# LAWS: Look Around and Warm-Start Natural Gradient Descent for Quantum Neural Networks

Zeyi Tao, Jindi Wu, Qi Xia and Qun Li

*Department of Computer Science*
*William & Mary*
Williamsburg, VA, USA
{ztao, jwu21, qxia01, liqun}@cs.wm.edu

*Abstract*—**Variational quantum algorithms (VQAs) have recently received significant attention from the research community due to their promising performance in Noisy Intermediate-Scale Quantum computers (NISQ). However, VQAs run on parameterized quantum circuits (PQC) with randomly initialized parameters are characterized by barren plateaus (BP) where the gradient vanishes exponentially in the number of qubits. In this paper, we first review quantum natural gradient (QNG), which is one of the most popular algorithms used in VQA, from the classical first-order optimization point of view. Then, we proposed a Look Around Warm-Start QNG (LAWS) algorithm to mitigate the widespread existing BP issues. LAWS is a combinatorial optimization strategy taking advantage of model parameter initialization and fast convergence of QNG. LAWS repeatedly reinitializes parameter search space for the next iteration parameter update. The reinitialized parameter search space is carefully chosen by sampling the gradient close to the current optimal. Moreover, we present a unified framework (WS-SGD) for integrating parameter initialization techniques into the optimizer. We provide the convergence proof of the proposed framework for both convex and non-convex objective functions based on Polyak-Lojasiewicz (PL) condition. Our experiment results show that the proposed algorithm could mitigate the BP and have better generalization ability in quantum classification problems.**

*Index Terms*—**Variational Quantum Algorithms, Optimization, Natural Gradient Descent**

## I. INTRODUCTION

The emergence of the Noisy Intermediate-Scale Quantum (NISQ) [1] technology has demonstrated its enormous potential in number factorization [2], quantum system simulation [3], or solving linear systems of equations [4]. Current state-of-the-art NISQ with 50-100 qubits may be able to perform tasks that outperform the capabilities of today's classical computers [5]. However, NISQ is limited by connectivity, qubit count, and gate fidelity, preventing the use of quantum error correction and making many quantum algorithms impractical [6]. To this end, one of the most promising computational models for using near-term quantum computers is proposed so-called Variational Quantum Algorithms (VQAs) [7].

The VQA is a quantum-classical hybrid algorithm. In VQA, a task of interest is prepared and evaluated via a parameterized quantum circuit (PQC) on a quantum computer, with variationally updating the parameters by a classical optimizer to find the optimum of some measurable cost function. The applications of VQA include the Variational Quantum Eigensolver

(VQE) [8], Quantum Approximate Optimization Algorithm (QAOA) [9], and Quantum Neural Networks (QNNs) [10]. The success of VQA is due to 1) VQA allows task-oriented programming making the design of quantum algorithms efficient [11]; 2) compared to the classical neural network, the expressibility of QNN is more significant even with shallow quantum circuits [12]; This low complexity in QNN mitigates the NISQ limitations.

Although it has been shown that the optimization task for minimizing the VQA cost function is, in general, an NP-hard problem [13], the effectiveness and efficiency of gradient-based optimizers are still charming. Many employ gradient-based optimizers as a backbone in VQA. For example, one uses gradient descent to reach the ground-state energy under a Hamiltonian in VQE study. Alternatively, for some variational classifiers, one uses stochastic gradient descent (SGD) to find the optimal PQC model. As a result, the performance of QNNs heavily depends on the power of such a classical optimizer. In particular, the quantum natural gradient (QNG) [14] has drawn much research attentions [15], [16] due to its extraordinary ability to discover the parameter space's geometric structure. Further, QNG has been proved in Ref. [17] that the VQE associated with QNG is equivalent to the imaginary time evolution [18] when the quantum Fubini-Study metric is applied to measure the geometric structure, making it more widespread.

However, a recently discovered phenomenon, so-called barren plateaus (BP) [19], where gradients of the cost functions vanish exponentially with the size of the system, dramatically limits the application of QNNs to practical problems. BP prevents PQC's parameter update from gradient changes when using gradient-based optimizers. To acquire the gradient information, exponential resources might be used for sampling errors in quantum measurements.

To address the BP issue, gradient rescaling [20], [21], PQC's parameter initialization [22], [23], and gradient-free optimizations [24] have been extensively studied. Our work is also motivated by addressing the BP issue. In this paper, we first review the gradient-based method, particularly QNG, in the view of mirror descent [25]. Then, we proposed a look around warm-start QNG (LAWS) algorithm as a primary instrument to mitigate the BP issue. The proposed algorithm is based on two observations: First, it has been reported in

Ref. [26] that the QNG can consistently find a global optimum and requires significantly fewer epochs than other optimizers. This outperformance holds even for large system sizes (40 qubits), indicating that using QNG to solve the QNN problem is suitable. Second, the success of applying parameter initialization in QNN reported in Ref. [22] demonstrates a potential direction for mitigating the BP issue, where it withstand the possible failure of using the gradient-based [27] or gradient-free [28] algorithm. In addition, many classical machine (deep) learning models benefit from parameter initialization strategies and gain performance improvement [29]. Based on the above, the intuition behind LAWS is that we repeatedly reinitialize the PQC's parameter while in training. We call this reinitialization during the training as warm-start. Essentially, we perform parameter initialization for every current optimum until some criteria meet (i.e., the value of the cost function is minimized). In this way, the fast convergence speed of QNG is adopted, and the BP could be mitigated via multiple parameter reinitializations. However, designing such an algorithm is non-trivial and challenging. The question, for example, of how to perform reinitialization while in training and how to reinitialization should be carefully treated. More motivations, discussions, and implementation details are present in section IV.

In summary, the contributions of this paper are following:

- First, we propose a new derivation of QNG by using a classical first-order optimization scheme known as mirror descent.
- Second, we proposed a new algorithm named LAWS for solving VQE and QNN in general. Our experiment results show that the proposed algorithm could mitigate the BP issue and have better generalization ability in quantum classification problems.
- Third, based on LAWS, we propose a unified framework WSSGD for the warm-start gradient descent algorithm that is easy to implement and compatible with the most current quantum learning libraries.
- Last, as a complementary part, we provide the convergence proof of the proposed framework for both convex and non-convex objective functions.

## II. RELATED WORK

In this section, we first introduce some notation used in the paper, then we summarize the recent studies in optimizing QPC and the presence of BP issue, which has been considered a major limitation in VQA.

### A. Notations

For $\theta, \mu \in R^d$, let $\sqrt{\theta}, \theta \odot \mu$, and $\theta/\mu$ denote the element-wise square root, multiplication, and division of the vectors. The $\|\theta\|_2^2$ is $l_2$-norm. We denote $\theta_k^t$ for parameter $\theta$ at $t$-th iteration $k$-th step.

### B. Optimization in VQA

Hybrid quantum-classical optimization performed via parameterized quantum circuits is a promsing approach for various new emerging quantum based application. A classical optimization scheme is utilized to update the parameters of such hybrid quantum-classical model. Two types of optimization are often used: gradient-based methods and gradient-free methods.

*Gradient-based Optimization* has been widely adopted in solving QPC such as stochastic gradient descent (SGD) [30]. SGD replaces the exact partial derivative with an unbiased gradient estimator at each optimization step. SGD is a promising method for almost all large-scale machine learning models, where it has been found to be efficient for gradient evaluation, fast convergence speed, etc. Properly choosing the step size of SGD is essential. In a recent study [20], the Vanilla SGD has been replaced with many modern adaptive learning rate methods such as AdaGrad [31], Adam [32], and AdaBelief [33] to achieve even faster convergence speed and better performance.

The natural gradient [34] (NG) automatically chooses gradient step size and moves in the steepest descent direction with respect to the Fisher information. The pioneering work [14] propose QNG as part of a general-purpose optimization framework for variational quantum algorithms. Later, [15] demonstrates some simple case studies for QNG via variational quantum eigensolver to reveal how the natural gradient optimizer uses the geometric property to change and improve the ordinary gradient method. It has been reported [16], [26] QNG could effectively avoid the local optimal and be stable of performance on all considered system sizes. QNG's computation is expensive; hence it becomes an obstacle for applying in both classical learning and VQA. A recent technical note [35] presents a time-efficient QNG method [36] to compute the inverse of the quantum Fisher information matrix.

*Gradient-free Optimization* such as the Nelder-Mead algorithm [37], Powell algorithm [38], and Constrained Optimization BY Linear Approximations is also welcome for optimizing PQC models. [24] proposes the first gradient-free quantum optimization for NISQ device as a detour solution for BP. However, [28] reports that gradient-free optimizers do not solve the BP problem due to the exponentially large resources demand.

### C. Mitigating the Barren Plateau

The barren plateau (BP) phenomenon in the cost function landscape was originally discovered in [19] where it was shown that deep (unstructured) parametrized quantum circuits exhibit BPs when randomly initialized. When a given cost function exhibits a BP, the magnitude of its partial derivatives will be, on average, exponentially vanishing with the system size [6]. Thus, BP has been recognized as a well-known bottleneck in VQA [27], especially when optimizing the QNNs using the gradient-based method. [27] demonstrates that even using high-order derivative information such as the Hessian, the exponential scaling associated with BP still exists.

Many works have been studied to mitigate BP, and they can be roughly categorized into two directions. The first type of approach uses problem-inspired ansatzes because problem-agnostic ansatzes, such as deep hardware efficient ansatzes,

could exhibit barren plateaus due to their high expressibility [11], [39]. The approach, for example [11], relaxes search space during the optimization to a smaller space that contains the solution to the problem or that at least contains a good approximation to the solution while maintaining a low expressibility.

Another line of study focus on QNN initialization. Parameter initialization has been proved to be helpful in classical machine learning. In [22], the proposed method uses the identity block strategy to limit the effective depth of the circuits used to calculate the first parameter update to avoid the QNN being stuck in a barren plateau at the start of training. Further, [23] proposes a parameter initialization strategy that transfers the small pre-trained layer blocks to the target model stacking by multiple identical basic blocks. This idea is based on transfer learning. The empirical results show that the gradient norm's variance is scaled and prove it is effective for mitigating BP where the trainability of QNNs is improved.

The method [21] uses adaptive learning rates induced from Gaussian kernels for the gradient method to avoid the gradient vanishing. Work [40] analyzes the existence of barren plateaus in various ansatzes, and sheds light on the role of the different initial states causing the presence or absence of barren plateaus. Then they provide an efficient framework for trainability-aware ansatz design strategies.

## III. BACKGROUND

### A. Variational Quantum Algorithms

The VQA is a quantum-classical hybrid algorithm that enables noisy quantum devices to work with the help of classical computers. The VQA runs a parameterized quantum circuit for ansatz preparation and expectation value measurement on a quantum computer. Meanwhile, VQA variationally updates the parameters by a classical optimizer to find a global minimum of the objective function. Generally, the PQC with unitary $U(\theta)$ having the form

$$U(\theta) = \prod_{l=1}^{L} U_l(\theta_l) \tag{1}$$

where

$$U_l(\theta_l) = \prod_{m=1}^{M} e^{-i\theta_{l,m} V_m/2} W_{l,m} \tag{2}$$

Here, the $l$ indicates the layer where $\theta_l = (\theta_{l,1}, \cdots, \theta_{l,M})$ contains the parameter of such layer and $\theta = \{\theta_l\}_{l=1}^{L}$. $V_m$ is Hermition operator that generates the unitary in the ansatz. In addition, $W_{l,m}$ is the unparameterized quantum gate. One of the widely used ansatzes in quantum chemistry, optimization, and quantum simulation is hamiltonian variational ansatz (HVA) which aims to prepare a trial ground state for a given Hamiltonian $H = \sum V_m$ (here $V_m$ are Hermitian operators, usual Pauli strings) by Trotterizing an adiabatic state preparation process [41]. The unitary of HVA, therefore, is given by $U(\theta) = \prod_l \left( \prod_m \exp(-i\theta_{l,m} V_{l,m}) \right)$ in the form of

Euqation 1. Without loss of generality, the cost of VQA can be expressed as

$$\mathcal{C}(\theta) = f(\{\rho_k\}, \{\mathcal{O}_k\}, U(\theta)) \tag{3}$$

where $f$ is some function, $\{\rho_k\}$ are input states, $\{\mathcal{O}_k\}$ is a set of observables, and $U(\theta)$ is parameterized unitary defined in Equation 1. To explicitly define function $f$, one could have cost in the form

$$\mathcal{C}(\theta) = \sum_{k=1}^{K} a_k \operatorname{Tr} \left[ U(\theta) \rho_k U^\dagger(\theta) \mathcal{O}_k \right] \tag{4}$$

where $a_k$ is the coefficients of the linear combination of expectation values. The subscript $k$ indicates the fact that one could work with different functions for each state in the training set. Particularly, in variational quantum eigensolver, one chooses $\mathcal{O} = H$, where $H$ is the physical Hamiltonian and lets training set with size $K = 1$. The cost function becomes the expectation value of the Hamiltonian $H$

$$\mathcal{C}(\theta) = \operatorname{Tr} \left[ U(\theta) \rho_0 U^\dagger(\theta) H \right] \tag{5}$$

where $\rho_0 = |0\rangle\langle 0|$ is the initial (pure) state. In the above Equation 5, let $|\psi(\theta)\rangle = U(\theta)|0\rangle$. For simplicity, we write $|\psi(\theta)\rangle$ as $|\psi\rangle$ henceforth. Given a $N$-dimensional complex Hilbert space $\mathbb{C}^N$, we define a projector $P_\psi$ as $|\psi\rangle\langle\psi| \in \mathbb{CP}^{N-1}$. Now we consider the following optimization problem

$$\min_{\theta} \mathcal{C}(\theta) \tag{6}$$

where

$$\mathcal{C}(\theta) = \operatorname{Tr}(P_\psi H) = \langle\psi|H|\psi\rangle \tag{7}$$

Note that $\psi$ is normalized since $U(\theta)$ is unitary.

### B. First-order Optimization

Globally optimizing the objective function $\mathcal{C}(\theta)$ is impractical due to the nonconvexity. To this end, practitioners search for local optima by solving the following dynamical system

$$\theta_{t+1} = \arg\min_{\theta} \left\{ \langle\theta, \nabla\mathcal{C}(\theta_t)\rangle + \frac{1}{2\eta} \|\theta - \theta_t\|_2^2 \right\} \tag{8}$$

which is equivalent to the gradient descent in the form $\theta_{t+1} = \theta_t - \eta\nabla\mathcal{C}(\theta_t)$. Notice, the stochastic gradient descent (SGD) is obtained when $g_t = \nabla\mathcal{C}(\theta_t, \xi)$ where $\xi$ is a sample drawn from dataset $\mathcal{D}$ such that $\mathbb{E}[g_t] = \nabla\mathcal{C}(\theta_t)$ is a unbaised estimator of $\nabla\mathcal{C}(\theta_t)$ and

$$\theta_{t+1} = \theta_t - \eta g_t \tag{9}$$

Optimization problem 8 or 9 is well-suited to assumptions regarding the objective function $\mathcal{C}$ which involve the Euclidean norm. Th intuition behind optimization in Equation 8 is objective function $\mathcal{C}$ is replaced by its linearization at $\theta_t$ plus a Euclidean distance term $\frac{1}{2\eta} \|\theta - \theta_t\|_2^2$, which prevents the next iterate $\theta_{t+1}$ from being too far from $\theta_t$.

Instead of using Vanilla SGD above, recent studies tackle the optimization problem 7 by using natural gradient descent, where we update the parameter as

$$\theta_{t+1} = \theta_t - \eta F(\theta)^{-1} g_t \tag{10}$$

Here, $F(\theta) = \Re[G(\theta)]$ is Fubini-Study metric tensor a $P \times P$ matrix recently identified as the (classical) Fisher information matrix. We define quantum geometric tensor $G(\theta)$ as

$$G_{i,j} = \left\langle \frac{\partial \psi}{\partial \theta_i}, \frac{\partial \psi}{\partial \theta_j} \right\rangle - \left\langle \frac{\partial \psi}{\partial \theta_i}, \psi \right\rangle \left\langle \psi, \frac{\partial \psi}{\partial \theta_j} \right\rangle \quad (11)$$

Seminar work [14] demonstrates the block-wise Fubini-Study metric tensor can be evaluated in terms of quantum expectation values of Hermitian observables which is thus experimentally realizable. To evaluate quantum geometric tensor, [42] reports that it requires total $\mathcal{O}(P^2)$ rounds of sampling. Each round involves a repeated but parallelizable evaluation of an ansatz circuit with $\mathcal{O}(P)$ gates. When the global phase of the ansatz state $|\psi\rangle$ is independent of the parameters, the second term of $G_{i,j}$ vanish and $G(\theta)$ becomes equivalent to the coefficient matrix in imaginary time evolution [18], [43]. The advantage of using QNG in VQE is that it can easily measure the distinguishability of the objective function in non-Euclidean parameter space.

### C. Barren Plateau Problem

The gradient ($\nabla C(\theta)$) or stochastic gradient ($g$) plays an essential role in the parameter optimization process via the gradient-based method. However, it has been shown recently that the cost function exhibits a BP where the cost function gradient vanishes exponentially with the system size (the number of qubits). This phenomenon is also identified as cost concentration in [44]. Here, we consider here the following generic definition of a barren plateau without loss of generality.

**Definition 1.** *(Barren Plateau). Consider the cost function $\mathcal{C}(\theta)$ defined in Eq 6. This cost exhibits a barren plateau if, for all $\theta_i \in \theta$, the expectation value of partial derivative $\partial_i \mathcal{C}_i(\theta) = \partial \mathcal{C}(\theta)/\partial \theta_i$ respect to the cost function is zero i.e., $\mathbb{E}[\partial_i \mathcal{C}_i(\theta)] = 0$. The variance of the above partial derivative vanishes exponentially with the number of qubits, i.e,*

$$Var_\theta[\partial_i \mathcal{C}_i(\theta)] \in \mathcal{O}(p^{-n}) \quad (12)$$

*for some $p > 1$.*

Notice, we have the following conclusion by using Chebyshev's inequality

$$P(|\partial_i \mathcal{C}_i(\theta)| \geq c) \leq \frac{\text{Var}_\theta[\partial_i \mathcal{C}_i(\theta)]}{c^2} \quad (13)$$

for some constant $c$. The definition and above inequality tells that the probability of finding a $\partial_i \mathcal{C}_i(\theta)$ that is larger than $c$ decreases exponentially when the variance of the partial derivative establishes an exponential decay. The presence of BPs exists in both deep unstructured PQC with randomly initialized parameters [19] and QNNs [22]. Ref. [22] theoretically analysis the BP based on the fact that when ansatzes become unitary 2-designs [45], the expected number of samples required to estimate $\partial C(\theta)$ is exponential in the system size which often refers the number of qubits $n$. A circuit forms unitary 2-design means that the distribution matches up to the second moment that of the uniform Haar distribution of
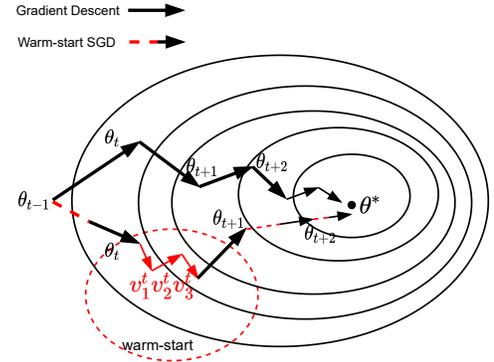


Fig. 1. A demonstration of gradient trajectory of gradient descent and warm-start QNG. Circle in red indicates the parameter re-initialization for next step parameter update.

unitaries, which will be used in the analysis of $\partial_i C(\theta)$). For example the partial derivative of $\theta_k$ in $l$-th layer respect to cost function in 7 can be explicitly described as

$$\partial_k \mathcal{C}(\theta) = i\langle 0|U_-^\dagger [V_k U_+^\dagger H U_+] U_-|0\rangle \quad (14)$$

where $U_- = \prod_{l=k-1}^1 U_l(\theta_l)$ and $U_+ = \prod_{l=L}^k U_l(\theta_l)$. BP is fatal in gradient-based optimization because it might halt the parameter update and quickly converge to some sub-optimal solution. In other words, one needs exponential resources to sample the gradient. Therefore, optimizing parameters in the BP region with the gradient-based method becomes hard.

## IV. MAIN RESULTS

In this section, we show that QNG corresponding to quantum probability space can be implemented as a classical first-order optimization known as mirror descent. Then we show the proposed LAWS algorithm and a general WSSGD framework.

### A. Quantum Information Geometry of Mirror Descent

In the seminar work [25], the Euclidean distance term $\frac{1}{2\eta}\|\theta - \theta_t\|_2^2$ in Equation 8 has been replaced with a general distance function $D_\Phi(\cdot, \cdot)$, i.e.,

$$D_\Phi(\theta_1, \theta_2) = \Phi(\theta_1) - \Phi(\theta_2) - \langle \nabla\Phi(\theta_2), \theta_1 - \theta_2 \rangle \quad (15)$$

where $\Phi(\cdot)$ is a carefully chosen continuously differentiable, strictly convex proximity function defined on some convex set. Notice, $D_\Phi(\theta_1, \theta_2) \geq 0$ with $D_\Phi(\theta_1, \theta_1) = 0$. $D_\Phi(\cdot, \cdot)$ defined above is also known as *Bregman divergence*, which is widely used in statistical inference, optimization, machine learning, and information geometry. As a result, a generalization of stochastic iterative optimization 8 has following

$$\theta_{t+1} = \arg\min_\theta \left\{ \langle \theta, g_t \rangle + \frac{1}{\eta} D_\Phi(\theta, \theta_t) \right\} \quad (16)$$

The above optimization is known as mirror descent (MD) [46] with proximity function $D_\Phi$. Note, if $\Phi(\theta) = \frac{1}{2}\|\theta\|_2^2$ convex, then $D_\Phi(\theta, \theta_t) = \frac{1}{2}\|\theta - \theta_t\|_2^2$ yields the standard gradient descent update 9. In addition, many modern machine

learning optimizations such as Vanilla SGD, AdaGrad and Adam [32] fall into MD 16 point view. For example, given Mahalanobis distance $\Phi(\theta) = \theta^\top A\theta$ where $A \succ 0$ is a positive (semi)definite matrix, i.e., $A = \sqrt{\sum_{i=1}^{t} g_i^2}$ a sum of all gradients for $t = 1$ to $t$. We have AdaGrad

$$\theta_{t+1} = \arg\min_{\theta} \left\{ \langle \theta, g_t \rangle + \frac{1}{2\eta}(\theta - \theta_t)^\top A(\theta - \theta_t) \right\} \quad (17)$$

which is equivalent to

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\sum_{i=1}^{t} g_i^2} + \epsilon} \odot g_t \quad (18)$$

where $\epsilon$ is a small number, typically set as $10^{-8}$, $\odot$ indicates the element-wise product. Moreover, if

$$A = \sqrt{(1 - \beta_2) \sum_{i=1}^{t} \beta_2^{t-i} g_i^2} \quad (19)$$

and set $m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t$ as exponential moving average (EMA) of stochastic gradient $g_t$ with $\beta_1, \beta_2 \in \mathbb{R}$ (typical values are $\beta_1 = 0.9$ $\beta_2 = 0.999$). We recover the Adam optimizer

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{(1 - \beta_2) \sum_{i=1}^{t} \beta_2^{t-i} g_i^2} + \epsilon} \odot m_t \quad (20)$$

In VQA, consider a parametric family of strictly positive probability distributions $p_\theta(x)$ parametrized by $\theta \in \mathbb{R}^d$ where $x \in [N]$ is a set of probability distributions on $N$ elements $[N] = \{1, \cdots, N\}$ and satisfies the normalization condition

$$\int p_\theta(x)dx = 1 \text{ for all } \theta \quad (21)$$

Assuming sufficient regularity, the derivatives of such densities satisfy the identity

$$\forall t > 0 \quad \int \frac{\partial^t p_\theta(x)}{\partial \theta^t} dx = \frac{\partial^t}{\partial \theta^t} \int p_\theta(x)dx = \frac{\partial^t 1}{\partial \theta^t} = 0 \quad (22)$$

To elucidate the geometry of the probability space $P$, we measure the density $p_\theta$ changes when one adds a small quantity $d\theta$ to its parameter. It can be achieved in a statistically meaningful way by using the Kullback-Leibler (KL) divergence [47]. Interestingly, KL-divergence is also a instance of Bregman divergence mentioned in Equation 15 by letting proximity function $\Phi(\theta) = \sum_i \theta_i \log(\theta_i)$ result in

$$D_\Phi(\theta, \theta + d\theta) = KL(\theta \| \theta + d\theta) = \mathbb{E}_{p_\theta} \left[ \log \left( \frac{p_\theta(x)}{p_{\theta+d\theta}(x)} \right) \right] \quad (23)$$

where $\mathbb{E}_{p_\theta}$ denotes the expectation with respect to the distribution $p_\theta$. Further, we can approximate the divergence with a second-order Taylor expansion such as

$$KL(\theta \| \theta + d\theta) = \mathbb{E}_{p_\theta} \left[ \log(p_\theta(x)) - \log(p_{\theta+d\theta}(x)) \right]$$
$$\approx -d\theta^\top \mathbb{E}_{p_\theta} \left[ \frac{\partial \log(p_\theta(x))}{\partial \theta} \right] + \frac{1}{2} d\theta^\top \mathbb{E}_{p_\theta} \left[ \frac{\partial^2 \log(p_\theta(x))}{\partial \theta^2} \right] d\theta \quad (24)$$

Applying the fact that first-order term is 0 shown in Equation 22, we have

$$D_\Phi(\theta, \theta + d\theta) = KL(\theta \| \theta + d\theta) \approx \frac{1}{2} d\theta^\top F(\theta) d\theta \quad (25)$$

$F(\theta)$ is defined by the Fisher information matrix (FIM)

$$F(\theta) = \mathbb{E}_{p_\theta} \left[ \frac{\partial^2 \log(p_\theta(x))}{\partial \theta^2} \right]$$
$$= \mathbb{E}_{p_\theta} \left[ \left( \frac{\partial \log(p_\theta(x))}{\partial \theta} \right) \left( \frac{\partial \log(p_\theta(x))}{\partial \theta} \right) \right] \quad (26)$$

We notice the second equality of $F(\theta)$ is often preferred because it makes clear that the $F(\theta)$ is symmetric and always positive semidefinite, though not necessarily positive definite. Finally, we plug the Bregman divergence defined on information entropy $\Phi(\theta) = \sum_i \theta_i \log(\theta_i)$ in MD optimization 16, and we have

$$\theta_{t+1} = \arg\min_{\theta} \left\{ \langle \theta, g_t \rangle + \frac{1}{2\eta}(\theta - \theta_t)^\top F(\theta - \theta_t) \right\} \quad (27)$$

The iterative solution of the above optimization problem 27 is

$$\theta_{t+1} = \theta_t - \eta F^{-1} g_t \quad (28)$$

where $F^{-1}$ is the pseudo-inverse of the Fisher information matrix, which recovers the natural gradient descent in 10. In the over-parameterized classical deep learning model, $F$ is singular. To make it invertible, one often adds a non-negative damping term $\delta$ such that $\theta_{t+1} = \theta_t - \eta(F + \delta I)^{-1} g_t$.

### B. Look Around Warm-start Natural Gradient

The presence of BP becomes one of the major bottlenecks in optimizing VQA, such as deep QNN. Notably, this does not preclude VQA, allowing for efficient gradient-based optimization. In section II-C, we discuss two mainstream techniques to mitigate BP. This work focuses on the optimization solution combined with the QNN parameter initialization strategy.

*1) Motivation:* Many studies in both classical and quantum regimes have shown that the parameter initialization could significantly improve the model performance and accelerate the training. Intuitively, a good parameter initialization (i.e., the distribution of initialized parameter close to optimal) requires a large amount of empirical studies, hyper-parameter tuning, and possibly human intervention, which is unproductive. Therefore, a natural question is raised: *can we perform efficient and effective parameter initialization for QNNs?* This is the primary motivation behind our approach. Second, the one-shot model initialization strategy initializes the model only at the beginning of the training process. However, as the training process proceeds, BP appears again when we use the gradient-based method to train the model. So, the questions, such as *can we design a hybrid method that effectively trains QNNs while adopting the superiority of parameter initialization*, or *can we periodically perform parameter initialization during the optimizing phase to mitigate the BP?* are also motivating us to explore various optimization strategies for QNNs. Besides finding a suitable initialization strategy, we also consider

**Algorithm 1:** Look Around and Warm-Start Natural Gradient Descent

---

**1** Input: Objective function $\mathcal{C}(\theta)$, learning data $\mathcal{D}$

**2** Initialization: $\theta_0$, learning rate $\eta_0$, warm-start learning rate $\mu_0$, warm-start iteration $K$ ($K = 5, 3,$ or $2$)

**3** **for** $t = 1, \cdots, T$ **do**

**4** $\quad$ $v_0^t = \theta_{t-1}$

**5** $\quad$ **for** $k = 1, \cdots, K$ **do**

**6** $\quad\quad$ Draw sample from batch data $\xi \sim \mathcal{D}_b$

**7** $\quad\quad$ $v_k^t = v_{k-1}^t - \mu_k \nabla C(\theta; \xi)$

**8** $\quad$ **end**

**9** $\quad$ $\theta_{\text{warm-start}}^t = v_k^t$

**10** $\quad$ Compute natural gradient $F_t = \text{FisherIM}(\theta_{\text{warm-start}}^t)$

**11** $\quad$ Compute new gradient $g_t = \theta_{\text{warm-start}}^t - \theta_{t-1}$

**12** $\quad$ Update parameter:

**13** $\quad$ $\theta_t = \theta_{\text{warm-start}}^t - \frac{\eta_t}{K} F_t^{-1} g_t$

**14** **end**

---

the algorithm's efficiency since computing quantum Fisher information in QNG is expensive, as discussed in section III-B.

*2) Proposed Method:* Here, we present our proposed algorithm shown in Algorithm 1. The key step in the proposed algorithm, in short, is that we perform the initialization after every parameter update instead of only initializing the PQC one at a time. The intuition behind this algorithm is that we try to warm-start the natural gradient descent for each iteration. Every time the optimizer finds a sub-optimal solution, say $\theta_t$, we utilize this $\theta_t$ and re-initialize the model around $\theta_t$ within a small region. In this way, we couple the parameter re-initialization and gradient descent method together, and the LAWS finds the (maybe) optimal solution in a "*look around*" manner. Later, we generalized the LAWS to accommodate all existing gradient-based methods in Algorithm 2.

There are two major advantages when using LAWS. First, LAWS could mitigate the BP issue by repeatedly performing the parameter re-initialization, where our empirical results also support this observation. Second, LAWS adopts a fast convergency speed, and it is more computationally efficient than the QNG. In the basic LAWS design, the parameter re-initialization is performed by another low-cost first-order optimization in fewer steps (usually $K = 5$). The expensive quantum Fisher information matrix evaluation is completed after parameter re-initialization. Assume we require total $T$ iterations for model training. We only evaluate qFIM for $T/K$ rounds. Third, we empirically find that LAWS achieves better generalization ability in the classification learning task.

*3) Implementation Details:* The implementation of LAWS is simple and is compatible with all existing gradient-based optimization frameworks. Therefore, how to effectively and efficiently perform warm-start (parameter re-initialization) is the key challenge in LAWS's design. Here, effective warm-

start means we expect the periodically re-initialized parameter to be close to the region where the possible optimal parameter resides. Furthermore, efficient warm-start means parameter re-initialization should be computational inexpensive and possibly without hyper-parameter tuning. To this end, the design of warm-start is based on a stochastic procedure, where the re-initialized parameter is sampled from a set of stochastic gradients. Fig 1 demonstrates the optimization trajectory of LAWS compared to the original QNG. We search gradients for fewer steps around the current optimal $\theta_t$ and then perform a natural gradient descent step ($F_t^{-1}$) on the accumulation of the previous gradient ($v_k^t - \theta_t$) at the re-initialized parameter point $v_k^t$.

We present two different warm-start strategies. The first one uses a K-step (K usually small, such as $K = 5$) inner loop (as the Algorithm 1 shows) to compute a set of K consecutive gradients such as $\mathcal{G}_K^t = \{g_1^t, g_2^t, \cdots, g_K^t\}$. Then, we compute a weighted average of gradients in $\mathcal{G}_K^t$ as a warm-start point of the next iteration

$$\theta_{\text{warm-start}} = \theta_{t-1} + \frac{1}{K} \sum_{k=1}^{K} g_k \tag{29}$$

for all $g_k^t \in \mathcal{G}_K^t$ such that

$$\mathcal{G}_K^t = \left\{ \nabla \mathcal{C}(v_k^t, \xi) | \xi \sim \mathcal{D}_b \right\} \tag{30}$$

where each $v_k^t$ is computed as line 9 in Algorithm 1. The second one uses a K-step inner loop to sample gradient candidates. But one significant difference compared to the first method is that sample gradient candidates are computed with respect to the same model parameter at current step $t-1$, say $\theta_{t-1}$. Mathematically, we have

$$\theta_{\text{warm-start}} = \theta_{t-1} + \frac{1}{K} \sum_{k=1}^{K} \nabla \mathcal{C}(\theta_{t-1}, \xi) \text{ where } \xi \sim \mathcal{D}_b \tag{31}$$

The two methods above ensure that the warm-start parameter is not far from the current optimal solution. The re-initialization is achieved using scholastic gradients sampled from randomly selected samples in (min) batch training data. It is worth mentioning that the FIM ($F$) in LAWS is evaluated on the warm-start parameter $\theta_{\text{warm-start}}$. In Algorithm 1, line 10 (or Algorithm 2, line 8), we explicitly show the $\theta_{\text{warm-start}} = v_k^t$ for clear presentation. In an actual implementation, we directly use $v_k^t$ as a warm-start parameter to save the memory usage.

Another discussion for Algorithm 1 is the gradient defined as $g_t = \theta_{\text{warm-start}}^t - \theta_{t-1}$ (line 11). LAWS simply uses the accumulation of all past gradients in $\mathcal{G}_K^t$ defined in 30. The past gradients' accumulation re-scales the gradient and may increase the magnitude of the new one compared to using a single gradient. We believe this is also a key point in making LAWS mitigate the BP. One could also perform this accumulation by using exponential moving average (EMA), similar to Adam settings, to apply more smooth gradient such

---

**Algorithm 2:** Generalized Warm-start Stochastic Gradient Descent (WSSGD)

---

**1** Input: Objective function $\mathcal{C}(\theta)$, learning data $\mathcal{D}$, warm-start optimizer $\mathcal{W}$, reparameterization coefficient function $\Delta$

**2** Initialization: $\theta_0$, warm-start learning rate $\mu_0$, warm-start iteration $K$

**3 for** $t = 1, \cdots, T$ **do**

**4**     $v_0^t = \theta_{t-1}$

**5**     **for** $k = 1, \cdots, K$ **do**

**6**        $v_k^t = \mathcal{W}(\mathcal{C}(\theta), v_0^t, \mu_0, \mathcal{D}_m)$

**7**     **end**

**8**     $\theta_{\text{warm-start}}^t = v_k^t$

**9**     Update parameter:

**10**     $\theta_{t+1} = \Delta_t \theta_t + (1 - \Delta_t)\theta_{\text{warm-start}}^t$

**11 end**

---

as

$$g_t = (1 - \beta) \sum_{i=1}^{K} \beta^{K-i} g_i^t \text{ for all } g_i^t \in \mathcal{G}_K^t \quad (32)$$

where we introduce a new hyper-parameter $\beta \in (0,1)$. Obviously, the process can be replaced with some gradient noise reduction techniques widely used in various classical machine learning applications. This is beyond the scope of this literature, and we leave this for future work. We empirically evaluate the above-mentioned warm-start strategies. More detailed results and analysis are shown in section V.

We notice that the proposed LAWS belongs to a certain first-order optimization in modern classical learning regime so-called *Lookahead optimizer*. Based on their extraordinary work, we propose a general warm-start framework for VQA in the next section.

*4) Generalized Warm-start Algorithm:* In the classical machine learning study, [48] proposed a new optimization algorithm named Lookahead. Lookahead is orthogonal to those aforementioned approaches [32] due to the different parameter update settings. The core idea of Lookahead is to maintain two kinds of model parameters, i.e., "fast parameter" $v_k^t$ and "slow parameter" $\theta_t$, and jointly update them. Specifically, the inner loop takes the slow weights ($\theta_{t-1}$) as initial point and updates the fast weights ($v_k^t$) $K$ times to receive $v_K^t$; while the outer loop updates the slow weights as

$$\theta_t = (1 - \alpha)\theta_{t-1} + \alpha v_K^t \quad \alpha \in (0,1) \quad (33)$$

Any standard optimizer, e.g., Vallina SGD, AdaGrad, and Adam, can serve as the inner-loop optimizer. In our speech, the inner-loop act as a warm-start initialization. In this way, the Lookahead optimizer achieves remarkable performance improvement over the standard optimizer. Further, due to its simplicity in implementation, negligible computation and memory cost, and compatibility with almost current ML libraries, Lookahead has been widely adopted.

Interestingly, we find LAWS also falls into this line of research. In algorithm 1, let $\lambda_t = \eta_t/K$, we compute the $\theta_t$ as

$$\begin{aligned} \theta_t &= \theta_{\text{warm-start}}^t - \lambda_t F_t^{-1} g_t \\ &= \theta_{\text{warm-start}}^t - \lambda_t F_t^{-1}(v_k^t - \theta_{t-1}) \quad (34) \\ &= \lambda_t F_t^{-1}\theta_{t-1} + (1 - \lambda_t F_t^{-1})\theta_{\text{warm-start}}^t \end{aligned}$$

the last equality is due to $v_k^t = \theta_{\text{warm-start}}^t$. From the above derivation, We see that the mathematical difference between LAWS and Lookahead is: we replace $\alpha$ in Lookahead 33 to some value such as

$$\alpha = 1 - \lambda_t F_t^{-1}$$

, where is not a fixed real coefficient but a Fisher information related quantity.

We can reveal a couple of insights from this observation. First, this modification makes the proposed algorithm LAWS significantly different from the Lookahead in terms of parameter update and design logic. Second, when the Fisher information matrix degenerates to the identity matrix $I$, we recover the standard Lookahead optimizer. On the other hand, one could use the Hessian matrix to replace $F$, where we have a Newton-type optimizer. Thrid, essentially, the information on the curvature of the cost surface is encoded in FIM for the Riemannian manifold. FIM can be interpreted as some specific "step size" in first-order optimization. In LAWS, the FIM reparameterizes the model according to its second-order information. The coefficient $\alpha$, therefore, may not be a fixed real number now.

To this end, we propose our unified framework WSSGD for a warm-start stochastic gradient descent algorithm for QNNs as shown in Algorithm 2. We first employ a generalized optimizer $\mathcal{W}$ for the warm-start inner loop. The choice of such an optimizer heavily affects re-initialization and model performance. As reported in [48], WSSGD may benefit from a larger learning rate in the inner loop. In other words, we could use a larger step size $\mu$. We also propose the general form for reparameterization coefficient $\Delta_t$ as a function of gradient, for example

$$\begin{aligned} \text{Lookahead: } & \Delta_t = 1 - \alpha \\ \text{WS-SGD: } & \Delta_t = 1 - \lambda_t F_t \\ \text{Adam-like SGD: } & \Delta_t = 1 - \lambda_t \sqrt{\sum_k \nabla C(v_k^t, \xi)^2} \end{aligned} \quad (35)$$

Moreover, WSSGD is simple and compatible with many VQA libraries such as Pennylane, Tenoerflow Quantum, etc. Our empirical results are present in section V, we conclude that WSSGD achieves faster convergence rates (i.e., smaller optimization error), and enjoys smaller generalization error. We release our source code implemented by Pennylane at https://github.com/taozeyi1990/LAWS

## C. Convergence Analysis

In this section, we present the convergence and generalization analysis of the proposed algorithm. We first provide some useful definitions and assumptions which have been widely adopted in classical machine learning. We provide the analysis of the convergence for both convex and non-convex objective functions $\mathcal{C}(\theta)$. We start by showing the proof of convergence on the convex problem to give some intuition first and then give the proof on a more realistic non-convex problem.

*1) Assumption:* The assumptions we are making are

**Assumption 1.** *(Bounded gradient). The function $\mathcal{C}(\theta)$ has bounded (stochastic) gradients, i.e., for any $\theta \in R^d$ we have*

$$||\nabla\mathcal{C}(\theta,\xi)||_2 \leq G \text{ for all } \xi \sim \mathcal{D} \tag{36}$$

**Assumption 2.** *(L-Lipschitz smooth). The function $\mathcal{C}(\theta)$ is L-Lipschitz smooth i.e.*

$$||\nabla\mathcal{C}(\theta) - \nabla\mathcal{C}(\mu)||_2 \leq L||\theta - \mu|| \text{ for all } \theta, \mu \in R^d \tag{37}$$

**Assumption 3.** *(M-Lipschitz continuous). The function $\mathcal{C}(\theta)$ is L-Lipschitz continous i.e.*

$$||\mathcal{C}(\theta) - \mathcal{C}(\mu)||_2 \leq M||\theta - \mu|| \text{ for all } \theta, \mu \in R^d \tag{38}$$

We also define Polyak-Lojasiewicz (PL) condition as

**Definition 2.** *(PL Condition) Let $\theta^* \in \arg\min_\theta \mathcal{C}(\theta)$. We say a function $\mathcal{C}(\theta)$ satisfies $\sigma$-PL condition if*

$$\sigma(\mathcal{C}(\theta) - \mathcal{C}(\theta^*)) \leq ||\nabla\mathcal{C}(\theta)||^2 \tag{39}$$

*with some constant $\sigma$.*

The above assumptions and definition can be easily obtained and verified in VQE. Now, we show our theoretical results below in the Theorem 1 and Theorem 2.

*2) Convex Objective function:* Given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ where $(x_i, y_i)$ is is drawn from an unknown distribution, one often minimizes the empirical risk $\mathcal{L}(\theta) = \frac{1}{n}\sum_{i=1}^n \mathcal{C}(\theta, x_i, y_i)$ via a randomized algorithm, e.g. SGD, to find an estimated optimum $\theta_T \in \arg\min_\theta L(\theta)$. However, this empirical solution $\hat{\theta}$, differs from the desired optimum $\theta^*$ of the population risk

$$\theta^* \in \arg\min_\theta L(\theta, \mathcal{D}) = \mathbb{E}_{x,y\sim\mathcal{D}}[\mathcal{C}(\theta, x)] \tag{40}$$

To begin with, we first investigate the convergence performance of WSSGD when its warm-start optimizer $\mathcal{W}$ is SGD. We summarize our main results in Theorem 1 below.

**Theorem 1.** *(Convex) Suppose the objective function $\mathcal{C}(\theta)$ is gamma-strongly convex, M-Lipschitz continuous, and L-Lipschitz smooth w.r.t., $\theta$. Let $\theta^* = \arg\min_\theta \mathcal{C}(\theta)$. Let warstart learning rate $\mu_k^t = \frac{c_0}{((t-1)k+K+2)}, c_0 \in (0,1]$, the optimization error of the output $\theta_T$ of WS-SGD satisfies*

$$\mathbb{E}[\mathcal{C}(\theta_T) - \mathcal{C}\theta^*)] \leq \frac{e^{2\Delta}L(k+2)^{2\Delta}}{2((T+1)K+2)^{2\Delta}}||\theta_0 - \theta^*||^2 + \frac{16LG^2}{c_0^2((T+1)K+2)^{2(1-\Delta)}(2\Delta-1)} \tag{41}$$

*Proof.* (Proof sketch) Due to the page limitation, we present the proof sketch to show some core results. This proof mainly follows the proof of Theorem 1 in [49]. We first bound the inner loop $\mathbb{E}||v_K^t - \theta^*||$, where we have $\mathbb{E}[||v_K^t - \theta^*||^2] \leq (1 - \gamma\mu_{K-1}^t)\mathbb{E}[||v_{K-1}^t - \theta^*||^2] + (\mu_{K-1}^t G)^2$. Let $\mu_k^t = \frac{c_0}{((t-1)k+K+2)}$, for some constant $c_0 > 0$, we roll the above recurrence relation from $k = 1$ to $K$, we have $\mathbb{E}[||v_K^t - \theta^*||^2] \leq \left(\left(\frac{(t-1)k+2}{tk+2}\right)^2\mathbb{E}[||\theta_{t-1} - \theta^*||^2] + \frac{16c_0G^2}{c_0^2(tk+2)^2}\right)$. The outer loop parameter update follows $\mathbb{E}[||\theta_t - \theta^*||^2] \leq \Delta\mathbb{E}[||\theta_{t-1} - \theta^*||^2] + (1 - \Delta)\mathbb{E}[||v_K^t - \theta^*||^2]$ according to line 10 in Algorithm 2. We have $\mathbb{E}[||\theta_t - \theta^*||^2] \leq \left(\Delta + (1-\Delta)\left(\frac{(t-1)K+2}{tK+2}\right)^2\right)\mathbb{E}[||\theta_{t-1} - \theta^*||^2] + \frac{16c_0(1-\Delta)G^2}{c_0^2(tK+2)^2}$. Again, unwinding this recurrence relation from $t = 1$ to $T$, yields $\mathbb{E}[||\theta_T - \theta^*||^2] \leq \frac{e^{2\Delta}(k+2)^{2\Delta}}{((T+1)K+2)^{2\Delta}}||\theta_0 - \theta^*|| + \frac{32G^2}{c_0^2((T+1)K+2)^{2(1-\Delta)}(2\Delta-1)}$. Since $\mathcal{C}(\theta)$ is M-smooth and $\theta^*$ is its optimum, we have $\mathbb{E}[\mathcal{C}(\theta_T) - \mathcal{C}(\theta^*)] \leq \frac{M}{2}\mathbb{E}[||\theta_T - \theta^*||^2]$, yields the result in Theorem 1. $\square$

*3) Non-Convex Objective function:* To prove no-convex objective function, we use the Polyak-Lojasiewicz condition defined in 2.

**Theorem 2.** *(Non-Convex) Suppose the objective function $\mathcal{C}(\theta)$ is M-Lipschitz continuous, and L-Lipschitz smooth w.r.t., $\theta$. In addition, suppose $\mathcal{C}(\theta)$ satisfies $\sigma$-PL condition. Let $\mu_k^t = \frac{1}{tK+t+1}$, we have*

$$\mathbb{E}[\mathcal{C}(\theta_T) - \mathcal{C}(\theta^*)] \leq \frac{4}{(TK+1)^{2\Delta}}\mathbb{E}[\mathcal{C}(\theta_0) - \mathcal{C}(\theta^*)] + \frac{2\Delta MG^2C_0}{(TK+1)^{2\Delta-1}} \tag{42}$$

*where $C_0 = \Delta + (1-\Delta)(K-1)$.*

We again omit the complete proof of this theorem. The proof sketch is we first bound $\mathbb{E}[||\mathcal{C}(v_K^t) - \mathcal{C}(\theta^*)||]$, then we use the relation of $v_K^t$ and $\theta_t$ defined in the Algorithm 1 line 10 to derive the final bound of $\mathbb{E}[\mathcal{C}(\theta_T) - \mathcal{C}(\theta^*)]$. In the next section, we present the numerical simulation results of LAWS, WS-SGD, and their variants.

## V. NUMERICAL SIMULATIONS

To evaluate the performance of LAWS, WS-SGD, and their variants, we use the open-source library PennyLane [50] 0.22.2 built on Python 3.7. Most of the experiments follow the open-source tutorials from official PennyLane website. We first show that the LAWS could mitigate the BP issue with random PQC energy problems (in Figure 2). Then we evaluate LAWS in a more practical problem of estimating the ground state energy of the hydrogen molecule (in Figure 3). Further, we conduct experiments on variational quantum classifiers (in Figure 4 and Figure 5) such that the quantum circuits can be trained from labeled data to classify new data samples. The classification training data is public and can be downloaded from the PennyLane tutorial. All the experimental results and
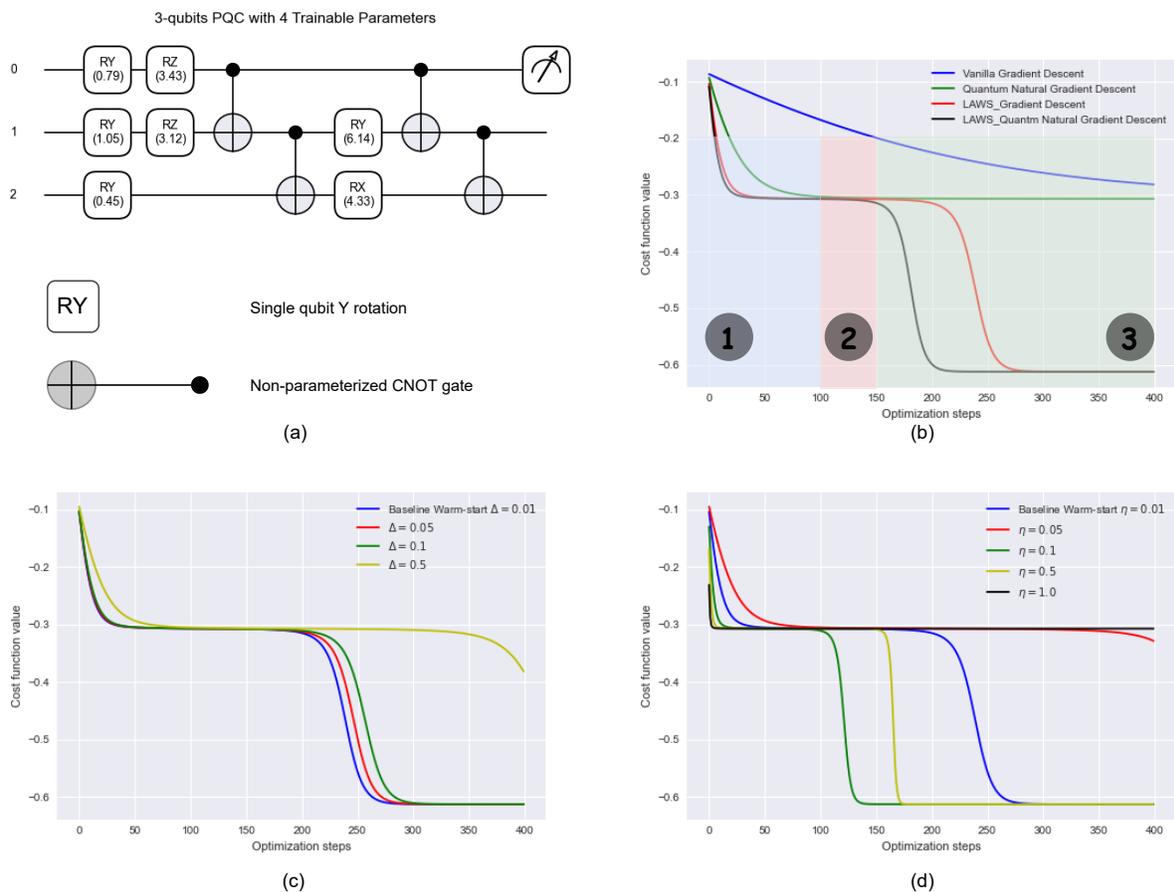
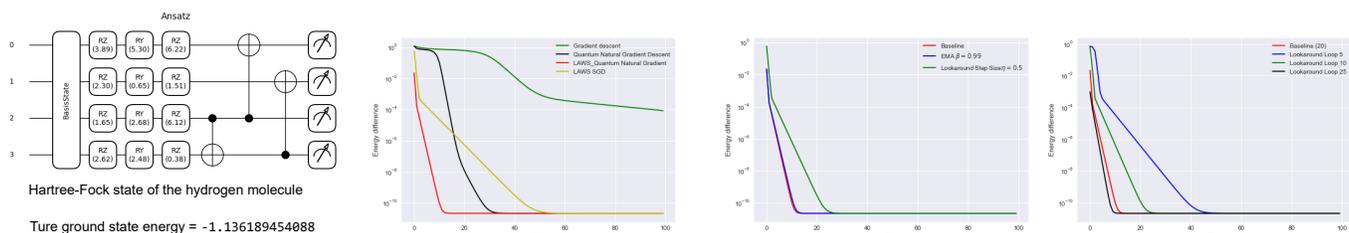Fig. 2. Evaluation on randomly designed PQC



Fig. 3. Performance on Hydrogen VQE

source code implementation can be found at https://github.com/taozeyi1990/LAWS.

### A. Evaluation on random PQC

Figure 2 above demonstrates the performance of LAWS compared to Vanilla SGD and Vanilla QNG in randomly designed PQC. The random PQC is shown on (a), where it contains 3 qubits and 4 trainable parameters. We perform the optimization for a total of 400 iterations. It is undeniable that LAWS outperforms Vanilia QNG and Vanilla SGD, as shown on (b). We divide the cost value region into 3 parts marked as blue ❶, red ❷, and green ❸. In region ❶, all methods are quickly converging, and LAWS is even faster. In region ❷,

methods are trapped because of BP, and there is no further cost value decreasing. Finally, in region ❸, LAWS and its variant mitigate the BP and find the optimal solutions.

Figure 2 (c) and (d) are showing the different configurations of LAWS for different $\Delta$'s and look around step size $\eta$'s. We see that the LAWS is sensitive in the change of $\eta$.

### B. Evaluation on Quantum Chemistry

Figure 3 shows the experimental results for a more practical problem of Hydrogen VQE. The primary purpose of this experiment is to approach ground state energy as close as we can. The exact value of the ground state energy defined by the above PQC is given by $-1.1361894$. The qubit register has
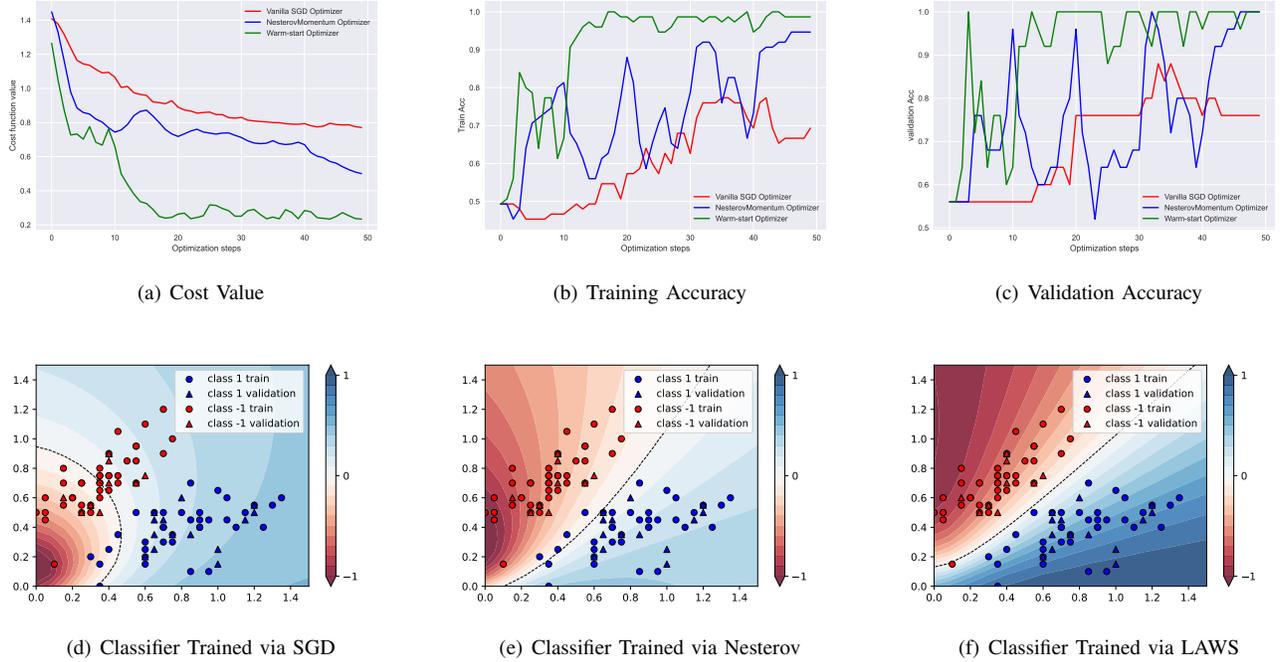
(a) Cost Value

(b) Training Accuracy

(c) Validation Accuracy



(d) Classifier Trained via SGD

(e) Classifier Trained via Nesterov

(f) Classifier Trained via LAWS

Fig. 4. Variational classifier for Iris classification task: SGD vs. Nesterov vs. LAWS



(a) costs of WS-SGD's variants

(b) Train Acc of WS-SGD's variants

(c) WS-SGD Warm-start with Random Sample

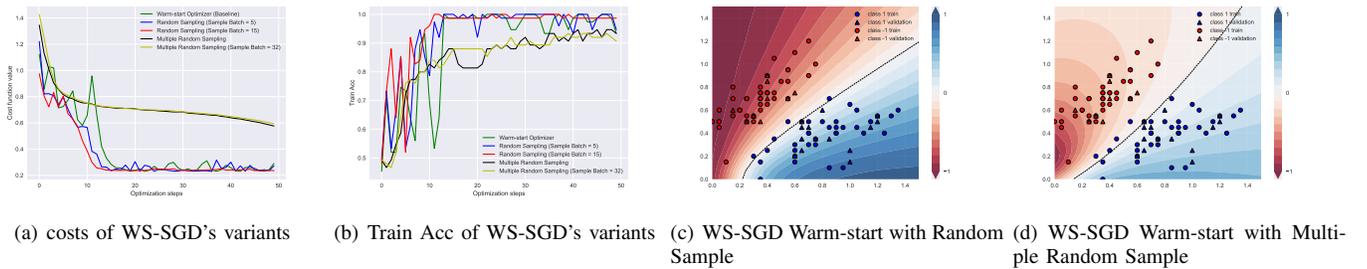(d) WS-SGD Warm-start with Multiple Random Sample

Fig. 5. Variational classifier for Iris classification task: Variants of LAWS

been initialized to $|1100\rangle$, which encodes for the Hartree-Fock state of the hydrogen molecule described on a minimal basis. We see the LAWS still outperforms than other state-of-the-art optimization methods.

*C. QNNs: Variational Classifier*

Last but least, we perform the binary Iris classification task, which is a simple but powerful QNN to show that the warm-start strategy has better generalization ability. The learning rate for SGD and Nesterov momentum optimizer is set to be 0.01. While the learning rate, look around rate and look around steps are 0.01, 0.5, and 5, respectively. We train QNN model within 50 iterations. Figure 4(a), 4(b), and 4(c) show the cost function value, training accuracy, and validation accuracy, respectively. As shown in each figure, the warm-start SGD in green demonstrates its superiority in this task. Figure 4(d), 4(e), and 4(f) indicate the decision boundaries of model trained with different optimizers. We observe that the

two classes in the train and validation dataset are perfectly being separated when using the warm-start optimizer, which indicates the WSSGD has stronger generalization ability. The result of the Nesterov optimizer seems to suffer from the under-fitted where the samples at the bottom left are mixed.

Moreover, we evaluate the different warm-start strategies discussed in section IV-B3 of WSSGD in Figure 5. We test the impacts on WSSGD of different configurations such as sampling batch size.

## VI. CONCLUSION

In this work, we propose an unified framework for QNG by using a classical first-order optimization scheme. The proposed new algorithm named WSSGD show its power in QVA learning. Our experiment results show that the proposed algorithm could mitigate the BP issue and have better generalization ability in quantum classification problems.

## REFERENCES

[1] J. Preskill, "Quantum computing in the nisq era and beyond," *Quantum*, vol. 2, p. 79, 2018.

[2] P. W. Shor, "Algorithms for quantum computation: discrete logarithms and factoring," in *Proceedings 35th annual symposium on foundations of computer science*. Ieee, 1994, pp. 124–134.

[3] S. Lloyd, "Universal quantum simulators," *Science*, vol. 273, no. 5278, pp. 1073–1078, 1996.

[4] A. W. Harrow, A. Hassidim, and S. Lloyd, "Quantum algorithm for linear systems of equations," *Physical review letters*, vol. 103, no. 15, p. 150502, 2009.

[5] 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.*, "Quantum supremacy using a programmable superconducting processor," *Nature*, vol. 574, no. 7779, pp. 505–510, 2019.

[6] M. Cerezo, A. Arrasmith, R. Babbush, S. C. Benjamin, S. Endo, K. Fujii, J. R. McClean, K. Mitarai, X. Yuan, L. Cincio *et al.*, "Variational quantum algorithms," *Nature Reviews Physics*, vol. 3, no. 9, pp. 625–644, 2021.

[7] J. R. McClean, J. Romero, R. Babbush, and A. Aspuru-Guzik, "The theory of variational hybrid quantum-classical algorithms," *New Journal of Physics*, vol. 18, no. 2, p. 023023, 2016.

[8] 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*, vol. 5, no. 1, pp. 1–7, 2014.

[9] E. Farhi, J. Goldstone, and S. Gutmann, "A quantum approximate optimization algorithm," *arXiv preprint arXiv:1411.4028*, 2014.

[10] E. Farhi and H. Neven, "Classification with quantum neural networks on near term processors," *arXiv preprint arXiv:1802.06002*, 2018.

[11] M. Cerezo, A. Sone, T. Volkoff, L. Cincio, and P. J. Coles, "Cost function dependent barren plateaus in shallow parametrized quantum circuits," *Nature communications*, vol. 12, no. 1, pp. 1–12, 2021.

[12] J. Tilly, H. Chen, S. Cao, D. Picozzi, K. Setia, Y. Li, E. Grant, L. Wossnig, I. Rungger, G. H. Booth *et al.*, "The variational quantum eigensolver: a review of methods and best practices," *arXiv preprint arXiv:2111.05176*, 2021.

[13] L. Bittel and M. Kliesch, "Training variational quantum algorithms is np-hard," *Physical Review Letters*, vol. 127, no. 12, p. 120502, 2021.

[14] J. Stokes, J. Izaac, N. Killoran, and G. Carleo, "Quantum natural gradient," *Quantum*, vol. 4, p. 269, 2020.

[15] N. Yamamoto, "On the natural gradient for variational quantum eigensolver," *arXiv preprint arXiv:1909.05074*, 2019.

[16] B. Koczor and S. C. Benjamin, "Quantum natural gradient generalised to non-unitary circuits," *arXiv preprint arXiv:1912.08660*, 2019.

[17] G. Carleo and M. Troyer, "Solving the quantum many-body problem with artificial neural networks," *Science*, vol. 355, no. 6325, pp. 602–606, 2017.

[18] S. McArdle, S. Endo, Y. Li, S. Benjamin, and X. Yuan, "Variational quantum simulation of imaginary time evolution with applications in chemistry and beyond," *arXiv preprint arXiv:1804.03023*, 2018.

[19] J. R. McClean, S. Boixo, V. N. Smelyanskiy, R. Babbush, and H. Neven, "Barren plateaus in quantum neural network training landscapes," *Nature communications*, vol. 9, no. 1, pp. 1–6, 2018.

[20] Y. Suzuki, H. Yano, R. Raymond, and N. Yamamoto, "Normalized gradient descent for variational quantum algorithms," in *2021 IEEE International Conference on Quantum Computing and Engineering (QCE)*. IEEE, 2021, pp. 1–9.

[21] T. Haug and M. Kim, "Optimal training of variational quantum algorithms without barren plateaus," *arXiv preprint arXiv:2104.14543*, 2021.

[22] E. Grant, L. Wossnig, M. Ostaszewski, and M. Benedetti, "An initialization strategy for addressing barren plateaus in parametrized quantum circuits," *Quantum*, vol. 3, p. 214, 2019.

[23] 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 preprint arXiv:2112.10952*, 2021.

[24] L. Franken, B. Georgiev, S. Muecke, M. Wolter, N. Piatkowski, and C. Bauckhage, "Gradient-free quantum optimization on nisq devices," *arXiv preprint arXiv:2012.13453*, 2020.

[25] A. S. Nemirovskii, D. B. Yudin, D. B. Iudin, and D. B. Iudin, *Problem complexity and method efficiency in optimization*. Wiley, 1983.

[26] D. Wierichs, C. Gogolin, and M. Kastoryano, "Avoiding local minima in variational quantum eigensolvers with the natural gradient optimizer," *Physical Review Research*, vol. 2, no. 4, p. 043246, 2020.

[27] M. Cerezo and P. J. Coles, "Impact of barren plateaus on the hessian and higher order derivatives," *arXiv e-prints*, pp. arXiv–2008, 2020.

[28] A. Arrasmith, M. Cerezo, P. Czarnik, L. Cincio, and P. J. Coles, "Effect of barren plateaus on gradient-free optimization," *Quantum*, vol. 5, p. 558, 2021.

[29] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *International conference on machine learning*. PMLR, 2013, pp. 1139–1147.

[30] A. W. Harrow and J. C. Napp, "Low-depth gradient measurements can improve convergence in variational hybrid quantum-classical algorithms," *Physical Review Letters*, vol. 126, no. 14, p. 140502, 2021.

[31] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization." *Journal of machine learning research*, vol. 12, no. 7, 2011.

[32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[33] J. Zhuang, T. Tang, Y. Ding, S. C. Tatikonda, N. Dvornek, X. Papademetris, and J. Duncan, "Adabelief optimizer: Adapting stepsizes by the belief in observed gradients," *Advances in neural information processing systems*, vol. 33, pp. 18 795–18 806, 2020.

[34] S.-I. Amari, "Natural gradient works efficiently in learning," *Neural computation*, vol. 10, no. 2, pp. 251–276, 1998.

[35] T. Jones and J. Gacon, "Efficient calculation of gradients in classical simulations of variational quantum algorithms," *arXiv preprint arXiv:2009.02823*, 2020.

[36] S. Soori, B. Can, B. Mu, M. Gürbüzbalaban, and M. M. Dehnavi, "Tengrad: Time-efficient natural gradient descent with exact fisher-block inversion," *arXiv preprint arXiv:2106.03947*, 2021.

[37] J. A. Nelder and R. Mead, "A simplex method for function minimization," *The computer journal*, vol. 7, no. 4, pp. 308–313, 1965.

[38] M. J. Powell, "An efficient method for finding the minimum of a function of several variables without calculating derivatives," *The computer journal*, vol. 7, no. 2, pp. 155–162, 1964.

[39] Z. Holmes, K. Sharma, M. Cerezo, and P. J. Coles, "Connecting ansatz expressibility to gradient magnitudes and barren plateaus," *PRX Quantum*, vol. 3, no. 1, p. 010313, 2022.

[40] M. Larocca, P. Czarnik, K. Sharma, G. Muraleedharan, P. J. Coles, and M. Cerezo, "Diagnosing barren plateaus with tools from quantum optimal control," *arXiv preprint arXiv:2105.14377*, 2021.

[41] D. Wecker, M. B. Hastings, and M. Troyer, "Progress towards practical quantum variational algorithms," *Physical Review A*, vol. 92, no. 4, p. 042303, 2015.

[42] B. van Straaten and B. Koczor, "Measurement cost of metric-aware variational quantum algorithms," *PRX Quantum*, vol. 2, no. 3, p. 030324, 2021.

[43] X. Yuan, S. Endo, Q. Zhao, Y. Li, and S. C. Benjamin, "Theory of variational quantum simulation," *Quantum*, vol. 3, p. 191, 2019.

[44] A. Arrasmith, Z. Holmes, M. Cerezo, and P. J. Coles, "Equivalence of quantum barren plateaus to cost concentration and narrow gorges," *arXiv preprint arXiv:2104.05868*, 2021.

[45] A. W. Harrow and R. A. Low, "Random quantum circuits are approximate 2-designs," *Communications in Mathematical Physics*, vol. 291, no. 1, pp. 257–302, 2009.

[46] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, "Robust stochastic approximation approach to stochastic programming," *SIAM Journal on optimization*, vol. 19, no. 4, pp. 1574–1609, 2009.

[47] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning," *Siam Review*, vol. 60, no. 2, pp. 223–311, 2018.

[48] M. Zhang, J. Lucas, J. Ba, and G. E. Hinton, "Lookahead optimizer: k steps forward, 1 step back," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[49] P. Zhou, H. Yan, X. Yuan, J. Feng, and S. Yan, "Towards understanding why lookahead generalizes better than sgd and beyond," *Advances in Neural Information Processing Systems*, vol. 34, 2021.

[50] V. Bergholm, J. Izaac, M. Schuld, C. Gogolin, M. S. Alam, S. Ahmed, J. M. Arrazola, C. Blank, A. Delgado, S. Jahangiri *et al.*, "Pennylane: Automatic differentiation of hybrid quantum-classical computations," *arXiv preprint arXiv:1811.04968*, 2018.