–1221 (2014).
- [35] Erik Torrontegui and Juan José García-Ripoll, “Unitary quantum perceptron as efficient universal approximator,” *EPL (Europhysics Letters)* **125**, 30004 (2019).
- [36] Francesco Tacchino, Chiara Macchiavello, Dario Gerace, and Daniele Bajoni, “An artificial neuron implemented on an actual quantum processor,” *npj Quantum Information* **5**, 1–8 (2019).
- [37] Kerstin Beer, Dmytro Bondarenko, Terry Farrelly, Tobias J Osborne, Robert Salzmann, Daniel Scheiermann, and Ramona Wolf, “Training deep quantum neural networks,” *Nature Communications* **11**, 1–6 (2020).
- [38] Dmytro Bondarenko and Polina Feldmann, “Quantum autoencoders to denoise quantum data,” *Physical Review Letters* **124**, 130502 (2020).
- [39] Kyle Poland, Kerstin Beer, and Tobias J Osborne, “No free lunch for quantum machine learning,” *arXiv preprint arXiv:2003.14103* (2020).
- [40] Vojtěch Havlíček, Antonio D Córcoles, Kristan Temme, Aram W Harrow, Abhinav Kandala, Jerry M Chow, and Jay M Gambetta, “Supervised learning with quantum-enhanced feature spaces,” *Nature* **567**, 209–212 (2019).

- [41] Kunal Sharma, Sumeet Khatri, Marco Cerezo, and Patrick J Coles, “Noise resilience of variational quantum compiling,” *New Journal of Physics* **22**, 043006 (2020).
- [42] Hsin-Yuan Huang, Richard Kueng, and John Preskill, “Predicting many properties of a quantum system from very few measurements,” *Nature Physics* **16**, 1050–1057 (2020).
- [43] S. Khatri, R. LaRose, A. Poremba, L. Cincio, A. T. Sornborger, and P. J. Coles, “Quantum-assisted quantum compiling,” *Quantum* **3**, 140 (2019).
- [44] Yuxuan Du, Min-Hsiu Hsieh, Tongliang Liu, and Dacheng Tao, “The expressive power of parameterized quantum circuits,” *Phys. Rev. Research* **2**, 033125 (2020).
- [45] Gian Giacomo Guerreschi and Mikhail Smelyanskiy, “Practical optimization for hybrid quantum-classical algorithms,” *arXiv preprint arXiv:1701.01450* (2017).
- [46] K. Mitarai, M. Negoro, M. Kitagawa, and K. Fujii, “Quantum circuit learning,” *Phys. Rev. A* **98**, 032309 (2018).
- [47] Maria Schuld, Ville Bergholm, Christian Gogolin, Josh Izaac, and Nathan Killoran, “Evaluating analytic gradients on quantum hardware,” *Physical Review A* **99**, 032331 (2019).
- [48] Christoph Dankert, Richard Cleve, Joseph Emerson, and Etera Livine, “Exact and approximate unitary 2-designs and their application to fidelity estimation,” *Physical Review A* **80**, 012304 (2009).
- [49] Fernando GSL Brandao, Aram W Harrow, and Michał Horodecki, “Local random quantum circuits are approximate polynomial-designs,” *Communications in Mathematical Physics* **346**, 397–434 (2016).
- [50] Aram Harrow and Saeed Mehraban, “Approximate unitary t -designs by short random quantum circuits using nearest-neighbor and long-range gates,” *arXiv preprint arXiv:1809.06957* (2018).
- [51] 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).
- [52] Carlos Ortiz Marrero, Mária Kieferová, and Nathan Wiebe, “Entanglement induced barren plateaus,” *PRX Quantum* **2**, 040316 (2021).
- [53] Neena Aloysius and M Geetha, “A review on deep convolutional neural networks,” in *2017 International Conference on Communication and Signal Processing (ICCSP)* (IEEE, 2017) pp. 0588–0592.
- [54] James Stokes, Josh Izaac, Nathan Killoran, and Giuseppe Carleo, “Quantum natural gradient,” *Quantum* **4**, 269 (2020).
- [55] Jonas M. Kübler, Andrew Arrasmith, Lukasz Cincio, and Patrick J. Coles, “An Adaptive Optimizer for Measurement-Frugal Variational Algorithms,” *Quantum* **4**, 263 (2020).
- [56] Bálint Koczor and Simon C Benjamin, “Quantum natural gradient generalised to non-unitary circuits,” *arXiv preprint arXiv:1912.08660* (2019).
- [57] Andrew Arrasmith, Lukasz Cincio, Rolando D Somma, and Patrick J Coles, “Operator sampling for shot-frugal optimization in variational algorithms,” *arXiv preprint arXiv:2004.06252* (2020).
- [58] Ryan LaRose, Arkin Tikku, Étude O’Neel-Judy, Lukasz Cincio, and Patrick J Coles, “Variational quantum state diagonalization,” *npj Quantum Information* **5**, 57 (2019).
- [59] Edward Grant, Leonard Wossnig, Mateusz Ostaszewski, and Marcello Benedetti, “An initialization strategy for addressing barren plateaus in parametrized quantum circuits,” *Quantum* **3**, 214 (2019).
- [60] Guillaume Verdon, Michael Broughton, Jarrod R McClean, Kevin J Sung, Ryan Babbush, Zhang Jiang, Hartmut Neven, and Masoud Mohseni, “Learning to learn with quantum neural networks via classical neural networks,” *arXiv preprint arXiv:1907.05415* (2019).
- [61] Tyler Volkoff and Patrick J Coles, “Large gradients via correlation in random parameterized quantum circuits,” *Quantum Science and Technology* **6**, 025008 (2021).
- [62] Adrien Boles and Markus Heyl, “Reinforcement learning for digital quantum simulation,” *Phys. Rev. Lett.* **127**, 110502 (2021).
- [63] Zbigniew Puchała and Jarosław Adam Miszczyk, “Symbolic integration with respect to the Haar measure on the unitary groups,” *Bulletin of the Polish Academy of Sciences Technical Sciences* **65**, 21–27 (2017).
- [64] Motohisa Fukuda, Robert König, and Ion Nechita, “RTNI—a symbolic integrator for Haar-random tensor networks,” *Journal of Physics A: Mathematical and Theoretical* **52**, 425303 (2019).
- [65] Arthur Pesah, M Cerezo, Samson Wang, Tyler Volkoff, Andrew T Sornborger, and Patrick J Coles, “Absence of barren plateaus in quantum convolutional neural networks,” *Phys. Rev. X* **11**, 041011 (2021).
- [66] Kaining Zhang, Min-Hsiu Hsieh, Liu Liu, and Dacheng Tao, “Toward trainability of quantum neural networks,” *arXiv preprint arXiv:2011.06258* (2020).
- [67] Chen Zhao and Xiao-Shan Gao, “Analyzing the barren plateau phenomenon in training quantum neural networks with the ZX-calculus,” *Quantum* **5**, 466 (2021).
- [68] Samson Wang, Enrico Fontana, Marco Cerezo, Kunal Sharma, Akira Sone, Lukasz Cincio, and Patrick J Coles, “Noise-induced barren plateaus in variational quantum algorithms,” *Nature communications* **12**, 1–11 (2021).
- [69] Amira Abbas, David Sutter, Christa Zoufal, Aurélien Lucchi, Alessio Figalli, and Stefan Woerner, “The power of quantum neural networks,” *Nature Computational Science* **1**, 403–409 (2021).
- [70] Taylor L Patti, Khadijeh Najafi, Xun Gao, and Susanne F Yelin, “Entanglement devised barren plateau mitigation,” *Physical Review Research* **3**, 033090 (2021).

Supplemental Material for *Trainability of Dissipative Perceptron-Based Quantum Neural Networks*

Here we provide proofs for the main results and theorems of the manuscript *Trainability of Dissipative Perceptron-Based Quantum Neural Networks*. In Section **A** we first present useful definitions and lemmas that will be employed to derive the main results. Then in Section **C** we show that $\langle \partial C \rangle = 0$ for the Dissipative Quantum Neural Networks (DQNN) considered in the main text. Finally, in Sections **E** and **D** we provide proofs for Theorem **1** and Theorem **2**, respectively. We note that we first provide a proof for Theorem **2**, since the proof of Theorem **1** can be built from the latter.

We also note that for our proofs we will assume a DQNN such that each perceptron only acts on one output qubit, and where there are no hidden layers. We generalize to DQNNs acting on m output qubits and with L hidden layers in Sections **F** and **G**, respectively.

A. Preliminaries

Properties of the Haar measure. Let $d\mu_H(V) \equiv d\mu(V)$ be the volume element of the Haar measure, with $V \in U(d)$, and where $U(d)$ denotes the unitary group of degree d . Then the following properties hold:

- The volume of the Haar measure is finite:

$$\int_{U(d)} d\mu(V) < \infty. \quad (\text{A1})$$

- The Haar measure is left- and right-invariant under the action of the unitary group of degree d . That is, for any integrable function $g(V)$ and for any $W \in U(d)$ we have

$$\int_{U(d)} d\mu(V)g(WV) = \int_{U(d)} d\mu(V)g(VW) = \int_{U(d)} d\mu(V)g(V). \quad (\text{A2})$$

- The Haar measure is uniquely defined up to a multiplicative constant factor. Let $d\omega(V)$ be an invariant measure, then, there exists a constant c such that

$$d\omega(V) = c \cdot d\mu(V). \quad (\text{A3})$$

Definition: t -design. Let be $\{V_y\}_{y \in Y}$, of size $|Y|$, be a set of unitaries V_y acting on a d -dimensional Hilbert space. In addition, let $P_t(W)$ be a polynomial of degree at most t in the matrix elements of V , and at most t in those of V^\dagger . Then $\{V_y \in U(d)\}_{y \in Y}$ is a unitary t -design if for every $P_t(W)$, the following holds [48]:

$$\frac{1}{|Y|} \cdot \sum_{y \in Y} P_t(V_y) = \int d\mu(V) P_t(V), \quad (\text{A4})$$

where it is implicit that the integral is over $U(d)$.

Symbolic integration. Here we recall formulas which allow for the symbolical integration with respect to the Haar measure on a unitary group [63]. For any $V \in U(d)$ the following expressions are valid for the first two moments:

$$\begin{aligned} \int d\mu(V) v_{ij} v_{pk}^* &= \frac{\delta_{ip} \delta_{jk}}{d}, \\ \int d\mu(V) v_{i_1 j_1} v_{i_2 j_2} v_{i'_1 j'_1}^* v_{i'_2 j'_2}^* &= \frac{\delta_{i_1 i'_1} \delta_{i_2 i'_2} \delta_{j_1 j'_1} \delta_{j_2 j'_2} + \delta_{i_1 i'_2} \delta_{i_2 i'_1} \delta_{j_1 j'_2} \delta_{j_2 j'_1}}{d^2 - 1} - \frac{\delta_{i_1 i'_1} \delta_{i_2 i'_2} \delta_{j_1 j'_2} \delta_{j_2 j'_1} + \delta_{i_1 i'_2} \delta_{i_2 i'_1} \delta_{j_1 j'_1} \delta_{j_2 j'_2}}{d(d^2 - 1)}, \end{aligned} \quad (\text{A5})$$

where v_{ij} are the matrix elements of V . Assuming $d = 2^n$, we use the notation $\mathbf{i} = (i_1, \dots, i_n)$ to denote a bitstring of length n such that $i_1, i_2, \dots, i_n \in \{0, 1\}$.

Useful identities. From Eqs. (A5), the following identities can be readily derived [19]

$$\int d\mu(V) \text{Tr} [V A V^\dagger B] = \frac{\text{Tr} [A] \text{Tr} [B]}{d} \quad (\text{A6})$$

$$\int d\mu(V) \text{Tr} [V A V^\dagger B V C V^\dagger D] = \frac{\text{Tr} [A] \text{Tr} [C] \text{Tr} [B D] + \text{Tr} [A C] \text{Tr} [B] \text{Tr} [D]}{d^2 - 1} - \frac{\text{Tr} [A C] \text{Tr} [B D] + \text{Tr} [A] \text{Tr} [B] \text{Tr} [C] \text{Tr} [D]}{d(d^2 - 1)} \quad (\text{A7})$$

$$\int d\mu(V) \text{Tr} [V A V^\dagger B] \text{Tr} [V C V^\dagger D] = \frac{\text{Tr} [A] \text{Tr} [B] \text{Tr} [C] \text{Tr} [D] + \text{Tr} [A C] \text{Tr} [B D]}{d^2 - 1} - \frac{\text{Tr} [A C] \text{Tr} [B] \text{Tr} [D] + \text{Tr} [A] \text{Tr} [C] \text{Tr} [B D]}{d(d^2 - 1)} \quad (\text{A8})$$

where A, B, C, D are linear operators on a d -dimensional Hilbert space. We point readers to [19] for detailed proofs of the aforementioned identities.

In addition, let us consider a bipartite Hilbert space $\mathcal{H} \equiv \mathcal{H}_1 \otimes \mathcal{H}_2$ of dimension $\dim(\mathcal{H}) = d_1 d_2$. Let $A, B : \mathcal{H} \rightarrow \mathcal{H}$ be linear operators. Then the following equality holds:

$$\int d\mu(V)(\mathbb{1}_1 \otimes V)A(\mathbb{1}_1 \otimes V^\dagger)B = \frac{\text{Tr}_2[A] \otimes \mathbb{1}_2}{d_2} B, \quad (\text{A9})$$

where Tr_2 indicates the partial trace over \mathcal{H}_2 . We note that Eq. (A9) follows by employing (A5).

Lemma 1. *Let $\mathcal{H} \equiv \mathcal{H}_1 \otimes \mathcal{H}_2$ denote a bipartite Hilbert space of dimension $\dim(\mathcal{H}) = d_1 d_2$, such that $d_1 = 2^n$ and $d_2 = 2^{n'}$. Let $A, B : \mathcal{H} \rightarrow \mathcal{H}$ be linear operators, and let $V \in U(d)$. Then we have that*

$$\text{Tr}[(\mathbb{1}_1 \otimes V)A(\mathbb{1}_1 \otimes V^\dagger)B] = \sum_{\mathbf{p}, \mathbf{q}} \text{Tr}[V A_{\mathbf{qp}} V^\dagger B_{\mathbf{pq}}]. \quad (\text{A10})$$

Where the summation runs over all bitstrings of length n' , and where we define

$$A_{\mathbf{qp}} = \text{Tr}_1[(|\mathbf{p}\rangle\langle\mathbf{q}| \otimes \mathbb{1}_2)A], \quad B_{\mathbf{pq}} = \text{Tr}_1[(|\mathbf{q}\rangle\langle\mathbf{p}| \otimes \mathbb{1}_2)B]. \quad (\text{A11})$$

Here $\mathbb{1}_1$ ($\mathbb{1}_2$) is the identity operator over \mathcal{H}_1 (\mathcal{H}_2).

Equation (A10) can be derived by expanding the operators in the computational basis as derived in [19].

Lemma 2. *Let $\mathcal{H} \equiv \mathcal{H}_1 \otimes \mathcal{H}_2 \otimes \mathcal{H}_3 \otimes \mathcal{H}_4$ denote a Hilbert space of dimension $\dim(\mathcal{H}) = d_1 d_2 d_3 d_4$. Consider the following linear operators $H : \mathcal{H}_1 \otimes \mathcal{H}_2 \rightarrow \mathcal{H}_1 \otimes \mathcal{H}_2$, $K : \mathcal{H}_1 \otimes \mathcal{H}_4 \rightarrow \mathcal{H}_1 \otimes \mathcal{H}_4$, $S, S' : \mathcal{H} \rightarrow \mathcal{H}$, $P, P' : \mathcal{H}_3 \otimes \mathcal{H}_4 \rightarrow \mathcal{H}_3 \otimes \mathcal{H}_4$, $V : \mathcal{H}_1 \otimes \mathcal{H}_3 \rightarrow \mathcal{H}_1 \otimes \mathcal{H}_3$, and $U : \mathcal{H}_1 \otimes \mathcal{H}_4 \rightarrow \mathcal{H}_1 \otimes \mathcal{H}_4$. Then, let us define the following operators acting on $\mathcal{H}_1 \otimes \mathcal{H}_2$:*

$$M = \text{Tr}_{34}[P[VUSU^\dagger V^\dagger], H], \quad B = \text{Tr}_{34}[P'VU[S', K]U^\dagger V^\dagger], \quad (\text{A12})$$

where Tr_{34} indicates the trace over subsystems $\mathcal{H}_3 \otimes \mathcal{H}_4$. The following equality holds:

$$\int d\mu(V)\text{Tr}_{12}[MB] = 0. \quad (\text{A13})$$

Lemma 3. *Let $\mathcal{H} \equiv \mathcal{H}_1 \otimes \mathcal{H}_2 \otimes \mathcal{H}_3 \otimes \mathcal{H}_4$ denote a Hilbert space of dimension $\dim(\mathcal{H}) = d_1 d_2 d_3 d_4$. Consider the following linear operators $V : \mathcal{H}_1 \otimes \mathcal{H}_2 \rightarrow \mathcal{H}_1 \otimes \mathcal{H}_2$, $U : \mathcal{H}_1 \otimes \mathcal{H}_3 \otimes \mathcal{H}_4 \rightarrow \mathcal{H}_1 \otimes \mathcal{H}_3 \otimes \mathcal{H}_4$, $P, P' : \mathcal{H}_2 \otimes \mathcal{H}_3 \otimes \mathcal{H}_4 \rightarrow \mathcal{H}_2 \otimes \mathcal{H}_3 \otimes \mathcal{H}_4$, and $Z, Z' : \mathcal{H}_4 \rightarrow \mathcal{H}_4$. Let P and P' be tensor product operators of the form $P = \Pi \otimes \tilde{P}$, and $P' = \Pi' \otimes \tilde{P}'$, where $\tilde{P}, \tilde{P}' : \mathcal{H}_3 \otimes \mathcal{H}_4 \rightarrow \mathcal{H}_3 \otimes \mathcal{H}_4$, and where $\Pi, \Pi' : \mathcal{H}_2 \rightarrow \mathcal{H}_2$ are rank one projector such that $\Pi^2 = \Pi$, $\text{Tr}[\Pi] = 1$, $(\Pi')^2 = \Pi'$, and $\text{Tr}[\Pi'] = 1$. Then, consider the following operators acting on subsystem \mathcal{H}_1 :*

$$\Omega = \text{Tr}_{2,3,4}[PVUZU^\dagger V^\dagger], \quad \tilde{\Omega} = \text{Tr}_{3,4}[\tilde{P}UZU^\dagger], \quad (\text{A14})$$

$$\Omega' = \text{Tr}_{2,3,4}[P'VUZ'U^\dagger V^\dagger], \quad \tilde{\Omega}' = \text{Tr}_{3,4}[\tilde{P}'UZ'U^\dagger], \quad (\text{A15})$$

where $\text{Tr}_{2,3,4}$ ($\text{Tr}_{3,4}$) indicates the trace over $\mathcal{H}_2 \otimes \mathcal{H}_3 \otimes \mathcal{H}_4$ ($\mathcal{H}_3 \otimes \mathcal{H}_4$). Then, the following equalities hold

$$\int d\mu(V)\text{Tr}_1[\Omega\Omega'] = \frac{d_1^2 d_2}{d_1^2 d_2^2 - 1} \left(\text{Tr}[\Pi\Pi'] - \frac{1}{d_1 d_2} \right) \text{Tr}_1[\tilde{\Omega}\tilde{\Omega}'] + \frac{d_1 d_2^2}{d_1^2 d_2^2 - 1} \left(1 - \frac{\text{Tr}[\Pi\Pi']}{d_2} \right) \text{Tr}_1[\tilde{\Omega}]\text{Tr}[\tilde{\Omega}']. \quad (\text{A16})$$

$$\int d\mu(V)\text{Tr}_1[\Omega]\text{Tr}[\Omega'] = \frac{d_1 d_2}{d_1^2 d_2^2 - 1} \left(\text{Tr}[\Pi\Pi'] - \frac{1}{d_2} \right) \text{Tr}_1[\tilde{\Omega}\tilde{\Omega}'] + \frac{d_1^2 d_2^2}{d_1^2 d_2^2 - 1} \left(1 - \frac{\text{Tr}[\Pi\Pi']}{d_1 d_2} \right) \text{Tr}_1[\tilde{\Omega}]\text{Tr}[\tilde{\Omega}']. \quad (\text{A17})$$

Lemmas 2, and 3 can be derived by explicitly integrating over V using (A5). In addition, one can also use the RTNI package for symbolic integrator over Haar-random tensor networks of Ref. [64].

Remark— We note that for simplicity we derive the proofs of our theorems for the case when the output states $|\phi_x^{\text{out}}\rangle$ in the training set are tensor product of computational basis states over n qubits. The proofs can then be trivially generalized for arbitrary tensor-product states over n qubits. Again, for simplicity, we provide rigorous proofs for the case when there are no hidden layers and later argue in Section G that our results hold for DQNNs with hidden layers. We also note that DQNNs considered in our work have a simple structure where unitaries act on $n + 1$ qubits. In Section F we generalize our results to the case when unitaries act on $n + m$ qubits.

B. Generalization of our results to quantum machine learning task

In the main text we have stated our main results in terms of state preparation, where the training set is of the form $\{|\phi_x^{\text{in}}\rangle, |\phi_x^{\text{out}}\rangle\}$. As we show here, this result can be used to generalize our result for other supervised learning quantum machine learning tasks.

Consider now the case when the training set is of the form $\{|\psi_x^{\text{in}}\rangle, y_x\}$. Here y_x are labels associated with each input quantum state $|\psi_x^{\text{in}}\rangle$, and the goal of the DQNN is to predict a label \tilde{y}_x that matches the true label. For simplicity we here consider the case of binary classification $y_x \in \{-1, 1\}$.

Following the scheme introduced in [40], the predicted labels \tilde{y}_x can be obtained by performing a binary measurement M_y on the states ρ_x^{out} . Without loss of generality, this measurement is taken to be on the z -basis, so that the measurement outcomes are bitstrings \mathbf{z} of length k , where k is the number of measured qubits. The measurement operator is then given by

$$M_y = \frac{\mathbb{1} + y\mathbf{h}}{2}, \quad \text{with} \quad \mathbf{h} = \sum_{\mathbf{z}} h(\mathbf{z})|\mathbf{z}\rangle\langle\mathbf{z}|. \quad (\text{B1})$$

Thus, the probability of obtaining label y from input state $|\psi_x^{\text{in}}\rangle$ is

$$p_y(|\psi_x^{\text{in}}\rangle) = \sum_{\mathbf{z}} h(\mathbf{z})\text{Tr}[O_{\mathbf{z}}\rho_x^{\text{out}}], \quad (\text{B2})$$

where

$$O_{\mathbf{z}} = \mathbb{1} \otimes |\mathbf{z}\rangle\langle\mathbf{z}|. \quad (\text{B3})$$

From these probabilities, the assigned labels are $\tilde{y}_x = p_y(|\psi_x^{\text{in}}\rangle) \geq p_{-y}(|\psi_x^{\text{in}}\rangle)$.

Note that Eq. (B3) is precisely of the form considered in the main text. Thus, the formalism in the main text can be directly extended for other supervised learning tasks.

C. Proof of $\langle\partial C\rangle = 0$

In this section, we prove that the average value of the partial derivative of the cost function $\langle\partial C\rangle$ is not biased towards any particular value. In particular, we show that $\langle\partial C\rangle = 0$. We prove this result for both cases: 1) random parameterized quantum circuits and 2) parameter matrix multiplication method.

Random parametrized quantum circuits. Let us first consider the case when the perceptrons are random parametrized quantum circuits of the form

$$V_j^l(\boldsymbol{\theta}_j^l) = \prod_{k=1}^{\eta_j^l} R_k(\theta^k)W_k, \quad (\text{C1})$$

with W_k an unparametrized unitary, $R_k(\theta^k) = e^{-(i/2)\theta^k\Gamma_k}$, and where Γ_k is a Hermitian operator with $\text{Tr}[\Gamma_k^2] \leq 2^{n+1}$.

We recall from the main text that the training set consisting of input and output pure quantum states: $\{|\phi_x^{\text{in}}\rangle, |\phi_x^{\text{out}}\rangle\}_{x=1}^N$. As mentioned in the Remark in the previous section, we now assume that $|\phi_x^{\text{out}}\rangle$ are computational basis states, i.e., $|\phi_x^{\text{out}}\rangle \equiv |\mathbf{z}^x\rangle = |z_1^x z_2^x \dots z_n^x\rangle$. The DQNN cost function is defined as

$$C = \frac{1}{N} \sum_{x=1}^N C_x, \quad \text{with} \quad C_x = \text{Tr}[O_x \rho_x^{\text{out}}], \quad (\text{C2})$$

where O_x is given by the global observable

$$O_x^G \equiv \mathbb{1} - |\phi_x^{\text{out}}\rangle\langle\phi_x^{\text{out}}| \quad (\text{C3})$$

or by the local operator

$$O_x^L = \mathbb{1} - \frac{1}{n} \sum_{i=1}^n |z_i^x\rangle\langle z_i^x| \otimes \mathbb{1}_{\bar{i}}, \quad (\text{C4})$$

where $\mathbb{1}_{\bar{j}}$ indicates identity on all qubits in the output layer except for qubit i . Then the partial derivative of C with respect to parameter θ^ν in a perceptron V_j^l is given by

$$\partial_\nu C = \frac{i}{2N} \sum_{x=1}^N \text{Tr} \left[A_j^l \tilde{\rho}_x^{\text{in}}(A_j^l)^\dagger [\mathbb{1}_{\bar{j}} \otimes \Gamma_k, (B_j^l)^\dagger \tilde{O}_x B_j^l] \right], \quad (\text{C5})$$

$$= -\frac{i}{2N} \sum_{x=1}^N \text{Tr} \left[(B_j^l)^\dagger \tilde{O}_x B_j^l [\mathbb{1}_{\bar{j}} \otimes \Gamma_k, A_j^l \tilde{\rho}_x^{\text{in}}(A_j^l)^\dagger] \right], \quad (\text{C6})$$

where the trace is taken all qubits in the QNN, and where $\mathbb{1}_{\bar{j}}$ indicates the identity on all qubits on which V_j^l does not act on. Moreover, we have defined

$$A_j^l = \mathbb{1}_{\bar{j}} \otimes \prod_{k=\nu}^{\eta_j^l} R_k(\theta^k) W_k, \quad (\text{C7})$$

$$B_j^l = \mathbb{1}_{\bar{j}} \otimes \prod_{k=1}^{\nu-1} R_k(\theta^k) W_k, \quad (\text{C8})$$

$$\tilde{\rho}_x^{\text{in}} = V_{j-1}^l \dots V_1^l (\rho_x^{\text{in}} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{\text{hid,out}}) (V_1^l)^\dagger \dots (V_{j-1}^l)^\dagger, \quad (\text{C9})$$

$$\tilde{O}_x = (V_{j+1}^l)^\dagger \dots (V_{n_{\text{out}}}^{\text{out}})^\dagger (\mathbb{1}_{\text{in,hid}} \otimes O_x) V_{n_{\text{out}}}^{\text{out}} \dots V_{j+1}^l. \quad (\text{C10})$$

Let us first consider the case when B_j^l is a 1-design. Assuming that all the perceptrons are independent, and that A_j^l and B_j^l are also independent, we can express

$$\langle \dots \rangle = \langle \dots \rangle_{V_1^l, \dots, A_j^l, B_j^l, \dots, V_{n_{\text{out}}}^{\text{out}}} = \langle \langle \dots \rangle_{B_j^l} \rangle_{V_1^l, \dots, A_j^l, \dots, V_{n_{\text{out}}}^{\text{out}}}. \quad (\text{C11})$$

Hence, we can first compute the average of $\partial_\nu C$ over B_j^l as

$$\langle \partial_\nu C \rangle_{B_j^l} = \frac{i}{2N} \sum_{x=1}^N \text{Tr} \left[A_j^l \tilde{\rho}_x^{\text{in}}(A_j^l)^\dagger [\mathbb{1}_{\bar{j}} \otimes \Gamma_k, \int d\mu(B_j^l) (B_j^l)^\dagger \tilde{O}_x B_j^l] \right] \quad (\text{C12})$$

$$= \frac{i}{2N} \sum_{x=1}^N \text{Tr} \left[A_j^l \tilde{\rho}_x^{\text{in}}(A_j^l)^\dagger [\mathbb{1}_{\bar{j}} \otimes \Gamma_k, \frac{1}{2^{m+1}} \text{Tr}_{jl}[\tilde{O}_x] \otimes \mathbb{1}_j^l] \right] \\ = 0. \quad (\text{C13})$$

Here, Tr_{jl} indicates the trace over the $m+1$ qubits on which B_j^l acts, and $\mathbb{1}_j^l$ ($\mathbb{1}_{\bar{j}}^l$) is the identity operator over the qubits on which B_j^l acts (does not act). For the second equality we used (A9), and in the third equality we used the fact that a commutator inside the trace is always zero. With a similar argument it is straightforward to show from (C6) that $\langle \partial_\nu C \rangle_{A_j^l} = 0$ if A_j^l is a 1-design. Finally, from Eq. (C11) we know that $\langle \partial_\nu C \rangle_{A_j^l} = 0$ ($\langle \partial_\nu C \rangle_{B_j^l} = 0$) leads to $\langle \partial_\nu C \rangle = 0$.

Parameter matrix multiplication. In this case the perceptrons are updated with the parameter multiplication matrix method of Ref. [37]. We recall that the perceptrons $V_j^l(0)$ are randomly initialized, and then at each step s they are updated via

$$V_j^l(s + \varepsilon) = e^{i\varepsilon H_j^l(s)} V_j^l(s), \quad (\text{C14})$$

where $H_j^l(s)$ is a Hermitian operator such that $\text{Tr}[(H_j^l)^2] \leq 2^{n+1}$.

As shown in [37], the derivative of $C(s)$ with respect to s is given by,

$$\partial_s C(s) = \lim_{\varepsilon \rightarrow 0} \frac{C(s + \varepsilon) - C(s)}{\varepsilon} \\ = \frac{i}{N} \sum_{x=1}^N \left[\sum_{l=1}^{\text{out}} \sum_{j=1}^{n_l} \text{Tr} \left[\mathbb{1}_{\bar{j}} \otimes H_j^l(s) [V_j^l(s) \tilde{\rho}_x^{\text{in}}(s) (V_j^l(s))^\dagger, \tilde{O}_x(s)] \right] \right] \quad (\text{C15})$$

$$= \frac{i}{N} \sum_{x=1}^N \left[\sum_{l=1}^{\text{out}} \sum_{j=1}^{n_l} \text{Tr} \left[\tilde{O}_x(s) [\mathbb{1}_{\bar{j}} \otimes H_j^l(s), V_j^l(s) \tilde{\rho}_x^{\text{in}}(s) (V_j^l(s))^\dagger] \right] \right], \quad (\text{C16})$$

where now

$$\tilde{\rho}_x^{\text{in}}(s) = V_{j-1}^l(s) \dots V_1^l(s) (\rho_x^{\text{in}} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{\text{hid,out}}) (V_1^l(s))^\dagger \dots (V_{j-1}^l(s))^\dagger, \quad (\text{C17})$$

$$\tilde{O}_x(s) = (V_{j+1}^l(s))^\dagger \dots (V_{n_{\text{out}}}^{\text{out}}(s))^\dagger (\mathbb{1}_{\text{in,hid}} \otimes O_x) V_{n_{\text{out}}}^{\text{out}}(s) \dots V_{j+1}^l(s). \quad (\text{C18})$$

Let us now consider time-step $s = 0$. If we assume that all the perceptrons are independently initialized we have

$$\langle \dots \rangle = \langle \dots \rangle_{V_1^l(0), \dots, V_{n_{\text{out}}}^{\text{out}}(0)} = \langle \langle \dots \rangle_{V_j^l(0)} \rangle_{V_1^l(0), \dots, V_{j-1}^l(0), V_{j+1}^l(0), \dots, V_{n_{\text{out}}}^{\text{out}}(0)}. \quad (\text{C19})$$

Therefore, if $V_j^l(0)$ is a 1-design, from (C12)–(C13) it follows that $\langle \partial_s C \rangle_{V_j^l(0)} = \langle \partial_s C \rangle = 0$.

D. Proof of Theorem 2

In this section, we provide a proof of Theorem 2. In what follows we will assume a DQNN with no hidden layers, and such that each perceptron only acts on one output qubit. We generalize to DQNNs acting on m output qubits and with L hidden layers in sections F and G, respectively.

We first recall Theorem 2 for convenience:

Theorem 2. *Consider a DQNN with deep global perceptrons, which are updated via the parameter matrix multiplication of (C14), and such that $V_j^l(0)$ form independent 2-designs over $n + 1$ qubits. Then the variance of the partial derivative of the cost function with respect to the time-step parameter s is upper bounded as*

$$\text{Var}[\partial_s C] \leq f(n), \quad \text{with } f(n) \in \mathcal{O}(1/2^n), \quad (\text{D1})$$

when O_x is the global operator as in (C3) or the local operator as in (C4).

Proof. We divide our proof in several subsections consisting of different cases. We also analyze the global and local cost functions separately. For simplicity we consider a DQNN with no hidden layers, and where both input and output layers consist of n qubits. Moreover, we denote the randomly initialized perceptrons at time step $s = 0$ as V_j , where we also henceforth drop the superscript index that denotes the layer.

Let

$$\sigma_x^{\text{in}} = |\phi_x^{\text{in}}\rangle\langle\phi_x^{\text{in}}| \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{\text{out}}, \quad \sigma_x^{\text{out}} = \mathbb{1}_{\text{in}} \otimes O_x, \quad (\text{D2})$$

$$G_j^x = [V_j \dots V_1 \sigma_x^{\text{in}} V_1^\dagger \dots V_j^\dagger, V_{j+1}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+1}]. \quad (\text{D3})$$

Then (C15) can be rewritten as

$$\partial_s C = \frac{i}{N} \sum_{x=1}^N \left[\sum_{j=1}^n \text{Tr}[G_j^x H_j] \right]. \quad (\text{D4})$$

Note that from the cyclicity of the trace, each term $\text{Tr}[G_j^x H_j]$ can also be expressed as

$$\text{Tr}[G_j^x H_j] = \text{Tr} \left[V_j \dots V_1 \sigma_x^{\text{in}} V_1^\dagger \dots V_j^\dagger \left[V_{j+1}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+1}, H_j \right] \right]. \quad (\text{D5})$$

The proof of Theorem 2 is constructed as follows. We first note that $\langle \partial C \rangle = 0$ and thus, $\text{Var}[\partial_s C]$ only depends on the second moment of partial derivatives. Moreover, $\langle (\partial_s C)^2 \rangle$ depends on terms of the following form $\text{Tr}[G_j^x H_j] \text{Tr}[G_{j'}^{x'} H_{j'}]$. We first consider a single term in the summation over j and x in (D4), and show that $\langle (\text{Tr}[G_j^x H_j])^2 \rangle \leq f(n)$ with $f(n)$ as in (D1). Then we prove that the cross terms in x and j also satisfy $\langle \text{Tr}[G_j^x H_j] \text{Tr}[G_{j'}^{x'} H_{j'}] \rangle \leq f(n)$.

1. Global Cost

a. Fixed j and fixed x

Let us first consider the case when the cost function is defined in terms of the global operator in (C3). As previously discussed, we analyze the scaling of the variance of a single term in (D4) with fixed x and j . By invoking Lemma 1,

$\text{Tr}[G_j^x H_j]$ can be expressed as

$$\text{Tr}[G_j^x H_j] = \sum_{\mathbf{p}, \mathbf{q}} \text{Tr}[V_j A_{\mathbf{qp}}^{(x,j)} V_j^\dagger B_{\mathbf{pq}}^{(x,j)}], \quad (\text{D6})$$

where

$$A_{\mathbf{qp}}^{(x,j)} = \text{Tr}_{\bar{j}}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{p}\rangle\langle\mathbf{q}|) A^{(x,j)}], \quad A^{(x,j)} = (V_{j-1} \dots V_1) \sigma_x^{\text{in}} (V_1^\dagger \dots V_{j-1}^\dagger), \quad (\text{D7})$$

$$B_{\mathbf{pq}}^{(x,j)} = \text{Tr}_{\bar{j}}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{q}\rangle\langle\mathbf{p}|) B^{(x,j)}], \quad B^{(x,j)} = [V_{j+1}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+1}, H_j]. \quad (\text{D8})$$

Here $\text{Tr}_{\bar{j}}$ indicates the trace over all qubits in the output layer except for qubit j . Moreover, we remark that the summation in (D6) runs over all bitstrings \mathbf{p} , and \mathbf{q} of length $n-1$, and we recall that the operator $|\mathbf{q}\rangle\langle\mathbf{p}|$ acts on all qubits in the output layer except on qubit j . From the definition of σ_x^{in} in (D2) and from (D7), it follows that $A_{\mathbf{qp}}^{(x,j)}$ is nonzero when

$$p_k = q_k = 0, \forall k \in \{j+1, \dots, n\}. \quad (\text{D9})$$

Similarly, from σ_x^{out} and (D8), it follows that $B_{\mathbf{pq}}^{(x,j)}$ is nonzero when

$$p_k = q_k = z_k^x, \forall k \in \{1, \dots, j-1\}. \quad (\text{D10})$$

Equations (D9) and (D10) follow from the fact that each V_j acts on $n+1$ number of qubits, and these equations imply that there is a single nonzero term in the summation (D6). This term can be identified by defining the following bitstring of length $n-1$:

$$\mathbf{r}^{(x,j)} \equiv (z_1^x, z_2^x, \dots, z_{j-1}^x, 0, \dots, 0). \quad (\text{D11})$$

Then, by invoking Lemma (A8) we get

$$\langle (\text{Tr}[G_j^x H_j])^2 \rangle_{V_j} = \int d\mu(V_j) \text{Tr}[V_j A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} V_j^\dagger B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}] \text{Tr}[V_j A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} V_j^\dagger B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}] \quad (\text{D12})$$

$$= \frac{1}{2^{2(n+1)} - 1} \left(\text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] - \frac{1}{2^{n+1}} \text{Tr}[A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}]^2 \right) \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \quad (\text{D13})$$

$$\leq \frac{1}{2^{2(n+1)} - 1} \text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2]. \quad (\text{D14})$$

where in the inequality we used the fact that $\text{Tr}[A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}] > 0$ as $A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}$ is a positive semidefinite operator, and therefore, we can drop the term with the negative sign.

Since $A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}$ and $B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}$ are functions of different preceptrons V_i , we can compute an upper bound on the expectation value of $(\text{Tr}[G_j^x H_j])^2$ as follows:

$$\langle \text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \rangle = \langle \text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \rangle_{V_1, \dots, V_{j-1}} \langle \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \rangle_{V_{j+1}, \dots, V_n}. \quad (\text{D15})$$

We note that since H_j only acts on all input qubits and on output qubit j , then $B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}$ can be expressed in the following compact form

$$B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} = s_{\mathbf{r}^{(x,j)}}^{(x,j)} [\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}, H_j], \quad (\text{D16})$$

where

$$\omega_{\mathbf{r}^{(x,j)}}^{(x,j)} = \frac{1}{s_{\mathbf{r}^{(x,j)}}^{(x,j)}} \text{Tr}_{\bar{j}}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|) V_{j+1}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+1}], \quad (\text{D17})$$

$$s_{\mathbf{r}^{(x,j)}}^{(x,j)} = \text{Tr}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|) V_{j+1}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+1}]. \quad (\text{D18})$$

It is straightforward to note that $\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}$ is a quantum state on all qubits in the input layer plus qubit j in the output layer. Let us now consider the following chain of inequalities:

$$\mathrm{Tr}[(\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}, H_j)^2] = 2 \left(\mathrm{Tr}[\omega_{\mathbf{r}^{(x,j)}}^{(x,j)} H_j \omega_{\mathbf{r}^{(x,j)}}^{(x,j)} H_j] - \mathrm{Tr}[(\omega_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 (H_j)^2] \right) \quad (\text{D19})$$

$$\begin{aligned} &\leq 2 \mathrm{Tr}[\omega_{\mathbf{r}^{(x,j)}}^{(x,j)} H_j \omega_{\mathbf{r}^{(x,j)}}^{(x,j)} H_j] \\ &\leq 2 \mathrm{Tr}[H_j \omega_{\mathbf{r}^{(x,j)}}^{(x,j)} H_j] \\ &= 2 \mathrm{Tr}[\omega_{\mathbf{r}^{(x,j)}}^{(x,j)} (H_j)^2] \\ &\leq 2 \mathrm{Tr}[(H_j)^2] \\ &\leq 2^{n+2}. \end{aligned} \quad (\text{D20})$$

The first inequality follows from the fact that $\omega_{\mathbf{r}^{(x,j)}}^{(x,j)} (H_j)^2 \omega_{\mathbf{r}^{(x,j)}}^{(x,j)}$ is a positive semidefinite operator, so that $\mathrm{Tr}[(\omega_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 (H_j)^2] \geq 0$. The second inequality follows by noting that $\omega_{\mathbf{r}^{(x,j)}}^{(x,j)} \leq \mathbb{1}$ and $H_j \omega_{\mathbf{r}^{(x,j)}}^{(x,j)} H_j \geq 0$. The third inequality follows from the fact that $\omega_{\mathbf{r}^{(x,j)}}^{(x,j)} \leq \mathbb{1}$, and that $(H_j)^2 \geq 0$. The last inequality holds from the assumption that $\mathrm{Tr}[(H_j)^2] \leq 2^{n+1}$. Finally, by combining Eqs. (F6) and (D20), we get that

$$\begin{aligned} \mathrm{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(j)})^2] &= (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \mathrm{Tr}[(\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}, H_j)^2] \\ &\leq 2^{n+2} (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2. \end{aligned} \quad (\text{D21})$$

Let us now evaluate the term $(s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2$. By invoking Lemma 1, we get

$$s_{\mathbf{r}^{(x,j)}}^{(x,j)} = \mathrm{Tr}[V_{j+1}(\mathbb{1}_{\mathrm{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|) V_{j+1}^\dagger V_{j+2}^\dagger \cdots V_n^\dagger \sigma_x^{\mathrm{out}} V_n \cdots V_{j+2}] \quad (\text{D22})$$

$$= \sum_{\mathbf{p}' \mathbf{q}'} \mathrm{Tr}[V_{j+1} C_{\mathbf{q}' \mathbf{p}'}^{(x,j+1)} V_{j+1}^\dagger D_{\mathbf{p}' \mathbf{q}'}^{(x,j+1)}]. \quad (\text{D23})$$

Here the summation is over all bitstrings \mathbf{p}' and \mathbf{q}' of length $n+1$, and

$$C_{\mathbf{q}' \mathbf{p}'}^{(x,j+1)} = \mathrm{Tr}_{\overline{j+1}}[(\mathbb{1}_{\mathrm{in},j+1} \otimes |\mathbf{p}'\rangle\langle\mathbf{q}'|)(\mathbb{1}_{\mathrm{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|)], \quad (\text{D24})$$

$$D_{\mathbf{p}' \mathbf{q}'}^{(x,j+1)} = \mathrm{Tr}_{\overline{j+1}}[(\mathbb{1}_{\mathrm{in},j+1} \otimes |\mathbf{q}'\rangle\langle\mathbf{p}'|) V_{j+2}^\dagger \cdots V_n^\dagger \sigma_x^{\mathrm{out}} V_n \cdots V_{j+2}], \quad (\text{D25})$$

where $\mathrm{Tr}_{\overline{j+1}}$ indicates the trace over all qubits in the output layer except qubit $j+1$.

Then from arguments similar to those used to deriving (D9) and (D10), we find

$$q'_j = p'_j = z_j^x, \quad (\text{D26})$$

$$q'_k = p'_k = r_k^{(x,j)}, \forall k \in \{1, 2, \dots, j-1, j+2, \dots, n\}. \quad (\text{D27})$$

We now point at a recursive relation. Let

$$\mathbf{r}^{(x,j+1)} \equiv (r_1^{(x,j)}, \dots, r_{j-1}^{(x,j)}, z_j^x, r_{j+2}^{(x,j)}, \dots, r_n^{(x,j)}). \quad (\text{D28})$$

Then $s_{\mathbf{r}^{(x,j)}}^{(x,j)}$ further simplifies to

$$s_{\mathbf{r}^{(x,j)}}^{(x,j)} = \mathrm{Tr} \left[V_{j+1} (\mathbb{1}_{\mathrm{in}} \otimes |r_{j+1}^{(x,j)}\rangle\langle r_{j+1}^{(x,j)}|) V_{j+1}^\dagger \mathrm{Tr}_{\overline{j+1}} \left(\mathbb{1}_{\mathrm{in},j+1} \otimes (|\mathbf{r}^{(x,j+1)}\rangle\langle\mathbf{r}^{(x,j+1)}|) V_{j+2}^\dagger \cdots V_n^\dagger \sigma_x^{\mathrm{out}} V_n \cdots V_{j+2} \right) \right] \quad (\text{D29})$$

$$= s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} \mathrm{Tr} \left[V_{j+1} (\mathbb{1}_{\mathrm{in}} \otimes |r_{j+1}^{(x,j)}\rangle\langle r_{j+1}^{(x,j)}|) V_{j+1}^\dagger \omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} \right]. \quad (\text{D30})$$

Here we defined

$$s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} = \mathrm{Tr} \left[(\mathbb{1}_{\mathrm{in},j+1} \otimes |\mathbf{r}^{(x,j+1)}\rangle\langle\mathbf{r}^{(x,j+1)}|) V_{j+2}^\dagger \cdots V_n^\dagger \sigma_x^{\mathrm{out}} V_n \cdots V_{j+2} \right], \quad (\text{D31})$$

$$\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} = \frac{1}{s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)}} \mathrm{Tr}_{\overline{j+1}} \left[(\mathbb{1}_{\mathrm{in},j+1} \otimes |\mathbf{r}^{(x,j+1)}\rangle\langle\mathbf{r}^{(x,j+1)}|) V_{j+2}^\dagger \cdots V_n^\dagger \sigma_x^{\mathrm{out}} V_n \cdots V_{j+2} \right]. \quad (\text{D32})$$

An upper bound on the average of $(s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2$ over V_{j+1} can be obtained as follows, provided that V_{j+1} forms a 2-design:

$$\begin{aligned} & \langle (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \rangle_{V_{j+1}} \\ &= \int d\mu(V_{j+1}) (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \end{aligned} \quad (\text{D33})$$

$$\begin{aligned} &= (s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)})^2 \int d\mu(V_{j+1}) \text{Tr} \left[V_{j+1} (\mathbb{1}_{\text{in}} \otimes |r_{j+1}^{(x,j)}\rangle\langle r_{j+1}^{(x,j)}|) V_{j+1}^\dagger \omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} \right] \text{Tr} \left[V_{j+1} (\mathbb{1}_{\text{in}} \otimes |r_{j+1}^{(x,j)}\rangle\langle r_{j+1}^{(x,j)}|) V_{j+1}^\dagger \omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} \right] \\ &= \frac{(s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)})^2}{2^{2(n+1)} - 1} \left(2^{2n} + 2^n \text{Tr}[(\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)})^2] - \frac{1}{2^{n+1}} (2^n + 2^{2n} \text{Tr}[(\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)})^2]) \right) \end{aligned} \quad (\text{D34})$$

$$\leq \frac{2^n (2^n + 1/2) (s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)})^2}{2^{2(n+1)} - 1}, \quad (\text{D35})$$

where we employed (A8) and used the fact that $\text{Tr}[(\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)})^2] \leq 1$ as $\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)}$ is a quantum state.

Here we remark that from $s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)}$ we can always define an operator $s_{\mathbf{r}^{(x,j+2)}}^{(x,j+2)}$ according to Eqs. (D22)–(D30). Moreover, by using the assumption that all randomly initialized perceptrons form 2-designs, we can recursively average over V_{j+2}, \dots, V_n . Therefore, from (D21), we get

$$\langle \text{Tr}[(B_{\mathbf{r}^{(x,j)}, \mathbf{r}^{(x,j)}}^{(x,j)})^2] \rangle_{V_{j+1}, \dots, V_n} \leq (s_{\mathbf{r}^{(x,n)}}^{(x,n)})^2 2^{n+1} \left(\frac{2^n (2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-j} \quad (\text{D36})$$

$$= 2^{3n+1} \left(\frac{2^n (2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-j}, \quad (\text{D37})$$

where we used that fact that $s_{\mathbf{r}^{(x,n)}}^{(x,n)} = \text{Tr}[\sigma_x^{\text{out}}] = 2^n$.

We now compute the average of $\text{Tr}[(A_{\mathbf{r}^{(x,j)}, \mathbf{r}^{(x,j)}}^{(x,j)})^2]$ over V_1, \dots, V_{j-1} . By following a similar procedure to the one previously employed, and from (D7), we have

$$A_{\mathbf{r}^{(x,j)}, \mathbf{r}^{(x,j)}}^{(x,j)} = \text{Tr}_{\bar{j}} \left[(\mathbb{1}_{\text{in}, j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle \mathbf{r}^{(x,j)}|) (V_{j-1} \dots V_1) \sigma_x^{\text{in}} (V_1^\dagger \dots V_{j-1}^\dagger) \right] \quad (\text{D38})$$

$$= q_{\mathbf{r}^{(x,j)}}^{(x,j)} \varphi_{\mathbf{r}^{(x,j)}}^{(x,j)}, \quad (\text{D39})$$

where

$$q_{\mathbf{r}^{(x,j)}}^{(x,j)} = \text{Tr} \left[(\mathbb{1}_{\text{in}, j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle \mathbf{r}^{(x,j)}|) (V_{j-1} \dots V_1) \sigma_x^{\text{in}} (V_1^\dagger \dots V_{j-1}^\dagger) \right], \quad (\text{D40})$$

$$\varphi_{\mathbf{r}^{(x,j)}}^{(x,j)} = \frac{1}{q_{\mathbf{r}^{(x,j)}}^{(x,j)}} \text{Tr}_{\bar{j}} \left[(\mathbb{1}_{\text{in}, j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle \mathbf{r}^{(x,j)}|) (V_{j-1} \dots V_1) \sigma_x^{\text{in}} (V_1^\dagger \dots V_{j-1}^\dagger) \right]. \quad (\text{D41})$$

Moreover, if V_{j-1} forms a 2-design, we can compute the expectation value of $\text{Tr}[(A_{\mathbf{r}^{(x,j)}, \mathbf{r}^{(x,j)}}^{(x,j)})^2]$ with respect V_{j-1} as

$$\int d\mu(V_{j-1}) \text{Tr}[(A_{\mathbf{r}^{(x,j)}, \mathbf{r}^{(x,j)}}^{(x,j)})^2] = \int d\mu(V_{j-1}) (q_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \text{Tr}[(\varphi_{\mathbf{r}^{(x,j)}}^{(x,j)})^2] \quad (\text{D42})$$

$$\leq \int d\mu(V_{j-1}) (q_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \quad (\text{D43})$$

$$\begin{aligned} &= \int d\mu(V_{j-1}) \left(\text{Tr} \left[V_{j-1}^\dagger (\mathbb{1}_{\text{in}, j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle \mathbf{r}^{(x,j)}|) V_{j-1} (V_{j-2} \dots V_1) \sigma_x^{\text{in}} (V_1^\dagger \dots V_{j-2}^\dagger) \right] \right)^2 \\ &\leq \frac{2^n (2^n + 1/2) (q_{\hat{\mathbf{r}}^{(x,j-1)}}^{(x,j-1)})^2}{2^{2(n+1)} - 1}, \end{aligned} \quad (\text{D44})$$

where we used arguments similar to those used in deriving (D22)–(D35). Here,

$$q_{\hat{\mathbf{r}}^{(x,j-1)}}^{(x,j-1)} = \text{Tr} \left[(\mathbb{1}_{\text{in}, j-1} \otimes |\hat{\mathbf{r}}^{(x,j-1)}\rangle\langle \hat{\mathbf{r}}^{(x,j-1)}|) (V_{j-2} \dots V_1) \sigma_x^{\text{in}} (V_1^\dagger \dots V_{j-2}^\dagger) \right], \quad (\text{D45})$$

$$\hat{\mathbf{r}}^{(x,j-1)} = (r_1^{(x,j)}, r_2^{(x,j)}, \dots, r_{j-2}^{(x,j)}, 0, r_{j+1}^{(x,j)}, \dots, r_n^{(x,j)}), \quad (\text{D46})$$

where $\hat{\mathbf{r}}^{(x,j-1)}$ denotes a bitstring of length $n-1$, and where $j-1$ in the superscript implies that $|\hat{\mathbf{r}}^{(x,j-1)}\rangle$ is a state on all qubits in the output layer, except the $(j-1)$ -th qubit. Then, since all randomly initialized perceptrons form 2-designs, we can recursively compute the average over V_{j-2}, \dots, V_1 . We get

$$\langle (\text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(j)})^2])^2 \rangle_{V_1, \dots, V_{j-1}} \leq (q_{\hat{\mathbf{r}}^{(x,1)}}^{(x,1)})^2 \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{j-1} \quad (\text{D47})$$

$$= \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{j-1}, \quad (\text{D48})$$

where we used the fact that $q_{\hat{\mathbf{r}}^{(x,1)}}^{(x,1)} = \text{Tr}[\sigma_x^{\text{in}}] = 1$. Then from (D14), (D37), and (D48), it follows that

$$\langle (\text{Tr}[G_j^x H_j])^2 \rangle \leq \frac{1}{2^{2(n+1)} - 1} \langle \text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \rangle \quad (\text{D49})$$

$$\leq \frac{2^{3n+2}}{2^{2(n+1)} - 1} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-1} \quad (\text{D50})$$

$$\leq f(n) \in \mathcal{O}(1/2^n). \quad (\text{D51})$$

b. Fixed j and different x

We now establish an upper bound on the cross terms with equal j but different x , i.e., on terms of the form: $\langle \text{Tr}[G_j^x H_j] \text{Tr}[G_j^{x'} H_j] \rangle$. Following (D6)–(D14), we find that

$$\langle \text{Tr}[G_j^x H_j] \text{Tr}[G_j^{x'} H_j] \rangle_{V_j} \leq \frac{1}{2^{2(n+1)} - 1} \left| \left(\Delta(A^{(x,x',j)})_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}} \right) \right| \left| \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} B_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)})] \right|. \quad (\text{D52})$$

Then, from (F6)–(D21), and by invoking the Cauchy-Schwarz inequality, we find that

$$\left| \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} B_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)})] \right| \leq \sqrt{\text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2]} \sqrt{\text{Tr}[(B_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)})^2]} \quad (\text{D53})$$

$$\leq 2^{n+2} s_{\mathbf{r}^{(x,j)}}^{(x,j)} s_{\mathbf{r}^{(x',j)}}^{(x',j)}. \quad (\text{D54})$$

From (A8), we compute the average of $s_{\mathbf{r}^{(x,j)}}^{(x,j)} s_{\mathbf{r}^{(x',j)}}^{(x',j)}$ with respect to V_{j+1} as follows:

$$\begin{aligned} & \langle s_{\mathbf{r}^{(x,j)}}^{(x,j)} s_{\mathbf{r}^{(x',j)}}^{(x',j)} \rangle_{V_{j+1}} \\ &= \int d\mu(V_{j+1}) s_{\mathbf{r}^{(x,j)}}^{(x,j)} s_{\mathbf{r}^{(x',j)}}^{(x',j)} \end{aligned} \quad (\text{D55})$$

$$\begin{aligned} &= s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} s_{\mathbf{r}^{(x',j+1)}}^{(x',j+1)} \int d\mu(V_{j+1}) \text{Tr} \left[V_{j+1} (\mathbb{1}_{\text{in}} \otimes |r_{j+1}^{(x,j)}\rangle \langle r_{j+1}^{(x,j)}|) V_{j+1}^\dagger \omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} \right] \text{Tr} \left[V_{j+1} (\mathbb{1}_{\text{in}} \otimes |r_{j+1}^{(x',j)}\rangle \langle r_{j+1}^{(x',j)}|) V_{j+1}^\dagger \omega_{\mathbf{r}^{(x',j+1)}}^{(x',j+1)} \right] \\ &= \frac{s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} s_{\mathbf{r}^{(x',j+1)}}^{(x',j+1)}}{2^{2(n+1)} - 1} \left(2^{2n} + 2^n \text{Tr}[\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} \omega_{\mathbf{r}^{(x',j+1)}}^{(x',j+1)}] - \frac{1}{2^{n+1}} (2^n + 2^{2n} \text{Tr}[\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} \omega_{\mathbf{r}^{(x',j+1)}}^{(x',j+1)}]) \right) \end{aligned} \quad (\text{D56})$$

$$\leq \frac{2^n(2^n + 1/2) s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} s_{\mathbf{r}^{(x',j+1)}}^{(x',j+1)}}{2^{2(n+1)} - 1}, \quad (\text{D57})$$

where we used the fact that $\text{Tr}[\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)} \omega_{\mathbf{r}^{(x',j+1)}}^{(x',j+1)}] \leq 1$ for quantum states $\omega_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)}$ and $\omega_{\mathbf{r}^{(x',j+1)}}^{(x',j+1)}$ defined as in (D17). Then by recursively integrating over each perceptron we find

$$\left\langle \left| \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} B_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)})] \right| \right\rangle_{V_{j+1}, \dots, V_n} \leq 2^{3n+2} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-j}. \quad (\text{D58})$$

We now establish an upper bound on $\left| \left(\Delta(A^{(x,x',j)})_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}} \right) \right|$. Consider that

$$\left| \left(\Delta(A^{(x,x',j)})_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}} \right) \right| = \left| \text{Tr}[A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} A_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)}] - \frac{1}{2^{n+1}} \text{Tr}[A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}] \text{Tr}[A_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)}] \right| \quad (\text{D59})$$

$$= \left| q_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}}^{(x,j)} q_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)} \text{Tr}[\varphi_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}}^{(x,j)} \varphi_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)}] - \frac{q_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}}^{(x,j)} q_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)}}{2^{n+1}} \text{Tr}[\varphi_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}}^{(x,j)}] \text{Tr}[\varphi_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)}] \right| \quad (\text{D60})$$

$$\leq q_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}}^{(x,j)} q_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)}. \quad (\text{D61})$$

The second equality follows from (D39). The inequality follows from the fact that $\text{Tr}[\varphi_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}}^{(x,j)} \varphi_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)}] \leq 1$ and $\text{Tr}[\varphi_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}}^{(x,j)}] = 1$. Then by following arguments similar to (D55)–(D57) we find that

$$\left\langle \left| \left(\Delta(A^{(x,x',j)})_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}} \right) \right| \right\rangle_{V_1, \dots, V_{j-1}} \leq \langle q_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}}^{(x,j)} q_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)} \rangle \quad (\text{D62})$$

$$\leq q_{\mathbf{r}^{(x,1)} \mathbf{r}^{(x',1)}}^{(x,1)} q_{\mathbf{r}^{(x',1)} \mathbf{r}^{(x',1)}}^{(x',1)} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{j-1} \quad (\text{D63})$$

$$= \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{j-1}. \quad (\text{D64})$$

Therefore, combining (D58) and (D64) leads to

$$\langle \text{Tr}[G_j^x H_j] \text{Tr}[G_j^{x'} H_j] \rangle \leq \frac{2^{3n+2}}{2^{2(n+1)} - 1} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-1} \quad (\text{D65})$$

$$\leq f(n) \in \mathcal{O}(1/2^n). \quad (\text{D66})$$

c. Different j and different x

In this subsection, we establish an upper bound on the average of cross terms of the form $\langle \text{Tr}[G_j^x H_j] \text{Tr}[G_k^{x'} H_k] \rangle$. Without loss of generality we assume that $j < k$. From (D5) we get

$$\text{Tr}[G_k^{x'} H_k] = \text{Tr} \left[V_j \cdots V_1 \sigma_x^{\text{in}} V_1^\dagger \cdots V_j^\dagger \cdot \left(V_{j+1}^\dagger \cdots V_k^\dagger \left[V_{k+1}^\dagger \cdots V_n^\dagger \sigma_{x'}^{\text{out}} V_n \cdots V_{k+1}, H_k \right] V_k \cdots V_{j+1} \right) \right]. \quad (\text{D67})$$

Then by following (D6)–(D14), we find that

$$\langle \text{Tr}[G_j^x H_j] \text{Tr}[G_k^{x'} H_k] \rangle_{V_j} \leq \frac{1}{2^{2(n+1)} - 1} \left(\Delta(A^{(x,x',j)})_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}} \right) \left(\text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} M_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)})] \right), \quad (\text{D68})$$

where

$$M_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)} = \text{Tr}_{\bar{j}}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x',j)}\rangle\langle\mathbf{r}^{(x',j)}|) M^{(x',j)}], \quad (\text{D69})$$

$$M^{(x',j)} = V_{j+1}^\dagger \cdots V_k^\dagger \left[V_{k+1}^\dagger \cdots V_n^\dagger \sigma_{x'}^{\text{out}} V_n \cdots V_{k+1}, H_k \right] V_k \cdots V_{j+1}. \quad (\text{D70})$$

Moreover, from (D8) it follows that

$$B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} = \text{Tr}_{\bar{j}}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|) [V_{j+1}^\dagger \cdots V_n^\dagger \sigma_x^{\text{out}} V_n \cdots V_{j+1}, H_j]]. \quad (\text{D71})$$

We now argue using Lemma 2 that the average of $\text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} M_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)})]$ is zero. Here, H_j in (D71) and H_k in (D70) correspond to H and K in Lemma 2, respectively. Moreover, $(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|)$ and $(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x',j)}\rangle\langle\mathbf{r}^{(x',j)}|)$ correspond to P and P' , respectively. Furthermore, V_{j+1}^\dagger corresponds to V , while $V_{j+2}^\dagger \cdots V_k^\dagger$ corresponds to U . Finally, $V_{k+1}^\dagger \cdots V_n^\dagger \sigma_x^{\text{out}} V_n \cdots V_{k+1}$ and $V_{k+1}^\dagger \cdots V_n^\dagger \sigma_{x'}^{\text{out}} V_n \cdots V_{k+1}$ correspond to S and S' , respectively in Lemma 2. Hence from Lemma 2 it follows that

$$\left\langle \left(\text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} M_{\mathbf{r}^{(x',j)} \mathbf{r}^{(x',j)}}^{(x',j)})] \right) \right\rangle_{V_{j+1}} = 0, \quad (\text{D72})$$

which implies

$$\langle \text{Tr}[G_j^x H_j] \text{Tr}[G_k^{x'} H_k] \rangle \leq 0. \quad (\text{D73})$$

Therefore, by combining results from Sections **D1 a**–**D1 c**, it follows that

$$\langle (\partial_s C)^2 \rangle \leq f(n), \quad (\text{D74})$$

with $f(n)$ as in **(D1)**.

2. Local Cost

We now estimate the scaling of the variance of the partial derivative of the local cost function. We first note that for a local cost function, **(D4)** gets transformed as follows:

$$\partial_s C^L = \frac{i}{N} \sum_{x=1}^N \left[\frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^n \text{Tr}[G_{(i,j)}^x H_j] \right) \right], \quad (\text{D75})$$

where

$$\sigma_{(x,i)}^{\text{out}} = \mathbb{1}_{\text{in}, \bar{i}} \otimes |z_i^x\rangle\langle z_i^x|, \quad (\text{D76})$$

$$G_{(i,j)}^x = [V_j \cdots V_1 \sigma_x^{\text{in}} V_1^\dagger \cdots V_j^\dagger, V_{j+1}^\dagger \cdots V_n^\dagger \sigma_{(x,i)}^{\text{out}} V_n \cdots V_{j+1}]. \quad (\text{D77})$$

Moreover, from the cyclicity of trace, each term $\text{Tr}[G_{(i,j)}^x H_j]$ can always be expressed as follows:

$$\text{Tr}[G_{(i,j)}^x H_j] = \text{Tr} \left[V_j \cdots V_1 \sigma_x^{\text{in}} V_1^\dagger \cdots V_j^\dagger \left[V_{j+1}^\dagger \cdots V_n^\dagger \sigma_{(x,i)}^{\text{out}} V_n \cdots V_{j+1}, H_j \right] \right]. \quad (\text{D78})$$

We now provide a proof of Theorem **2** for local cost functions in the following subsections. Similarly to the proof for the global cost in Section **D1**, here we individually consider all different cases than can arise from the tripple summation in **(D75)**.

a. Fixed j , fixed i , and fixed x

Let us consider first the case $i < j$. From **(D8)** it follows that

$$B^{(x,i,j)} = [\mathbb{1}_{\text{in}, \bar{i}} \otimes |z_i^x\rangle\langle z_i^x|, \mathbb{1}_{\bar{j}} \otimes H_j] = 0, \quad (\text{D79})$$

which can be combined with **(D6)** to imply that $\text{Tr}[G_{i,j}^x H_j] = 0$.

Let us consider the case when $i = j$. From **(D8)** it follows that

$$B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,i,j)} = [\mathbb{1}_{\text{in}} \otimes |z_j^x\rangle\langle z_j^x|, H_j], \quad (\text{D80})$$

which further implies that

$$\text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,i,j)})^2] \leq 2^{n+2}, \quad (\text{D81})$$

where we used arguments similar to **(D20)**. Then, combining the previous inequality with **(D14)** and **(D48)**, we get

$$\langle (\text{Tr}[G_{i,j}^x H_j])^2 \rangle \leq \frac{2^{n+2}}{2^{2(n+1)} - 1} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{j-1} \quad (\text{D82})$$

$$\leq f(n) \in \mathcal{O}(1/2^n). \quad (\text{D83})$$

We now consider the case when $i > j$. By following **(D8)** again, we get

$$B^{(x,i,j)} = [V_{j+1}^\dagger \cdots V_i^\dagger (\mathbb{1}_{\text{in}, j+1, \dots, i-1} \otimes |z_i^x\rangle\langle z_i^x|) V_i \cdots V_{j+1} \otimes \mathbb{1}_j, H_j] \otimes \mathbb{1}_{1, \dots, j-1, i+1, i+2, \dots, n}, \quad (\text{D84})$$

$$B_{\mathbf{r}^{(x,j)}}^{(x,i,j)} = [v_j, H_j], \quad (\text{D85})$$

where

$$v_j = \text{Tr}_{\text{in},j}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{j+1,\dots,i})\Upsilon], \quad (\text{D86})$$

$$\Upsilon = (V_{j+1}^\dagger \dots V_i^\dagger (\mathbb{1}_{\text{in},j,j+1,\dots,i-1} \otimes |z_i^x\rangle\langle z_i^x|) V_i \dots V_{j+1}), \quad (\text{D87})$$

with v_j acting on all input qubits and on the j -th output qubit, and where the subscript in $|\mathbf{0}\rangle\langle\mathbf{0}|_{j+1,\dots,i}$ indicates that the projector acts on qubits $j+1, \dots, i$ in the output layer.

We note that $\Upsilon \leq \mathbb{1}_{\text{in},j+1,\dots,i}$ since $(\mathbb{1}_{\text{in},j,j+1,\dots,i-1} \otimes |z_i^x\rangle\langle z_i^x|) \leq \mathbb{1}_{\text{in},j,\dots,i}$, and since unitary transformations do not change the spectrum of an operator. Then, the following inequality holds:

$$\text{Tr}_{\text{in},j}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{j+1,\dots,i})(\mathbb{1}_{\text{in},j,\dots,i} - \Upsilon)(\mathbb{1}_{\text{in},j} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{j+1,\dots,i})] \geq 0, \quad (\text{D88})$$

which implies that $v_j \leq \mathbb{1}_{\text{in},j}$. Using arguments similar to the ones used in deriving (D19)–(D20), we find

$$\text{Tr}[(B_{\mathbf{r}(x,j)\mathbf{r}(x,j)}^{(x,i,j)})^2] \leq 2^{n+2}. \quad (\text{D89})$$

Again by combining this with (D14) and (D48), we get

$$\langle (\text{Tr}[G_{(i,j)}^x H_j])^2 \rangle \leq \frac{2^{n+2}}{2^{2(n+1)} - 1} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{j-1} \quad (\text{D90})$$

$$\leq f(n) \in \mathcal{O}(1/2^n). \quad (\text{D91})$$

b. Fixed j , different i , and different x

In this subsection, we establish a bound on $\langle (\text{Tr}[G_{(i,j)}^x H_j] \text{Tr}[G_{(i',j)}^{x'} H_j]) \rangle$. We first note that if either of i or i' are smaller than j , then from (D79) we have

$$\langle (\text{Tr}[G_{(i,j)}^x H_j] \text{Tr}[G_{(i',j)}^{x'} H_j]) \rangle = 0. \quad (\text{D92})$$

Let $i = j$ and $i' > j$. By means of the Cauchy-Schwarz inequality and invoking both (D81) and (D89), we find

$$|\text{Tr}[B_{\mathbf{r}(x,j)\mathbf{r}(x,j)}^{(x,i,j)} B_{\mathbf{r}(x',j)\mathbf{r}(x',j)}^{(x',i',j)}]| \leq \sqrt{\text{Tr}[(B_{\mathbf{r}(x,j)\mathbf{r}(x,j)}^{(x,i,j)})^2]} \sqrt{(\text{Tr}[B_{\mathbf{r}(x',j)\mathbf{r}(x',j)}^{(x',i',j)}])^2} \quad (\text{D93})$$

$$= 2^{n+2}. \quad (\text{D94})$$

Again by combining this with (D64), we get

$$\langle \text{Tr}[G_{(i,j)}^x H_j] \text{Tr}[G_{(i',j)}^{x'} H_j] \rangle \leq f(n) \in \mathcal{O}(1/2^n). \quad (\text{D95})$$

c. Different j , different i , and different x

In this subsection, we find an upper bound on $\langle (\text{Tr}[G_{(i,j)}^x H_j] \text{Tr}[G_{(i',k)}^{x'} H_k]) \rangle$. We first note that $\text{Tr}[G_{i',k}^{x'} H_k]$ can be expressed as (D65), where $\sigma_{x',i'}^{\text{out}}$ is replaced by $\sigma_{(x',i')}^{\text{out}}$ as in (D76). Without loss of generality we assume that $j < k$. Since the proof in Section D1c holds for any form of σ_x^{out} and $\sigma_{x'}^{\text{out}}$, it follows that

$$\langle (\text{Tr}[G_{(i,j)}^x H_j] \text{Tr}[G_{(i',k)}^{x'} H_k]) \rangle_{V_{j+1}} = 0, \quad (\text{D96})$$

and therefore, by combining results from Sections D2a–D2c, we find that

$$\langle (\partial_s C^L)^2 \rangle \leq f(n), \quad (\text{D97})$$

with $f(n)$ as in (D1). □

E. Proof of Theorem 1

In this section, we provide a proof of Theorem 1, which we recall for convenience.

Theorem 1. *Consider a DQNN with deep global perceptrons parametrized as in (C1), such that A_1^1, B_1^1 in (C7)–(C8), and V_j^1 ($\forall j$) form independent 2-designs over $n+1$ qubits. Then, the variance of a partial derivative of the cost function with respect to θ^ν is upper bounded as*

$$\text{Var}[\partial_\nu C^G] \leq g(n), \quad \text{with } g(n) \in \mathcal{O}(1/2^{2n}), \quad (\text{E1})$$

if O_x is the global operator of (C3), and upper bounded as

$$\text{Var}[\partial_\nu C^L] \leq h(n), \quad \text{with } h(n) \in \mathcal{O}(1/2^n), \quad (\text{E2})$$

when O_x is the local operator in Eq. (C4).

Proof. Here we first analyze the global cost function, and then we consider the case of local cost functions. Similarly to the proofs of the previous sections, we divide our derivations in several subsections consisting of different cases.

1. Global cost

Similar to Section D, we consider a DQNN with n input and n output qubits, and with no hidden layer. Then, we recall that we can compute the partial derivative of the cost function with respect to a parameter θ^ν in a given V_j , i.e., $\partial C/\partial\theta^\nu \equiv \partial_\nu C$, as

$$\partial_\nu C = \frac{i}{2N} \sum_{x=1}^N \partial_\nu C_x, \quad \text{with } \partial_\nu C_x \equiv \text{Tr} \left[A_j^1 \tilde{\sigma}_x^{\text{in}} (A_j^1)^\dagger [\mathbb{1}_{\bar{j}} \otimes \Gamma_k, (B_j^1)^\dagger \tilde{\sigma}_x^{\text{out}} B_j^1] \right], \quad (\text{E3})$$

with $B_j^1 = \mathbb{1}_{\bar{j}} \otimes \prod_{\nu=1}^{k-1} R_k(\theta^k) W_\nu$, $A_j^1 = \mathbb{1}_{\bar{j}} \otimes \prod_{\nu=k}^{n_j} R_k(\theta^k) W_\nu$, and where

$$\tilde{\sigma}_x^{\text{in}} = V_{j-1} \dots V_1 \sigma_x^{\text{in}} V_1^\dagger \dots V_{j-1}^\dagger, \quad (\text{E4})$$

$$\tilde{\sigma}_x^{\text{out}} = V_{j+1}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+1}. \quad (\text{E5})$$

a. Fixed x

We now establish an upper bound on a single term in (E3). That is, we consider a term $\langle (\partial_\nu C_x)^2 \rangle$ with fixed x . From (A8) it follows that

$$\langle (\partial_\nu C_x)^2 \rangle_{A_j^1, B_j^1} = \frac{2^n \text{Tr}[\Gamma_k^2]}{(2^{2n+2} - 1)^2} \sum_{\substack{\mathbf{p}\mathbf{q} \\ \mathbf{p}'\mathbf{q}'}} \Delta(\Omega^{(x,j)})_{\mathbf{q}\mathbf{p}}^{\mathbf{q}'\mathbf{p}'} \Delta(\Psi^{(x,j)})_{\mathbf{p}\mathbf{q}'}^{\mathbf{p}'\mathbf{q}'}, \quad (\text{E6})$$

where the summation runs over all bitstrings $\mathbf{p}, \mathbf{q}, \mathbf{p}', \mathbf{q}'$ of length $n-1$. In addition, we defined

$$\Delta(\Omega^{(x,j)})_{\mathbf{q}\mathbf{p}}^{\mathbf{q}'\mathbf{p}'} = \text{Tr}[\Omega_{\mathbf{q}\mathbf{p}}^{(x,j)} \Omega_{\mathbf{q}'\mathbf{p}'}^{(x,j)}] - \frac{\text{Tr}[\Omega_{\mathbf{q}\mathbf{p}}^{(x,j)}] \text{Tr}[\Omega_{\mathbf{q}'\mathbf{p}'}^{(x,j)}]}{2^{n+1}}, \quad (\text{E7})$$

$$\Delta(\Psi^{(x,j)})_{\mathbf{p}\mathbf{q}'}^{\mathbf{p}'\mathbf{q}'} = \text{Tr}[\Psi_{\mathbf{p}\mathbf{q}'}^{(x,j)} \Psi_{\mathbf{p}'\mathbf{q}'}^{(x,j)}] - \frac{\text{Tr}[\Psi_{\mathbf{p}\mathbf{q}'}^{(x,j)}] \text{Tr}[\Psi_{\mathbf{p}'\mathbf{q}'}^{(x,j)}]}{2^{n+1}}, \quad (\text{E8})$$

where $\Omega_{\mathbf{q}\mathbf{p}}^{(x,j)}$ and $\Psi_{\mathbf{p}\mathbf{q}'}^{(x,j)}$ are operators on $n+1$ qubits (all qubits in the input layer plus the j -th qubit in the output layer) defined as

$$\Omega_{\mathbf{q}\mathbf{p}}^{(x,j)} = \text{Tr}_{\bar{j}} \left[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{p}\rangle\langle \mathbf{q}|_{\bar{j}}) \tilde{\sigma}_x^{\text{out}} \right] \quad (\text{E9})$$

$$\Psi_{\mathbf{p}\mathbf{q}'}^{(x,j)} = \text{Tr}_{\bar{j}} \left[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{q}\rangle\langle \mathbf{p}|_{\bar{j}}) \tilde{\sigma}_x^{\text{in}} \right], \quad (\text{E10})$$

and where $\text{Tr}_{\bar{j}}$ indicates the trace over the subsystem of all qubits in the output layer except for the j -th qubit.

Similar to Section D, we assume that the output state is in the computational basis $|\phi_x^{\text{out}}\rangle \equiv |\mathbf{z}^x\rangle = |z_1^x z_2^x \dots z_n^x\rangle$. Then following arguments similar to the ones employed in Eq. (D9) and (D10), we find that

$$p_k = q_k = 0, \forall k \in \{j+1, \dots, n\}, \quad (\text{E11})$$

$$p_k = q_k = z_k^x, \forall k \in \{1, \dots, j-1\}, \quad (\text{E12})$$

which leads to a bitstring $\mathbf{r}^{(x,j)}$ as in (D11).

By recursively integrating over each randomly initialized perceptron as in (D39)–(D48), we find that

$$\langle \Delta(\Psi^{(x,j)})_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}} \rangle_{V_1, \dots, V_{j-1}} \leq \langle \text{Tr}[(\Psi_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \rangle \quad (\text{E13})$$

$$\leq \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{j-1}. \quad (\text{E14})$$

On the other hand, using (D17) and (D18), the following inequality holds:

$$\Delta(\Omega^{(x,j)})_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}} \leq (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \text{Tr}[(\omega_{\mathbf{r}^{(x,j)}}^{(x,j)})^2] \quad (\text{E15})$$

$$\leq (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \quad (\text{E16})$$

where $s_{\mathbf{r}^{(x,j)}}^{(x,j)}$ is given by (D18), and where we again used the fact that $\text{Tr}[(\omega_{\mathbf{r}^{(x,j)}}^{(x,j)})^2] \leq 1$ as $\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}$ is a quantum state. Finally, following (D22)–(D37), we find

$$\langle (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \rangle_{V_{j+1}, \dots, V_n} \leq 2^{2n} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-j}. \quad (\text{E17})$$

Therefore, combining (E10), (E17), and (E8), leads to

$$\langle (\partial_\nu C_x)^2 \rangle \leq \frac{2^{3n} \text{Tr}[\Gamma_k^2]}{(2^{2n+2} - 1)^2} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-1} \quad (\text{E18})$$

$$\leq \frac{2^{4n+1}}{(2^{2n+2} - 1)^2} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-1} \quad (\text{E19})$$

$$\leq g(n) \in \mathcal{O}(1/2^{2n}). \quad (\text{E20})$$

b. Different x

We now establish a bound on cross terms with different x , i.e., on terms of the form $\langle \partial_\nu C_x \partial_\nu C_{x'} \rangle$. From (A8), it follows that

$$\langle \partial_\nu C_x \partial_\nu C_{x'} \rangle_{A_j^1, B_j^1} = \frac{2^n \text{Tr}[\Gamma_k^2]}{(2^{2n+2} - 1)^2} \sum_{\substack{\mathbf{p}\mathbf{q} \\ \mathbf{p}'\mathbf{q}'}} \Delta(\Omega^{(x,x',j)})_{\mathbf{q}\mathbf{p}}^{\mathbf{q}'\mathbf{p}'} \Delta(\Psi^{(x,x',j)})_{\mathbf{p}\mathbf{q}}^{\mathbf{p}'\mathbf{q}'}, \quad (\text{E21})$$

where

$$\Delta(\Omega^{(x,x',j)})_{\mathbf{q}\mathbf{p}}^{\mathbf{q}'\mathbf{p}'} = \text{Tr}[\Omega_{\mathbf{q}\mathbf{p}}^{(x,j)} \Omega_{\mathbf{q}'\mathbf{p}'}^{(x',j)}] - \frac{\text{Tr}[\Omega_{\mathbf{q}\mathbf{p}}^{(x,j)}] \text{Tr}[\Omega_{\mathbf{q}'\mathbf{p}'}^{(x',j)}]}{2^{n+1}}, \quad (\text{E22})$$

$$\Delta(\Psi^{(x,x',j)})_{\mathbf{p}\mathbf{q}}^{\mathbf{p}'\mathbf{q}'} = \text{Tr}[\Psi_{\mathbf{p}\mathbf{q}}^{(x,j)} \Psi_{\mathbf{p}'\mathbf{q}'}^{(x',j)}] - \frac{\text{Tr}[\Psi_{\mathbf{p}\mathbf{q}}^{(x,j)}] \text{Tr}[\Psi_{\mathbf{p}'\mathbf{q}'}^{(x',j)}]}{2^{n+1}}, \quad (\text{E23})$$

and where $\Omega_{\mathbf{q}\mathbf{p}}^{(x,j)}$, and $\Psi_{\mathbf{q}\mathbf{p}}^{(x,j)}$ are defined according to Eqs. (E9), and (E10), respectively.

From arguments similar to those used in deriving (D61)–(D64), we find that

$$\left\langle \Delta(\Psi^{(x,x',j)})_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x',j)}} \right\rangle_{V_1, \dots, V_{j-1}} \leq \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{j-1}. \quad (\text{E24})$$

Similarly, from arguments similar to (F6)–(D18), we obtain

$$\Omega_{\mathbf{r}^{(x,j)}\mathbf{r}^{(x,j)}}^{(x,j)} = s_{\mathbf{r}^{(x,j)}}^{(x,j)} \omega_{\mathbf{r}^{(x,j)}}^{(x,j)}, \quad (\text{E25})$$

where $\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}$ and $s_{\mathbf{r}^{(x,j)}}^{(x,j)}$ are given by (D18) and (D17), respectively. Then, it is straightforward to show that

$$\left\langle \Delta(\Omega_{\mathbf{r}^{(x,x',j)}\mathbf{r}^{(x,j)}\mathbf{r}^{(x',j)}}^{(x,x',j)}) \right\rangle_{V_{j+1}, \dots, V_n} \leq \langle s_{\mathbf{r}^{(x,j)}}^{(x,j)} s_{\mathbf{r}^{(x',j)}}^{(x',j)} \rangle_{V_{j+1}, \dots, V_n} \quad (\text{E26})$$

$$\leq 2^{2n} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-j}. \quad (\text{E27})$$

Therefore, by combing (E21)–(E27), we find

$$\langle \partial_\nu C_x \partial_\nu C_{x'} \rangle \leq \frac{2^{4n+1}}{(2^{2n+2} - 1)^2} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+1)} - 1} \right)^{n-1} \quad (\text{E28})$$

$$\leq g(n) \in \mathcal{O}(1/2^{2n}). \quad (\text{E29})$$

Hence, recalling that $\partial_\nu C = (i/2N) \sum_{x=1}^N \partial_\nu C_x$, we get

$$\langle (\partial_\nu C)^2 \rangle \leq g(n) \in \mathcal{O}\left(\frac{1}{2^{2n}}\right). \quad (\text{E30})$$

2. Local Cost

We now estimate the scaling of the variance of the partial derivative of the local cost function. We define the local cost function as follows:

$$C^L = \frac{1}{nN} \sum_{x=1}^N \sum_{i=1}^n C_{x,i}^L, \quad \text{with} \quad C_{x,i}^L = \text{Tr} \left[\sigma_{(x,i)}^{\text{out}} V_n \dots V_1 \sigma_x^{\text{in}} V_1^\dagger \dots U_n^\dagger \right], \quad (\text{E31})$$

and where $\sigma_{(x,i)}^{\text{out}} = \mathbb{1}_{\text{in}, \bar{i}} \otimes |z_i^x\rangle\langle z_i^x|$. Here $\mathbb{1}_{\text{in}, \bar{i}}$ indicates the identity over all qubits in the input layer plus all the qubits in the output layer except for qubit i .

In what follows we first consider a single term in the summation over x in (E31). Here we have to consider the three following cases $i < j$, $i > j$, and $i = j$. Moreover, we remark that (E4)–(E10) remain the same, except for the fact that σ_x^{out} is replaced by $\sigma_{(x,i)}^{\text{out}}$.

a. Fixed x and $i < j$.

We first consider the case $i < j$. From (E9) and (E10), we find that

$$p_k = q_k = 0, \forall k \in \{j+1, \dots, n\}, \quad (\text{E32})$$

$$p_i = q_i = z_i^x, \quad (\text{E33})$$

$$p_k = q_k, \forall k \in \{1, \dots, i-1, i+1, \dots, j-1\}. \quad (\text{E34})$$

Then, we define the following set of bitstrings of length $n-1$:

$$\mathbf{r}^{(x,i,j,\mathbf{p})} = (p_1^{(x,i,j)}, p_2^{(x,i,j)}, \dots, p_{i-1}^{(x,i,j)}, z_i^x, p_{i+1}^{(x,i,j)}, \dots, p_{j-1}^{(x,i,j)}, 0, \dots, 0), \quad (\text{E35})$$

where j in the superscript implies that $\mathbf{r}^{(x,i,j,\mathbf{p})}$ is a bitstring over all qubits in the output layer, except the j -th qubit. The bold notation \mathbf{p} in (E35) indicates that each $p_k^{(x,i,j)} \in \{0, 1\}$. Then from (E9), we find that

$$\Omega_{\mathbf{r}^{(x,i,j,\mathbf{p})}\mathbf{r}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} = \text{Tr}_{\bar{j}} \left[(\mathbb{1}_{\text{in}, j} \otimes |\mathbf{r}^{(x,i,j,\mathbf{p})}\rangle\langle \mathbf{r}^{(x,i,j,\mathbf{p})}|) (\mathbb{1}_{\text{in}, \bar{i}} \otimes |z_i^x\rangle\langle z_i^x|) \right] \quad (\text{E36})$$

$$= \mathbb{1}_{\text{in}, j}, \quad (\text{E37})$$

which further implies

$$\Delta(\Omega^{(x,i,j,\mathbf{p},\mathbf{p}')})_{\mathbf{r}^{(x,i,j,\mathbf{p})} \mathbf{r}^{(x,i,j,\mathbf{p}')}} = \text{Tr}[\mathbb{1}_{\text{in},j} \mathbb{1}_{\text{in},j}] - \frac{\text{Tr}[\mathbb{1}_{\text{in},j}] \text{Tr}[\mathbb{1}_{\text{in},j}]}{2^{n+1}} \quad (\text{E38})$$

$$= 2^{n+1} - \frac{2^{2(n+1)}}{2^{n+1}} \quad (\text{E39})$$

$$= 0. \quad (\text{E40})$$

Therefore, from (G8) we get

$$\langle (\partial_\nu C_x)^2 \rangle_{A_j^1, B_j^1} = 0. \quad (\text{E41})$$

b. Fixed x and $i > j$.

Let us consider the case when $i > j$. We note that this case is different from the one studied in the previous section due to the fact that the perceptron unitaries do not commute with each other. We now have

$$p_k = q_k = 0, \forall k \in \{j+1, \dots, n\}, \quad (\text{E42})$$

$$p_k = q_k, \forall k \in \{1, \dots, j-1\}. \quad (\text{E43})$$

Similarly to (E35), we define here the bitstrings

$$\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} = (p_1^{(x,i,j)}, p_2^{(x,i,j)}, \dots, p_{j-1}^{(x,i,j)}, 0, \dots, 0). \quad (\text{E44})$$

In this case, from (E9) we obtain the operator

$$\Omega^{(x,i,j)} = \text{Tr}_{j+1, \dots, i} [(\mathbb{1}_{\text{in}} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{j+1, \dots, i}) V_{j+1}^\dagger \dots V_i^\dagger (\mathbb{1}_{\text{in}, j+1, \dots, i-1} \otimes |z_i^x\rangle\langle z_i^x|) V_i \dots V_{j+1}] \otimes \mathbb{1}_j \quad (\text{E45})$$

where $\mathbb{1}_j$ is the identity over qubit j in the output layer, and where $|\mathbf{0}\rangle\langle\mathbf{0}|_{j+1, \dots, i}$ is the projector onto the all-zero state on qubits $j+1, \dots, i$ in the output layer. Note that now $\Omega^{(x,i,j)}$ in (E45), and $\Delta(\Omega^{(x,i,j)})$ in (E7) are independent of the bitstring \mathbf{p} .

Similarly, by using (E8) we define

$$\Psi_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} = \text{Tr}_j [(\mathbb{1}_{\text{in},j} \otimes |\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}\rangle\langle\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}|) V_{j-1} \dots V_1 (\sigma_x^{\text{in}}) V_1^\dagger \dots V_{j-1}^\dagger] \quad (\text{E46})$$

$$= q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} \varphi_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})}, \quad (\text{E47})$$

where

$$q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} = \text{Tr} [(\mathbb{1}_{\text{in},j} \otimes |\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}\rangle\langle\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}|) V_{j-1} \dots V_1 (\sigma_x^{\text{in}}) V_1^\dagger \dots V_{j-1}^\dagger], \quad (\text{E48})$$

$$\varphi_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} = \frac{1}{q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})}} \text{Tr}_j [(\mathbb{1}_{\text{in},j} \otimes |\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}\rangle\langle\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}|) V_{j-1} \dots V_1 (\sigma_x^{\text{in}}) V_1^\dagger \dots V_{j-1}^\dagger]. \quad (\text{E49})$$

We now note that

$$\Delta(\Psi_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})})_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}} \leq q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p}')} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p}')}}^{(x,i,j,\mathbf{p}')}, \quad (\text{E50})$$

which follows from the fact that $\text{Tr}[(\varphi_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})})^2] \leq 1$ as $\varphi_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})}$ is a quantum state.

Then by combining (E45) and (E50), we get the following inequality:

$$\langle (\partial_\nu C_{x,i}^L)^2 \rangle_{A_j^1, B_j^1} \leq \frac{2^{2n+1}}{(2^{2n+2} - 1)^2} \text{Tr}[(\Omega^{(x,i,j)})^2] \sum_{\mathbf{p}, \mathbf{p}'} q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p}')} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p}')}}^{(x,i,j,\mathbf{p}')} \quad (\text{E51})$$

$$= \frac{2^{2n+1}}{(2^{2n+2} - 1)^2} \text{Tr}[(\Omega^{(x,i,j)})^2] \left(\sum_{\mathbf{p}} q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} \right) \left(\sum_{\mathbf{p}'} q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p}')} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p}')}}^{(x,i,j,\mathbf{p}')} \right) \quad (\text{E52})$$

$$= \frac{2^{2n+1}}{(2^{2n+2} - 1)^2} \text{Tr}[(\Omega^{(x,i,j)})^2], \quad (\text{E53})$$

where we used the fact that

$$\sum_{\mathbf{p}} q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} = \text{Tr}[(\mathbb{1}_{\text{in},1,\dots,j-1} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{j+1,\dots,n}) V_{j-1} \dots V_1 (\sigma_x^{\text{in}}) V_1^\dagger \dots V_{j-1}^\dagger] = 1. \quad (\text{E54})$$

Finally, by recursively applying Lemma 3, we find that

$$\langle (\partial_\nu C_{x,i}^L)^2 \rangle \leq \frac{2^{2n+1}}{(2^{2n+2} - 1)^2} \left\langle \text{Tr}[(\Omega^{(x,i,j)})^2] \right\rangle_{V_{j+1}, \dots, V_i} \quad (\text{E55})$$

$$\leq h(n) \in \mathcal{O}(1/2^n). \quad (\text{E56})$$

c. Fixed x and $i = j$.

In this case it can be easily shown that

$$\Omega^{(x,j,j)} = \mathbb{1}_{\text{in}} \otimes |z_j^x\rangle\langle z_j^x|, \quad (\text{E57})$$

and hence from (E7) we have

$$\Delta(\Omega^{(x,j)}) \leq 2^n. \quad (\text{E58})$$

Then, following a similar procedure to the one employed in the previous section (see Eqs.(E46)–(E54)), we obtain

$$\langle (\partial_\nu C_{x,i}^L)^2 \rangle \leq \frac{2^{3n+1}}{(2^{2n+2} - 1)^2} \quad (\text{E59})$$

$$\leq h(n) \in \mathcal{O}(1/2^n). \quad (\text{E60})$$

d. Different x and different i .

In this section we consider the cross terms of the form: $\langle \partial_\nu C_{x,i}^L \partial_\nu C_{x',i'}^L \rangle$. We first consider the case when either of i or i' is smaller than j . Then from arguments similar to those used in deriving (E36)–(E40), it can be shown that

$$\Delta(\Omega^{(x,x',i,i',j,\mathbf{p},\mathbf{p}')})_{\mathbf{r}^{(x',i',j,\mathbf{p}')} \mathbf{r}^{(x,i,j,\mathbf{p})}} = 0, \quad (\text{E61})$$

and therefore,

$$\langle (\partial_\nu C_{x,i}^L \partial_\nu C_{x',i'}^L) \rangle = 0. \quad (\text{E62})$$

Let us now consider the case when $i = j$ and $i' > j$. Then from (E57) and (E45) it follows that

$$\Omega^{(x,i,j)} = \mathbb{1}_{\text{in}} \otimes |z_j^x\rangle\langle z_j^x|, \quad (\text{E63})$$

$$\Omega^{(x',i',j)} = \text{Tr}_{j+1,\dots,i'}[(\mathbb{1}_{\text{in}} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{j+1,\dots,i'}) V_{j+1}^\dagger \dots V_{i'}^\dagger (\mathbb{1}_{\text{in},j+1,\dots,i'-1} \otimes |z_{i'}^x\rangle\langle z_{i'}^x|) V_{i'} \dots V_{j+1}] \otimes \mathbb{1}_j, \quad (\text{E64})$$

which further implies that $\Delta(\Omega^{(x,x',i,i',j)})$ in (E7) is independent of the bitstrings \mathbf{p} and \mathbf{p}' , and hence

$$\Delta(\Omega^{(x,x',i,i',j)}) \leq \text{Tr}[\tilde{\Omega}^{(x',i',j)}] \text{Tr}[|z_j^x\rangle\langle z_j^x|] \quad (\text{E65})$$

$$= \text{Tr}[\tilde{\Omega}^{(x',i',j)}], \quad (\text{E66})$$

where

$$\tilde{\Omega}^{(x',i',j)} = \text{Tr}_{j+1,\dots,i'}[(\mathbb{1}_{\text{in}} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{j+1,\dots,i'}) V_{j+1}^\dagger \dots V_{i'}^\dagger (\mathbb{1}_{\text{in},j+1,\dots,i'-1} \otimes |z_{i'}^x\rangle\langle z_{i'}^x|) V_{i'} \dots V_{j+1}]. \quad (\text{E67})$$

Finally, by using arguments similar to those used in deriving (E54), and by recursively invoking Lemma 1, the following bound holds:

$$\langle (\partial_\nu C_{x,i}^L \partial_\nu C_{x',i'}^L) \rangle \leq \frac{2^{3n+1}}{(2^{2n+2} - 1)^2} \quad (\text{E68})$$

$$\leq h(n) \in \mathcal{O}(1/2^n). \quad (\text{E69})$$

Let us finally consider the final case when both i and i' are greater than j and $i < i'$. By following the proof in [E 2 b](#), we find that

$$\begin{aligned} & \langle (\partial_\nu C_{x,i}^L \partial_k C_{x',i'}^L) \rangle_{A_j^1, B_j^1} \\ & \leq \frac{2^{2n+1}}{(2^{2n+2} - 1)^2} \left(\text{Tr}[\Omega^{(x,i,j)} \Omega^{(x',i',j)}] \right) \left(\sum_{\mathbf{p}} q_{\hat{\mathbf{r}}^{(x,i,j,\mathbf{p})} \hat{\mathbf{r}}^{(x,i,j,\mathbf{p})}}^{(x,i,j,\mathbf{p})} \right) \left(\sum_{\mathbf{p}'} q_{\hat{\mathbf{r}}^{(x',i',j,\mathbf{p}')} \hat{\mathbf{r}}^{(x',i',j,\mathbf{p}')}}^{(x',i',j,\mathbf{p}')} \right) \end{aligned} \quad (\text{E70})$$

$$= \frac{2^{2n+1}}{(2^{2n+2} - 1)^2} \left(\text{Tr}[\Omega^{(x,i,j)} \Omega^{(x',i',j)}] \right). \quad (\text{E71})$$

Then by following arguments as the one used in deriving [\(E55\)](#) and [\(E68\)](#), we get

$$\langle (\partial_k C_{x,i}^L \partial_\nu C_{x',i'}^L) \rangle \leq h(n) \in \mathcal{O}(1/2^n). \quad (\text{E72})$$

Therefore, by combining results from Sections [E 2 a–E 2 c](#), it follows that

$$\langle (\partial_\nu C^L)^2 \rangle \leq h(n) \in \mathcal{O}(1/2^n). \quad (\text{E73})$$

□

F. DQNNs with unitaries acting on $n + m$ qubits

In this section, we generalize our results for the case when unitaries in a DQNN acts on $n + m$ qubits. Here, n denotes the number of qubits in the layer l and m qubits are from the layer $l + 1$. In Sections [D](#) and [E](#), we proved our results for the case when $m = 1$. For instance, in [Fig. 4](#), we show a DQNN with two layers where the initial layer has 2 nodes and the final layer has 4 nodes. Here, $n = m = 2$, and the action of perceptrons can be described in terms of two unitaries, each acting on four qubits (2 in the input layer, and 2 in the output layer). Note that, we henceforth assume that each qubit in the $(l + 1)$ -th layers is acted upon by only one perceptron.

Since our theorem statements for this case are generalizable from the existing proofs, we only provide a brief sketch of our proof. As discussed previously, [Theorem 1](#) and [Theorem 2](#) for two different methods of updating the perceptrons. Below, we guide readers to a derivation similar to that of [Section D 1](#) and then state a new theorem for the general case.

Let us consider the case when the cost function is defined in terms of the global operator in [\(C3\)](#). We consider a particular example where the output state is on mn qubits. In general, it does not have to depend on n .

Following the arguments similar to [\(D6\)–\(D8\)](#), we get

$$p_k = q_k = 0, \forall k \in \{mj + 1, \dots, mn\}, p_k = q_k = z_k^x, \forall k \in \{1, \dots, m(j - 1)\}. \quad (\text{F1})$$

To write [\(F1\)](#) compactly, we define the following bitstring of length $m(n - 1)$:

$$\mathbf{r}^{(x,j)} \equiv (z_1^x, z_2^x, \dots, z_{m(j-1)}^x, 0, \dots, 0). \quad (\text{F2})$$

Using $\mathbf{r}^{(x,j)}$ we define $A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}$ and $B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}$, as in [\(D7\)](#) and [\(D8\)](#), respectively. Then by invoking [Lemma A8](#) we get

$$\langle (\text{Tr}[G_j^x H_j])^2 \rangle_{V_j} = \int d\mu(V_j) \text{Tr}[V_j A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} V_j^\dagger B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}] \text{Tr}[V_j A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} V_j^\dagger B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}] \quad (\text{F3})$$

$$= \frac{1}{2^{2(n+m)} - 1} \left(\text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] - \frac{1}{2^{n+m}} \text{Tr}[A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}]^2 \right) \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \quad (\text{F4})$$

$$\leq \frac{1}{2^{2(n+m)} - 1} \text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2], \quad (\text{F5})$$

where for getting the inequality we used the fact that $\text{Tr}[A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}] > 0$ as $A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}$ is a positive semidefinite operator.

We follow [\(D15\)](#) to compute the upper bound on the expectation value of $(\text{Tr}[G_j^x H_j])^2$. Let us note that since H_j only acts on all input qubits and on output qubits $m(j - 1) + 1, \dots, mj$, i.e., m output qubits in total, then $B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)}$ can be expressed in the following compact form

$$B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)} = s_{\mathbf{r}^{(x,j)}}^{(x,j)} [\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}, H_j], \quad (\text{F6})$$

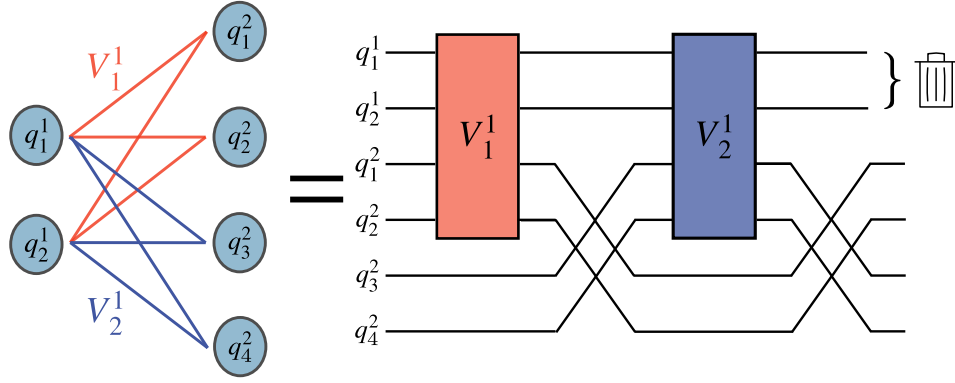


FIG. 4. Schematic diagram of a DQNN where each perceptron acts non trivially on n qubits in the l -th layer, and m qubits in the $(l + 1)$ -th layer. Shown is the case when $n = m = 2$.

where

$$\omega_{\mathbf{r}^{(x,j)}}^{(x,j)} = \frac{1}{s_{\mathbf{r}^{(x,j)}}^{(x,j)}} \text{Tr}_{\bar{j}}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|) V_{j+1}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+1}], \quad (\text{F7})$$

$$s_{\mathbf{r}^{(x,j)}}^{(x,j)} = \text{Tr}[(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|) V_{j+1}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+1}]. \quad (\text{F8})$$

Here, the subscript \bar{j} implies all output qubits besides $\{m(j-1) + 1, m(j-1) + 2, \dots, mj\}$. Similarly, $\mathbb{1}_{\text{in},j}$ implies an identity acting on all qubits in the input layer and on qubits $\{m(j-1) + 1, m(j-1) + 2, \dots, mj\}$ in the output layer. We henceforth follow this notation. Moreover, note that $\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}$ is a quantum state on all qubits in the input layer plus j -th block of m qubits (i.e., $\{m(j-1) + 1, m(j-1) + 2, \dots, mj\}$ qubits) in the output layer. Then from arguments similar to those used in deriving (D19)–(D20), we get

$$\text{Tr}[(\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}, H_j)^2] \leq 2^{n+m+1}, \quad (\text{F9})$$

where we used the assumption that $\text{Tr}[(H_j)^2] \leq 2^{n+m}$.

Finally, by combining Eqs. (F6) and (F9), we get that

$$\begin{aligned} \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(j)})^2] &= (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2 \text{Tr}[(\omega_{\mathbf{r}^{(x,j)}}^{(x,j)}, H_j)^2] \\ &\leq 2^{n+m+1} (s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2. \end{aligned} \quad (\text{F10})$$

Let us now evaluate the term $(s_{\mathbf{r}^{(x,j)}}^{(x,j)})^2$. By invoking Lemma 1, we get

$$s_{\mathbf{r}^{(x,j)}}^{(x,j)} = \text{Tr}[V_{j+1}(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|) V_{j+1}^\dagger V_{j+2}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+2}] \quad (\text{F11})$$

$$= \sum_{\mathbf{p}' \mathbf{q}'} \text{Tr}[V_{j+1} C_{\mathbf{q}' \mathbf{p}'}^{(x,j+1)} V_{j+1}^\dagger D_{\mathbf{p}' \mathbf{q}'}^{(x,j+1)}]. \quad (\text{F12})$$

Here the summation is over all bitstrings \mathbf{p}' and \mathbf{q}' of length $m(n-1)$, and

$$C_{\mathbf{q}' \mathbf{p}'}^{(x,j+1)} = \text{Tr}_{\bar{j+1}}[(\mathbb{1}_{\text{in},j+1} \otimes |\mathbf{p}'\rangle\langle\mathbf{q}'|)(\mathbb{1}_{\text{in},j} \otimes |\mathbf{r}^{(x,j)}\rangle\langle\mathbf{r}^{(x,j)}|)], \quad (\text{F13})$$

$$D_{\mathbf{p}' \mathbf{q}'}^{(x,j+1)} = \text{Tr}_{\bar{j+1}}[(\mathbb{1}_{\text{in},j+1} \otimes |\mathbf{q}'\rangle\langle\mathbf{p}'|) V_{j+2}^\dagger \dots V_n^\dagger \sigma_x^{\text{out}} V_n \dots V_{j+2}], \quad (\text{F14})$$

where $\text{Tr}_{\bar{j+1}}$ indicates the trace over all qubits in the output layer except qubits $\{mj+1, \dots, m(j+1)\}$.

Then from arguments similar to those used to deriving (D9) and (D10), we find

$$q'_k = p'_k = z_k^x, \forall k \in \{m(j-1) + 1, \dots, mj\} \quad (\text{F15})$$

$$q'_k = p'_k = r_k^{(x,j)}, \forall k \in \{1, 2, \dots, m(j-1), m(j+1) + 1, \dots, mn\}. \quad (\text{F16})$$

Let

$$\mathbf{r}^{(x,j+1)} \equiv (r_1^{(x,j)}, \dots, r_{m(j-1)}^{(x,j)}, z_{m(j-1)+1}^x, \dots, z_{mj}^x, r_{m(j+1)+1}^{(x,j)}, \dots, r_n^{(x,j)}). \quad (\text{F17})$$

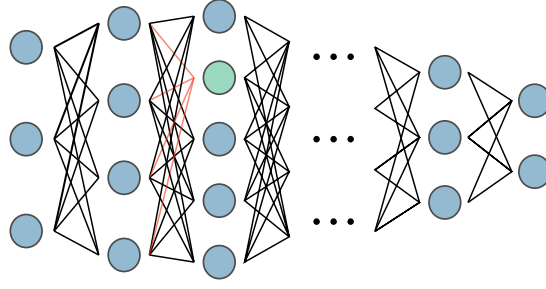


FIG. 5. Schematic diagram for generalizing the results to a DQNN with an arbitrary number of layers. We now consider the case when the partial derivative is taken with respect to a parameter in V_j^l , i.e., in a perceptron of the l -th layer acting on the j -th qubit of the $(l + 1)$ -th layer.

Then the average of $s_{\mathbf{r}^{(x,j)}}^{(x,j)}$ over V_{j+1} can be upper bounded as follows:

$$s_{\mathbf{r}^{(x,j)}}^{(x,j)} \leq \frac{2^n(2^n + 1/2)(s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)})^2}{2^{2(n+m)} - 1}, \quad (\text{F18})$$

where we employed (A8) and used arguments similar to those used in deriving (D35).

Here we remark that from $s_{\mathbf{r}^{(x,j+1)}}^{(x,j+1)}$ we can always define an operator $s_{\mathbf{r}^{(x,j+2)}}^{(x,j+2)}$ according to Eqs. (D22)–(D30). Moreover, by using the assumption that all randomly initialized perceptrons form 2-designs, we can recursively average over V_{j+2}, \dots, V_n . Therefore, from (F10), we get

$$\langle \text{Tr}[(B_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2] \rangle_{V_{j+1}, \dots, V_n} \leq (s_{\mathbf{r}^{(x,n)}}^{(x,n)})^2 2^{n+m+1} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+m)} - 1} \right)^{n-j} \quad (\text{F19})$$

$$= 2^{3n+m+1} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+m)} - 1} \right)^{n-j}, \quad (\text{F20})$$

where we used that fact that $s_{\mathbf{r}^{(x,n)}}^{(x,n)} = \text{Tr}[\sigma_x^{\text{out}}] = 2^n$.

We now compute the average of $\text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(x,j)})^2]$ over V_1, \dots, V_{j-1} . By following a similar procedure to the one previously employed, and we get

$$\langle (\text{Tr}[(A_{\mathbf{r}^{(x,j)} \mathbf{r}^{(x,j)}}^{(j)})^2])^2 \rangle_{V_1, \dots, V_{j-1}} \leq \left(\frac{2^n(2^n + 1/2)}{2^{2(n+m)} - 1} \right)^{j-1}. \quad (\text{F21})$$

Then from (D14), (D37), and (D48), it follows that

$$\langle (\text{Tr}[G_j^x H_j])^2 \rangle \leq \frac{2^{3n+m+1}}{2^{2(n+m)} - 1} \left(\frac{2^n(2^n + 1/2)}{2^{2(n+m)} - 1} \right)^{n-1} \quad (\text{F22})$$

$$\leq f(n, m) \in \mathcal{O}(1/2^{(2m-1)n}). \quad (\text{F23})$$

Thus, using the aforementioned result, and from a similar analysis of Sections D1b–D2c, we find that for both global cost functions

$$\langle (\partial_s C^G)^2 \rangle \leq f(n, m) \in \mathcal{O}(1/2^{(2m-1)n}). \quad (\text{F24})$$

Similarly, one can generalize the proof for local cost functions.

G. DQNNs with hidden layers

In this section we show how the results obtained in Sections D and E can be generalized to the case when the DQNN has L hidden layers.

First, let us follow the proof in Section E for a single input state labeled x , and note that here we can redefine the following quantities from Eq. (E4):

$$\sigma_x^{\text{in}} = |\phi_x^{\text{in}}\rangle\langle\phi_x^{\text{in}}| \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{\text{hid}} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{\text{out}}, \quad (\text{G1})$$

$$\sigma_x^{\text{out}} = \mathbb{1}_{\text{in}} \otimes \mathbb{1}_{\text{hid}} \otimes O_x. \quad (\text{G2})$$

Here, $|\mathbf{0}\rangle_{\text{hid}}$ denotes the input state to the hidden layers, and $\mathbb{1}_{\text{hid}}$ denotes the identity in the Hilbert space associated with the qubits in the hidden layers.

As shown in Fig. (5), we now consider the case when the partial derivative is taken with respect to a parameter in V_j^l , i.e., in a perceptron of the l -th layer acting on the j -th qubit of the $(l+1)$ -th layer. For convenience of notation, let us here define the following sets of qubit indexes:

- Set \mathcal{S}_1 : composed of all the indexes for qubits in the input layer, and all those for qubits in the first $(l-2)$ hidden layers.
- Set \mathcal{S}_2 : composed of all the indexes for qubits in the $(l-1)$ -th hidden layers.
- Set \mathcal{S}_3 : composed of all the indexes for qubits in the l -th layer with indexes smaller than j .
- Set \mathcal{S}_4 : composed of all the indexes for qubits in the l -th layer with indexes larger than j , and all those for the qubits in remaining hidden layers with indexes larger than l .

Therefore, the union of these sets along with the index of j -th qubit in the l -th layer describe all the qubits in the DQNN.

Here we note that the action of the unitaries prior to V_j^l , make it so that the quantum state in Eq. (G1) can be expressed as

$$\sigma_x^{(i,j-)} = |\phi_x^{(i,j)}\rangle\langle\phi_x^{(i,j)}|_{\mathcal{S}_1, \mathcal{S}_2} \otimes |\mathbf{0}\rangle\langle\mathbf{0}|_{j, \mathcal{S}_3}, \quad (\text{G3})$$

where $|\phi_x^{(i,j)}\rangle_{\mathcal{S}_1, \mathcal{S}_2}$ is the joint state of all the qubits with indexes in the sets \mathcal{S}_1 and \mathcal{S}_2 .

From the definition in (E9) and from (E10), it follows that

$$p_k = q_k = 0, \quad \text{for all qubits with index in } \mathcal{S}_4. \quad (\text{G4})$$

$$p_k = q_k = z_k^x, \quad \text{for all qubits with index in } \mathcal{S}_2 \text{ and } \mathcal{S}_3, \quad (\text{G5})$$

$$p_k = q_k, \quad \text{for all qubits with index in } \mathcal{S}_1. \quad (\text{G6})$$

Then, let us recall the definition $\Omega_{\mathbf{qp}}^{(x,j)} = \text{Tr}_{\mathcal{S}_1, \mathcal{S}_3, \mathcal{S}_4} [(\mathbb{1}_{\mathcal{S}_2, j} \otimes |\mathbf{p}\rangle\langle\mathbf{q}|_{\mathcal{S}_1, \mathcal{S}_3, \mathcal{S}_4}) \tilde{\sigma}_x^{\text{out}}]$. Here, since $\tilde{\sigma}_x^{\text{out}}$ acts trivially on all qubits in \mathcal{S}_1 , we have that

$$\text{Tr}_{\mathcal{S}_1, \mathcal{S}_3, \mathcal{S}_4} [(\mathbb{1}_{\mathcal{S}_2, j} \otimes |\mathbf{p}\rangle\langle\mathbf{q}|_{\mathcal{S}_1, \mathcal{S}_3, \mathcal{S}_4}) \tilde{\sigma}_x^{\text{out}}] = \text{Tr}_{\mathcal{S}_3, \mathcal{S}_4} [(\mathbb{1}_{\mathcal{S}_2, j} \otimes |\mathbf{p}\rangle\langle\mathbf{q}|_{\mathcal{S}_3, \mathcal{S}_4}) (\mathbb{1}_{\mathcal{S}_2, j, \mathcal{S}_3, \mathcal{S}_4} \otimes O_x)]. \quad (\text{G7})$$

and hence $\Omega_{\mathbf{qp}}^{(x,j)}$ is independent of the \mathbf{q} and \mathbf{p} indexes in \mathcal{S}_1 .

From Eq. (G8), we can now take the summation over \mathbf{q} and \mathbf{p} indexes in \mathcal{S}_1 to note that

$$\sum_{\mathbf{pq} \in \mathcal{S}_1} \Psi_{\mathbf{qp}}^{(x,j)} = \text{Tr}_{\mathcal{S}_1, \mathcal{S}_3, \mathcal{S}_4} [(\mathbb{1}_{\mathcal{S}_2, j} \otimes |\mathbf{p}\rangle\langle\mathbf{q}|_{\mathcal{S}_1, \mathcal{S}_3, \mathcal{S}_4}) \tilde{\sigma}_x^{\text{in}}] = \text{Tr}_{\mathcal{S}_1, \mathcal{S}_3, \mathcal{S}_4} [(\mathbb{1}_{\mathcal{S}_1, \mathcal{S}_2, j} \otimes |\mathbf{p}\rangle\langle\mathbf{q}|_{\mathcal{S}_3, \mathcal{S}_4}) \tilde{\sigma}_x^{\text{in}}], \quad (\text{G8})$$

where we have used the fact that $\sum_{\mathbf{pq} \in \mathcal{S}_1} |\mathbf{p}\rangle\langle\mathbf{q}|_{\mathcal{S}_1} = \sum_{\mathbf{pq} \in \mathcal{S}_1} |\mathbf{p}\rangle\langle\mathbf{p}|_{\mathcal{S}_1} = \mathbb{1}_{\mathcal{S}_1}$ (here $p_k = q_k$ from Eq. (G6)).

The latter allows us to get rid of all qubit indexes in \mathcal{S}_1 . Here, the remaining sets of equations

$$p_k = q_k = 0, \quad \text{for all qubits with index in } \mathcal{S}_4. \quad (\text{G9})$$

$$p_k = q_k = z_k^x, \quad \text{for all qubits with index in } \mathcal{S}_2 \text{ and } \mathcal{S}_3, \quad (\text{G10})$$

are exactly like those in Eqs. (E11) and (E12) of Section E. One can now follow steps similar to those used in deriving Eq. (E16) and Eq. (E17) to find that for a global cost function

$$\langle(\partial_\nu C_x)^2\rangle \leq g(n) \in \mathcal{O}(1/2^{2n}). \quad (\text{G11})$$

A similar analysis can be done for cross terms in x to show that if $\partial_\nu C = (i/2N) \sum_{x=1}^N \partial_\nu C_x$, then

$$\langle(\partial_\nu C)^2\rangle \leq g(n) \in \mathcal{O}\left(\frac{1}{2^{2n}}\right). \quad (\text{G12})$$

Similarly, one can generalize the proof for local cost functions with no hidden layers, to the case when the DQNN has L hidden layers.