

# RESPONSE LETTER

We would like to gratefully thank the Reviewers for their thorough study and valuable feedback which have been very helpful for us to improve the quality of our manuscript. Below is an itemized reply to each point raised by the Reviewers. The comments are reproduced verbatim in **black**, and the responses are presented in **blue**. We hope each query is addressed to the Reviewers' satisfaction.

## Reply to Reviewer 1 (Hao Wang):

This manuscript proposes distributing the computation of sub-gradients of the loss function of VQAs to multiple QPUs to speed up the wall clock time. As a straightforward approach applied extensively in classical machine learning, this idea has not been tested in QML.

The authors show the anticipated wall-clock time speedup with simulated gate noises on a straightforward classification problem using a shallow hardware-efficient ansatz. Also, the authors prove the convergence speed of the gradient-based optimizer in the distributed setting has the same upper bound as in the traditional VQAs.

Moreover, it is empirically illustrated that when the noise distributions differ significantly across local QPUs, the gradient-based optimizer's convergence speed is hampered quite a bit, for which the authors propose a rotation rule for averaging out the discrepancy of the noises across different QPUs.

Reply: We gratefully acknowledge the Reviewer for his comprehensive summary of the manuscript.

1. I think the theoretical part is not very satisfactory, primarily based on Ref. [44], as you pointed out in the paper. It is crucial to look at the structure of the noise closely, i.e., the covariance of the estimated partial derivatives, i.e.,  $Cov\left(\frac{\partial \hat{L}}{\partial \theta_i}, \frac{\partial \hat{L}}{\partial \theta_j}\right)$  ("hat" means it is estimated with parameter-shift rule), due to quantum noises. In the traditional, centralized setup, this covariance is often used to mitigate errors (at least measurement errors) for the gradient-based methods. In the distributed setup, I presume that the covariance admits a block-diagonal structure (since local QPUs are independent), which should play a role in your analysis.

Reply: We thank the Reviewer for raising this very important question.

The covariance of the estimated partial derivatives does play an important role in mitigating error on traditional field. In particular, this method can usually be used for suppressing the error caused by using only part of the data instead of all of it during the

mini-batch learning [arXiv:2002.03979 (2020)], or enhancing the generalization ability of data-parallel training [PMLR 162:18347-18377 (2022), arXiv:2205.09305 (2022)]. However, the parameter-parallel strategy proposed in our work does not split the data and therefore does not encounter these problems.

Moreover, the error considered in our article is quantum noise, which is caused by factors such as decoherence of quantum devices. We know that characterizing and modeling quantum noise is extremely complex, and we choose the “worst-case” noise channel, *i.e.*, the depolarizing channels in our work, which not only describes circuit noise, but also portrays measurement error. In Theorem 1, we discuss the convergence of PPD-VQA under depolarizing noise. This theorem is the theoretical guarantee that PPD-VQA can accelerate the training of VQAs, not a tool for error mitigation.

As we mentioned in the Conclusion section, the PPD-VQA could be combined naturally with data-parallel training and the mini-batch learning. Thus, Reviewer 1' suggestions can indeed be considered further in these future combinations. We have included some discussion in the Conclusion section of the revised version for clarifying and providing the reader some suggestions for the future works. If the Reviewer has further comments (perhaps provides clearer references), we are happy to continue to revise to further improve the quality of the manuscript.

2. It is unclear what you meant by "convergence test." Is it nothing more than selecting a QPU u.a.r. and getting the accuracy score on the classification task?

**Reply:** A quick answer is yes. The convergence test refers to selecting a QPU and getting the accuracy score on the classification task.

The convergence test is necessary for two reasons: 1) It provides a criterion to determine when PPD-VQA can be stopped. In our example, the training of PPD-VQA terminates when the classification accuracy on training set reaches a certain percentage. This criterion could be changed according to specific machine learning task. 2) The ultimate goal of PPD-VQA is to obtain a high-performance trained model, which is composed of a quantum hardware and a trained parameterized quantum circuit. By implementing the convergence test, we can monitor the performance of the trained parameterized quantum circuit on the chosen QPU (as each QPU has different noise, it seems best to choose a fixed QPU) to ensure that the final parameterized quantum circuit will perform well on the chosen QPU.

To make this part clear, we have added some explanations in **step 2** on page 2.

3. IMHO, the gradient compression part is superfluous, which is indeed a bottleneck in deep learning, where we have to face millions of parameters. However, in VQAs, we cannot afford that, correct? due to the current limitation of the hardware implementation and also to the theory that says barren plateaus will kick us out of the game if we have  $O(\text{poly}(n))$  layers ( $n$  is the number of qubits). Therefore, I do not think gradient compression is needed. (Yes, in table I, you show some compression results; but I think you are saving the communication cost from a small overhead).

**Reply:** We thank the Reviewer for raising this very important point which indeed merits further discussions. We still maintain the view that gradient compression is useful for PPD-VQA, for the following two reasons, especially the second one:

First, we agree with the Reviewer that the current distributed VQA is indeed immune with the problem of communication bottlenecks. However, for the future deep parameterized quantum circuit and large-scale datasets, the number of parameters may explode, leading to the role of parameter compression coming to the fore. Thus, the gradient compression strategy proposed here is in response to potential future possibilities.

Second, gradient compression may help to mitigate errors in the experimental implementation of the PPD-VQA, due to the following two reasons: 1) For most quantum computing systems, it is not easy to implement the tiny angular rotations of single-qubit quantum operations with high precision. Gradient compression can avoid updates of tiny angles and thus potentially improve experimental accuracy. 2) As the frequency of updating parameters (especially those with small gradient changes) is reduced, the number of gate operations that need to be changed by the quantum device is consequently reduced and the accumulation of quantum operation errors is naturally suppressed.

We added these discussions to the top right corner of page 8 of the revised version.

4. The most interesting aspect to me is Eq. (3), which writes down the biased of the estimated gradient on each local node/QPU. I think the author should investigate the relationship between the noise level and the bias term, which can drastically change the gradient direction if the noise is high and the magnitude of the sub-gradient is small.

**Reply:** As shown in the newly added Appendix B in the revised version, the relationship between the noise level and the bias term is

$$\text{bias term} \leq (2 + 9\lambda\pi)\tilde{p}_i,$$

where  $\tilde{p}_i$  is the depolarizing probability of  $i$ -th QPU, which can be calculated according to Eq. (1).

It can be observed that the larger the noise, the larger the upper bound on the bias term.

5. Also, in light of the above question, I wonder if the distributed scheme converges at all. Does the bias term also scale down with the diminishing sub-gradient/partial derivative when approaching the critical points on the quantum loss landscape? Otherwise, if the partial derivatives go to zero while the bias does not, we could have a serious problem. Maybe I overthink this part, as perhaps the bias is so tiny that it can be ignored. Please comment on this.

**Reply:** This is an interesting point, and we are happy to discuss it with the Reviewer.

We should note that although we could simply interpret the estimated partial derivative  $[\nabla L]_{i,j}$  under noise case as  $[\nabla \bar{L}]_{i,j} = [\nabla L]_{i,j} + \text{bias term}$ , where  $[\nabla L]_{i,j}$  is the ideal partial derivative without noise, the training process of VQA is to optimize the parameterized quantum circuit to make  $[\nabla \bar{L}]_{i,j}$  converge instead of  $[\nabla L]_{i,j}$ . Thus, what we really care about is whether the whole composed of  $[\nabla L]_{i,j}$  and *bias term*,  $[\nabla L]_{i,j} + \text{bias term}$ , can converge. Thus, from this perspective, independent consideration of the bias term may give us some information, but it does not directly correlate to whether the estimated partial derivative converges or not. For example, if a single-qubit gate in the parameterized quantum circuit rotates one degree more each time, this is a coherent noise, the training process of VQA could adaptively adjust to make the single-qubit gate rotate one degree less to compensate for this noise. Thus, in this case, the bias term would always be present, but the training of VQA would not be affected by the bias term.

To conclude, we would like to point that the proved convergence of PPD-VQA in **Theorem 1** is sufficient to address the Reviewer's query.

6. Do you contemplate any error mitigation approach on each local node?

**Reply:** We do not contemplate error mitigation approach on each local node. The reason is that we expect to explore the effect of noise during the training of PPD-VQA, and the introduction of error mitigation would complicate the problem.

We know that error mitigation techniques could allow us to reduce the computational errors and then evaluate accurate results from noisy quantum circuits. The combination of error mitigation techniques and PPD-VQAs is a potential future work, which has been discussed in the conclusion section in the page 8.

7. Please improve some usage of the language/jargon.

\* “...the estimates of the gradients of each parameter...” → the gradients/partial derivatives of the observable w.r.t. to each parameter.

**Reply:** Corrected.

\* I don't like the notation used in the expression under “Step 1”. It is more standard in math to express it by  $\theta_i^{(t)} = (\theta_{1+(i-1)n}^{(t)}, \theta_{2+(i-1)n}^{(t)}, \dots, \theta_{in}^{(t)})$ ,  $n = d/M$ .

**Reply:** Agreed and corrected.

\* Also, please define the loss function/expectation  $L$  first (should simply take one sentence).

**Reply:** We have defined the loss function  $L$  before theorem 1 in the revised version.

\*  $\nabla L_i(\theta^{(t)})$  is confusing/non-standard to me, which immediately implies you have a sub-function  $L_i$  of  $L$ . I assume this means the loss function is defined on a batch of data sets, correct? Please either provide an explicit definition thereof or use a more understandable notation. Provided that I understand it correctly, the gradient step should also be divided by the batch size, right? See equation  $\theta^{(t+1)} = \theta^{(t)} - \eta \sum_{ij} \dots$  at the very bottom of page 3.

**Reply:** We thank the Referee for pointing this out.

$\nabla L_i(\theta^{(t)})$  represents the sub-gradient which is acquired on  $i - th$  local node.

Accordingly, we have used  $[\nabla L(\theta^{(t)})]_{i,j}$  to represent the  $j - th$  component of sub-gradient acquired on  $i - th$  local node in the revised manuscript.

\* Please clarify if the random variables on the same QPU are independent, which is essential for the upper bound derived in the paper.

**Reply:** We thank the Referee for pointing this out. The random variables in different expressions of gradient components are independent, we have clarified that in the newly added Appendix A.

\* “The similar convergence rate guarantees that PPD-VQA promises an intuitive linear runtime speedup concerning the increased number of local nodes” → you mean the speedup of the computation of the gradient, right? Please be very clear here.

**Reply:** Yes, we refer to the speedup of the computation of gradient. We have rewritten this statement on Page 4.

\* “Even if the average noise of each quantum processor is the same, the noise environment of the qubits executing quantum circuits in different processors is unlikely to be consistent” → What do you mean by this statement? The hardware noise? But, then, why is this different from bullet point (1) you mentioned before this sentence? Also, please always provide references to such important messages.

**Reply:** We thank the Referee for pointing this out.

In general, the error rate  $\varepsilon_i$  of each qubit on a quantum processor is different, and we let the **average error rate** of the processor be  $\bar{\varepsilon} = \frac{1}{N} \sum_i \varepsilon_i$ , where  $N$  is number of qubits.

Based on this definition, the point (1) refers to the case that different QPUs have different **average error rate**, while the point (2) refers to the case that different QPUs have the **same average error rate**, but the **error rate per qubit** in the processor may also be different.

We have rewritten this statement at the end of the Page 2.

In all, given what I have observed, I feel the paper is not mature enough for SciPost Physics, considering the acceptance criterion of the journal:

- Detail a groundbreaking theoretical/experimental/computational discovery;
- Present a breakthrough on a previously-identified and long-standing research stumbling block;
- Open a new pathway in an existing or a new research direction, with clear potential for multipronged follow-up work;
- Provide a novel and synergetic link between different research areas.

**Reply:** We would like to thank the Reviewer for her/his helpful comments, which have helped us to further improve the quality of the manuscript. We believe that our work passes the bar for the acceptance criteria of *SciPost* for the following reasons:

VQA is a promising near-term technique to explore practical quantum advantage on near-term devices. However, the inefficient parameter training process due to the incompatibility with backpropagation and the cost of a large number of measurements, posing a great challenge to the large-scale development of VQAs.

Parallel training is a natural potential solution to this bottleneck. However, **due to the presence of quantum noise, it remains unknown whether parallel training is effective, especially in the presence of differences in noise across QPUs**. We not only prove the convergence of PPD-VQA, but also propose an efficient strategy, alternate strategy, to suppress the performance degradation caused by the difference in

noise across QPUs. Our work makes parallel training of VQA in realistic noisy environments feasible, and thus **open a new pathway in an existing research direction**. Moreover, the efficient parallel training can lead to many applications in distributed scenarios, such as distributed quantum machine learning and federal quantum machine learning. Moreover, the PPD-VQA has good compatibility with other distributed strategies such as data-parallel and error mitigation techniques, to further improve the practicality of VQA. Thus, our work could provide **potential for multipronged follow-up work**.

We hope that our point-by-point responses and concomitant changes to the manuscript make a convincing case now that our manuscript is suitable for publication in *SciPost*.

### **Reply to Reviewer 2:**

The author brings forward an algorithm for variational quantum computing in which they parallelize the evaluation of gradients of the parameters over different quantum processing units (QPUs). The authors fairly and clearly state that the idea in itself is not groundbreaking, however, what makes this work interesting and useful is that they evaluate the performance of such parallelization in the presence of noise. It could be in fact possible that the idea would not work well when each QPU is noisy, and in its own way.

The authors show that even in presence of noise there is a significant speed up from parallelization, and actually the noisy systems scale as well as the ideal one.

For the case of QPUs with significantly different noise, the authors showed also that one can rotate through the QPUs used for different set of parameters, doing so "averaging out" the different noise scenarios for each set of parameters and thus performing in a more consistent way.

They also evaluate the performance of gradient clipping, which they use to remove gradients which are too small and thus reducing the communication bottleneck. In this case they show that even reductions of 60% result in almost no change in the speed-up.

**Reply:** We gratefully acknowledge the Reviewer for his/her comprehensive summary of the manuscript.

Some comments:

1) When introducing the issues faced by VQA, the authors do not mention barren plateaus, which are then mentioned in the conclusions. Is there a particular reason for this choice?

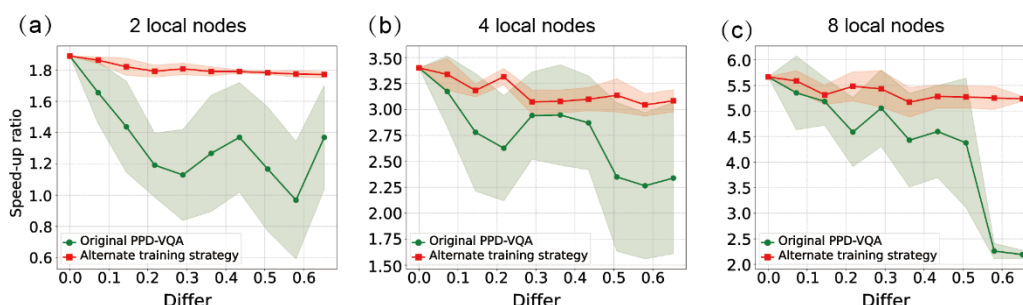
**Reply:** We thank the Referee for pointing this out.

The PPD-VQA we proposed is to accelerate the training process of VQA by using parallel training, but not to solve issue of the barren plateaus. To remove this

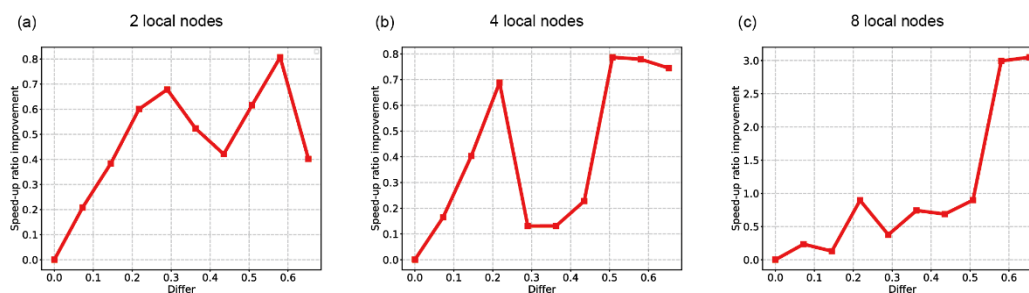
misinformation, we have removed the discussion of barren plateaus from the conclusion section.

2) Regarding the alternate training strategy, it seems to me that the more the setup is parallelized the closer becomes to the first training proposed. If this is true, it could imply that if one parallelizes more, he/she should not need the alternate training. What do the authors think about this?

**Reply:** The alternate training is to suppress the performance degradation caused by the difference in noise across QPUs. As shown in the Fig. 4 in the manuscript, as the number of nodes increases, the alternate training strategy does not show signs of failing to improve the PPD-VQA.



To make this point clearer, we subtract the speedup ratio obtained by not using alternate training from the speedup ratio obtained by using alternate training to get the following results:



It can be observed that the speedup ratio improvement from alternate training strategy does not decrease as the number of nodes increases.

3) Do the authors have an idea of why the speed-up seems to be unaffected by the presence of noise? The curves in Fig.3(b, c) are remarkably close to each other.

**Reply:** We are grateful to the Reviewer for reminding us that our Fig. 3(c) does have the potential to cause misunderstandings among readers.



Actually, the speed-up is unaffected by the presence of noise. As shown in the Fig. 3(b), the required number of iterations is increase, as the noise increases, while does not show variability as the number of local node changes, which may be caused by the proximity of noise levels of  $M$  local nodes.

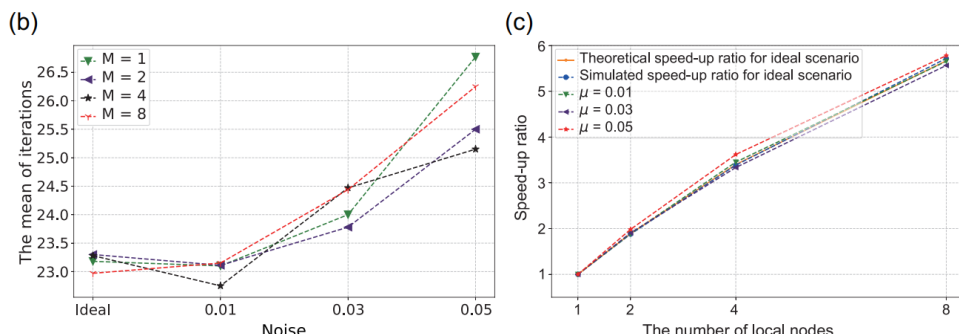


Fig. 3(c) is the result for speed-up ratio, which is ratio of the training speed of  $M$  nodes to that of 1 node for the same noise level. Based on the analysis for Fig.3(b) and the definition of speed-up ratio, PPD-VQA produces an approximate  $M$ -fold acceleration when using  $M$  local nodes in Fig.3 (c). The reason that the speed-up ratio is almost independent of noise is that the training speed of 1 nodes and  $M$  nodes slows down at the same time as the noise grows (as shown in Fig. 3(b)). We have added explanations to the revised version to remove this possible misunderstanding.

4) It is not clear to me how to compare the values of the noise used with those found on QPUs like Zuchongzhi. Can the authors say a bit more about this so that we can understand better when such study would be really applicable to available QPUs?

**Reply:** We thank the Referee for pointing this out.

It is known that characterizing and modeling quantum noise in the real quantum device is extremely complex, and we choose the “worst-case” noise channel—the depolarizing channels in our work. The noise value in the Fig. (3) is the depolarizing probability for single-qubit gate.

The average single-qubit gate Pauli error of *Zuchongzhi* is 0.14%, which is calibrated using XEB. The following formula [Nature, 574(7779), 505-510 (2019)] can be used to convert depolarizing probability to Pauli errors for single-qubit gate,

$$e_p = \frac{3}{4}p,$$

where  $e_p$  is single-qubit gate Pauli error. We have provided some discussions in the caption of Figure 3.

5) Regarding the classification task, I think that more information could be provided to the broad audience who could read this article. For instance, how are the data encoded classically? How is it converted to be used on the quantum machine? More clarity could

be given on the number of parameters used. There is some description in the caption, but I think that it could be expanded in the main text.

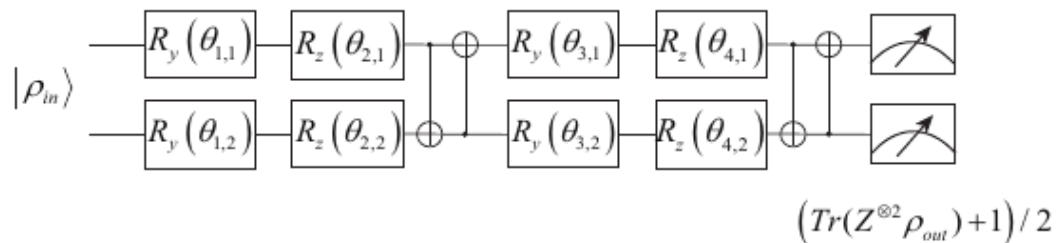
**Reply:** We are very grateful to the Reviewer for pointing this out, and it does need to be explained in more detail for the broad audience.

Each sample in IRIS dataset has 4 attributes. Assume that  $(a, b, c, d)$  represents a sample in IRIS, we can transform a sample into 2-qubits quantum states by amplitude encoding, i.e.,

$$|\varphi\rangle = \frac{a|00\rangle + b|01\rangle + c|10\rangle + d|11\rangle}{\sqrt{a^2 + b^2 + c^2 + d^2}}.$$

Regarding the number of parameters, we use a hardware-efficient ansatz with 2 layer (as shown in the Fig. 2(b), each layer has 2  $R_y$  parameterized gates,  $R_z$  parameterized gates), and thus there are 8 parameters in total.

(b)



We have provided these descriptions on page 4 of the revised version.

6) Fig.1. For  $M - th$  QPU, I think that the portion that is optimized should be color-coded in red, but it is in orange.

**Reply:** Corrected.

7) Fig.3(a). It is not clear to me what the various error bars, colored areas and circles mean. More detail and explanation could help.

**Reply:** We thank the Reviewer for raising this point which needs more explanations. The middle line of the boxes, which is the median of the data, represents the average of the sample data. The upper and lower limits of the box, which are the upper and lower quartiles of the data, respectively, which means that the box contains 50% examples of the data. Therefore, the width of the box reflects to some extent the degree of fluctuation of the data. Above and below the box, there are other lines each representing the maximum and minimum values and the circles represent outliers.

We have added this explanation in the caption of Fig. 3.

8) Fig.4. It is not clear to me what the shaded area exactly represents.

**Reply:** The shaded area represents statistical error caused by 100 independent experiments. We have added this explanation in the caption of Fig. 4.

9) The English would need a bit of polishing, as there are a lot of minor issues, but this is something that can be done at a later stage.

**Reply:** We have polished the English in our revised version.

To conclude, I do not think that the paper is ready to be published in Scipost, and I think that the authors may have to present a more convincing case for the paper to be published in Scipost Physics instead of Scipost Physics Core. In fact I am not convinced that the paper satisfy at least one of the 4 criteria

1. Detail a groundbreaking theoretical/experimental/computational discovery;
2. Present a breakthrough on a previously-identified and long-standing research stumbling block;
3. Open a new pathway in an existing or a new research direction, with clear potential for multipronged follow-up work;
4. Provide a novel and synergetic link between different research areas.

**Reply:** We would like to thank the Reviewer for her/his helpful comments, which have helped us to further improve the quality of the manuscript. We believe that our work passes the bar for the acceptance criteria of *SciPost* for the following reasons:

VQA is a promising near-term technique to explore practical quantum advantage on near-term devices. However, the inefficient parameter training process due to the incompatibility with backpropagation and the cost of a large number of measurements, posing a great challenge to the large-scale development of VQAs. Parallel training is a natural potential solution to this bottleneck. However, **due to the presence of quantum noise, it remains unknown whether parallel training is effective, especially in the presence of differences in noise across QPUs.** We not only prove the convergence of PPD-VQA, but also propose an efficient strategy, alternate strategy, to suppress the performance degradation caused by the difference in noise across QPUs. Our work makes parallel training of VQA in realistic noisy environments feasible, and thus **open a new pathway in an existing research direction.** Moreover, the efficient parallel training can lead to many applications in distributed scenarios, such as distributed quantum machine learning and federal quantum machine learning. Moreover, the PPD-VQA has good compatibility with other distributed strategies such as data-parallel and error mitigation techniques, to further improve the practicality of VQA. Thus, our work could provide **potential for multipronged follow-up work.**

We hope that our point-by-point responses and concomitant changes to the manuscript make a convincing case now that our manuscript is suitable for publication in *SciPost*.

# Parameter-Parallel Distributed Variational Quantum Algorithm

Yun-Fei Niu,<sup>1,\*</sup> Shuo Zhang,<sup>1,†</sup> Chen Ding,<sup>1</sup> Wan-Su Bao,<sup>1,‡</sup> and He-Liang Huang<sup>1,2,3,4,§</sup>

<sup>1</sup>Henan Key Laboratory of Quantum Information and Cryptography, Zhengzhou, Henan 450000, China

<sup>2</sup>Hefei National Research Center for Physical Sciences at the Microscale and School of Physical Sciences, University of Science and Technology of China, Hefei 230026, China

<sup>3</sup>Shanghai Research Center for Quantum Science and CAS Center for Excellence in Quantum Information and Quantum Physics, University of Science and Technology of China, Shanghai 201315, China

<sup>4</sup>Hefei National Laboratory, University of Science and Technology of China, Hefei 230088, China

(Dated: January 16, 2023)

Variational quantum algorithms (VQAs) have emerged as a promising near-term technique to explore practical quantum advantage on noisy intermediate-scale quantum (NISQ) devices. However, the inefficient parameter training process due to the incompatibility with backpropagation and the cost of a large number of measurements, posing a great challenge to the large-scale development of VQAs. Here, we propose a parameter-parallel distributed variational quantum algorithm (PPD-VQA), to accelerate the training process by parameter-parallel training with multiple quantum processors. To maintain the high performance of PPD-VQA in the realistic noise scenarios, an alternate training strategy is proposed to alleviate the acceleration attenuation caused by noise differences among multiple quantum processors, which is an unavoidable common problem of distributed VQA. Besides, the gradient compression is also employed to overcome the potential communication bottlenecks. The achieved results suggest that the PPD-VQA could provide a practical solution for coordinating multiple quantum processors to handle large-scale real-world applications.

## I. INTRODUCTION

Quantum computing holds the promise of solving certain problems that intractable for classical computers, such as factoring large numbers [1–3], database search [4, 5], solving linear systems of equations [6–8]. However, a universal fault-tolerant quantum computer that can solve efficiently the above problems would require millions of qubits with low error rates [9, 10], which is still a long way from current techniques and may take decades. Thus, we will be in the noisy intermediate-scale quantum (NISQ) era for a long time [11–16]. Variational quantum algorithms (VQAs) leverage a quantum device to minimize a specific cost function [17, 18], by employing a classical optimizer (e.g., Adam optimizer [19]) to train parameter quantum circuits (PQCs). Such algorithms were shown to have natural noise resilience [20] and even benefit from noise, making it particularly suitable for near-term quantum devices, and thus be considered the most promising path to quantum advantage on practical problems in NISQ era [18]. Previous studies have exhibited the application of VQAs on a variety of problems, including classification task [21–24] and generative task [25–27], combinatorial optimization [28–32], quantum many-body problem [33] and quantum chemistry [34–39].

The training process of VQAs is actually not very efficient compared to the classical neural network, due to the following two main reasons: 1) The quantum state of the intermediate process of the quantum circuit cannot be stored, making VQAs impossible to use the backpropagation to update the parameters as efficiently as the classical neural network; 2) A

large number of measurements is required for the result read-out of the quantum circuit, which is time-consuming. Therefore, the training of VQAs will face significant challenges, as the amount of data and trainable parameters increases.

To address the above issue, a distributed VQA based on data-parallel has been proposed by Du *et. al.* to accelerate the training of VQA [40]. In this work, a parameter-parallel distributed variational quantum algorithm (PPD-VQA) is proposed to further accelerate the training process by parameter-parallel training with multiple quantum processors. Although the idea of parallel training is not difficult to come up with, including data-parallel or parameter-parallel, it is worth investigating whether the approach works in the realistic scenario that the local quantum nodes will inevitably be affected by quantum noise, and the noise intensity of each node is different. We first proof the convergence of the PPD-VQA, even if each local node has different quantum noise. Further, we design an alternate training strategy to alleviate the acceleration attenuation caused by excessive noise differences among multiple quantum processors, and adopt the gradient compression to cut a large amount of communication bandwidth, to enhance the practicality and scalability of PPD-VQA.

## II. PPD-VQA

The conventional VQAs employ PQCs and update their parameters  $\theta$  via a classical optimization training procedure, to find the global minimum of the given loss functions  $L$ . Usually, in the training procedure, the gradients of each parameter is evaluated by the parameter-shift rule [41, 42]. The PPD-VQA leverages the fact that **the partial derivatives of the observable with respect to each parameter** are genuinely independent of one another at each iteration to accelerate the training of conventional VQA, by parallelizing the gradient estimation across multiple quantum processing unit (QPU) nodes.

\* These two authors contributed equally

† These two authors contributed equally; [shuoshuo19851115@163.com](mailto:shuoshuo19851115@163.com)

‡ [bws@qiclab.cn](mailto:bws@qiclab.cn)

§ [quanhhl@ustc.edu.cn](mailto:quanhhl@ustc.edu.cn)

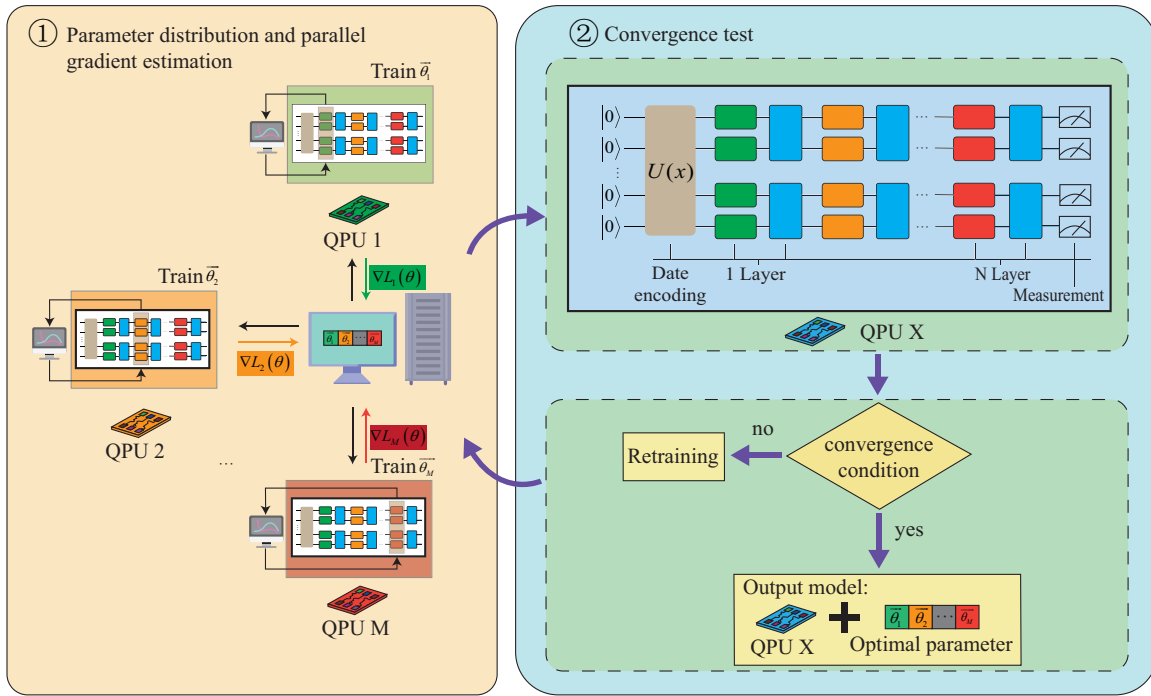


FIG. 1. **Schematic diagram of PPD-VQA.** The diagram illustrates two main steps of the PPD-VQA workflow. Firstly, the central parameter server allocates the trainable parameters to  $M$  local nodes, consisting of a QPU and a classic computer, for parallel training. Each local node only trains a part of the trainable parameters, and synchronizes the gradient information to the central parameter server. Secondly, a local node, named QPU X, is selected to verify that the convergence condition is met. If it does not converge, repeat Steps 1 and 2, otherwise, output the optimal parameters and selected local node as the trained model.

Conceptually, a classical central parameter server and  $M$  local nodes constitute the framework of PPD-VQA, where each local node consists of a QPU and a classical optimizer. As shown in Fig. 1 and Algorithm 1, at each iteration, the central parameter server divides the trainable parameters  $\theta$  into  $M$  parts, each of the  $M$  local nodes is tasked with computing the gradient of the parameters for a given component. Then, the complete gradient information is obtained through information sharing between local nodes and central parameter sever, which is used to update the trainable parameters as the initial parameters of next iteration. This process is repeated until the optimal parameters are found. The specific process can be divided into the following two steps:

**Step 1: Parameter distribution and parallel gradient estimation.** At the beginning of  $t$ -th iteration, the classical central server distribute the complete parameter  $\theta^{(t)}$  of PQC to each local node as the initial parameters, as well as instructions on which parameters the  $i$ -th local node is assigned for training. The default instruction is to divide the trainable parameters  $\theta^{(t)}$  into  $M$  equal parts

$$\theta^{(t)} = [\theta_1^{(t)}, \dots, \theta_M^{(t)}], \theta_i^{(t)} = (\theta_{i,1}^{(t)}, \theta_{i,2}^{(t)}, \dots, \theta_{i,n}^{(t)}), n = \frac{d}{M},$$

and the  $i$ -th local node is responsible for estimating the corresponding component of gradient  $\nabla L_i(\theta^{(t)})$ . After the training on each node, the local node synchronizes  $\{\nabla L_i(\theta^{(t)})\}_{i=1}^M$  to the central parameter server, and the central parameter server combines the information from each lo-

cal node into a completed gradient  $\nabla L(\theta^{(t)})$  used to update  $\theta^{(t)}$  to  $\theta^{(t+1)}$ .

**Step 2: Convergence test.** Choose a **fixed** local node from the  $M$  local nodes, and substitute the parameters  $\theta^{(t+1)}$  into this local node. After that, employ this model to the test dataset to to determine whether the convergence condition has been met. The setting of the convergence condition depends on the machine learning task. For example, for the classification task, the convergence condition might be set as a certain classification accuracy threshold. **The convergence test refers to selecting a QPU and getting the accuracy score on the classification task.** If the convergence condition has not been met, return to Step 1 for the subsequent training iteration; otherwise, output the final parameters and chosen local node as the trained model and terminate the training procedure. **By implementing the convergence test, we can monitor the performance of the trained PQC on the chosen QPU to ensure that the final PQC will perform well on the chosen QPU.**

The core idea of PPD-VQA is simple and natural. However, distributed quantum machine learning faces different challenges than its classical counterpart, the main one being that the quantum processors on different local nodes are not identical, due to the inevitable quantum noise. **In general, the error rate  $\varepsilon_i$  of each qubit on a quantum processor is different, and we let the average error rate of the processor be  $\bar{\varepsilon} = \frac{1}{N} \sum_i \varepsilon_i$ , where  $N$  is number of qubits. Thus, the non-uniformity mentioned above manifests itself in two ways: 1) The average error**

ror rates of each quantum processors are different. For example, some processors have lower noise and some have higher noise; 2) Even if the average error rate of each quantum processor is the same, the error rate of each qubit in these processors is unlikely to be consistent. In such a realistic scenario, it remains to be verified whether the parameter-parallel training is still effective, and whether the convergence conditions can be achieved. This important issue is directly related to the practical utility of our scheme and will be discussed in the next section.

---

**Algorithm 1** The pseudocode of PPD-VQA.

---

**Require:**  $\theta \in [0, 2\pi)^d$ : the parameters of ansatz;  $L$ : loss function;  $M$ : the number of local nodes, and we donate  $M_i$  as the  $i$ -th local node;

**Ensure:** optimal parameters  $\theta^*$

- 1: **while** convergence condition is not satisfied **do**
- 2: The central parameter server divides the parameter  $\theta$  into  $M$  parts and allocates  $\theta$  to  $M$  local nodes
- 3: **for** Local nodes  $M_i, \forall i \in \{1, \dots, M\}$  in parallel **do**
- 4: Calculate gradient component  $[\nabla L(\theta)]_i$
- 5: **end for**
- 6: Synchronize  $\nabla L(\theta)$  by merging  $\{[\nabla L(\theta)]_i\}_{i=1}^M$
- 7: Update  $\theta$  with a classical optimizer, such as ADAM
- 8: Choose a local node from  $\{M\}_{i=1}^M$  for convergence test
- 9: **if** convergence condition is satisfied **then**
- 10: break
- 11: **end if**
- 12: **end while**

---

### III. PERFORMANCE ANALYSIS AND ERROR MITIGATION STRATEGY IN THE REALISTIC NOISE SCENARIO

Gradient represents the optimization direction during the training procedure of VQA, which plays an decisive role in the process of finding the global minima of loss function. Thus, by examining the gradient, we analyze how noise affects convergence of PPD-VQA in the realistic scenario that noise varies for each quantum processor. Furthermore, we will propose a strategy to mitigate the negative consequences that may be caused by this realistic scenario.

#### A. Convergence and acceleration

We apply the ‘‘worst-case’’ noise channel—the depolarizing channels [43] for the following research. According to the Lemma 6 in Ref. [44] all noisy channels  $\varepsilon(\cdot)$ , which are separately applied to each layer of the ansatz, can be merged together and represented by a new depolarizing channel acting on the whole ansatz, i.e.,

$$\tilde{\varepsilon}(\rho) = (1 - \tilde{p})\rho + \tilde{p}\frac{\mathbb{I}}{2^n} \quad (1)$$

where  $\tilde{p} = 1 - (1 - p)^N$ ,  $p$  is the depolarizing probability in  $\varepsilon(\cdot)$ , and  $N$  refers to depth of ansatz. To facilitate understanding, we donate  $\tilde{p}_i$  as the noise level of the  $i$ -th QPU. Ob-

viously, the depolarizing noise turns the quantum state into a maximally mixed state with a certain probability, which could make the gradients obtained by parameter-shift-rule in the experiment deviate from that of the ideal environment without noise.

We firstly simplify some notations and introduce basic concepts in optimization theory for ease of subsequent discussion. Donate  $D = \{x_k, y_k\}$  as the training dataset, where  $x_k \in R^{2^n}$  and  $y_i \in R$  refer to example and the corresponding label respectively. We define  $L$  as the loss function,  $\nabla L(\theta^{(t)})$  as the gradient of loss function  $L$ . Here we employ the mean square error (MSE) loss function, i.e.,

$$L = \frac{1}{2N_D} \sum_k (\hat{y}_k - y_k)^2 + \frac{\lambda}{2} \|\theta^{(t)}\|^2, \quad (2)$$

where  $\hat{y}_k = \langle O \rangle$  is the predicted label with  $\langle O \rangle$  being the outcome of the observable  $O$  by  $K$  measurements,  $N_D$  is the number of the data, and  $\lambda \geq 0$  is the regularizer coefficient.

According to parameter-shift-rule,  $\nabla_j L_i(\theta^{(t)})$  ( $j$ -th component of parameters vector  $\theta^{(t)}$ ) satisfies

$$\frac{1}{N_D} \left[ \sum_k (\hat{y}_k^{(t)} - y_k) \frac{\hat{y}_k^{(t, \pm j)} - \hat{y}_k^{(t, -j)}}{2} + \lambda \theta_{i,j}^{(t)} \right], \quad (3)$$

where  $\hat{y}_k^{(t, \pm j)}$  denotes the output of PQC with shifted parameter  $\theta^{(t)} \pm \frac{\pi}{2} e_{i,j}$ , and  $e_{i,j}$  denotes the unit vector. Thus, for each data, the local node should implement  $1 + 2d/M$  quantum circuits for the gradient estimation, where  $d/M$  is the number of parameter in each local node.

Now we quantify the convergence of PPD-VQA with multiple local nodes that have different performance, by using the following utility metric [44]:

$$R_1(\theta^{(T)}) = \mathbb{E}[\|\nabla L(\theta^{(T)})\|^2] \quad (4)$$

where  $T$  is the number of iterations and the expectation  $\mathbb{E}[\bullet]$  is taken over the random variables associated with depolarizing noise. This metric evaluates how far the result is away from the stationary point. The upper bounds of  $R_1(\theta^{(T)})$  when implementing PPD-VQA with multiple non-identical processors are summarized in the following theorem.

**Theorem 1** Suppose that  $M$  noisy local nodes of PPD-VQA have different depolarizing noise with depolarizing probability  $\{\tilde{p}_i\}_{i=1}^M$ , the metric  $R_1(\theta^{(T)})$  has following upper bound

$$R_1 \leq \frac{1 + 9\pi^2 \lambda d}{2T(1 - \tilde{p}_{max})^2} + \frac{2G + d}{(1 - \tilde{p}_{max})^2} (2 - \tilde{p}_{max}) \tilde{p}_{max} (1 + 10\lambda)^2 + \frac{2dK + d}{2N_D K^2} \frac{1}{(1 - \tilde{p}_{max})^2},$$

where loss function  $L$  is  $S$ -smooth with  $S = (3/2 + \lambda)d^2$ ,  $G$ -Lipschitz with  $G = d(1 + 3\pi\lambda)$ , and  $\tilde{p}_{max} = \max\{\tilde{p}_i\}_{i=1}^M$ .

The proof of Theorem 1 is essentially similar with conventional VQA, for both of them acquire the complete gradient



information only once in an iteration. Therefore, one can obtain the upper bound of  $R_1(\theta^{(T)})$  of PPD-VQA in noise scenario by following a similar proof procedure of Theorem 1 in Ref. [44]. We briefly sketch our proof as follows.

The first step is to establish the relation between **true gradient component**  $[\nabla L(\theta^{(t)})]_{i,j}$  (unbiased) and that in the **estimated gradient**  $[\nabla \bar{L}(\theta^{(t)})]_{i,j}$  (biased) that is evaluated from QPU  $i$  (see Appendix A for the detailed derivation),

$$[\nabla \bar{L}(\theta^{(t)})]_{i,j} = (1 - \tilde{p}_i)^2 [\nabla L(\theta^{(t)})]_{i,j} + C_{i,j}^{(t)} + \varsigma_{i,j}^{(t)}, \quad (5)$$

where  $C_{i,j}^{(t)}$  originates from the depolarizing noise, and  $\varsigma_{i,j}^{(t)}$  is a item related to random variables, which has zero mean.

Then one can further utilize the  $S$ -smooth and  $G$ -Lipschitz of the  $L$  to calculate the loss difference, i.e.,

$$\begin{aligned} & L(\theta^{(t+1)}) - L(\theta^{(t)}) \\ & \leq \langle \nabla \bar{L}(\theta^{(t)}), \theta^{(t+1)} - \theta^{(t)} \rangle + \frac{S}{2} \|\theta^{(t+1)} - \theta^{(t)}\|_2^2 \end{aligned} \quad (6)$$

Substitute Eq.(5) and  $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla \bar{L}(\theta^{(t)})$  (we set the learning rate  $\eta = 1/S$ ) into Eq.(6) and take the expectation over the random variable  $\varsigma_{i,j}^{(t)}$ , one have

$$\begin{aligned} & \mathbb{E}_{\varsigma_{i,j}^{(t)}} [L(\theta^{(t+1)}) - L(\theta^{(t)})] \\ & \leq \sum_{i,j} \left[ -\frac{1}{2S} (1 - \tilde{p}_i)^2 \left( [\nabla L(\theta^{(t)})]_{i,j} \right)^2 \right. \\ & \quad \left. + \frac{2G/d + 1}{2S} (2 - \tilde{p}_i) \tilde{p}_i (1 + 10\lambda)^2 \right] \\ & \quad + \frac{2dK + d}{4SN_D K^2} \end{aligned} \quad (7)$$

Note that  $-\sum_{i,j} (1 - \tilde{p}_i)^2 \left( [\nabla L(\theta^{(t)})]_{i,j} \right)^2 \leq -(1 - \tilde{p}_{max})^2 \|\nabla L(\theta^{(t)})\|^2$ , and  $(2 - \tilde{p}_i) \tilde{p}_i \leq (2 - \tilde{p}_{max}) \tilde{p}_{max}$ , we obtain

$$\begin{aligned} & \|\nabla L(\theta^{(t)})\|^2 \\ & \leq 2S \frac{L(\theta^{(t)}) - \mathbb{E}_{\varsigma_{i,j}^{(t)}} L(\theta^{(t+1)})}{(1 - \tilde{p}_{max})^2} \\ & \quad + \frac{2G + d}{(1 - \tilde{p}_{max})^2} (2 - \tilde{p}_{max}) \tilde{p}_{max} (1 + 10\lambda)^2 \\ & \quad + \frac{2dK + d}{4SN_D K^2} \frac{1}{(1 - \tilde{p}_{max})^2}. \end{aligned} \quad (8)$$

Finally, by summing over  $t = 0, 1, \dots, T$ , the upper bound of  $R_1(\theta^{(T)})$  is achieved.

From Theorem 1 above and Theorem 1 in Ref. [44], we can observe that the convergence rate between conventional VQA and PPD-VQA is similar, i.e., both of them scale with  $O(1/\sqrt{T})$  [44], since the second term and the third term are constant in above inequality when  $\{\tilde{p}\}_{i=1}^M$  is fixed. The similar convergence rate guarantees that PPD-VQA promises a intuitive linear runtime speedup of the computation of the gradient with respect to the increased number of local nodes  $M$ .

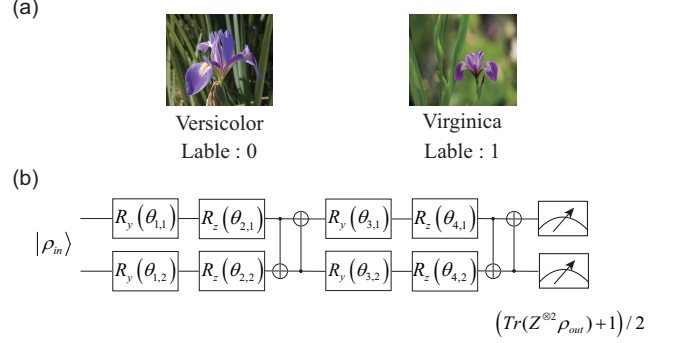


FIG. 2. **Classification task on Iris dataset and the ansatz in numerical simulation.** (a) A visualization of training examples sampled from Iris dataset. We choose the data of iris versicolor (label 0) and iris virginica (label 1) for binary classification. (b) Ansatz of PPD-VQA for the classification.

Next, we perform numerical experiment to study the performance of PPD-VQA in the realistic noise scenario.

In our simulations, we apply PPD-VQA to the binary classification task, by employing the Iris dataset and ansatz shown in Fig. 2. We choose 100 examples from Iris dataset with 50 versicolors (label 0) and 50 virginicas (label 1), where 75% examples are randomly selected as the training set and the remaining 25% as the test set. **We encode the classical example  $x_k$  in Iris dataset into the a two-qubit quantum state  $\rho_k$  by amplitude encoding, i.e.,**

$$|\psi_k\rangle = \frac{\alpha_{k1} |00\rangle + \alpha_{k2} |01\rangle + \alpha_{k3} |10\rangle + \alpha_{k4} |11\rangle}{\sqrt{|\alpha_{k1}|^2 + |\alpha_{k2}|^2 + |\alpha_{k3}|^2 + |\alpha_{k4}|^2}}, \quad (9)$$

where  $(\alpha_{k1}, \alpha_{k2}, \alpha_{k3}, \alpha_{k4})$  is the feature of  $x_k$ . Then a hardware-efficient PQC with 8 trainable single-qubit gates, as shown in Fig. 2(b), is employed for the training. After the quantum state  $\rho_k$  has evolved, we perform  $K$  global measurements on the final quantum state. We then derive the expectation of observable and map it to  $[0, 1]$  by linear mapping, where the observable is set as  $(Z^{\otimes 2} + I)/2$ .

We implement the task using the PPD-VQA with  $M = 1$  (conventional VQA), 2, 4, 8 local nodes, respectively. For each type of PPD-VQA, we also set different noise parameters separately. Specifically, for each node the PPD-VQA, the depolarizing probability  $p_i$  for single-qubit gate is set by sampling from a Gaussian distribution i.e.,  $p_i \sim N(\mu, \sigma^2)$ , where the mean  $\mu$  varies from 0.01 to 0.05 with step 0.02 and  $\sigma = \mu/9$ . The depolarizing probability of two-qubit gate is set as  $4p_i$  refer to the performance of SOTA quantum processor *Zuchongzhi* [15]. Each local node's noise will be somewhat different as a result of such random sampling. A total of 100 independent experiments were run for each setting, and in each experiment, the measurement shots is set to 8192, batch-size is set to 5, and the convergence condition is that the classification accuracy on the training set exceeds 96%.

As shown in Fig. 3(a, b), for both conventional VQA and PPD-VQAs with 2, 4, and 8 local nodes, the number of iterations required to achieve a preset training accuracy increases

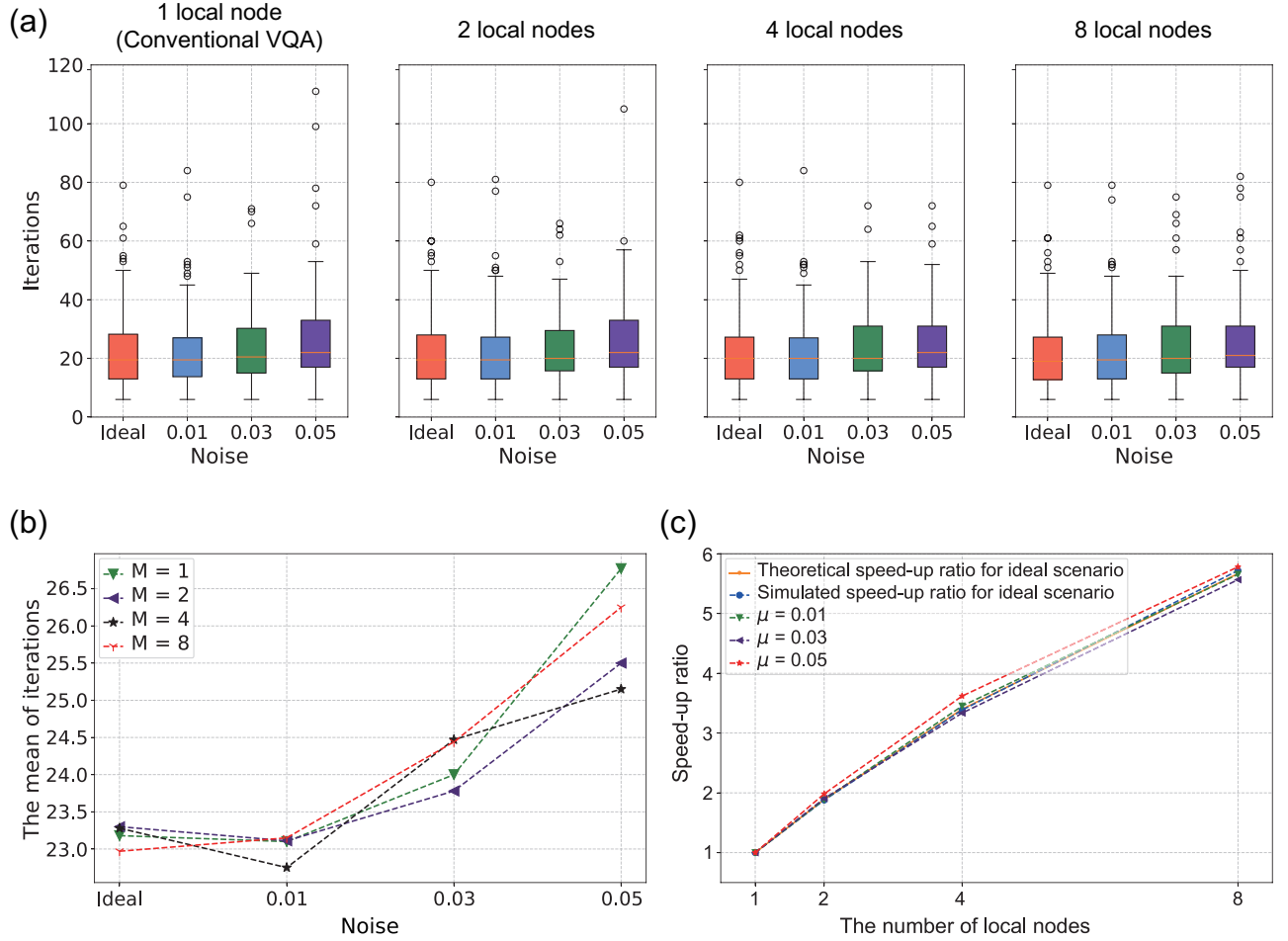


FIG. 3. **Simulation results of PPD-VQA with  $M$  local nodes under noise scenario for Iris dataset classification.** (a) Boxplots count the iterations of PPD-VQA with  $M$  local nodes, where  $M = 1, 2, 4, 8$  from left to right, when achieving a predefined training accuracy. The depolarizing probability  $p_i$  for single-qubit gate is set by sampling from a Gaussian distribution i.e.,  $p_i \sim N(\mu, \sigma^2)$ , where the mean  $\mu = 0$  (ideal case), 0.01, 0.03, 0.05, and  $\sigma = \mu/9$ . The depolarizing probability can be converted to Pauli errors  $e_p$  for single-qubit gate by using  $e_p = \frac{3}{4}p_i$  [13]. As a reference, the single-qubit gate pauli error of the Zuchongzhi processor is 0.14% [15]. The middle line of the boxes, which is the median of the data, represents the average of the number of iterations. The upper and lower limits of the box, which are the upper and lower quartiles of that, respectively, which means that the box contains 50% examples. Above and below the box, there are other lines each representing the maximum and minimum values and the circles represent outliers. (b) Scaling behavior of the mean of the iterations in (a) for increasing noise ( $\mu$ ). The results of PPD-VQA with  $M = 1, 2, 4, 8$  local nodes are shown. (c) Scaling behavior of speed-up ratio in clock-time for increasing number of local nodes  $M$ . The results of different depolarizing probabilities are shown.

with the mean of noise  $\mu$ , and the PPD-VQA with multiple local nodes has a similar convergence speed as conventional VQA (see Fig. 3(b)), which is consistent with Theorem 1. **Meanwhile, the number of iterations is not sensitive to changes of the number of local nodes, which may be caused by the proximity of noise levels of  $M$  local nodes.**

We further introduce a metric, i.e.  $R_S = T_1/T_m$  to evaluate the speed-up ratio of the PPD-VQA with  $M = m > 1$  local nodes compared to the conventional VQA with just  $M = 1$  local node, where  $T_1$  and  $T_m$  are the time consuming of conventional VQA and PPD-VQA from the start of training to meeting the convergence conditions, respectively. **Noticeably,  $T_1$  and  $T_m$  are obtained on the same noise scenario.** Assuming that the time consumption of implementing each quantum circuit is the same (since the number of measurements is the

same, and only the rotation angle of the single-qubit gate will be changed each time the circuit is executed), the formula of the speedup ratio  $R_S$  can be further rewritten as

$$R_S = \frac{(1 + 2d) \times N_D \times N_I^1}{(1 + \frac{2d}{M}) \times N_D \times N_I^M} = \frac{(1 + 2d) \times N_I^1}{(1 + \frac{2d}{M}) \times N_I^M}, \quad (10)$$

where  $d$  is the number of parameters,  $N_I^1$  and  $N_I^M$  are the total number of iterations for the conventional VQA and PPD-VQA, respectively. In the ideal scenario of noiseless,  $N_I^1 = N_I^M$ , thus  $R_S = \frac{1+2d}{1+\frac{2d}{M}}$  in ideal scenario.

Figure 3(c) shows that speed-up ratio for the PPD-VQA with 1 (conventional VQA), 2, 4, 8 local nodes under a variety of noise scenarios. No matter how the mean of noise  $\mu$  changes, the speed-up ratio of PPD-VQA is almost only re-



lated to the number of local nodes and is extremely close to the ideal case. **The reason that the speed-up ratio is almost independent of noise is that the training speed of 1 nodes and  $M$  nodes slows down at the same time as the noise grows.** This result strongly supports that PPD-VQA can achieve a very good acceleration in realistic scenarios.

### B. Alternate training strategy for mitigating the negative effects of large noise differences between different local nodes

In the previous subsection, the difference in the noise of the quantum processors of each node is not particularly large, because the noise is set by sampling from a Gaussian distribution  $N(\mu, \sigma^2)$ , where  $\sigma = \mu/9$ . In this subsection, we will study the performance of PPD-VQA in cases where the noise difference is more pronounced.

We first monitor the performance of PPD-VQA when the noise difference of different local nodes changes from small to large. To quantify the noise differences of local nodes, we introduce a metric, named *Differ*, which is defined as

$$Differ = D_{KL}(P(p) \parallel P_{\text{Uniform}}),$$

where  $D_{KL}$  is Kullback-Leibler (K-L) divergence [45],  $P_{\text{Uniform}}$  refers to the uniform distribution,  $P(p)$  is the normalized distribution of depolarization probability of each local node, where  $P(p)_k = p_k / \sum_{i=1}^M p_i$ , and  $p_k$  is the depolarization probability of the  $k$ -th local node. With this metric, a noise setting with a resulting distribution that corresponds to a higher K-L divergence with respect to uniform distribution would mean greater noise variance between local nodes. Besides, we set another constraint that the mean of  $\{p_i\}_{i=1}^M$  is 0.04. For each PPD-VQA with  $M \in [1, 2, 4, 8]$  local nodes, *Differ* varies from 0 to 0.625, we generate 10 instances of noise setting for each *Differ*, and for each instance 50 experiments with different initial parameters are implemented. As shown in Fig. 4, the speed-up ratio tends to become smaller as the *Differ* increases, indicating that the advantage of PPD-VQA in terms of speedup is diminished in extreme cases where the noise difference between local nodes is significant.

To suppress acceleration decay of PPD-VQA caused by excessive noise difference between local nodes, we propose a simple but effective approach named as alternate training strategy, whose core idea is decoupling the trainable parameter groups and corresponding quantum processors. The process of this alternate training strategy is as follows: Suppose that at the first iteration, the  $i$ -th local node is scheduled to train the parameters  $\theta_i$ . We denote this process as  $\{\theta_i : \text{QPU}_i\}_{i=1}^M$ . Then in the next iteration, The corresponding relationship between trainable parameters and local nodes becomes  $\{\theta_M : \text{QPU}_1\} \cup \{\theta_i : \text{QPU}_{i+1}\}_{i=1}^{M-1}$ , that is, we perform a cyclic shift on the correspondence between the trainable parameters and local nodes. The alternate training strategy is repeated with the training process, which makes each parameter group  $\theta_i$  be trained in turn by all quantum processors throughout the whole training process.

The numerical simulation results of PPD-VQA with alternate training strategy are shown in Fig. 4. An immediate observation is that when the noise difference between local nodes increases, PPD-VQA performance degrades relatively little thanks to the alternate training strategy. Besides, the performance of PPD-VQA becomes more stable as the variance of the mean of different experiments is significantly smaller. These two benefits suggest that this strategy can be effectively employed for mitigating the negative effects of large noise differences between different local nodes.

## IV. GRADIENT COMPRESSION

Another challenge of distributed machine learning is the large amount of communication bandwidth for gradient exchange [46]. With the development of quantum computing hardware, this problem may also arise in large-scale distributed quantum machine learning. To overcome this potential problem, we adopt the technique of gradient compression [47] widely used in the classical community to PPD-VQA, to reduce the communication bandwidth for distributed training. The pseudocode of PPD-VQA with gradient compression for local node  $i$  in PPD-VQA is as follows.

---

### Algorithm 2 The pseudocode of PPD-VQA with gradient compression.

---

**Require:**  $\theta \in [0, 2\pi)^d$ : the parameters of ansatz;  $L$ : loss function;  $M$ : the number of local nodes, and we denote  $M_i$  as the  $i$ -th local node;  $Mask = (0, \dots, 0)$  has the same dimension with  $\theta_i$  defined in section *PPD-VQA*, and  $\odot$  is hadamada product, i.e.,  $\mathbf{a} \odot \mathbf{b} = (a_1 b_1, \dots, a_n b_n)$

**Ensure:** optimal parameters  $\theta^*$

- 1: Calibrate threshold *thr*
- 2:  $[\nabla L(\theta)]_i = \mathbf{0}$  for  $i \in [1, 2, \dots, M]$
- 3: **while** convergence condition is not satisfied **do**
- 4:   The central parameter server divides the parameter  $\theta$  into  $M$  parts and allocates  $\theta$  to  $M$  local nodes
- 5:   **for** Local nodes  $M_i, \forall i \in [M]$  in parallel **do**
- 6:     Calculate gradient component  $G_i(\theta)$
- 7:      $\nabla L_i(\theta) = \nabla L_i(\theta) + G_i(\theta)$
- 8:     **for**  $j = 1 + (i-1) \frac{d}{M}, \dots, 1 + i \frac{d}{M}$  **do**
- 9:       **if**  $|\nabla L(\theta)_{i,j}| > thr$  **then**
- 10:           $Mask[j] = 1$
- 11:       **end if**
- 12:     **end for**
- 13:      $g_i(\theta) = [\nabla L(\theta)]_i \odot Mask$
- 14:      $[\nabla L(\theta)]_i = [\nabla L(\theta)]_i \odot \neg Mask$
- 15:   **end for**
- 16:   Synchronize  $g_i(\theta)$  by merging  $\{g_i(\theta)\}_{i=1}^M$ ;
- 17:   Update  $\theta$  with a classical optimizer, such as ADAM;
- 18:   Choose a local node from  $\{M\}_{i=1}^M$  for convergence test.
- 19:   **if** convergence condition is satisfied **then**
- 20:     break
- 21:   **end if**
- 22: **end while**

---

The idea of gradient compression is gradient clipping, which makes the gradient sparse by comparing its individual components with a threshold *thr*. Only the components of the

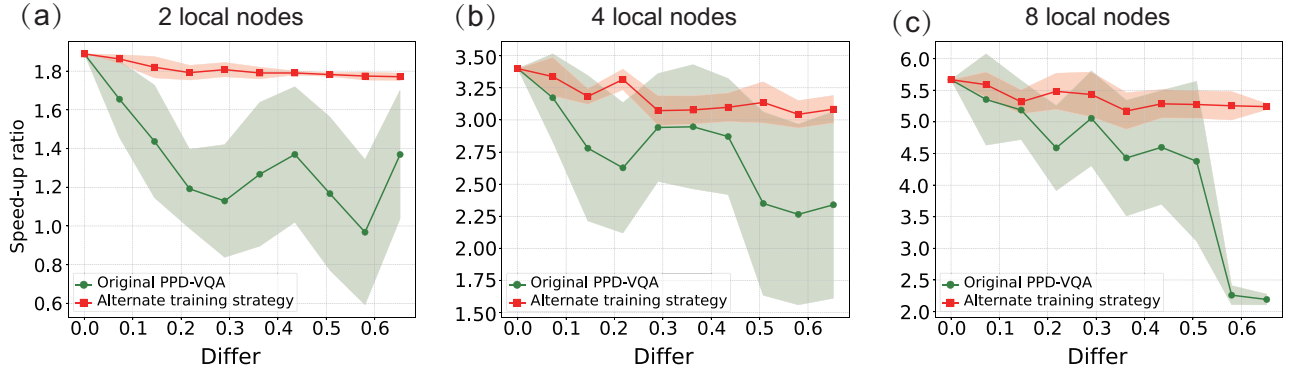


FIG. 4. **Simulation results of PPD-VQA with  $M$  local nodes in cases where the noise difference is more pronounced.** (a) The average speed-up ratio as a function of  $Differ$  (a metric for quantifying the noise differences of local nodes) for the PPD-VQA with  $M = 2$ (a), 4(b), 8(c) local nodes. 100 independent experiments are implemented for each setting. The green and red solid lines present the results for original PPD-VQA and PPD-VQA with alternate training strategy, respectively. **The shaded area represents statistical error caused by 100 independent experiments.** It is obvious that the red curves have noticeable larger values and smaller variances than the green curves in all three cases.

M	noise setting	without gradient compression		with gradient compression		result	
		iterations	communication volume	iterations	communication volume	compression ratio	speed-up ratio
2	$\mu = 0.016$	2324	$2324 \times 8$	2259	7016	62.3%	$1.87 \rightarrow 1.92$
	$\mu = 0.064$	2680	$2680 \times 8$	2780	8565	60.1%	$2.04 \rightarrow 1.97$
4	$\mu = 0.016$	2324	$2324 \times 8$	2444	3630	80.0%	$3.36 \rightarrow 3.20$
	$\mu = 0.064$	2712	$2712 \times 8$	3295	2862	86.9%	$3.64 \rightarrow 2.99$
8	$\mu = 0.016$	2282	$2282 \times 8$	2463	3634	80.0%	$5.71 \rightarrow 5.29$
	$\mu = 0.064$	2702	$2702 \times 8$	3200	2818	87%	$6.09 \rightarrow 5.14$
2	$\mu = 0.016$	2324	$2324 \times 8$	5659	701	96.3%	$1.87 \rightarrow 0.77$
	$\mu = 0.064$	2680	$2680 \times 8$	6463	787	96.3%	$2.04 \rightarrow 0.84$
4	$\mu = 0.016$	2324	$2324 \times 8$	6338	623	96.7%	$3.36 \rightarrow 1.23$
	$\mu = 0.064$	2712	$2712 \times 8$	6410	791	96.4%	$3.64 \rightarrow 1.54$
8	$\mu = 0.016$	2282	$2282 \times 8$	6043	607	96.7%	$5.71 \rightarrow 2.15$
	$\mu = 0.064$	2702	$2702 \times 8$	6322	781	96.4%	$6.09 \rightarrow 2.60$

TABLE I. **The comparison of the performance between PPD-VQA without gradient compression and with gradient compression.** The results of PPD-VQA under different number of local nodes ( $M = 2$ ,  $M = 4$  and  $M = 8$ ) and noise settings ( $\mu = 0.016$  and  $\mu = 0.064$ ) are presented. In the table we count total number of iterations for all 100 instances in each setting. Communication volume  $CV$  is defined as the total number of gradient components after clipping transmitting between the central parameter server and multiple local nodes, and compression ratio is  $1 - CV_{with}/CV_{without}$ , where  $CV_{with}$  ( $CV_{without}$ ) is the communication volume for PPD-VQA with (without) gradient compression. The symbol  $\rightarrow$  indicates the change of speed-up ratio from left (PPD-VQA without gradient compression) to right (PPD-VQA with gradient compression). Two typical results help us explore the relationship between acceleration of PPD-VQA and compression ratio: (top) The acceleration of PPD-VQA with gradient compression has only a slight decay when the gradient compression ratio is over 60%. (bottom) The acceleration of PPD-VQA with gradient compression decreases significantly when the gradient compression ratio is too high (over 96%).

gradient with larger absolute values compared with  $thr$  can be synchronized to the central parameter server, thus ensuring that the general direction of the parameter update remains correct. The remaining components smaller than  $thr$  are still retained in corresponding local node and counted as a part of new gradient in next iteration. Thus we obtain the uncropped original  $\nabla L_i(\theta)$  in local node  $i$ . This method greatly reduces the actual communication bandwidth required in PPD-VQA. However, due to the existence of quantum noise, it is also unknown whether gradient compression works on PPD-VQA,

so next we will perform numerical simulations to address this concern.

We test the gradient compression on a PPD-VQA with  $M = 4$  local nodes, where the noise  $p_i$  in each node is set by sampling from the Gaussian distribution  $N(\mu, \sigma^2)$  with  $\mu = 0.016$  and  $\sigma = \mu/9$ . In our simulation, the threshold value  $thr$  varies from 0 to 0.7 with step 0.1, and we still implement 100 independent experiments for each setting. As shown in Fig. 5, by setting a reasonable compression threshold, we can greatly reduce the communication cost. It can

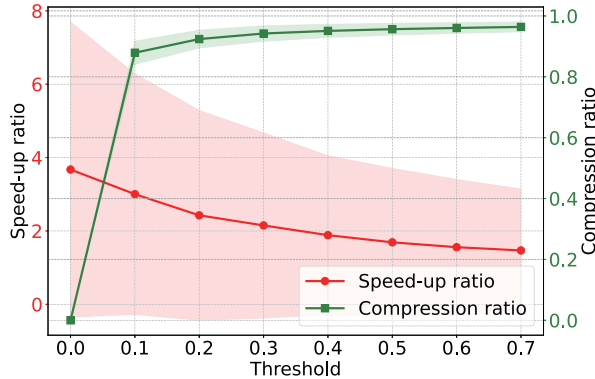


FIG. 5. The speed-up ratio and compression ratio as a function of threshold value for the PPD-VQA with  $M = 4$  local nodes. The depolarizing noise  $p_i$  in each node is set by sampling from the gaussian distribution  $N(\mu, \sigma^2)$  with  $\mu = 0.016$  and  $\sigma = \mu/9$ .

be also observed that the increase of the gradient compression ratio leads to the decay of the acceleration of PPD-VQA. When  $thr > 0.1$ , the growth of gradient compression ratio becomes very slow, while speed-up ratio is still decreasing rapidly. Thus, we need to find a balance between the decay of acceleration advantage and reducing the communication volume. When the threshold value is 0.1, we can achieve a relatively high gradient compression ratio ( $> 80\%$ ) without losing too much acceleration advantage ( $R_S > 2.7$ ).

In Table I, we further show two types of typical results for the PPD-VQA with  $M = 2, 4, 8$  local nodes, and noise level  $\mu = 0.016, 0.064$ . For each setting, 100 independent experiments are implemented. In the first typical result, we set a reasonable compression ratio, so that the speed-up ratio is almost not lost compared to the uncompressed case. However, in this scenario, the compression ratio can still be higher than 60%, or even up to 87%, indicating that we can solve the problem on communication bottleneck without losing too much acceleration advantage of PPD-VQA. In the second typical result, to achieve a more aggressive compression ratio (above 96%), the speedup of PPD-VQA is significantly reduced. Possible reason is that fewer trainable parameters make the gradient not have the sparsity compared with deep neural networks, which leads to a significant increase in the number of iteration when we apply the gradient compression algorithm to PPD-VQA. Anyway, our experiments demonstrate that gradient compression is very suitable for PPD-VQA, even in the realistic noise

scenario.

In addition to reducing communication bandwidth, gradient compression may help to mitigate errors in the experimental implementation of the PPD-VQA, due to the following two reasons: 1) For most quantum computing systems, it is not easy to implement the tiny angular rotations of single-qubit quantum operations with high precision. Gradient compression can avoid updates of tiny angles and thus potentially improve experimental accuracy. 2) As the frequency of updating parameters (especially those with small gradient changes) is reduced, the number of gate operations that need to be changed by the quantum device is consequently reduced and the accumulation of quantum operation errors is naturally suppressed.

## V. CONCLUSION

Our results show that PPD-VQA is highly promising as it achieves approximately linear acceleration over the training process of conventional VQA, both in theory and simulation results. The PPD-VQA exhibits good resilience to the excessive noise differences among local nodes, by employing the alternate training strategy. Furthermore, by adopting the gradient compression strategy, potential communication bottlenecks can also be addressed to support the future scalability of PPD-VQA.

The PPD-VQA is naturally compatible with the data-parallel distributed VQA proposed in [40], so the combination of the two approaches could enable a stronger acceleration for the training of VQA. When doing such a combination, some methods [48–50] can be employed to enhance the generalization ability of data-parallel training [49, 50]. Besides, error mitigation techniques [51, 52] have the potential to further improve the capability of PPD-VQA on near-term quantum devices. Some more complex application scenarios, such as privacy-preserving distributed VQA, requires more in-depth discussions in the future works.

## ACKNOWLEDGMENTS

H.-L. H. acknowledges support from the Youth Talent Lifting Project (Grant No. 2020-JCJQ-QT-030), National Natural Science Foundation of China (Grants No. 11905294), China Postdoctoral Science Foundation, and the Open Research Fund from State Key Laboratory of High Performance Computing of China (Grant No. 201901-01).

[1] P. W. Shor, *SIAM Rev.* **41**, 303 (1999).  
[2] C.-Y. Lu, D. E. Browne, T. Yang, and J.-W. Pan, *Phys. Rev. Lett.* **99**, 250504 (2007).  
[3] H.-L. Huang, Q. Zhao, X. Ma, C. Liu, Z.-E. Su, X.-L. Wang, L. Li, N.-L. Liu, B. C. Sanders, C.-Y. Lu, *et al.*, *Phys. Rev. Lett.* **119**, 050503 (2017).

[4] L. K. Grover, in *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing* (1996) pp. 212–219.  
[5] G. L. Long, Y. S. Li, W. L. Zhang, and L. Niu, *Phys. Lett. A* **262**, 27 (1999).  
[6] A. W. Harrow, A. Hassidim, and S. Lloyd, *Phys. Rev. Lett.* **103**, 150502 (2009).

- [7] X.-D. Cai, C. Weedbrook, Z.-E. Su, M.-C. Chen, M. Gu, M.-J. Zhu, L. Li, N.-L. Liu, C.-Y. Lu, and J.-W. Pan, *Phys. Rev. Lett.* **110**, 230501 (2013).
- [8] H.-L. Huang, Y.-W. Zhao, T. Li, F.-G. Li, Y.-T. Du, X.-Q. Fu, S. Zhang, X. Wang, and W.-S. Bao, *Front. Phys.* **12**, 120305 (2017).
- [9] C. Gidney and M. Ekerå, *Quantum* **5**, 433 (2021).
- [10] Y. Zhao, Y. Ye, H.-L. Huang, Y. Zhang, D. Wu, H. Guan, Q. Zhu, Z. Wei, T. He, S. Cao, *et al.*, *Phys. Rev. Lett.* **129**, 030501 (2022).
- [11] J. Preskill, *Quantum* **2**, 79 (2018).
- [12] H.-L. Huang, D. Wu, D. Fan, and X. Zhu, *Sci. China Inf. Sci.* **63**, 180501 (2020).
- [13] F. Arute, K. Arya, R. Babbush, D. Bacon, J. C. Bardin, R. Barends, R. Biswas, S. Boixo, F. G. Brandao, D. A. Buell, *et al.*, *Nature* **574**, 505 (2019).
- [14] H.-S. Zhong, H. Wang, Y.-H. Deng, M.-C. Chen, L.-C. Peng, Y.-H. Luo, J. Qin, D. Wu, X. Ding, Y. Hu, *et al.*, *Science* **370**, 1460 (2020).
- [15] Y. Wu, W.-S. Bao, S. Cao, F. Chen, M.-C. Chen, X. Chen, T.-H. Chung, H. Deng, Y. Du, D. Fan, *et al.*, *Phys. Rev. Lett.* **127**, 180501 (2021).
- [16] Q. Zhu, S. Cao, F. Chen, M.-C. Chen, X. Chen, T.-H. Chung, H. Deng, Y. Du, D. Fan, M. Gong, *et al.*, *Sci. Bull.* **67**, 240 (2022).
- [17] K. Bharti, A. Cervera-Lierta, T. H. Kyaw, T. Haug, S. Alperin-Lea, A. Anand, M. Degroote, H. Heimonen, J. S. Kottmann, T. Menke, *et al.*, *Rev. Mod. Phys.* **94**, 015004 (2022).
- [18] M. Cerezo, A. Arrasmith, R. Babbush, S. C. Benjamin, S. Endo, K. Fujii, J. R. McClean, K. Mitarai, X. Yuan, L. Cincio, *et al.*, *Nat. Rev. Phys.* **3**, 625 (2021).
- [19] D. P. Kingma and J. Ba, *arXiv:1412.6980* (2014), <https://doi.org/10.48550/arXiv.1412.6980>.
- [20] J. R. McClean, J. Romero, R. Babbush, and A. Aspuru-Guzik, *New J. Phys.* **18**, 023023 (2016).
- [21] I. Cong, S. Choi, and M. D. Lukin, *Nat. Phys.* **15**, 1273 (2019).
- [22] V. Havlíček, A. D. Córcoles, K. Temme, A. W. Harrow, A. Kandala, J. M. Chow, and J. M. Gambetta, *Nature* **567**, 209 (2019).
- [23] J. Liu, K. H. Lim, K. L. Wood, W. Huang, C. Guo, and H.-L. Huang, *Sci. China Phys. Mech. Astron.* **64**, 290311 (2021).
- [24] H.-L. Huang, X.-L. Wang, P. P. Rohde, Y.-H. Luo, Y.-W. Zhao, C. Liu, L. Li, N.-L. Liu, C.-Y. Lu, and J.-W. Pan, *Optica* **5**, 193 (2018).
- [25] S. Lloyd and C. Weedbrook, *Phys. Rev. Lett.* **121**, 040502 (2018).
- [26] H.-L. Huang, Y. Du, M. Gong, Y. Zhao, Y. Wu, C. Wang, S. Li, F. Liang, J. Lin, Y. Xu, *et al.*, *Phys. Rev. Appl.* **16**, 024051 (2021).
- [27] M. S. Rudolph, N. B. Toussaint, A. Katarawa, S. Johri, B. Peropadre, and A. Perdomo-Ortiz, *Phys. Rev. X* (2020), <https://doi.org/10.1103/PhysRevX.12.031010>.
- [28] S. Ebadi, A. Keesling, M. Cain, T. T. Wang, H. Levine, D. Bluvstein, G. Semeghini, A. Omran, J.-G. Liu, R. Samajdar, *et al.*, *Science*, eab6587 (2022).
- [29] M. P. Harrigan, K. J. Sung, M. Neeley, K. J. Satzinger, F. Arute, K. Arya, J. Atalaya, J. C. Bardin, R. Barends, S. Boixo, *et al.*, *Nat. Phys.* **17**, 332 (2021).
- [30] G. E. Crooks, *arXiv:1811.08419* (2018), <https://doi.org/10.22331/q-2022-07-07-759>.
- [31] E. Farhi and A. W. Harrow, *arXiv:1602.07674* (2016), <https://doi.org/10.48550/arXiv.1602.07674>.
- [32] L. Zhou, S.-T. Wang, S. Choi, H. Pichler, and M. D. Lukin, *Phys. Rev. X* **10**, 021067 (2020).
- [33] M. Gong, H.-L. Huang, S. Wang, C. Guo, S. Li, Y. Wu, Q. Zhu, Y. Zhao, S. Guo, H. Qian, *et al.*, *arXiv:2201.05957* (2022), <https://doi.org/10.48550/arXiv.2201.05957>.
- [34] A. Kandala, A. Mezzacapo, K. Temme, M. Takita, M. Brink, J. M. Chow, and J. M. Gambetta, *Nature* **549**, 242 (2017).
- [35] G. A. Quantum, Collaborators, F. Arute, K. Arya, R. Babbush, D. Bacon, J. C. Bardin, R. Barends, S. Boixo, M. Broughton, B. B. Buckley, *et al.*, *Science* **369**, 1084 (2020).
- [36] J. Romero, R. Babbush, J. R. McClean, C. Hempel, P. J. Love, and A. Aspuru-Guzik, *Quantum Sci. Technol.* **4**, 014008 (2018).
- [37] C. Cade, L. Mineh, A. Montanaro, and S. Stanisic, *Phys. Rev. B* **102**, 235122 (2020).
- [38] C. Hempel, C. Maier, J. Romero, J. McClean, T. Monz, H. Shen, P. Jurcevic, B. P. Lanyon, P. Love, R. Babbush, *et al.*, *Phys. Rev. X* **8**, 031022 (2018), doi: 10.1103/PhysRevX.8.031022.
- [39] Y. Nam, J.-S. Chen, N. C. Panti, K. Wright, C. Delaney, D. Maslov, K. R. Brown, S. Allen, J. M. Amini, J. Apisdorf, *et al.*, *npj Quantum Inf.* **6**, 33 (2020).
- [40] Y. Du, Y. Qian, X. Wu, and D. Tao, *IEEE Trans. Quantum Eng.* **3**, 1 (2022).
- [41] K. Mitarai, M. Negoro, M. Kitagawa, and K. Fujii, *Phys. Rev. A* **98**, 032309 (2018).
- [42] M. Schuld, V. Bergholm, C. Gogolin, J. Izaac, and N. Killoran, *Phys. Rev. A* **99**, 032331 (2019).
- [43] M. M. Wilde, *Quantum Information Theory* (Cambridge University Press, Cambridge, 2013) pp. 128–129.
- [44] Y. Du, M.-H. Hsieh, T. Liu, S. You, and D. Tao, *PRX Quantum* **2**, 040337 (2021).
- [45] D. A. Meyer and N. R. Wallach, *J. Math. Phys.* **43**, 4273 (2002).
- [46] J. Verbracken, M. Wolting, J. Katzy, J. Kloppenburg, T. Verbeelen, and J. S. Rellermeier, *ACM Comput. Surv.* **53**, 1 (2020).
- [47] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, *arXiv:1712.01887* (2017), <https://doi.org/10.48550/arXiv.1712.01887>.
- [48] W. Zhu, X. Chen, and W. B. Wu, *Journal of the American Statistical Association*, 1 (2021).
- [49] M. H. Zhu, L. N. Ezzine, D. Liu, and Y. Bengio, *arXiv preprint arXiv:2205.09305* (2022), <https://doi.org/10.48550/arXiv.2205.09305>.
- [50] A. Rame, C. Dancette, and M. Cord, in *International Conference on Machine Learning* (PMLR, 2022) pp. 18347–18377.
- [51] H.-L. Huang, X.-Y. Xu, C. Guo, G. Tian, S.-J. Wei, X. Sun, W.-S. Bao, and G.-L. Long, *arXiv:2211.08737* (2022), <https://doi.org/10.48550/arXiv.2211.08737>.
- [52] Z. Cai, R. Babbush, S. C. Benjamin, S. Endo, W. J. Huggins, Y. Li, J. R. McClean, and T. E. O’Brien, *arXiv:2210.00921* (2022).

## VI. APPENDIX

### A. The derivation of estimated gradient

According to parameter-shift-rule, the  $j$ -th component of the analytical gradient  $\nabla_j L(\theta^{(t)})$  at  $t$ -th iteration satisfies

$$\begin{aligned} & \nabla_j L(\theta^{(t)}) \\ &= \frac{1}{N_D} \left[ \sum_k (\hat{y}_k^{(t)} - y_k) \frac{\hat{y}_k^{(t,+j)} - \hat{y}_k^{(t,-j)}}{2} + \lambda \theta_{i,j}^{(t)} \right]. \end{aligned}$$

Here, the label  $\hat{y}_k^{(t)} = \text{Tr}[U(\boldsymbol{\theta}^{(t)})\rho_k U^\dagger(\boldsymbol{\theta}^{(t)})O]$  is defined by the expected output of noiseless PQC  $U(\boldsymbol{\theta}^{(t)})$  with the observable  $O$  and input state  $\rho_k$ ;  $\hat{y}_k^{(t,\pm j)} = \text{Tr}[U(\boldsymbol{\theta}^{(t)} \pm \frac{\pi}{2}\mathbf{e}_j)\rho_k U^\dagger(\boldsymbol{\theta}^{(t)} \pm \frac{\pi}{2}\mathbf{e}_j)O]$  denote the output of PQC with shifted parameter  $\boldsymbol{\theta}^{(t)} \pm \frac{\pi}{2}\mathbf{e}_j$  where  $\mathbf{e}_j$  denotes the unit vector with its  $j$ -th component equals to one. Thus, for each data, the local node should implement  $1 + 2d/M$  quantum circuits for the gradient estimation, where  $d/M$  is the number of parameter in each local node.

In the realistic scenario, the  $j$ -th component of estimated gradient calculated on  $i$ -th QPU  $[\nabla \bar{L}(\boldsymbol{\theta}^{(t)})]_{i,j}$  is obtained by the estimated values of  $\bar{y}_k^{(t)}$  and  $\bar{y}_k^{(t,\pm j)}$ , *i.e.*,

$$\begin{aligned} & [\nabla \bar{L}(\boldsymbol{\theta}^{(t)})]_{i,j} \\ &= \frac{1}{N_D} \sum_k \left[ (\bar{y}_k^{(t)} - y_k) \frac{\bar{y}_k^{(t,+j)} - \bar{y}_k^{(t,-j)}}{2} + \lambda \theta_{i,j}^{(t)} \right]. \end{aligned}$$

According to the definition and notation of the depolarizing noise model in III.A, the estimated label  $\bar{y}_k^{(t)}$  has the mean

value  $v_{k,i}^{(t)}$  and variance  $(\sigma_{k,i}^{(t)})^2$  after  $K$  measurements:

$$\begin{aligned} v_{k,i}^{(t)} &= (1 - \tilde{p}_i) \hat{y}_{k,i}^{(t)} + \tilde{p}_i \frac{\text{Tr}[O]}{2^n}, \\ (\sigma_{k,i}^{(t)})^2 &= \frac{1}{K} (1 - \tilde{p}_i) \tilde{p}_i (\hat{y}_{k,i}^{(t)} - \frac{\text{Tr}[O]}{2^n})^2. \end{aligned}$$

We further introduce the random variable  $\xi_{k,i}^{(t)}$  with zero mean and variance  $(\sigma_{k,i}^{(t)})^2$  to describe the output of PQC on QPU  $i$ , *i.e.*,

$$\bar{y}_k^{(t)} = v_{k,i}^{(t)} + \xi_{k,i}^{(t)}.$$

Similarly, we can define  $v_{k,i}^{(t,\pm j)}$  and the random variable  $\xi_{k,i}^{(t,\pm j)}$  to describe the output of QPU  $i$  with shifted-parameter in the same way.

Therefore, a formulaic description of the relation between the estimated partial derivative  $[\nabla \bar{L}(\boldsymbol{\theta}^{(t)})]_{i,j}$  and the analytic gradients  $\nabla_j L(\boldsymbol{\theta}^{(t)})$  is as follows,

$$\begin{aligned} & [\nabla \bar{L}(\boldsymbol{\theta}^{(t)})]_{i,j} \\ &= \frac{1}{N_D} \sum_k \left[ (v_{k,i}^{(t)} + \xi_{k,i}^{(t)} - y_k) \left( \frac{v_{k,i}^{(t,+j)} - v_{k,i}^{(t,-j)} + \xi_{k,i}^{(t,+j)} - \xi_{k,i}^{(t,-j)}}{2} + \lambda \theta_{i,j}^{(t)} \right) \right] \\ &= \frac{1}{N_D} \sum_k \left[ (1 - \tilde{p}_i)^2 [\nabla L(\boldsymbol{\theta}^{(t)}, \rho_k)]_{i,j} \right] + \frac{1}{N_D} \sum_k \left[ (1 - \tilde{p}_i) \tilde{p}_i \left( \frac{\text{Tr}[O]}{2^n} - y_k \right) \frac{\hat{Y}_{k,i}^{(t,+j)} - \hat{Y}_{k,i}^{(t,-j)}}{2} + (2\tilde{p}_i - \tilde{p}_i^2) \lambda \theta_{i,j}^{(t)} \right] \\ &\quad + \frac{1}{N_D} \sum_k \left[ (v_{k,i}^{(t)} - y_k) (\xi_{k,i}^{(t,+j)} - \xi_{k,i}^{(t,-j)}) + \frac{\hat{Y}_{k,i}^{(t,+j)} - \hat{Y}_{k,i}^{(t,-j)}}{2} \xi_{k,i}^{(t)} + \xi_{k,i}^{(t)} (\xi_{k,i}^{(t,+j)} - \xi_{k,i}^{(t,-j)}) \right] \\ &= (1 - \tilde{p}_i)^2 \nabla_j L(\boldsymbol{\theta}^{(t)}) + C_{i,j}^{(t)} + \varsigma_{i,j}^{(t)}, \end{aligned}$$

where

$$\begin{aligned} C_{i,j}^{(t)} &= \frac{1}{N_D} \sum_k \left[ (1 - \tilde{p}_i) \tilde{p}_i \left( \frac{\text{Tr}[O]}{2^n} - y_k \right) \frac{\hat{Y}_{k,i}^{(t,+j)} - \hat{Y}_{k,i}^{(t,-j)}}{2} + (2\tilde{p}_i - \tilde{p}_i^2) \lambda \theta_{i,j}^{(t)} \right], \\ \varsigma_{i,j}^{(t)} &= \frac{1}{N_D} \sum_k \left[ (v_{k,i}^{(t)} - y_k) (\xi_{k,i}^{(t,+j)} - \xi_{k,i}^{(t,-j)}) + \frac{\hat{Y}_{k,i}^{(t,+j)} - \hat{Y}_{k,i}^{(t,-j)}}{2} \xi_{k,i}^{(t)} + \xi_{k,i}^{(t)} (\xi_{k,i}^{(t,+j)} - \xi_{k,i}^{(t,-j)}) \right]. \end{aligned}$$

Obviously,  $\varsigma_{i,j}^{(t)}$  has the mean zero, where the random variables  $\xi_{k,i}^{(t)}$ ,  $\xi_{k,i}^{(t,\pm j)}$  in above formula are independent, because these random variables arise when we separately calculate the expectation of the observable with different shifted parameters.

## B. Bias term analysis

In this subsection, we give an error analysis on estimated gradient. By leveraging the explicit form of estimated gradient in Appendix A, estimated gradient has following average bias

term compared to the ideal gradient,

$$\begin{aligned} \text{bias term} &= |E_{\xi_{k,i}^{(t)}, \xi_{k,i}^{(t,\pm j)}} \left( [\nabla \bar{L}(\boldsymbol{\theta}^{(t)})]_{i,j} - \nabla_j L(\boldsymbol{\theta}^{(t)}) \right)| \\ &= |(\tilde{p}_i^2 - 2\tilde{p}_i) \nabla_j L(\boldsymbol{\theta}^{(t)}) + C_{i,j}^{(t)}| \\ &\leq |(\tilde{p}_i^2 - 2\tilde{p}_i) \nabla_j L(\boldsymbol{\theta}^{(t)})| + |C_{i,j}^{(t)}| \\ &\leq 2\tilde{p}_i \left( \frac{1}{2} + 3\lambda\pi \right) + (1 + 6\lambda\pi) \tilde{p}_i \\ &= (2 + 9\lambda\pi) \tilde{p}_i \end{aligned}$$



where the inequality uses the upper bound of  $\nabla_j L(\boldsymbol{\theta}^{(t)})$ , i.e.,  $\nabla_j L(\boldsymbol{\theta}^{(t)}) \leq 1 \times \frac{1}{2} + 3\lambda\pi$  when  $\theta \in [\pi, 3\pi]$  and  $\hat{g}_k^{(t)}, \hat{g}_k^{(t, \pm j)} \in [0, 1]$ , and the upper bound of  $C_{i,j}^{(t)}$ , i.e.,  $C_{i,j}^{(t)} \leq \tilde{p}_i + 2\tilde{p}_i \cdot \lambda \cdot 3\pi$ .

Therefore, the larger the noise, the larger the upper bound of the bias term, which visually demonstrates the effect of noise on the gradient calculation. However, in actual training process, we pay more attention to the estimated gradient  $[\nabla \bar{L}(\boldsymbol{\theta}^{(t)})]_{i,j}$  containing the bias term, because VQAs are typical adaptive noise approaches.