

# Classical Splitting of Parametrized Quantum Circuits

Cenk Tüysüz,<sup>1,2,\*</sup> Giuseppe Clemente,<sup>1</sup> Arianna Crippa,<sup>1,2</sup> Tobias Hartung,<sup>3,4</sup> Stefan Kühn,<sup>4</sup> and Karl Jansen<sup>1</sup>

<sup>1</sup>*Deutsches Elektronen-Synchrotron (DESY), Platanenallee 6, 15738 Zeuthen, Germany*

<sup>2</sup>*Institut für Physik, Humboldt-Universität zu Berlin, Newtonstr. 15, 12489 Berlin, Germany*

<sup>3</sup>*Department of Mathematical Sciences, 4 West, University of Bath, Claverton Down, Bath BA2 7AY, UK*

<sup>4</sup>*Computation-Based Science and Technology Research Center,  
The Cyprus Institute, 20 Kavafi Street, 2121 Nicosia, Cyprus*

Barren plateaus appear to be a major obstacle to using variational quantum algorithms to simulate large-scale quantum systems or replace traditional machine learning algorithms. They can be caused by multiple factors such as expressivity, entanglement, locality of observables, or even hardware noise. We propose classical splitting of ansätze or parametrized quantum circuits to avoid barren plateaus. Classical splitting is realized by splitting an  $N$  qubit ansatz to multiple ansätze that consists of  $\mathcal{O}(\log N)$  qubits. We show that such an ansatz can be used to avoid barren plateaus. We support our results with numerical experiments and perform binary classification on classical and quantum datasets. Then, we propose an extension of the ansatz that is compatible with variational quantum simulations. Finally, we discuss a speed-up for gradient-based optimization and hardware implementation, robustness against noise and parallelization, making classical splitting an ideal tool for noisy intermediate scale quantum (NISQ) applications.

## I. INTRODUCTION

Variational quantum algorithms (VQAs)[1] are promising tools to solve a wide range of problems, such as finding the ground state of a given hamiltonian via the variational quantum eigensolver (VQE)[2], solving combinatorial optimization problems with the quantum approximate optimization algorithm (QAOA)[3] or solving classification problems using quantum neural networks [4].

VQAs are suitable for noisy intermediate scale quantum (NISQ) [5] hardware as they can be implemented with a small number of layers and gates for simple tasks. However, a scalability problem arises with the increasing number of qubits, hindering a possible advantage. VQAs rely on a classical optimization loop that updates the parameters of the ansatz iteratively until a condition on the cost function is satisfied. Classical optimizers use the information on the parametrized cost landscape to find the minimum. The updates on the parameters move the ansatz to a lower point on the cost surface. In 2018, McClean et al. showed that the cost landscape flattens with the increasing number of qubits, making it exponentially harder to find the solution for the optimizer [6]. The flattening was first observed by looking at the distribution of gradients across the parameter space, and the problem was named barren plateaus (BPs). A VQA is said to have a BP if its gradients decay exponentially with respect to one of its hyper-parameters, such as the number of qubits or layers.

Since the discovery of the BP problem, there has been significant progress that improved our understanding of what causes BPs and several methods to avoid them have been proposed. It has been shown that noise [7], entanglement [8], and the locality of the observable [9] play

an essential role for determining whether an ansatz will exhibit BPs. It has also been shown that the choice of ansatz (e.g. expressivity) of the circuit is one of the decisive factors that impact BPs [10]. For instance, the absence of BPs has been shown for quantum convolutional neural networks (QCNN) [11, 12] and tree tensor networks (TTN) [13, 14]. On the other hand, the hardware efficient ansatz (HEA) [6, 14, 15] and matrix product states (MPS) [14] have been shown to have BPs.

One of the essential discoveries showed that BPs are equivalent to cost concentration and narrow gorges [16]. This implies that BPs are not only a result of the exponentially decaying gradient but also of the cost function itself, and they can be identified by analyzing random points on the cost surface. As a result, gradient-free optimizers are also prone to BPs and do not offer a way to circumvent this problem [17].

Many methods have been suggested to mitigate BPs in the literature. Some of these methods suggest to use different ansätze or cost functions [18, 19], determining a better initial point to start the optimization [20–23], determining the step size during the optimization based on the ansatz [24], correlating parameters of the ansatz (e.g., restricting the directions of rotation) [10, 25], or combining multiple methods [26, 27].

In this work, we propose a novel idea in which we claim that if any ansatz of  $N$  qubits is classically separated to a set of ansätze with  $\mathcal{O}(\log N)$  qubits, the new ansatz will not exhibit Barren Plateaus. This work is not the first proposal in the literature that considers partitioning an ansatz. However, our proposal is significantly different. Most work in the literature first considers an ansatz and then emulates the result of that ansatz through many ansätze (exponentially many in general) with less number of qubits (which increases the effective size of quantum simulations) using gate decompositions, entanglement forging, divide and conquer or other methods [28–35]. On the other hand, this work proposes us-

---

\* [cenk.tueysuez@desy.de](mailto:cenk.tueysuez@desy.de)

ing ansätze that are classically split, meaning that there are no two-qubit gate operations between the subcircuits before splitting. This way, there is no need for gate decompositions or other computational steps. Our results show that this approach provides many benefits such as better trainability, robustness against noise and faster implementation on NISQ devices.

In the remainder of the paper, we start by giving an analytical illustration of the method in Section II. Then, we provide numerical evidence for our claim in Section III and extend our results to practical use cases by comparing binary classification performance of classical splitting for classical and quantum data. Next, we propose an extension of the classical splitting ansatz and perform experiments to simulate the ground state of the transversal-field ising hamiltonian. Finally, we discuss the advantages of employing classical splitting, make comments on future directions in Section IV and give an outlook in Section V.

## II. AVOIDING BARREN PLATEAUS

Barren plateaus (BPs) can be identified by investigating how the gradients of an ansatz scale with respect to a parameter. Here, we will start with the notation of McClean et al. and extend it to classical splitting [6]. The ansatz is composed of consecutive parametrized ( $V$ ) and non-parametrized entangling ( $W$ ) layers. We define  $U_l(\theta_l) = \exp(-i\theta_l V_l)$ , where  $V_l$  is a Hermitian operator and  $W_l$  is a generic unitary operator. Then the ansatz can be expressed with a multiplication of layers,

$$U(\boldsymbol{\theta}) = \prod_{l=1}^L U_l(\theta_l) W_l. \quad (1)$$

Then, for an observable  $O$  and input state of  $\rho$ , the cost is given as

$$C(\boldsymbol{\theta}) = \text{Tr}[OU(\boldsymbol{\theta})\rho U^\dagger(\boldsymbol{\theta})]. \quad (2)$$

The ansatz can be separated into two parts to investigate a certain layer, such that  $U_- \equiv \prod_{l=1}^{j-1} U_l(\theta_l) W_l$  and  $U_+ \equiv \prod_{l=j}^L U_l(\theta_l) W_l$ . Then, the gradient of the  $j^{\text{th}}$  parameter can be expressed as

$$\partial_j C(\boldsymbol{\theta}) = \frac{\partial C(\boldsymbol{\theta})}{\partial \theta_j} = i \text{Tr}[[V_j, U_+^\dagger O U_+] U_- \rho U_-^\dagger]. \quad (3)$$

The expected value of the gradients can be computed using the Haar measure. Please see Appendix A for more details on the Haar measure, unitary t-designs and details of the proofs in this section. If we assume the ansatz  $U(\boldsymbol{\theta})$  forms a unitary 2-design, then this implies that  $\langle \partial_k C(\boldsymbol{\theta}) \rangle = 0$  [6]. Since the average value of the gradients

are centered around zero, the variance of the distribution, which is defined as,

$$\text{Var}[\partial_k C(\boldsymbol{\theta})] = \langle (\partial_k C(\boldsymbol{\theta}))^2 \rangle - \langle \partial_k C(\boldsymbol{\theta}) \rangle^2, \quad (4)$$

can inform us about the size of the gradients. The variance of the gradients of the  $j^{\text{th}}$  parameter of the ansatz, where  $U_-$  and  $U_+$  are both assumed to be unitary 2-designs, and the number of qubits is  $N$ , is given as [6, 10],

$$\text{Var}[\partial_j C(\boldsymbol{\theta})] \approx \mathcal{O}\left(\frac{1}{2^{6N}}\right). \quad (5)$$

This means that for a unitary 2-design the gradients of the ansatz vanish exponentially with respect to the number of qubits  $N$ . Details of this proof is provided in Appendix A. Now, let us consider the classical splitting (CS) case. We split the ansatz  $U(\boldsymbol{\theta})$  to  $k$  many  $m$ -qubit ansätze, where we assume without loss of generality that  $N = k \times m$ . Then, we introduce a new notation for each classically split layer,

$$U_l^i(\theta_l^i) = e^{-i\theta_l^i V_l^i} W_l^i, \quad (6)$$

where index  $l$  determines the layer and index  $i$  determines which sub-circuit it belongs to. This notation combines the parametrized and entangling gates under  $U_l^i$ . Then, the overall CS ansatz can be expressed as,

$$U(\boldsymbol{\theta}) = \prod_{l=1}^L \bigotimes_{i=1}^k U_l^i(\theta_l^i) = \bigotimes_{i=1}^k \prod_{l=1}^L U_l^i(\theta_l^i) = \bigotimes_{i=1}^k U^i(\boldsymbol{\theta}^i). \quad (7)$$

The CS ansatz can be seen in Fig. 1a. Next, we will assume the observable and the input state to be classically split, such that they both can be expressed as a tensor product of  $m$ -qubit observables or states. This assumption restricts our proof to be valid only for  $m$ -local quantum states and  $m$ -local observables. It is important to note here that we use a definition that is different from the literature throughout the paper. For this proof, an  $m$ -local observable is an observable such that there are no operators that act on overlapping groups of  $m$  qubits. A generic  $m$ -local observable can be expressed as,

$$O_{m\text{-local}} = \sum_{i=1}^k O_i \otimes \mathbb{1}_{\bar{i}} = \sum_{i=1}^k \bigotimes_{j=1}^k (O_i - \mathbb{1}) \delta_{i,j} + \mathbb{1}, \quad (8)$$

where  $O_i$  is an observable over the qubits  $\{(i-1)m+1, (i-1)m+2, \dots, im\}$ , and  $\bar{i}$  represents the remaining  $N-m$  qubits. Then, the cost function becomes;

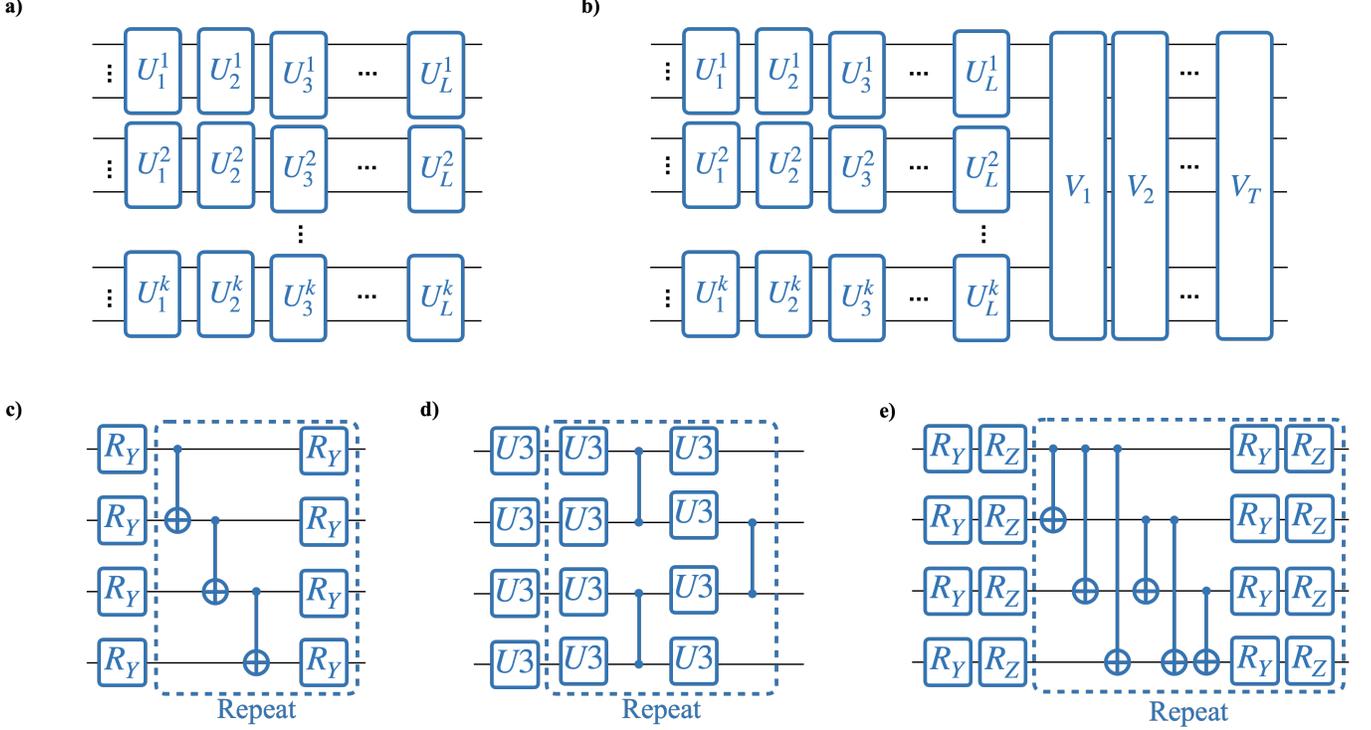


FIG. 1. All types of ansätze used in this work. (a) An  $N$ -qubit generic ansatz consisting of  $L$  layers of the parametrized unitary  $U$  are separated in to  $k = N/m$  many  $m$ -qubit ansätze. This ansatz will be referred to as the classically split (CS) ansatz. The standard ansatz can be recovered by setting  $m = N$ . (b) Extended classically split (ECS) ansatz. This is an extension to the CS ansatz. First  $L$  layers of the ansatz consists of  $k = N/m$  many  $m$  qubit  $U$  blocks. Then,  $T$  layers of  $N$  qubit  $V$  layers are applied. (c) A simple ansatz that consists of  $R_Y$  rotation gates and  $CX$  gates connected in a “ladder” layout. (d) Hardware Efficient Ansatz (HEA) that is used to produce the quantum dataset. Parameters of the first column of  $U3$  gates are sampled from a uniform distribution  $\in [-1, 1]$ , while the rest of the parameters are provided by the dataset [36]. (e) EfficientSU2 ansatz with “full” entangler layers [37].

$$\begin{aligned}
 C(\boldsymbol{\theta}) &= \sum_{i=1}^k \text{Tr} \left[ \bigotimes_{j=1}^k ((O_i - \mathbb{1}) \delta_{i,j} + \mathbb{1}) U^j(\boldsymbol{\theta}^j) \rho_j U^{j\dagger}(\boldsymbol{\theta}^j) \right] \\
 &= \sum_{i=1}^k \prod_{j=1}^k \text{Tr} \left[ ((O_i - \mathbb{1}) \delta_{i,j} + \mathbb{1}) U^j(\boldsymbol{\theta}^j) \rho_j U^{j\dagger}(\boldsymbol{\theta}^j) \right] \\
 &= \sum_{i=1}^k \text{Tr} [O_i U^i(\boldsymbol{\theta}^i) \rho_i U^{i\dagger}(\boldsymbol{\theta}^i)].
 \end{aligned} \tag{9}$$

This can be written as a simple sum,

$$C(\boldsymbol{\theta}) = \sum_{i=1}^k C^i(\boldsymbol{\theta}^i), \tag{10}$$

where,

$$C^i(\boldsymbol{\theta}^i) = \text{Tr} [O_i U^i(\boldsymbol{\theta}^i) \rho_i U^{i\dagger}(\boldsymbol{\theta}^i)]. \tag{11}$$

Then, the costs of each classically separated circuit are independent of each other. The gradient of  $j^{\text{th}}$  parameter of the  $i^{\text{th}}$  ansatz can be written as,

$$\begin{aligned}
 \partial_{i,j} C(\boldsymbol{\theta}) &= \partial_{i,j} C^i(\boldsymbol{\theta}^i) \\
 &= \partial_{i,j} (\text{Tr} [O_i U^i(\boldsymbol{\theta}^i) \rho_i U^{i\dagger}(\boldsymbol{\theta}^i)]).
 \end{aligned} \tag{12}$$

Now, let us consider each ansatz  $U^i(\boldsymbol{\theta}^i)$  to be a unitary 2-design. We want to choose the integer  $m$  such that it scales logarithmically in  $N$ . Hence, we choose  $\beta$  and  $\gamma$  appropriately, such that  $m = \beta \log_\gamma N$  holds. Then, if we combine Eq. (5) with Eq. (12), the variance of the gradient of  $j^{\text{th}}$  parameter can be expressed as

$$\text{Var}[\partial_j C(\boldsymbol{\theta})] \approx \mathcal{O} \left( \frac{1}{2^{(6m)}} \right) = \mathcal{O} \left( \frac{1}{N^{6\beta \log_\gamma 2}} \right). \tag{13}$$

Here, the dependence on  $i$  or  $j$  becomes irrelevant (a simpler choice for ansatz design would be to choose every new ansatz to be the same), so it can be dropped for a simpler notation. Similar to Eq. (5) the variance scales with the dimension of the hilbert space (e.g.  $\mathcal{O}(2^m)$ ). Then, the overall expression scales with,  $\mathcal{O}(N^{-6\beta \log_\gamma 2})$ , where  $\beta$  and  $\gamma$  are constant (e.g.  $\beta = 1$  and  $\gamma = 2$  results in  $m = \log_2 N$ ). As a result, the variance of the classically splitting ansatz scales with  $\mathcal{O}(\text{poly}(N)^{-1})$  instead of

$\mathcal{O}(\exp(N)^{-1})$ . Therefore, a CS ansatz, irrespective of its choice of gates or layout, can be used without leading to BPs.

### III. NUMERICAL EXPERIMENTS

In this section, we report results of four numerical experiments. We investigate the scaling of gradients under classical splitting by computing variances over many samples in Section III A. Then, we perform three experiments to observe how classical splitting affects performance of an ansatz. This task by itself leads to many questions as there are multitudes of metrics that one needs to compare and as many different problems one can consider. For this purpose, we consider problems well known in the literature, where trainability of ansätze plays a significant role.

First, we perform binary classification on a synthetic classical dataset in Section III B. The dataset contains two distributions that are called as classes. The goal is to predict the class of each sample. We perform the same task for distribution of quantum states in Section III C. Then, we give practical remarks in Section III D. Finally, we propose an extension to the CS ansatz and employ it for quantum simulating the ground state of the transverse field ising hamiltonian in Section III E.

For the first three experiments (Sections III A to III C), we consider the CS ansatz with layers that consists of  $R_Y$  rotation gates and CX entangling gates applied in a ladder formation for each layer. This layer can be seen in Fig. 1c. As the observable, we construct the 1-local observable defined in Eq. (14), where  $Z_i$  represents the Pauli-Z operator applied on the  $i^{\text{th}}$  qubit and  $\mathbb{1}_{\bar{i}}$  represents the identity operator applied on the rest of the qubits.

$$O = \frac{1}{N} \sum_{i=1}^N Z_i \otimes \mathbb{1}_{\bar{i}} \quad (14)$$

#### A. Barren Plateaus

Barren Plateaus are typically identified by looking at the variance of the first parameter over a set of random samples [6]. Recently, it has been shown that this is equivalent to looking at the variance of samples from the difference of two cost values evaluated at different random points of the parameter space [16]. Since the gradient-free optimization methods are also affected from BPs, the values of the cost become a more inclusive indicator [17]. For this reason, we will report our findings with respect to the cost, rather than the gradients to draw a broader picture. Results with respect to the gradient of the first parameter is presented in Appendix B for the sake of completeness.

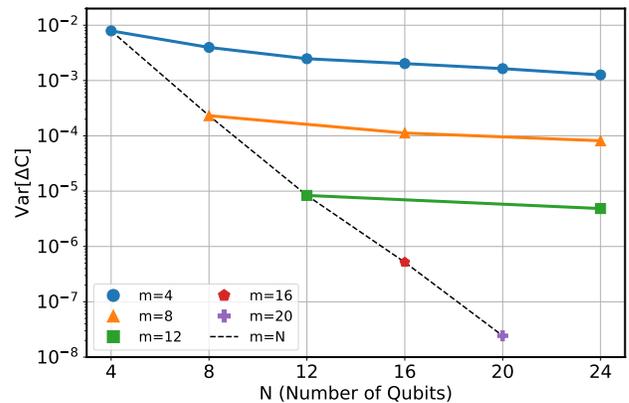


FIG. 2. The variance of the change in cost vs. the number of qubits for varying values of  $m$ . Each color/marker represents a certain value of  $m$  and data points of the standard ansatz ( $m = N$ ) is plotted with a dashed black line.

The experiments were performed using analytical gradients and expectation values, assuming a perfect quantum computer and infinite number of measurements, using PennyLane [38] and Pytorch [39]. Variances are computed over 2000 samples, where the values of the parameters are randomly drawn from a uniform distribution over  $[0, 2\pi]$ .

We start by presenting the variances over different values of  $m$  and  $N$  in Fig. 2. We fix the number of layers ( $L$ ) to  $N$ , so that the ansatz exhibits BPs in the no classical splitting setting ( $m = N$ ). The results indicate that a constant value of  $m$  resolves the exponential behaviour, as expected from Eq. (13). Furthermore, it is evident that larger values of  $m$  can allow the ansatz to escape BPs, given that  $m$  grows slow enough (e.g.  $\mathcal{O}(\log N)$ ).

Our theoretical findings illustrate that the classical splitting can be used to avoid BPs irrespective of the number of layers. In our first experiment, we numerically showed that this holds when we set  $L = N$ . Recent findings showed that, a transition to BPs happens at a depth of  $\mathcal{O}(\log N)$  for an ansatz with a local cost function [9]. Therefore, there is great importance in investigating the behaviour for larger values of  $L$ . For considerably low values of  $N$  (e.g.  $N < 32$ ), we can assume a constant value for  $m$  (e.g.  $m = 4$ ), such that  $m$  is approximately  $\mathcal{O}(\log N)$ . We present variances of two ansätze ( $m = 4$ ,  $m = N$ ) for up to 200 layers and 16 qubits in Fig. 3. For the standard ansatz, we see a clear transition to BPs with increasing number of layers, as expected [9]. On the other hand, the CS ansatz ( $m = 4$ ) shows a robust behavior from small to large number of layers.

These two experiments show the potential of the classical splitting in avoiding BPs. However, the question of whether this potential can be transferred in-to practice (e.g. binary classification performance or quantum simulation) still lacks an answer. Next, we will be addressing this question.

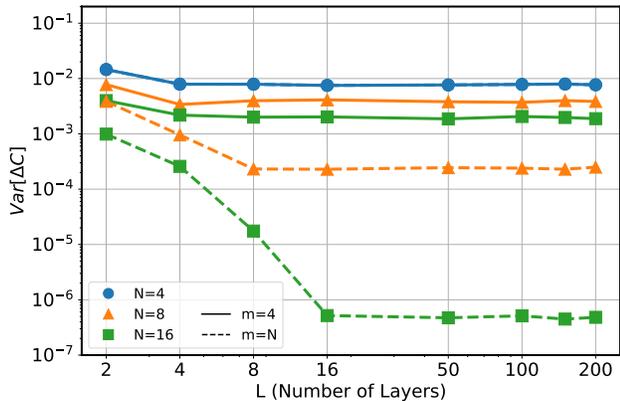


FIG. 3. The variance of the change in cost vs. the number of layers for  $m = 4$  (solid lines) and  $m = N$  (dashed lines) with varying number of qubits.

### B. Binary classification using a classical dataset

In this experiment, we will continue using the same ansatz with same assumptions to perform binary classification using a classical dataset. Our goal here is to compare performance of the CS ansatz to the standard case for increasing number of qubits. We need a dataset that can be scaled for this purpose. However, datasets are typically constant in dimension and do not offer an easy way to test the scalability in this sense. Therefore, we employ an ad-hoc dataset that can be produced with different number of features.

Three datasets ( $N = 4, 8$  and  $16$ ) were produced using the `make_classification` function of `scikit-learn`<sup>1</sup> [40]. This tool allows us to draw samples from an  $N$ -dimensional hypercube, where samples of each class are clustered around the vertices. Each dataset contains 420 training and 180 testing samples. Each of the data samples were encoded using one  $R_Y$  gate per qubit, such that each ansatz uses the same number of features of the given dataset. Please see Appendix C for more details on the production of the dataset and distributions of samples.

The binary classification was performed using the expectation value over the observable defined in Eq. (14) and the binary cross entropy function was used as the loss function during training, such that,

$$L(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y}), \quad (15)$$

where  $y$  (i.e.  $y \in \{0, 1\}$ ) is the class label of the given data sample and  $\hat{y}$  is the prediction (i.e.  $\hat{y} =$

<sup>1</sup> The classical dataset is produced for 600 data samples with a 420/180 train/test split, a class separation value of 1.0, 2.0% class assignment error and no redundant or repeated features.

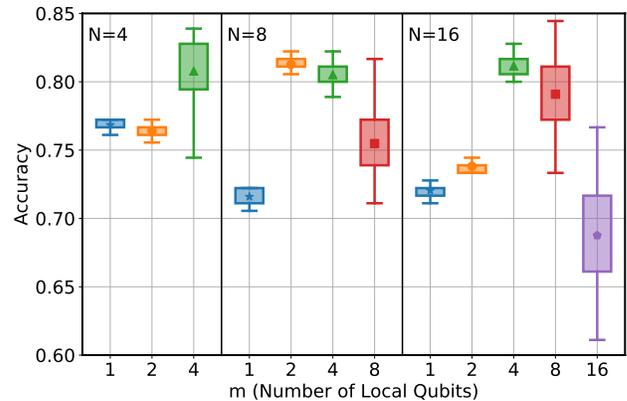


FIG. 4. Box plot of the best test accuracy obtained over 50 runs plotted with respect to the relevant local number of qubits ( $m$ ). Each column represents a problem with a different sample size (4, 8, 16). Each marker is placed on the median, boxes cover the range from the first to third quartiles and the error bars extend the quartiles by 3 times range. Each  $m$  value is plotted with a different marker and color.

$\text{Tr}[OU(\boldsymbol{\theta})\rho(x)U^\dagger(\boldsymbol{\theta})]$ , where  $x$  is the data sample)<sup>2</sup>. The ADAM optimizer [41] with a learning rate of 0.1 was used and all models are trained for 100 epochs using full batch size ( $\text{bs}=420$ )<sup>3</sup>. We report our results based on 50 runs for each setting.

Classification performance of ansätze for changing values of  $m$  using the three datasets are presented in Fig. 4. Here, the results show the distribution of accuracies over the test set. For the  $N = 4$  case, we see that the standard ( $m = N$ ) ansatz performs the best. However, this is not the case as we go to more qubits. For the 8 and 16 qubit cases, it is evident that  $m < N$  ansätze can match the performance of the standard ansatz. We can also see that the constant choice of  $m = 4$  can provide a robust performance with increasing number of qubits (at least up-to  $N = 16$ ), matching our expectations. Training curves of all settings are presented in Appendix D.

### C. Binary classification using a quantum dataset

The binary classification performance of the classical splitting over the classical datasets provides the first numerical evidence for their advantage against the standard ansätze. It is also important to investigate if they can be

<sup>2</sup> Here, the expectation value can have values between  $[-1, 1]$ , we scale it to be  $[0, 1]$  to compensate for the discrepancy between the class labels.

<sup>3</sup> In the case of  $N = m = L = 16$  full batch size was not possible due to vast memory requirement. Therefore,  $\text{bs}=60$  was used only for this case. In Appendix D, we show that using a smaller batch size does not affect the performance of the model significantly.

extended to problems where the data consists of quantum states. Our proof in Section II assumed the input states to be tensor product states. Now, we remove this constraint and use a quantum dataset.

For this experiment, we use the NTangled dataset [36]. NTangled dataset provides parameters to produce distributions of quantum states that are centered around different Concentrable Entanglement (CE) [42] values. CE is a measure of entanglement, which is defined as follows,

$$\text{CE}(|\Psi\rangle) = 1 - \frac{1}{2^N} \sum_{\alpha \in Q} \text{Tr}[\rho_\alpha^2], \quad (16)$$

where  $Q$  is the power set of the set  $\{1, 2, \dots, N\}$ , and  $\rho_\alpha$  is the reduced state of subsystems labeled by the elements of  $\alpha$  associated to  $|\Psi\rangle$ . The NTangled dataset provides three ansätze trained for different CE values for  $N=3, 4$  and  $8$ . We choose the Hardware Efficient Ansatz (Fig. 1d) with depth=5, such that the parameters of the first layer of  $U3$  gates are sampled from a unitary distribution  $\in [-1, 1]$  and the others are provided by the dataset. Then, we apply the same CS ansatz used in Section III B and perform binary classification such that the CE values are the labels of classes. The CE distributions of the produced quantum states are presented in Appendix E.

For the binary classification task, the same training settings are used as in Section III B, except this time models are trained until 50 epochs, as most models were able to reach 100% test accuracy. We report our results using different pairs of distributions in Table I. In the case of  $N = 4$ , we observed that classical splitting can perform at similar accuracy, even if the ansatz do not have any entangling gates ( $m = 1$ ). We see that entangling gates are needed for better performance if the problem gets harder (e.g. 0.25 vs. 0.35 case). If we go to a problem with more qubits, we can safely say that the CS ansatz can match the performance of the standard ansatz and converge faster.

#### D. Practical remarks on classical splitting

The efficacy of classical splitting relies on the parts of the circuit before and after the set of gates that undergo classical splitting. This can be seen most clearly if we set  $m = 1$  and apply classical splitting to the entire circuit after a possible initialization. In this case, we only perform single qubit operations after initialization. Hence, if the initialization produces a tensor product state, then the circuit subject to classical splitting with  $m = 1$  can no longer generate any entanglement. Similarly, if we initialize with the HEA (Fig. 1d) and apply classical splitting with  $m = 1$  to the remaining circuit, then no tensor product state can be found.

More generally,  $m = 1$  produces a circuit that cannot change the amount of entanglement. For other choices

of  $m$ , the picture becomes more complicated but, generally, the set of states that can be generated by the quantum circuit before classical splitting will be reduced to a subset based on the characteristics of the remaining initialization.

A naïve implementation of classical splitting therefore requires knowledge of the correct initialization such that the final solution can still be reached with the classically split circuit. In generic applications, this knowledge is likely not available. Hence, an adaptive approach to classical splitting should be considered.

One adaptive approach would be to increase  $m$  to check for improvements. After we observe no further training improvement with  $m = 1$ , we could move to  $m = 2$ . This enlarges the set of states the quantum circuit can reach, and thus may lead to further training improvements, at the cost of possibly stronger BP effects. However, if  $m = 1$  has already converged fairly well, then the state is already fairly close to the  $m = 2$  solution and it is unlikely to find a BP. With  $m = 2$  converged, we can then move to  $m = 4$  and continue the process by doubling  $m$  one step at a time.

If, for example, we consider the  $N = 4$  “0.25 vs. 0.3” case of Table I, we may start training with  $m = 1$ . This training converges to about 90% accuracy. Increasing  $m$  to  $m = 2$  will lead to further improvements that converge to about 98% accuracy. Finally, we can further improve the 98% to 100% accuracy by going to  $m = 4$ .

In this way, we utilize the efficiency of classical splitting to obtain an approximate solution which we then refine by trading efficiency for circuit expressivity through increasing  $m$ . At this point, the efficiency reduction should no longer lead to insurmountable complications as we already are close to the optimal solution for the current  $m$  value.

Another adaptive approach would be to use classical splitting to check and bypass plateaus. For example, if a VQE appears to be converged, it may also just be stuck in a plateau. Applying classical splitting at this point would reduce the effect of the plateau. Thus, if the VQE continues optimizing after classically splitting a seemingly converged circuit, we can conclude that this was in fact a plateau. After a suitable number of updates using the classically split circuit, we can then return to the full circuit in the hopes of having passed the plateau.

Unfortunately, this approach cannot be used to positively distinguish between true local optima and plateaus since the classical splitting reduces expressivity and thus introduces artificial constraints. Hence, if the set of states expressible by the classically split circuit is orthogonal to the gradient in the cost function landscape, then a plateau will be replaced with a local optimum and, thus, no improvements will be obtained. In this case, we therefore cannot conclude that the VQE has converged simply because classical splitting shows no improvements. However, experimenting with different implementations of classical splitting may result in cases that do not replace the plateau with an artificial local optimum.

TABLE I. Binary classification performance of ansätze with different values of  $m$  over different distributions of quantum states from the NTangled dataset [36]. Average of 50 runs are presented with errors showing the difference to maximum and minimum observed values. Best average value of each metric for the given task is printed in bold.

N	Task [CE Values]	L	m	Train Accuracy (%)	Avg. epochs to reach 90% Train Accuracy	Avg. epochs to reach 100% Train Accuracy	Test Accuracy (%)	Avg. epochs to reach 90% Test Accuracy	Avg. epochs to reach 100% Test Accuracy
4	0.05 vs. 0.35	4	1	$94.6^{+2.5}_{-1.7}$	$6.7^{+11.3}_{-5.7}$	N/A	$94.6^{+3.8}_{-1.8}$	$6.1^{+11.9}_{-5.1}$	N/A
			2	<b><math>100.0^{+0.0}_{-0.5}</math></b>	<b><math>4.9^{+12.1}_{-3.9}</math></b>	N/A	$100.0^{+0.0}_{-0.0}$	<b><math>3.9^{+11.1}_{-2.9}</math></b>	$10.8^{+26.2}_{-9.8}$
			4	$99.9^{+0.1}_{-1.6}$	$5.4^{+6.6}_{-4.4}$	N/A	$100.0^{+0.0}_{-1.1}$	$4.1^{+8.9}_{-3.1}$	N/A
4	0.25 vs. 0.35	4	1	$90.4^{+4.1}_{-3.5}$	N/A	N/A	$86.4^{+6.9}_{-5.9}$	N/A	N/A
			2	$98.2^{+1.5}_{-1.3}$	$7.7^{+25.3}_{-5.7}$	N/A	$97.1^{+2.3}_{-1.6}$	$7.9^{+27.1}_{-6.9}$	N/A
			4	<b><math>100.0^{+0.0}_{-0.4}</math></b>	<b><math>5.1^{+9.9}_{-4.1}</math></b>	N/A	<b><math>100.0^{+0.0}_{-1.1}</math></b>	<b><math>4.5^{+11.5}_{-3.5}</math></b>	N/A
8	0.15 vs. 0.45	8	1	$99.9^{+0.1}_{-0.2}$	$3.3^{+3.7}_{-2.3}$	N/A	$100.0^{+0.0}_{-0.0}$	$2.4^{+2.6}_{-1.4}$	$6.1^{+10.9}_{-5.1}$
			2	$100.0^{+0.0}_{-0.0}$	$2.5^{+2.5}_{-1.5}$	$7.1^{+11.9}_{-6.1}$	$100.0^{+0.0}_{-0.0}$	$1.5^{+2.5}_{-0.5}$	$3.2^{+6.8}_{-2.2}$
			4	$100.0^{+0.0}_{-0.0}$	<b><math>2.4^{+1.6}_{-1.4}</math></b>	<b><math>4.6^{+4.4}_{-3.6}</math></b>	$100.0^{+0.0}_{-0.0}$	<b><math>1.4^{+1.6}_{-0.4}</math></b>	<b><math>2.9^{+7.1}_{-1.9}</math></b>
			8	$100.0^{+0.0}_{-0.0}$	$2.8^{+2.2}_{-0.8}$	$7.8^{+11.2}_{-5.8}$	$100.0^{+0.0}_{-0.0}$	$1.8^{+2.2}_{-0.8}$	$4.8^{+8.2}_{-3.8}$
8	0.40 vs. 0.45	8	1	$99.9^{+0.1}_{-0.4}$	$3.1^{+2.9}_{-2.1}$	N/A	$99.6^{+0.4}_{-0.7}$	$2.2^{+2.8}_{-1.2}$	N/A
			2	$100.0^{+0.0}_{-0.0}$	$2.8^{+4.2}_{-1.8}$	$9.2^{+11.8}_{-8.2}$	$100.0^{+0.0}_{-0.0}$	$1.9^{+4.1}_{-0.9}$	$5.2^{+5.8}_{-4.2}$
			4	$100.0^{+0.0}_{-0.0}$	<b><math>2.4^{+1.6}_{-1.4}</math></b>	<b><math>5.3^{+12.7}_{-4.3}</math></b>	$100.0^{+0.0}_{-0.0}$	<b><math>1.5^{+1.5}_{-0.5}</math></b>	<b><math>3.2^{+9.8}_{-2.2}</math></b>
			8	$100.0^{+0.0}_{-0.0}$	$2.9^{+3.1}_{-0.9}$	$8.2^{+9.8}_{-6.2}$	$100.0^{+0.0}_{-0.0}$	$1.9^{+3.1}_{-0.9}$	$5.7^{+5.3}_{-4.7}$

### E. Extending classical splitting to VQE

Until now, we have investigated using classical splitting for binary classification problems. It succeeded by showing an overall better training performance in Section III B and a competitive performance and faster convergence in Section III C. In this section, we consider simulating the ground state of the transverse-field Ising Hamiltonian (TFIH) on a 1D chain. The TFIH with periodic boundary conditions can be defined as;

$$H = -J \sum_{i=1}^N Z_i Z_{i+1} - h \sum_{i=1}^N X_i, \quad (17)$$

for  $N$  lattice sites, where  $J$  determines the strength of interactions and  $h$  determines the strength of the external field. Simulating the TFIH on a 1D chain requires connectivity of qubits on the 1D chain. This contradicts with the assumption we made, when we proved absence of BPs for classically split ansätze in Section II, since the TFIH does not fit the definition we had for an  $m$ -local observable in Eq. (8). Therefore, we need to rely on the numerical experiments to talk about BPs under the new constraints.

The CS ansätze can only produce local entangled states, for this reason we need an extension of the ansatz in Fig. 1a. We propose to extend the classically split ansatz by adding standard layers at the end. The reason for adding them at the end is to keep the base of light cones<sup>4</sup> produced by the classically split layers constant.

Then, when we add the standard layers, the light cones will grow at a pace that is determined by the newly-added part<sup>5</sup>. This way, the overall ansatz can still escape BPs as long as the newly-added part does not exhibit BPs.

We define the extended classically split (ECS) ansatz with two types of layers. First  $L$  layers consist of classically split  $m$  qubit gate blocks. Then, there are  $T$  layers of any no-BP ansatz (see Fig. 1b). Since the first  $L$  layers can only produce  $m$ -local product states (i.e.  $m < \mathcal{O}(\log N)$ ), the existence of BPs depends only on the remaining  $T$  layers. This way we can choose very large  $L$ , but need to keep  $T$  small as standard ansätze reach BPs rather rapidly (e.g.  $\mathcal{O}(\log N)$  depth for a ladder connected ansatz [9]). We provide numerical evidence for avoiding BPs with the ECS ansatz in Appendix F.

For the experiment, we consider the Hamiltonian defined in Eq. (17) with  $J = 1, h = 1$ . Then, we implement the ECS ansatz with  $m = 4$  for total depth of 2, 4, 6 and 8. Each side of the ansatz consists of EfficientSU2 layers [37] (see Fig. 1e). The first  $L$  layers are classically split to subcircuits of  $m$  qubits, while the next  $T$  layers do not have any splitting. Total depth ( $D$ ) corresponds to  $L + T$ , where  $T = 0$  is equivalent to the CS ansatz,  $T = D$  is equivalent to the standard EfficientSU2 ansatz and other values explore hybrid use cases of the ECS ansatz. We report the energy error, which is the absolute difference between the final energy measurement and the exact ground state energy in Fig. 5. Results of 10 runs are averaged and plotted with their minimum and maximum values as the error bars. Experiments are performed under no noise assumption using 10k shots. The SPSA optimizer [44] is used with 10k iterations. Results

<sup>4</sup> A light cone or a causal cone of an ansatz is an abstract concept that illustrates how information spreads as more gates are applied. The types of gates and their connectivity determines the opening angle of the cone. The evidence from the literature suggests that there is a correspondence between the opening angle of the cone, BPs and quantum circuit complexity [9, 43].

<sup>5</sup> It also depends on the choice of  $m$ , but since we already have a constraint on  $m$  (i.e.  $m = \mathcal{O}(\log N)$ ) the newly-added ansatz will be the dominant component.

with  $m = 2$  and training curves of all runs are presented in Appendix G and H.

The upper panel shows that the mean error increases with increasing total depth in the no classical splitting setting ( $T = D$ ). This is mainly due to the flattening of the cost landscape, which makes the optimization process harder. On the other hand, setting  $T$  (e.g.  $T = 1$ ) to a low number provides a better error, since it preserves trainability despite the increasing total depth. This is a clear indication that the classical splitting allows deeper ansätze.

The lower panel shows the best error obtained in all the runs for two settings. Here, we observe that both settings achieve better errors with increasing depth initially. Then, the no CS setting shows rapidly increasing errors as it loses trainability rather quickly, compared to the ECS ansatz.

In this experiment, the best error was achieved with the fully classically split ansatz ( $T = 0$ ). This is mainly due to the employed EfficientSU2 ansatz not being a very good choice for this particular problem. This means that by employing other ansätze, the observed behaviour might change, making a larger value of  $T$  perform the best. Nevertheless, the results are still a good indication of how the trainability of the ansatz is affected by the choice of  $L$  and  $T$ . We plan to draw a more detailed picture of the tradeoff between values of  $L$  and  $T$  in a future work.

Simulating larger size systems requires a deep ansatz (linear or larger in system size) in general [1]. Although a problem-agnostic ansatz can perform well at small sizes, BPs forbid the scalability. Our results show that the ECS can help circumvent this issue and allow deeper ansätze. Here, we haven't investigated the potential of classical splitting to obtain the exact ground state energy of the model, but focused on the trainability aspect. Such a study is left as future work. Our goal here is to show that classical splitting can allow one to build wide and deep ansätze without exhibiting BPs. Typically, faster convergence or a better final energy might be achieved with a different ansatz or an optimizer, but this is out of scope of this work.

#### IV. DISCUSSION

In this work, we showed that the classical splitting of the ansätze can be used to escape BPs both analytically and numerically. Then, we investigated if the classical splitting hinders the learning capacity of the ansatz. Our experiments showed that this is not the case, and the classically split ansatz can match the performance at low number of qubits and is potentially superior at larger number of qubits.

In general the benefits of classical splitting comes from the reducing the effective Hilbert Space that the CS ansatz can explore. Classical splitting only allows the ansatz to produce  $m$ -qubit tensor product states, if the

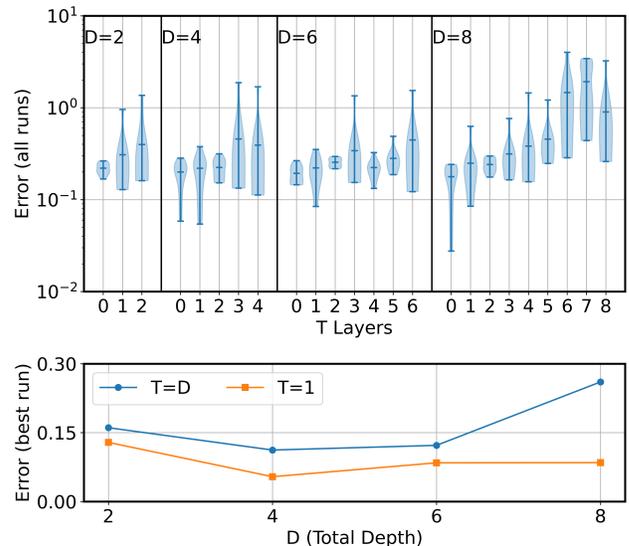


FIG. 5. Energy errors of ansätze with increasing total depth for  $N = 12$  TFIH using the extended classical splitting (ECS) ansatz with EfficientSU2 subblocks (see Fig. 1b) and  $m = 4$ . Total depth ( $D$ ) corresponds to  $L + T$ , where  $T = 0$  is equivalent to the CS ansatz,  $T = D$  is equivalent to standard EfficientSU2 and other values explore hybrid use cases of the ECS ansatz. Final energy measurements of 10 runs are averaged and plotted with their minimum and maximum values as the error bars on the upper panel. The lower panel shows the best errors obtained for  $T = 1$  and  $T = D$  settings. Energy error is the absolute difference of the energy measurement and the exact ground state energy.

input state is also a tensor product state following our assumptions in Section II. This, as a result, reduces the expressivity of the ansatz. Nevertheless, this also allows the ansatz to avoid BPs [10] by limiting the scaling behavior to the more favorable case of  $m$ -qubit systems. In the case of the classical splitting, the exponential increase of the Hilbert Space dimension is prevented and instead a polynomial scaling is enforced. For the  $m$ -local CS ansatz, each local Hilbert Space have  $\dim(H_k) = 2^m = N^{\beta \log_2 \gamma}$ . Although the advantage of using classical splitting may look trivial, there are many benefits of employing such an ansatz besides the numerical experiments we performed in Section III.

In our binary classification experiments using a classical dataset, we relied on single qubit and single rotation gate data encoding. This meant that any classically split ansatz had less information in each group. This could in fact be improved with embedding methods such as data re-uploading, where one can encode all the data points to each single qubit independently, such that there are alternating layers of rotation gates that encode the data and parametrized gates that are to be optimized [45]. Data re-uploading ansätze showed great classification performance even for low number of qubits. Since the classical splitting doesn't have a limit on the amount of layers,

data re-uploading would potentially be great way to get a performance increase.

Classical splitting can provide faster training when used with gradient based optimizers. In general, the exact gradients of ansätze are computed with the well-known parameter shift rule [46, 47]. However, this requires 2 instances of the same circuit to be executed per parameter. This quickly results in a bottleneck for the optimization procedure. An ansatz with  $L = N$  layers, where each layer has  $N$  parameters, requires  $\mathcal{O}(N^2)$  circuit executions to compute gradients for a single data sample. On the other hand, classical splitting provides cost functions that are independent of each other, as it was shown in Eq. (11). This allows gradients to be computed simultaneously across different instances of the classically split ansatz. As a result, the classically split ansatz optimization requires  $\mathcal{O}(N \log N)$  circuit executions for  $m = \mathcal{O}(\log N)$ .

The bottleneck in optimization is only one of the challenges of implementing scalable VQAs. Another problem that is worth mentioning here is the amount of two-qubit gates. NISQ hardware provides limited connectivity of qubits. The topology of the devices plays an essential role in the efficient implementation of quantum circuits [48]. Typically, a quantum circuit compilation (or transpilation) procedure is required to adapt a given circuit to be able to be compatible with the capabilities of the devices (e.g. converting gates to native gates, applying SWAP gates to connect qubits which are not physically connected) [49].

Classical splitting provides a significant reduction in number of two qubit gates as it divides a large qubit to many circuits with less qubits. To show the scale of the reduction, we can construct a set of hypothetical devices that has a 2D grid topology (square lattice with no diagonal connections). We start by considering the CS ansatz that consists the ansätze in Fig. 1c and extend it to a fully entangled architecture. A linear entangled ansatz has  $\mathcal{O}(N)$  two qubit gates, while a fully entangled one has  $\mathcal{O}(N^2)$  per layer. Then, we use Qiskit’s transpiler<sup>6</sup> [37] to fit these ansätze to the hypothetical devices and report the two qubit gate counts in Table II.

The amount of gates are not only important to have a better implementation but also to have a more precise results, since NISQ devices come with noisy gates. We consider the CX gate errors reported by IBM for their devices, which can be taken as  $\mathcal{O}(10^{-2})$  on average<sup>7</sup>. Then, as a figure of merit, we can assume 50% to be the limit,

TABLE II. Two qubit gate counts of different ansätze transpiled for hypothetical devices that has a 2D grid topology (square lattice with no diagonal connections).

$m$	$L$	amount of two qubit gates					
		linear entanglement			full entanglement		
		$N = 4$	$N = 16$	$N = 36$	$N = 4$	$N = 16$	$N = 36$
$N$	2	6	33	121	24	696	3601
	$N$	12	240	1362	46	5372	65040
4	2	6	24	54	24	92	250
	$N$	12	192	978	46	964	4376
2	2	4	16	36	4	16	42
	$N$	8	128	654	8	134	648

in which we can still get meaningful results. This would allow us to use 50 CX gates at most. Now, the results from Table II implies that it is possible to construct a 36 qubit, 2 layer ansatz with linear entanglement, if we employ classical splitting. This would not be possible for the standard case as it comes with more than twice two qubit gates. The reduction only gets better if we consider a full entanglement case. Following the same logic, to implement a 36 qubit, 36 layer, fully entangled ansatz, a CX gate error of  $\mathcal{O}(10^{-6})$  is needed, while the classically split ansatz only requires a CX gate error of  $\mathcal{O}(10^{-4})$ . A similar reduction in noise is also possible for other types of circuit partitioning methods [50].

Classically splitting an ansatz further allows faster implementation on hardware. A generic ansatz consists of two-qubit gates that follow one and another, matching a certain layout. We mentioned some of these as ladder/linear or full. However, this means that the hardware implementation of such an ansatz requires execution of these gates sequentially, taking a significant amount of time. To overcome such obstacles, ansätze such as the HEA (see Fig. 1d) are widely used in the literature [15]. Classically splitting an ansatz can reduce the implementation time significantly since it allows simultaneous two-qubit gates across different local circuits. This can mean a speed-up of from  $\mathcal{O}(N/\log N)$  to  $\mathcal{O}((N/\log N)^2)$  depending on the connectivity of the original ansatz.

Finally, the formulation we used in Section III B allows the CS ansatz to be implemented on smaller quantum computers instead of a single large quantum computer. This means that for similar problems, there are many implementation options available. These include using one large device, using many small devices (e.g.,  $\mathcal{O}(N/\log N)$  many  $\mathcal{O}(\log N)$  qubit devices) and parallelizing the task or using one small device and performing all computation sequentially. All of these features makes the classical splitting an ideal approach for Quantum Machine Learning (QML) applications using NISQ devices.

<sup>6</sup> Qiskit’s transpiler algorithm is a stochastic algorithm, meaning that it is possible to get better values if the algorithm is executed many times. Here, we run the algorithm two times and take the best results using optimization level 3, and sabre-sabre layout and routing methods. Although, It is possible to obtain better gate counts with more runs or different transpilation algorithms, the best values obtained wouldn’t change our conclusions.

<sup>7</sup> This value is chosen after a survey of devices listed on [IBM Quantum Cloud](#).

## V. CONCLUSION

In this work, we presented some foundational ideas of applying classical splitting to generic ansätze. Our results indicate many benefits of using classical splitting, such as better trainability, faster hardware implementation, faster convergence, robustness against noise and parallelization under certain conditions. These suggest that classical splitting or variations of this idea might play an essential role in how we are designing ansätze for QML problems. We also presented an extension to the initial classical splitting idea so that these types of ansätze can be used in VQE. The initial results that we presented in this work suggest that classical splitting can help improve the trainability and reach better error values. However, it is still an open question to what extent VQE can benefit from classical splitting. Our results encourages employing approaches that are based upon classically splitting or partitioning parametrized quan-

tum circuits [28–35], as they are in general more robust against hardware noise. We consider in-depth analysis and applications with VQE and QAOA as future directions for this work.

## ACKNOWLEDGMENTS

C.T. and A.C. are supported in part by the Helmholtz Association - “Innopolis Project Variational Quantum Computer Simulations (VQCS)”. S.K. acknowledges financial support from the Cyprus Research and Innovation Foundation under project “Future-proofing Scientific Applications for the Supercomputers of Tomorrow (FAST)”, contract no. COMPLEMENTARY/0916/0048, and “Quantum Computing for Lattice Gauge Theories (QC4LGT)”, contract no. EXCELLENCE/0421/0019. We thank Lena Funcke for valuable discussions.

- 
- [1] M. Cerezo, A. Arrasmith, R. Babbush, S. C. Benjamin, S. Endo, K. Fujii, J. R. McClean, K. Mitarai, X. Yuan, L. Cincio, and P. J. Coles, Variational quantum algorithms, *Nature Reviews Physics*, 625 (2021).
  - [2] 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).
  - [3] E. Farhi, J. Goldstone, and S. Gutmann, A Quantum Approximate Optimization Algorithm, [arXiv:1411.4028](https://arxiv.org/abs/1411.4028) (2014).
  - [4] E. Farhi and H. Neven, Classification with Quantum Neural Networks on Near Term Processors, [arXiv:1802.06002](https://arxiv.org/abs/1802.06002) (2018).
  - [5] J. Preskill, Quantum computing in the NISQ era and beyond, *Quantum* 2, 1 (2018).
  - [6] J. R. McClean, S. Boixo, V. N. Smelyanskiy, R. Babbush, and H. Neven, Barren plateaus in quantum neural network training landscapes, *Nature Communications* 9, 4812 (2018).
  - [7] S. Wang, E. Fontana, M. Cerezo, K. Sharma, A. Sone, L. Cincio, and P. J. Coles, Noise-induced barren plateaus in variational quantum algorithms, *Nature Communications* 12, 6961 (2021).
  - [8] C. Ortiz Marrero, M. Kieferová, and N. Wiebe, Entanglement-Induced Barren Plateaus, *PRX Quantum* 2, 040316 (2021).
  - [9] M. Cerezo, A. Sone, T. Volkoff, L. Cincio, and P. J. Coles, Cost function dependent barren plateaus in shallow parametrized quantum circuits, *Nature Communications* 12, 1791 (2021).
  - [10] Z. Holmes, K. Sharma, M. Cerezo, and P. J. Coles, Connecting Ansatz Expressibility to Gradient Magnitudes and Barren Plateaus, *PRX Quantum* 3, 010313 (2022).
  - [11] I. Cong, S. Choi, and M. D. Lukin, Quantum convolutional neural networks, *Nature Physics* 15, 1273–1278 (2019).
  - [12] A. Pesah, M. Cerezo, S. Wang, T. Volkoff, A. T. Sornborger, and P. J. Coles, Absence of Barren Plateaus in Quantum Convolutional Neural Networks, *Physical Review X* 11, 041011 (2021).
  - [13] E. Grant, M. Benedetti, S. Cao, A. Hallam, J. Lockhart, V. Stojevic, A. G. Green, and S. Severini, Hierarchical quantum classifiers, *npj Quantum Information* 4, 17 (2018).
  - [14] C. Zhao and X.-S. Gao, Analyzing the barren plateau phenomenon in training quantum neural networks with the ZX-calculus, *Quantum* 5, 466 (2021).
  - [15] 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).
  - [16] A. Arrasmith, Z. Holmes, M. Cerezo, and P. J. Coles, Equivalence of quantum barren plateaus to cost concentration and narrow gorges, [arXiv:2104.05868](https://arxiv.org/abs/2104.05868) (2021).
  - [17] A. Arrasmith, M. Cerezo, P. Czarnik, L. Cincio, and P. J. Coles, Effect of barren plateaus on gradient-free optimization, *Quantum* 5, 558 (2021).
  - [18] A. Wu, G. Li, Y. Ding, and Y. Xie, Mitigating Noise-Induced Gradient Vanishing in Variational Quantum Algorithm Training, [arXiv:2111.13209](https://arxiv.org/abs/2111.13209) (2021).
  - [19] K. Zhang, M.-H. Hsieh, L. Liu, and D. Tao, Toward Trainability of Deep Quantum Neural Networks, [arXiv:2112.15002](https://arxiv.org/abs/2112.15002) (2021).
  - [20] E. Grant, M. Ostaszewski, L. Wossnig, and M. Benedetti, An initialization strategy for addressing barren plateaus in parametrized quantum circuits, *Quantum* 3, 214 (2019).
  - [21] H.-Y. Liu, T.-P. Sun, Y.-C. Wu, Y.-J. Han, and G.-P. Guo, A Parameter Initialization Method for Variational Quantum Algorithms to Mitigate Barren Plateaus Based on Transfer Learning, [arXiv:2112.10952](https://arxiv.org/abs/2112.10952) (2021).
  - [22] A. Rad, A. Seif, and N. M. Linke, Surviving The Barren Plateau in Variational Quantum Circuits with Bayesian

- Learning Initialization, [arXiv:2203.02464 \(2022\)](#).
- [23] K. Zhang, M.-H. Hsieh, L. Liu, and D. Tao, Gaussian initializations help deep variational quantum circuits escape from the barren plateau, [arXiv:2203.09376 \(2022\)](#).
- [24] S. H. Sack, R. A. Medina, A. A. Michailidis, R. Kueng, and M. Serbyn, Avoiding barren plateaus using classical shadows, [arXiv:2201.08194 \(2022\)](#).
- [25] T. Volkoff and P. J. Coles, Large gradients via correlation in random parameterized quantum circuits, *Quantum Science and Technology* **6**, 025008 (2021), arXiv: 2005.12200.
- [26] T. L. Patti, K. Najafi, X. Gao, and S. F. Yelin, Entanglement devised barren plateau mitigation, *Physical Review Research* **3**, 033090 (2021).
- [27] L. Broers and L. Mathey, Optimization of Quantum Algorithm Protocols without Barren Plateaus, [arXiv:2111.08085 \(2021\)](#).
- [28] S. Bravyi, G. Smith, and J. A. Smolin, Trading Classical and Quantum Computational Resources, *Physical Review X* **6**, 021043 (2016).
- [29] T. Peng, A. W. Harrow, M. Ozols, and X. Wu, Simulating Large Quantum Circuits on a Small Quantum Computer, *Physical Review Letters* **125**, 150504 (2020).
- [30] W. Tang, T. Tomesh, M. Suchara, J. Larson, and M. Martonosi, CutQC: Using Small Quantum Computers for Large Quantum Circuit Evaluations, *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, 473 (2021), arXiv: 2012.02333.
- [31] M. A. Perlin, Z. H. Saleem, M. Suchara, and J. C. Osborn, Quantum circuit cutting with maximum-likelihood tomography, *npj Quantum Information* **7**, 1 (2021).
- [32] A. Eddins, M. Motta, T. P. Gujarati, S. Bravyi, A. Mezzacapo, C. Hadfield, and S. Sheldon, Doubling the size of quantum simulators by entanglement forging, [arXiv:2104.10220 \(2021\)](#).
- [33] Z. H. Saleem, T. Tomesh, M. A. Perlin, P. Gokhale, and M. Suchara, Quantum Divide and Conquer for Combinatorial Optimization and Distributed Computing, [arXiv:2107.07532 \(2021\)](#).
- [34] K. Fujii, K. Mizuta, H. Ueda, K. Mitarai, W. Mizukami, and Y. O. Nakagawa, Deep Variational Quantum Eigensolver: A Divide-And-Conquer Method for Solving a Larger Problem with Smaller Size Quantum Computers, *PRX Quantum* **3**, 010346 (2022).
- [35] S. C. Marshall, C. Gyurik, and V. Dunjko, High Dimensional Quantum Learning With Small Quantum Computers, [arXiv:2203.13739 \(2022\)](#).
- [36] L. Schatzki, A. Arrasmith, P. J. Coles, and M. Cerezo, Entangled Datasets for Quantum Machine Learning, [arXiv:2109.03400 \(2021\)](#).
- [37] M. Treinish, J. Gambetta, P. Nation, P. Kassebaum, qiskit bot, D. M. Rodríguez, S. d. I. P. González, S. Hu, K. Krsulich, L. Zdanski, J. Garrison, J. Yu, J. Gacon, D. McKay, J. Gomez, L. Capelluto, Travis-S-IBM, M. Marques, A. Panigrahi, J. Lishman, lerongil, R. I. Rahman, S. Wood, L. Bello, T. Itoko, D. Singh, Drew, E. Arbel, J. Schwarm, and J. Daniel, *Qiskit: An Open-source Framework for Quantum Computing* (2022).
- [38] V. Bergholm, J. Izaac, M. Schuld, C. Gogolin, M. S. Alam, S. Ahmed, J. M. Arrazola, C. Blank, A. Delgado, S. Jahangiri, K. McKiernan, J. J. Meyer, Z. Niu, A. Száva, and N. Killoran, PennyLane: Automatic differentiation of hybrid quantum-classical computations, [arXiv:1811.04968 \(2020\)](#).
- [39] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, PyTorch: An Imperative Style, High-Performance Deep Learning Library, in *Advances in Neural Information Processing Systems*, Vol. 32, edited by H. Wallach, H. Larochelle, A. Beygelzimer, F. d. Alché-Buc, E. Fox, and R. Garnett (Curran Associates, Inc., 2019).
- [40] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research* **12**, 2825 (2011).
- [41] D. P. Kingma and J. Ba, Adam: A Method for Stochastic Optimization, [arXiv:1412.6980 \(2017\)](#).
- [42] J. L. Beckey, N. Gigena, P. J. Coles, and M. Cerezo, Computable and Operationally Meaningful Multipartite Entanglement Measures, *Physical Review Letters* **127**, 140501 (2021).
- [43] J. Haferkamp, P. Faist, N. B. T. Kothakonda, J. Eisert, and N. Y. Halpern, Linear growth of quantum circuit complexity, [arXiv:2106.05305 \(2021\)](#).
- [44] J. C. Spall, Overview of the simultaneous perturbation method for efficient optimization, *Johns Hopkins APL Technical Digest* **19**, 482 (1998).
- [45] A. Pérez-Salinas, A. Cervera-Lierta, E. Gil-Fuster, and J. I. Latorre, Data re-uploading for a universal quantum classifier, *Quantum* **4**, 226 (2020).
- [46] K. Mitarai, M. Negoro, M. Kitagawa, and K. Fujii, Quantum circuit learning, *Physical Review A* **98**, 032309 (2018).
- [47] M. Schuld, V. Bergholm, C. Gogolin, J. Izaac, and N. Killoran, Evaluating analytic gradients on quantum hardware, *Physical Review A* **99**, 1 (2019).
- [48] J. Weidenfeller, L. C. Valor, J. Gacon, C. Tornow, L. Bello, S. Woerner, and D. J. Egger, Scaling of the quantum approximate optimization algorithm on superconducting qubit based hardware, [arXiv:2202.03459 10.48550/arXiv.2202.03459 \(2022\)](#).
- [49] A. Botea, A. Kishimoto, and R. Marinescu, On the Complexity of Quantum Circuit Compilation, *Proceedings of the International Symposium on Combinatorial Search* **9**, 138 (2018).
- [50] S. Basu, A. Saha, A. Chakrabarti, and S. Sur-Kolay, *i*-QER: An Intelligent Approach towards Quantum Error Reduction, [arXiv:2110.06347 \(2022\)](#).
- [51] J. A. Miszczak and Z. Puchała, Symbolic integration with respect to the Haar measure on the unitary groups, *Bulletin of the Polish Academy of Sciences: Technical Sciences*; 2017; 65; No 1; 21-27 (2017).

### Appendix A:

When analyzing the size of the gradients of an ansatz we need tools that allows integration over all states allowed by the ansatz over the  $d$ -dimensional Hilbert Space. This can be achieved by using the Haar measure. Haar measure is an invariant measure over the  $SU(d)$  group. An ensemble of unitary operators  $U$  is called as a unitary  $t$ -design if they are equal to the Haar measure  $\mu(U)$  up-to polynomial order  $t$ . Then, the expectation of ensemble  $U$ , where unitary  $V_i$  can be sampled with probability  $p_i$  is given as,

$$\mathbb{E}_H^t(\rho) = \int U^{\otimes t} \rho (U^{\otimes t})^\dagger dU = \sum_i p_i V_i^{\otimes t} \rho (V_i^{\otimes t})^\dagger. \quad (\text{A1})$$

Then, to perform symbolic integration over the Haar measure we will need to use some properties of the measure [51]. For the first moment we have,

$$\int d\mu(U) U_{ij} U_{km}^* = \frac{\delta_{ik} \delta_{jm}}{d}, \quad (\text{A2})$$

where  $d$  is the dimension of the Unitary, such that  $d = 2^N$  and  $N$  is number of qubits. Then, for the second moment we have,

$$\begin{aligned} \int d\mu(U) U_{i_1 j_1} U_{i_2 j_2} U_{k_1 m_1}^* U_{k_2 m_2}^* &= \\ &= \frac{\delta_{i_1 k_1} \delta_{j_1 m_1} \delta_{i_2 k_2} \delta_{j_2 m_2} + \delta_{i_1 k_2} \delta_{i_2 k_1} \delta_{j_1 m_2} \delta_{j_2 m_1}}{d^2 + 1} - \frac{\delta_{i_1 k_1} \delta_{j_2 m_2} \delta_{j_1 m_2} \delta_{j_2 m_1} + \delta_{i_1 k_2} \delta_{i_2 k_1} \delta_{j_1 m_1} \delta_{j_2 m_2}}{d(d^2 + 1)} \end{aligned} \quad (\text{A3})$$

Then one can derive the following identities for integrals over the Haar measure [6, 9, 10],

$$\int d\mu(U) \text{Tr}[U A U^\dagger B] = \frac{\text{Tr}[A] \text{Tr}[B]}{d}. \quad (\text{A4})$$

We can extend this to the second moment to obtain the following identity,

$$\begin{aligned} \int d\mu(U) \text{Tr}[U A U^\dagger B U C U^\dagger D] &= \\ &= \frac{\text{Tr}[A] \text{Tr}[C] \text{Tr}[B D] + \text{Tr}[A C] \text{Tr}[B] \text{Tr}[D]}{d^2 - 1} - \frac{\text{Tr}[A C] \text{Tr}[B D] + \text{Tr}[A] \text{Tr}[B] \text{Tr}[C] \text{Tr}[D]}{d(d^2 - 1)}. \end{aligned} \quad (\text{A5})$$

We also have,

$$\begin{aligned} \int d\mu(U) \text{Tr}[U A U^\dagger B] \text{Tr}[U C U^\dagger D] &= \\ &= \frac{\text{Tr}[A C] \text{Tr}[B] \text{Tr}[D] + \text{Tr}[A C] \text{Tr}[B D]}{d^2 - 1} - \frac{\text{Tr}[A C] \text{Tr}[B] \text{Tr}[D] + \text{Tr}[A] \text{Tr}[C] \text{Tr}[B D]}{d(d^2 - 1)}. \end{aligned} \quad (\text{A6})$$

Now, we can use these identities to compute the average value of the gradients. Let's start by reminding ourselves the definitions we used before. The ansatz is composed of consecutive parametrized ( $V$ ) and non-parametrized entangling ( $W$ ) layers. We define  $U_l(\theta_l) = \exp(-i\theta_l V_l)$ , where  $V_l$  is a Hermitian operator and  $W_l$  is a generic unitary operator. Then, the curcuit ansatz can be expressed with a multiplication of layers,

$$U(\boldsymbol{\theta}) = \prod_{l=1}^L U_l(\theta_l) W_l \quad (\text{A7})$$

For an observable  $O$  and an input state  $\rho$ , the cost function is given as

$$C(\boldsymbol{\theta}) = \text{Tr}[OU(\boldsymbol{\theta})\rho U^\dagger(\boldsymbol{\theta})] \quad (\text{A8})$$

The ansatz can be separated into two parts to investigate a certain layer, such that  $U_- \equiv \prod_{l=1}^{j-1} U_l(\theta_l)W_l$  and  $U_+ \equiv \prod_{l=j}^L U_l(\theta_l)W_l$ . Then, the gradient of the  $j^{\text{th}}$  parameter can be expressed as [6]

$$\partial_j C(\boldsymbol{\theta}) = \frac{\partial C(\boldsymbol{\theta})}{\partial \theta_j} = i \text{Tr}[[V_j, U_+^\dagger O U_+] U_- \rho U_-^\dagger] \quad (\text{A9})$$

Then the expected value of the gradient can be computed by using the Haar integral such that,

$$\langle \partial_j C(\boldsymbol{\theta}) \rangle = i \int d\mu(U_-) d\mu(U_+) \text{Tr}[[V_j, U_+^\dagger O U_+] U_- \rho U_-^\dagger] \quad (\text{A10})$$

$$= \frac{i \text{Tr}[\rho]}{d} \int d\mu(U_+) \text{Tr}[[V_j, U_+^\dagger O U_+] ] = 0, \quad (\text{A11})$$

where we use Eq. (A4) to obtain (A10) and use the fact that trace of the commutator is zero in (A11). This proves that the gradients are centered around zero. Then, the variance of the gradient can inform us about the size of the gradients. The variance is defined as,

$$\text{Var}[\partial_j C(\boldsymbol{\theta})] = \langle (\partial_j C(\boldsymbol{\theta}))^2 \rangle - \langle \partial_j C(\boldsymbol{\theta}) \rangle^2 = \langle (\partial_j C(\boldsymbol{\theta}))^2 \rangle \quad (\text{A12})$$

We can compute the expected value of the variance using the same logic. Then we have,

$$\text{Var}[\partial_j C(\boldsymbol{\theta})] = - \int d\mu(U_-) d\mu(U_+) \text{Tr}[[V_j, U_+^\dagger O U_+] U_- \rho U_-^\dagger]^2 \quad (\text{A13})$$

$$= - \frac{1}{d^2 - 1} \left( \int d\mu(U_+) \text{Tr}[\rho^2] \text{Tr}[[V_j, U_+^\dagger O U_+]^2] + \text{Tr}[\rho^2] \text{Tr}[[V_j, U_+^\dagger O U_+]^2] \right) \quad (\text{A14})$$

$$+ \frac{1}{d(d^2 - 1)} \left( \int d\mu(U_+) \text{Tr}[\rho^2] \text{Tr}[[V_j, U_+^\dagger O U_+]^2] + \text{Tr}[\rho^2] \text{Tr}[[V_j, U_+^\dagger O U_+]^2] \right) \quad (\text{A15})$$

$$= - \left( \text{Tr}[\rho^2] - \frac{1}{d} \right) \frac{1}{d^2 - 1} \int d\mu(U_+) \text{Tr}[[V_j, U_+^\dagger O U_+]^2] \quad (\text{A16})$$

We use Eq. (A6) to obtain Eq. (A14). Then, use the fact that commutator being traceless to obtain Eq. (A16). To compute the integral of Eq. (A16) we need another identity such that [10],

$$\text{Tr}[[V_j, U_+^\dagger O U_+]^2] = 2\text{Tr}[U_+ V_j U_+^\dagger O U_+ V_j U_+^\dagger O] - 2\text{Tr}[U_+ V_j^2 U_+^\dagger O^2]. \quad (\text{A17})$$

Then, the variance becomes,

$$\text{Var}[\partial_j C(\boldsymbol{\theta})] = - \left( \text{Tr}[\rho^2] - \frac{1}{d} \right) \frac{2}{d^2 - 1} \left( \int d\mu(U_+) \text{Tr}[U_+ V_j U_+^\dagger O U_+ V_j U_+^\dagger O] + \int d\mu(U_+) \text{Tr}[U_+ V_j^2 U_+^\dagger O^2] \right). \quad (\text{A18})$$

The first integral can be computed using Eq. (A5) and the second can be computed using Eq. (A4). Then we obtain,

$$\begin{aligned} \text{Var}[\partial_j C(\boldsymbol{\theta})] &= - \left( \text{Tr}[\rho^2] - \frac{1}{d} \right) \frac{2}{d^2 - 1} \left( \frac{1}{d^2 - 1} (\text{Tr}[V]^2 \text{Tr}[O]^2 + \text{Tr}[V^2] \text{Tr}[O^2]) \right. \\ &\quad \left. - \frac{1}{d(d^2 - 1)} (\text{Tr}[V]^2 \text{Tr}[O]^2 + \text{Tr}[V^2] \text{Tr}[O^2]) - \frac{1}{d} \text{Tr}[V^2] \text{Tr}[O^2] \right) \\ &= - \left( \text{Tr}[\rho^2] - \frac{1}{d} \right) \frac{2 \text{Tr}[V^2] \text{Tr}[O^2]}{d^2 - 1} \left( \frac{d - 1}{d(d^2 - 1)} (1 + \text{Tr}[V]^2 \text{Tr}[O]^2) - \frac{1}{d} \right). \quad (\text{A19}) \end{aligned}$$

Finally, the asymptotic behaviour of the variance can be expressed as

$$\text{Var}[\partial_j C(\boldsymbol{\theta})] \approx \mathcal{O}\left(\frac{1}{d^6}\right) \approx \mathcal{O}\left(\frac{1}{2^{6N}}\right), \tag{A20}$$

where  $d = 2^N$ . Thus, the variance vanishes exponentially with respect to N.

## Appendix B

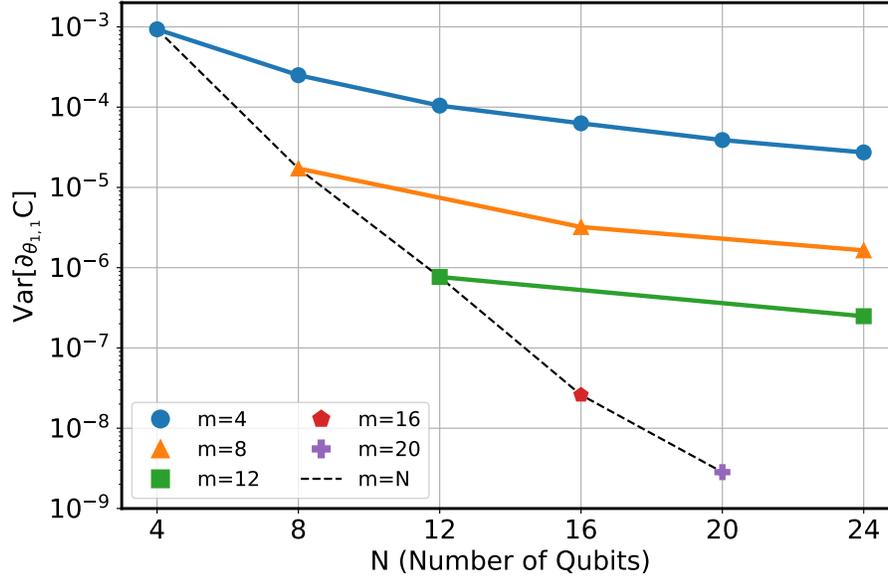


FIG. 6. The variance of the gradients of the first parameter of the ansatz as a function of the number of qubits for varying values of  $m$ . Each color/marker represents a certain value of  $m$  and data points of the standard ansatz ( $m = N$ ) is plotted with a dashed black line.

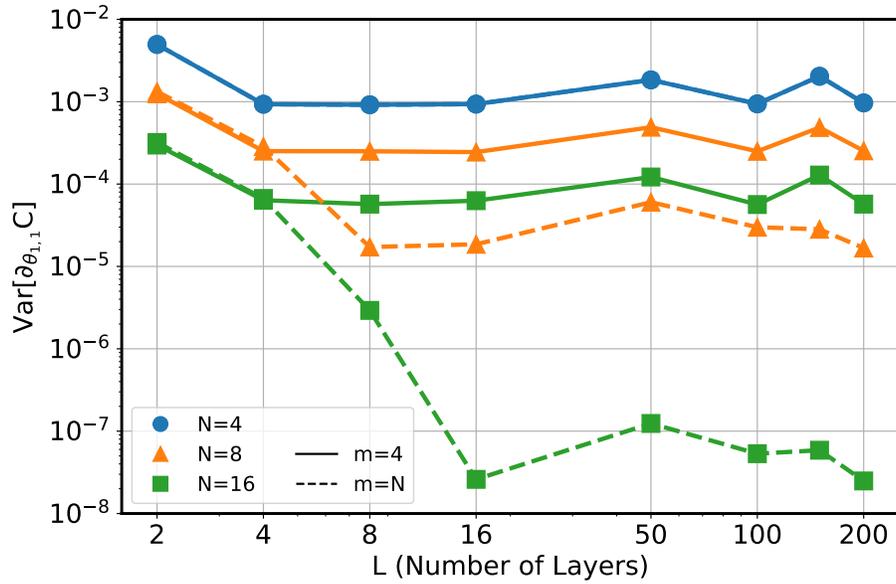


FIG. 7. The log plot of variance of the gradients of the first parameter of the ansatz vs. number of layers for  $m = 4$  (solid lines) and  $m = N$  (dashed lines) with varying number of qubits.

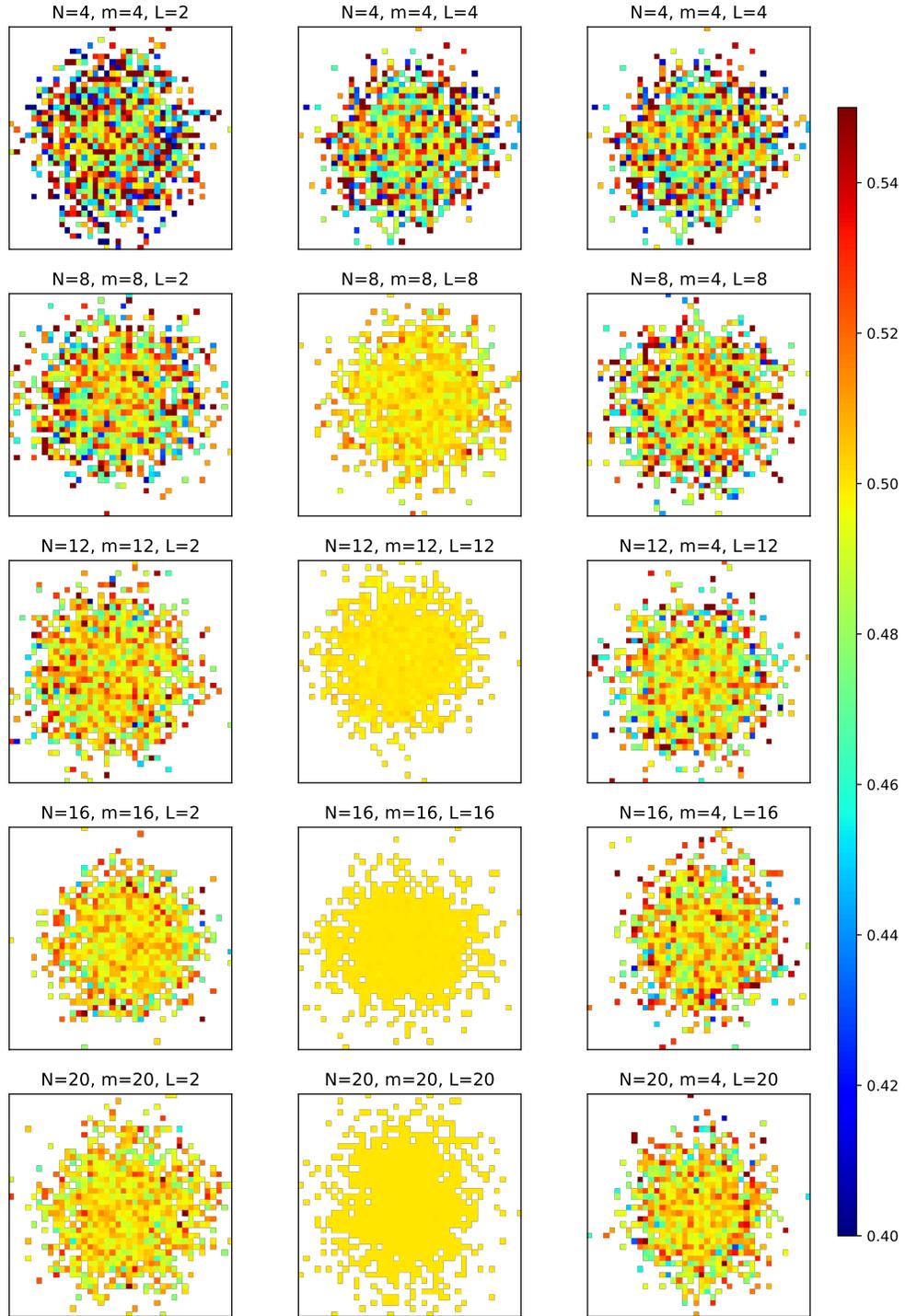


FIG. 8. Cost landscapes of ansätze with different settings. Parameters of the ansatz are reduced down to two using PCA and the  $x$  and  $y$  axis of the plots represents the PCA variables in same scale but with arbitrary units. The cost values (shown with the colormap) are obtained using the definitions in Section III B. First column shows cost values of an  $L = 2$  standard ansatz for increasing number of qubits. Second column shows the results for the same ansatz but with  $L = N$  layers. As, expected the landscape flattens with more qubits and we see a single color for  $N > 12$ . Third column shows results for splitting (for  $m = 4$ ) of the ansatz in the case of  $L = N$  layers. We see that the landscape does not become flatter with more qubits.

## Appendix C

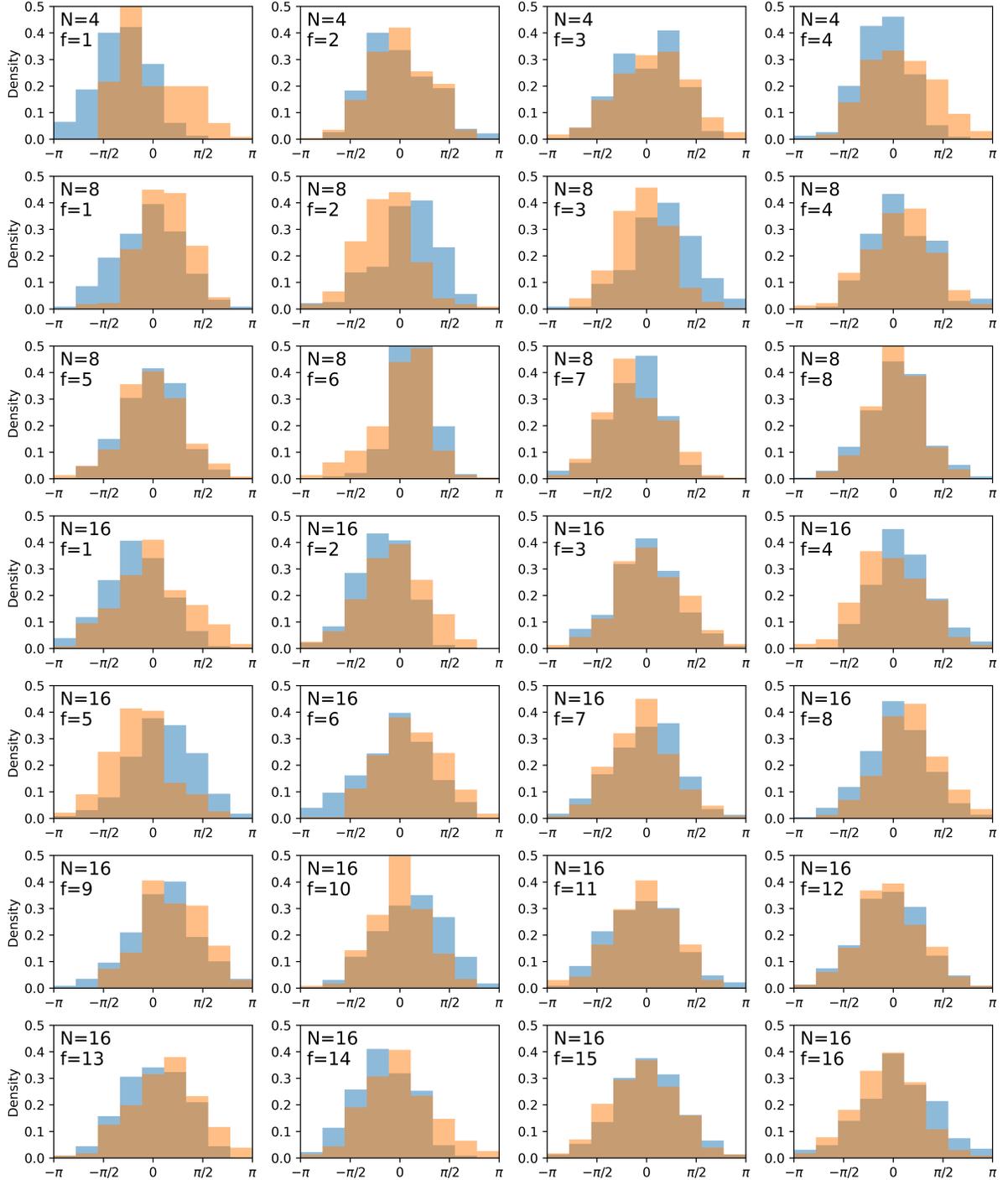


FIG. 9. Distributions of the ad-hoc dataset used in Section III B. Each panel shows distribution of a single feature from one of three datasets.  $N$  denotes the size of the dataset (number of features), while  $f$  denotes the feature number. There exists 600 samples of  $N$  features for a size  $N$  dataset. Colors represent two classes. During training, data samples are divided with a 420/180 train/test ratio. The dataset is produced using `make_classification` function of scikit-learn [40] with a class separation value of 1.0, 2% class assignment error and no redundant or repeated features.

## Appendix D

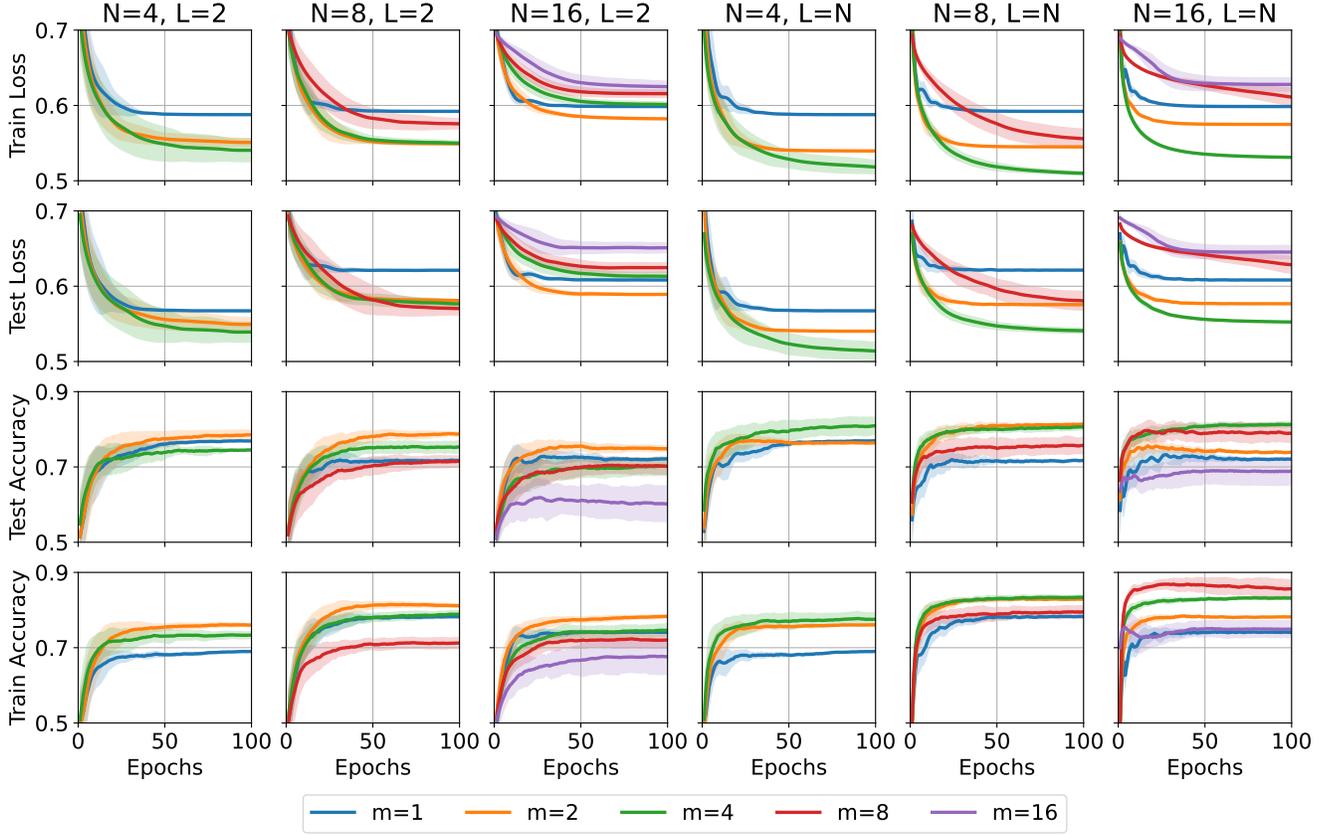


FIG. 10. Training curves showing four different metrics for the problem described in Section III B. Panels of each row show a different metric. First three columns show training results from  $L = 2$  ansätze, the last three columns show training results from  $L = N$  ansätze for  $N \in \{4, 8, 16\}$ . Each value of  $m$  is plotted with a different color. Lines are obtained by averaging 50 runs and their standard deviation is shown with shades.

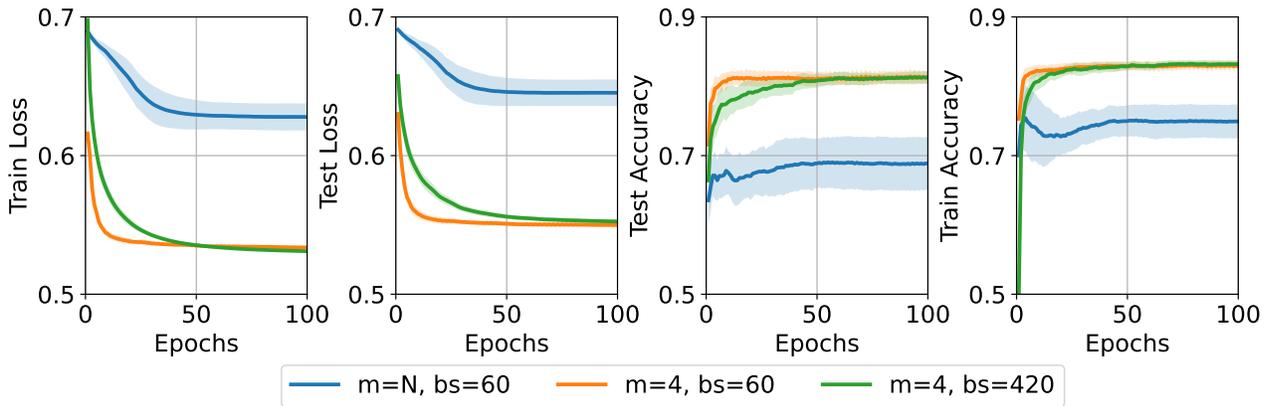


FIG. 11. Batch size comparison for the training of  $N = 16$ ,  $m = 16$  and  $m = 4$ . Training the  $N = L = 16$  model requires vast computational resources, especially memory. This restricted us from using a full batch size during the training of  $N = m = L = 16$  setting. Therefore, we presented results from a training that used a batch size of 60 instead of 420 (full). Here, we show training curves for  $m = 4$  on addition to  $m = 16$  for two different batch size (bs). Behaviour of the curves show that the gain in performance has nothing to do with the batch size difference.

## Appendix E

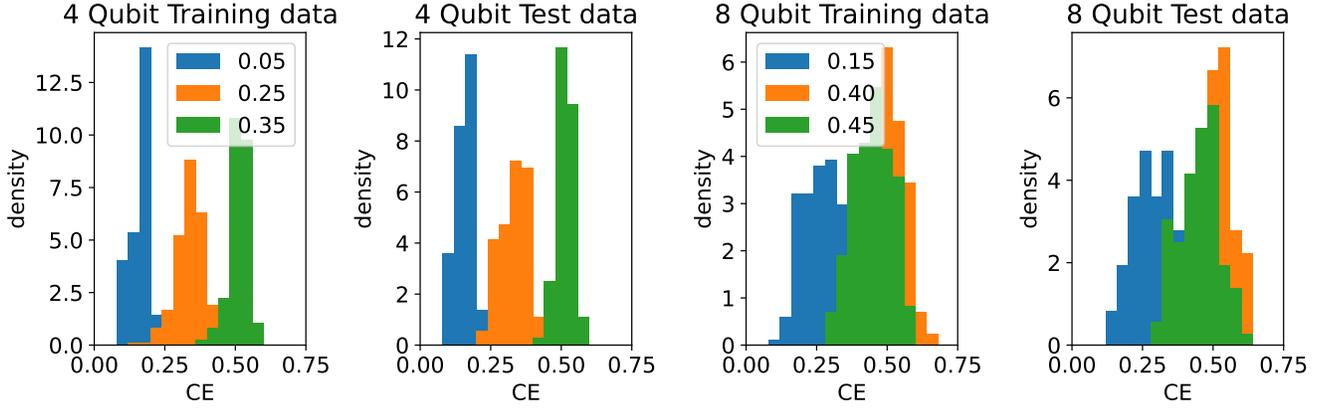


FIG. 12. Distributions of the NTangled [36] dataset with respect to the  $CE$  values described in Section III C. The HEA ansatz (Fig. 1d) is used to produce the distributions. Each training set has 420 and each test set has 180 data samples. We see a mismatch for  $CE \in \{0.40, 0.45\}$  in the 8 qubit case. We are not sure what causes this, but it is not an issue for our problem as we are not interested in the  $CE$  values themselves but the quantum states as a whole. So, they are valid quantum state distributions as long as they can be separated with a given metric for our problem. Our results show that this is in fact true.

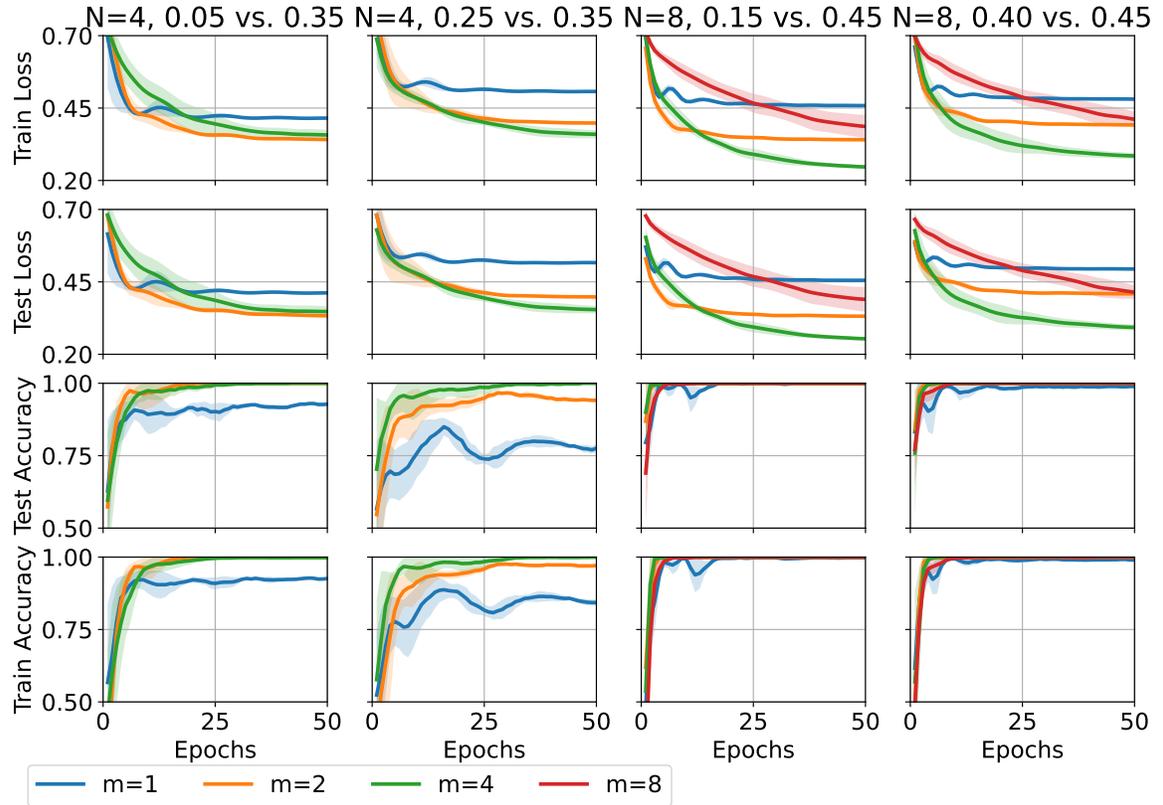


FIG. 13. Training curves showing four different metrics for the problem described in Section III C. Panels of each row show a different metric. Each column presents a different task, where  $N$  determines the problem size and the  $CE$  values are the labels of the classes. Each value of  $m$  is plotted with a different color. Lines are obtained by averaging 50 runs and their standard deviation is shown with shades.

## Appendix F

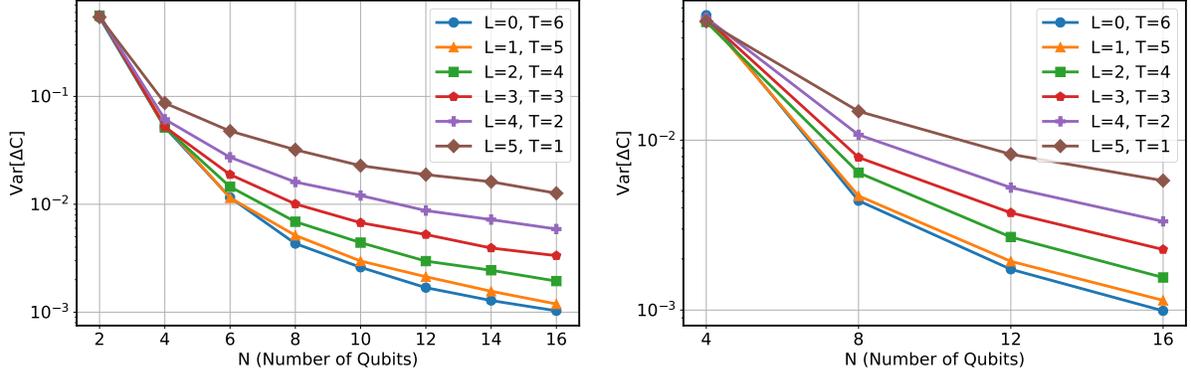


FIG. 14. The variance of the change in cost as a function of the number of qubits for varying values of  $L$  and  $T$  for  $L+T = D = 6$ . The cost function is the TFIH hamiltonian defined in Eq. 17 and the ansatz is the ECS ansatz with EfficientSU2 subblocks (see Fig. 1). Each line depicts a different value of  $T$ . The left panels show results for  $m = 2$  and the right panel shows results for  $m = 4$ . Variances are obtained over 2000 cost samples, where  $\Delta C$  is the difference of any arbitrary two cost samples.

## Appendix G

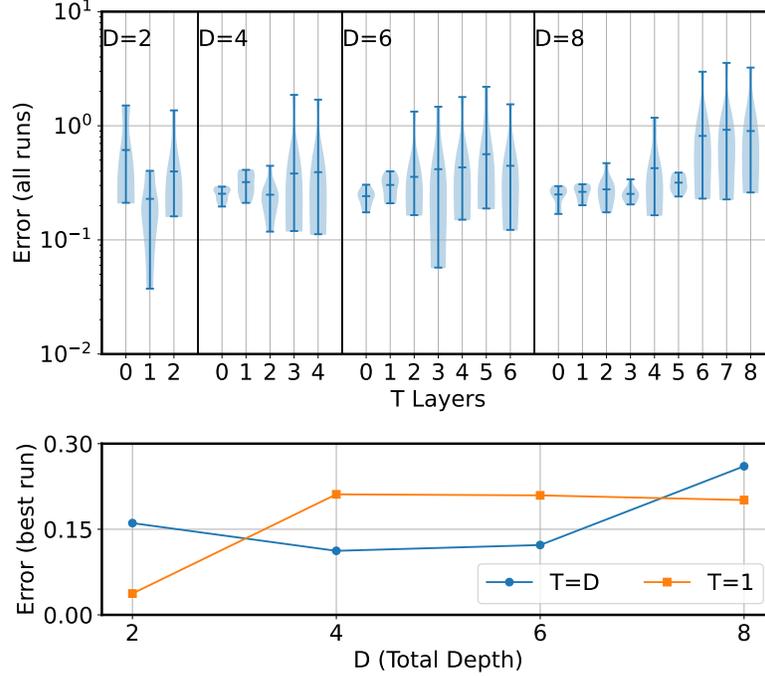


FIG. 15. Energy errors of ansätze with increasing total depth for  $N = 12$  TFIH using the extended classical splitting (ECS) ansatz with EfficientSU2 subblocks (see Fig. 1b) and  $m = 2$ . Total depth ( $D$ ) corresponds to  $L+T$ , where  $T = 0$  is equivalent to the CS ansatz,  $T = D$  is equivalent to standard EfficientSU2 and other values explore hybrid use cases of the ECS ansatz. Final energy measurements of 10 runs are averaged and plotted with their minimum and maximum values as the error bars on the upper panel. The lower panel shows the best errors obtained for  $T = 1$  and  $T = D$  settings. Energy error is the absolute difference of the energy measurement and the exact ground state energy.

## Appendix H

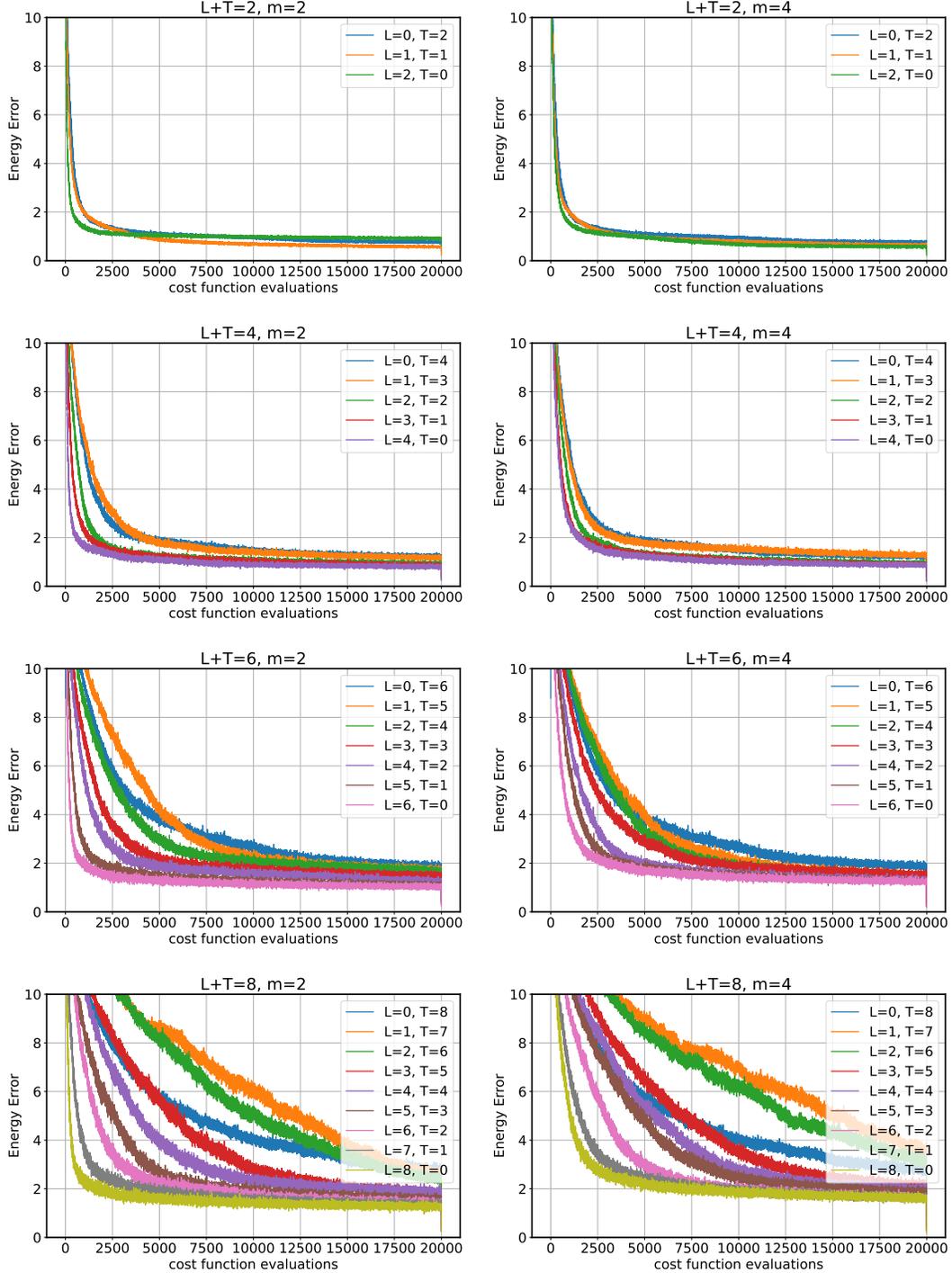


FIG. 16. Energy error curves of ansätze TFIH for  $N = 12$  using the extended classical splitting (ECS) ansatz with EfficientSU2 subblocks (see Fig. 1b). Columns corresponds to ansätze with  $m = 2$  and  $m = 4$  respectively. Each row shows results with increasing total depth ( $D$ ), such that  $L + T = D$ . Energy errors of 10 runs are averaged and their mean is presented. Energy error is the absolute difference of the energy measurement and the exact ground state energy. It becomes harder to optimize an ansatz with no classical splitting ( $L = 0$ ) as depth increases. However, we see that the optimization does not get as hard if we set  $T$  to a small value, e.g. to 1, and employ classical splitting. We observe similar conclusions with  $m = 2$  and  $m = 4$ .