# Normalized Gradient Descent for Variational Quantum Algorithms

Yudai Suzuki*, Hiroshi Yano†, Rudy Raymond‡§, and Naoki Yamamoto†§

* Department of Mechanical Engineering, Keio University, Hiyoshi 3-14-1, Kohoku, Yokohama 223-8522, Japan
† Department of Applied Physics and Physico-Informatics, Keio University, Hiyoshi 3-14-1, Kohoku, Yokohama 223-8522, Japan
‡ IBM Quantum, IBM Japan, 19-21 Nihonbashi Chuo-ku, Tokyo 103-8510, Japan
§ Quantum Computing Center, Keio University, Hiyoshi 3-14-1, Kohoku, Yokohama 223-8522, Japan

*Abstract*—**Variational quantum algorithms (VQAs) are promising methods that leverage noisy quantum computers and classical computing techniques for practical applications. In VQAs, the classical optimizers such as gradient-based optimizers are utilized to adjust the parameters of the quantum circuit so that the objective function is minimized. However, they often suffer from the so-called vanishing gradient or barren plateau issue. On the other hand, the normalized gradient descent (NGD) method, which employs the normalized gradient vector to update the parameters, has been successfully utilized in several optimization problems. Here, we study the performance of the NGD methods in the optimization of VQAs for the first time. Our goal is twofold. The first is to examine the effectiveness of NGD and its variants for overcoming the vanishing gradient problems. The second is to propose a new NGD that can attain the faster convergence than the ordinary NGD. We performed numerical simulations of these gradient-based optimizers in the context of quantum chemistry where VQAs are used to find the ground state of a given Hamiltonian. The results show the effective convergence property of the NGD methods in VQAs, compared to the relevant optimizers without normalization. Moreover, we make use of some normalized gradient vectors at the past iteration steps to propose the novel *historical NGD* that has a theoretical guarantee to accelerate the convergence speed, which is observed in the numerical experiments as well.**

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

## I. INTRODUCTION

Along with the recent rapid advances in quantum information processing devices, the increasing attention has been paid to the possibility that quantum computers can outperform the classical (conventional) computers in various research fields such as machine learning, chemistry, and finance. However, since the currently available quantum computers are not fault-tolerant, their capabilities are limited [1]. This has led the necessity to develop hybrid quantum-classical approaches that enable those noisy quantum devices to work, with the help of classical computers. The variational quantum algorithm (VQA) is one such strategy [2]–[4]. The VQA runs a parameterized quantum circuit (PQC) on a quantum computer, with variationaly updating the parameters by a classical optimizer to find a global minimum of the objective function. This approach is expected to show some quantum advantages, because PQCs with even short depth potentially have bigger expressibility than classical models such as the neural network. To date, a variety of VQAs has been proposed; the variational quantum eigensolver (VQE) for quantum chemistry [5]–[7], the quantum approximate optimization algorithm (QAOA) for combinatorial optimizations [8]–[10], and quantum circuit learning algorithms for machine learning problems [11]–[13].

Of course the performance of VQAs heavily depends on the power of classical optimizer, particularly the convergence speed of the optimization. Actually various optimizers have been tested in VQAs; the gradient-based optimizers such as adaptive moment estimation (ADAM) [14], the conjugate gradient (CG) [15], the simultaneous perturbation stochastic approximation (SPSA) [16] and the natural gradient [17], as well as the gradient-free optimizers such as Nelder-Mead [18] and COBYLA [19]. In this work, we study the gradient-based approach, with special attention to the serious issue recently recognized in this method. That is, it has been demonstrated that the VQA has the so-called vanishing gradient or the barren plateau issue, where the gradient vector of the objective function becomes exponentially small with the increase of the number of qubits [20], [21]. This makes the gradient-based optimizer inefficient for the optimization of VQAs. So far, several circumventing approaches have been proposed to remedy this issue, such as the initialization-engineering technique and the tailored PQCs [22]–[26].

Note that the similar vanishing gradient issue can generally occur for non-convex optimization problems. Additionally, in such non-convex optimization problems, the so-called exploding gradient issue is also observed, meaning that the norm of the gradient vector takes a huge value and as a result the training becomes unstable. These detrimental issues, however, can be circumvented via rather a simple method that uses the normalized gradient vector to update the parameter through the learning process. This method is called the *normalized gradient descent (NGD)* [27]. Though the vanishing or exploding gradient issues never happen, the NGD could be disadvantageous in view of the convergence properties because it does not have a norm-dependent flexibility to search the optimal parameter. Nonetheless, under some conditions, the NGD is proven to evade the saddle points faster than the ordinary gradient descent method in the setting of continuous-

time dynamics [28].

Here, we study the performance of the NGD method in the optimization of VQAs. To the best of our knowledge, this work is the first to investigate the NGD methods in regards to the VQAs. This paper focuses on two objectives. The first objective is to examine the effectiveness of NGD and its variants for resolving the vanishing gradient issues. We applied those optimizers to several VQE problems, where the ground state of a given Hamiltonian is variationally sought. The numerical simulations show the good convergence property of the NGD methods, in comparison with the relevant optimizers that do not use the normalized gradient vector. As the second objective, we exploit the normalized gradient vectors at the past iteration steps to propose the *historical NGD*. In particular, we derive a set of proper learning rates of the historical normalized gradient vectors, under the assumption of the strictly-locally-quasi-convex (SLQC) objective functions. Hence the historical NGD is guaranteed to show the faster convergence than the ordinary NGD, which is also demonstrated in the numerical experiments.

The rest of the paper is organized as follows. In Section II, we show some gradient-based optimizers studied in this work. We then describe our historical NGD method that utilizes the normalized gradient vectors at the previous iteration step. Subsequently, Section IV demonstrates the numerical simulation of these optimizers for several VQE problems. At last, we conclude the paper in Section V.

## II. PRELIMINARIES

The gradient-based optimizers use the gradient vector of a objective function to update the parameters, for finding the global minimum of the function. The straightforward method is the gradient descent (GD), which is expressed as follows:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta g_t, \tag{1}$$

where $\mathbf{x}_t$ represents the set of parameters, $g_t = \nabla f(\mathbf{x}_t)$ is the gradient vector of the differentiable function $f(\mathbf{x}_t)$, and $\eta \in \mathbb{R}$ is the learning rate. GD can effectively find the global minimum when the objective function is convex. However, when the objective function is non-convex, GD often poorly perform owing to the trainability problems, typically the vanishing gradient issue caused by the plateau landscape of the objective function and the exploding gradient one due to the steep cliffs.

To date, numerous variants of GD have been proposed to circumvent those drawbacks. In this section, we briefly describe some gradient-based optimizers used in this paper.

### A. Normalized Gradient Descent

NGD is a simple method that resolves the aforementioned vanishing or exploding gradient issues by normalizing the gradient vector. NGD updates the parameters according to the following formula:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \hat{g}_t, \tag{2}$$

where $\hat{g}_t = \nabla f(\mathbf{x}_t)/\|\nabla f(\mathbf{x}_t)\|$ is the normalized gradient vector. Because $\|\hat{g}_t\| = 1$ for all $t$, NGD of course neither vanishes nor explodes. Hence we can expect NGD will show good convergence property as well as stable learning. Actually, it has been proven in the framework of continuous-time dynamics that NGD can escape the saddle point faster than GD [28]. Also, Ref. [27] proved that NGD can converge to a global minimum for a wider class of functions, assuming the strictly-locally-quasi-convex (SLQC) property of the objective function, which we will explain later on.

### B. Nesterov's Accelerated Gradient Method

Nesterov's Accelerated Gradient (NAG) method [29] is a simple modification of the momentum method [30], which is also a variant of GD. In the momentum method, a moving average of the past gradients are taken into account to realize faster convergence and alleviate the oscillation along the ridges of the canyon in the landscape of the objective function. The momentum method updates the parameters in the following way;

$$m_t = \beta m_{t-1} - \eta \nabla f(\mathbf{x}_t), \quad \mathbf{x}_{t+1} = \mathbf{x}_t + m_t,$$

where $\beta$ and $\eta$ are positive scalars.

As for NAG method, the update rule is represented as

$$\mathbf{x}_{t+1} = \mathbf{y}_t - \eta \nabla f(\mathbf{y}_t), \tag{3}$$

$$\mathbf{y}_t = \mathbf{x}_t + \gamma_t(\mathbf{x}_t - \mathbf{x}_{t-1}), \tag{4}$$

$$\gamma_t = \frac{\rho_{t-1} - 1}{\rho_t}, \tag{5}$$

$$\rho_t = \frac{1 + \sqrt{1 + 4\rho_{t-1}^2}}{2}. \tag{6}$$

Here, $\rho_t$ and $\gamma_t$ at each iteration step are recursively calculated so that the optimal convergence rate for the smooth convex function is achieved [29]. In our experiment, similar to [31] we set the initial value of $\rho$ as 1, i.e. $\rho_0 = 1$.

Note that the NAG method is rewritten as

$$m_t = \beta m_{t-1} - \eta \nabla f(\mathbf{x}_{t-1} + \beta m_{t-1}), \quad \mathbf{x}_{t+1} = \mathbf{x}_t + m_t,$$

which shows that the only difference between the momentum method and NAG is the point at which we calculate the gradient [31]. The NAG method recently attracts much attention in the convex optimization community due to its good convergence property in some situations.

### C. Adaptive Moment Estimation

ADAM [14] is a widely-used optimizer in the field of machine learning, particularly for deep neural networks. In a broad sense, ADAM utilizes the advantages of both the momentum method and RMSprops [32]; ADAM not only keeps the moving average of the past gradients, but also

computes the adaptive learning rate to reduce the oscillation. The update rule of ADAM is described as follows;

$$\bar{m}_t = \beta_1 \bar{m}_{t-1} - (1 - \beta_1) \nabla f(\mathbf{x}_t), \tag{7}$$

$$\bar{v}_t = \beta_2 \bar{v}_{t-1} - (1 - \beta_2) \nabla f(\mathbf{x}_t)^2, \tag{8}$$

$$m_t = \frac{\bar{m}_t}{1 - \beta_1^{t+1}}, \tag{9}$$

$$v_t = \frac{\bar{v}_t}{1 - \beta_2^{t+1}}, \tag{10}$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{m_t}{\sqrt{v_t} + \epsilon}, \tag{11}$$

where $\epsilon$ is a small constant introduced for numerical stability. In [14], the hyperparameters are set as $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$, which we also choose in our experiments.

Note that we can take another type of gradient-based optimizer, which utilizes higher-order derivative information such as the Hessian of an objective function. Typically used are Newton's method and quasi Newton method (e.g., the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm [33] and a sequential least squares programming (SLSQP) algorithm [34]). While they are computationally expensive compared to the variants of GD, these methods show faster convergence. However, this work focuses on the effectiveness of the optimizers with the normalized gradient vectors for VQAs, and thus we will not consider the higher-order gradient-based optimizer.

## III. Historical NGD

As mentioned before, the possible drawback of NGD is that it could be slower compared to GD due to the restricted norm condition. The historical NGD described here may be used to mitigate this issue; note that this method can also be applied to general optimization problem other than VQAs. The basics of this method is the provable convergence property of NGD under the assumption of the strictly-locally-quasi-convex (SLQC) objective function [27]. A particularly useful result is that, according to [27], NGD can converge to an $\epsilon$-optimal minimum of the SLQC objective function with the rate $\mathcal{O}(1/\epsilon^2)$; in fact we will make use of the lemmas related to this fact, to derive the proper learning rates of NGD with historical gradients.

In this section, we firstly introduce the relationship between the SLQC and NGD shown in [27], which is followed by deriving lemmas used to prove our result. Then we show that the one-step historical NGD can indeed accelerate the convergence speed. This method is further generalized to the historical NGD based on arbitrary $m$ normalized gradient vectors used in the past iteration steps.

### A. NGD for the Strictly-Locally-Quasi-Convex Function

In a broad sense, the SLQC function is the generalization of unimodal (or quasi-convex) functions with multi-dimensions, which can take even a plateau shape. The definition of the SLQC is as follows [27].

**Definition III.1.** *(Local-Quasi-Convexity) Let $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$ with the dimension $d$. Also let $\kappa$ and $\epsilon$ be positive constants. We consider a differentiable function $f : \mathbb{R}^d \mapsto \mathbb{R}$. Then $f$ is called $(\epsilon, \kappa, \mathbf{z})$-Strictly-Locally-Quasi-Convex (SLQC) at $\mathbf{x}$, if at least one of the following conditions holds:*

1) $f(\mathbf{x}) - f(\mathbf{z}) \leq \epsilon$.
2) $\|\nabla f(\mathbf{x})\| > 0$, and for every $\mathbf{y} \in \mathbb{B}(\mathbf{z}, \epsilon/\kappa)$ it holds that $\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq 0$,

*where $\mathbb{B}(\mathbf{x}, r)$ denotes a ball of radius $r$ around $\mathbf{x}$.*

In [27], the authors proved that NGD converges to the global minimum of a SLQC objective function, by taking the learning rate $\eta = \epsilon/\kappa$ in Eq. (2). Note that the SLQC condition is not too strict; some intriguing SLQC functions are studied in [27], such as the generalized linear model with certain setups.

Now we consider a SLQC function $f(\mathbf{x})$, and assume that it has a local or global minimum point $\mathbf{x}^*$. Also recall that $\hat{g}_t = \nabla f(\mathbf{x}_t)/\|\nabla f(\mathbf{x}_t)\|$ is a normalized gradient vector, and the NGD policy is given by $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \hat{g}_t$ with the learning rate $\eta = \epsilon/\kappa$. Then the following three lemmas hold [27]. Note that, the first lemma is proved in [27] and we obtain the remaining lemmas by following a similar proof.

**Lemma III.1.**
$$\langle \hat{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \geq \epsilon/\kappa.$$

Using this lemma, Ref. [27] showed that at every update of $\mathbf{x}_t$ in NGD, the distance between $\mathbf{x}_t$ and $\mathbf{x}^*$ is reduced by at least $(\epsilon/\kappa)^2$.

**Lemma III.2.** *Define $\delta = \langle \hat{g}_{t+1}, \hat{g}_t \rangle$. Then,*
$$\langle \hat{g}_{t+1}, \mathbf{x}_t - \mathbf{x}^* \rangle \geq (\epsilon/\kappa)(1 + \delta).$$

*Proof.* Using Lemma III.1, we have

$$\langle \hat{g}_{t+1}, \mathbf{x}_{t+1} - \mathbf{x}^* \rangle \geq (\epsilon/\kappa)$$
$$\Leftrightarrow \langle \hat{g}_{t+1}, \mathbf{x}_t - (\epsilon/\kappa)\hat{g}_t - \mathbf{x}^* \rangle \geq (\epsilon/\kappa)$$
$$\Leftrightarrow \langle \hat{g}_{t+1}, \mathbf{x}_t - \mathbf{x}^* \rangle - (\epsilon/\kappa)\langle \hat{g}_{t+1}, \hat{g}_t \rangle \geq (\epsilon/\kappa)$$
$$\Leftrightarrow \langle \hat{g}_{t+1}, \mathbf{x}_t - \mathbf{x}^* \rangle \geq (1 + \delta)(\epsilon/\kappa).$$

$\square$

**Lemma III.3.** *Define $\delta_{a,b} = \langle \hat{g}_{t+a}, \hat{g}_{t+b} \rangle$. Then, for a natural number $m \in \mathbb{N}$,*

$$\langle \hat{g}_{t+m}, \mathbf{x}_t - \mathbf{x}^* \rangle \geq (\epsilon/\kappa)\left(1 + \sum_{i=0}^{m-1} \delta_{m,i}\right).$$

*Proof.* Using Lemma III.1, we have

$$\langle \hat{g}_{t+m}, \mathbf{x}_{t+m} - \mathbf{x}^* \rangle \geq (\epsilon/\kappa)$$
$$\Leftrightarrow \langle \hat{g}_{t+m}, \mathbf{x}_t - (\epsilon/\kappa)\sum_{i=0}^{m-1}\hat{g}_{t+i} - \mathbf{x}^* \rangle \geq (\epsilon/\kappa)$$
$$\Leftrightarrow \langle \hat{g}_{t+m}, \mathbf{x}_t - \mathbf{x}^* \rangle - (\epsilon/\kappa)\sum_{i=0}^{m-1}\langle \hat{g}_{t+m}, \hat{g}_{t+i} \rangle \geq (\epsilon/\kappa)$$
$$\Leftrightarrow \langle \hat{g}_{t+m}, \mathbf{x}_t - \mathbf{x}^* \rangle \geq \left(1 + \sum_{i=0}^{m-1} \delta_{m,i}\right)(\epsilon/\kappa).$$

$\square$

## B. One step historical NGD

Here we show that the historical NGD, which updates the variable using $\hat{g}_{t+1}$ in addition to $\hat{g}_t$, helps accelerating the convergent speed of NGD compared to the ordinary NGD. Note that the inner product $\delta_t = \langle \hat{g}_t, \hat{g}_{t+1} \rangle$ used in the following result satisfies $-1 < \delta_t \leq 1$ because $\hat{g}_t$ and $\hat{g}_{t+1}$ are normalized.

**Corollary III.1.** *Let $\delta_t = \langle \hat{g}_t, \hat{g}_{t+1} \rangle$, satisfying $-1 < \delta_t \leq 1$. Consider the following update policy:*

$$\mathbf{x}_{t+2} = \mathbf{x}_t + \eta_1 \hat{g}_t + \eta_2 \hat{g}_{t+1},$$

*where the learning rate $\eta_1, \eta_2$ are determined as follows:*
- *When $-1 < \delta_t \leq (\sqrt{5} - 1)/2$,*

$$\eta_1 = -\left(\frac{\epsilon}{\kappa}\right) \frac{1 - \delta_t - \delta_t^2}{1 - \delta_t^2}, \quad \eta_2 = -\left(\frac{\epsilon}{\kappa}\right) \frac{1}{1 - \delta_t^2}.$$

- *When $(\sqrt{5} - 1)/2 < \delta_t \leq 1$,*

$$\eta_1 = 0, \quad \eta_2 = -\left(\frac{\epsilon}{\kappa}\right)(1 + \delta_t).$$

*Then it holds that*

$$\|\mathbf{x}_{t+2} - \mathbf{x}^*\|^2 < \|\mathbf{x}_t - \mathbf{x}^*\|^2 - c\left(\epsilon/\kappa\right)^2,$$

*for $c \geq 2$.*

*Proof.* By simple algebra, we can obtain

$$
\begin{aligned}
\|\mathbf{x}_{t+2} - \mathbf{x}^*\|^2 &= \langle \mathbf{x}_{t+2} - \mathbf{x}^*, \mathbf{x}_{t+2} - \mathbf{x}^* \rangle \\
&= \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \|\eta_1 \hat{g}_t + \eta_2 \hat{g}_{t+1}\|^2 \\
&\quad + 2\langle \mathbf{x}_t - \mathbf{x}^*, \eta_1 \hat{g}_t + \eta_2 \hat{g}_{t+1} \rangle \\
&= \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \eta_1^2 + \eta_2^2 + 2\eta_1\eta_2\delta_t \\
&\quad + 2\eta_1 \langle \hat{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle + 2\eta_2 \langle \hat{g}_{t+1}, \mathbf{x}_t - \mathbf{x}^* \rangle \\
&\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \eta_1^2 + \eta_2^2 + 2\eta_1\eta_2\delta_t \\
&\quad + 2\eta_1(\epsilon/\kappa) + 2\eta_2(\epsilon/\kappa)(1 + \delta_t).
\end{aligned}
$$

(12)

Here, we substitute $\mathbf{x}_t + \eta_1 \hat{g}_t + \eta_2 \hat{g}_{t+1}$ for $\mathbf{x}_{t+2}$ in the second equality. Also, the last inequality is due to Lemmas III.1 and III.2 as well as $\eta_1, \eta_2 \leq 0$. The proof follows by substituting the values of $\eta_1$ and $\eta_2$ to obtain:
- When $-1 < \delta_t \leq (\sqrt{5} - 1)/2$,

$$\|\mathbf{x}_{t+2} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{2 - \delta_t^2}{(1 - \delta_t^2)}\left(\epsilon/\kappa\right)^2.$$

- When $(\sqrt{5} - 1)/2 < \delta_t \leq 1$,

$$\|\mathbf{x}_{t+2} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - (1 + \delta_t)^2\left(\epsilon/\kappa\right)^2.$$

$\square$

Note that, when the ordinary NGD is successively applied twice, we have

$$\|\mathbf{x}_{t+2} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\left(\epsilon/\kappa\right)^2,$$

meaning that the historical NGD can show better convergence than the ordinary NGD.

## C. Multi-step historical NGD

In the previous subsection, we consider the historical NGD based on the gradient information at one past iteration step. Here we generalize the idea to NGD with arbitrary $m$ normalized gradient vectors in the past iteration steps.

Note that, in the previous case, the key to have the result is the proper choice of learning rates $\eta_1$ and $\eta_2$, and this can be done by quadratic programming. Actually, the learning rates can be obtained by minimizing the second and subsequent terms in the right hand side of the last inequality in Eq. (12), which can be expressed as the quadratic function (denoted as $h(\eta_1, \eta_2)$) in the following way;

$$
\begin{aligned}
h(\eta_1, \eta_2) &= \eta_1^2 + \eta_2^2 + 2\eta_1\eta_2\delta_t + 2\eta_1\left(\epsilon/\kappa\right) + 2\eta_2(\epsilon/\kappa)(1 + \delta_t) \\
&= \begin{bmatrix} \eta_1 & \eta_2 \end{bmatrix} \begin{bmatrix} 1 & \delta_t \\ \delta_t & 1 \end{bmatrix} \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \\
&\quad + \begin{bmatrix} 2\epsilon/\kappa & 2\epsilon/\kappa(1 + \delta_t) \end{bmatrix} \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \\
&= \mathbf{y}^{\mathrm{T}} A \mathbf{y} + C\mathbf{y}.
\end{aligned}
$$

(13)

The learning rates are obtained by solving the minimum of the above quadratic function, under the constraint $\eta_1, \eta_2 \leq 0$. We can apply this strategy to the case using $m$ previous gradient vectors for the historical NGD, which we call NGD$m$. The update rule of the NGD$m$ is given by

$$\mathbf{x}_{t+m} = \mathbf{x}_t + \sum_{i=1}^{m} \eta_i \hat{g}_{t+i-1},$$

(14)

where $\hat{g}_{t+m}$ is the normalized gradient vector at $t+m$ iteration step. Then we determine the learning rates $\{\eta_k\}_{k=1}^m$ so that $\mathbf{x}_{t+m}$ is closer to $\mathbf{x}^*$ than that via the ordinary NGD. More precisely, we derive the inequality connecting $\|\mathbf{x}_{t+m} - \mathbf{x}^*\|^2$ to $\|\mathbf{x}_t - \mathbf{x}^*\|^2$, like Eq. (12). Consequently, the problem for determining the proper learning rates reduces to the one to find the minimum of the following quadratic function $h(\mathbf{y})$ with $\mathbf{y} = [\eta_1, \ldots, \eta_m]$:

$$
\begin{aligned}
h(\mathbf{y}) &= \mathbf{y}^{\mathrm{T}} A \mathbf{y} + C\mathbf{y}, \\
A &= \begin{bmatrix} 1 & \delta_{0,1} & \ldots & \delta_{0,m-1} \\ \delta_{0,1} & 1 & \ldots & \delta_{1,m-1} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_{0,m-1} & \delta_{1,m-1} & \ldots & 1 \end{bmatrix}, \\
C &= 2\epsilon/\kappa \begin{bmatrix} 1 & 1 + \delta_{0,1} & \ldots & 1 + \sum_{i=0}^{m-2}\delta_{i,m-1} \end{bmatrix},
\end{aligned}
$$

(15)

under the constraint $\mathbf{y} \leq 0$, where $\delta_{i,j} = \langle \hat{g}_{t+i}, \hat{g}_{t+j} \rangle$. Here we utilize Lemma III.3 to derive $C$ in $h(\mathbf{y})$. Since the quadratic programming can be efficiently solved, our method can compute the proper learning rates at each iteration step quickly. However there is a caveat for this method in terms of numerical calculation; that is, the computation process can become unstable. For instance, when we use NGD2 (which is exactly the same as the one-step historical case), $\delta = -1$ results in the divergence of $\eta_1$ and $\eta_2$. This is because the

minimization of the function $h(\eta_1, \eta_2)$ in Eq. (13) is reduced to lowering $\eta_2$ as small as possible, as shown below;

$$h(\eta_1, \eta_2) = \eta_1^2 + \eta_2^2 - 2\eta_1\eta_2 + 2\eta_1 (\epsilon/\kappa)$$
$$= (\eta_1 - \eta_2 + \epsilon/\kappa)^2 + 2\eta_2(\epsilon/\kappa) - (\epsilon/\kappa)^2.$$

Hence, in the numerical simulation shown in the next section, we add extra constraint $\mathbf{y} \geq k$ with a negative constant $k$ to avoid the computational instability.

Lastly note that the proposed NGD is somehow similar to the momentum method in the sense that both use the past gradient information. However, our proposal differs with respect to the update rule, which is derived based on SLQC property where NGD can converge to the global minimum.

## IV. Numerical simulations

In this section, we show numerical simulations for several optimization problems appearing quantum chemistry, to test the performance of NGD. The goal is to find the ground state of a given Hamiltonian $H$, using VQE. The basic procedure of VQE is as follows; given a PQC $U(\theta)$ with $n_p$ tunable parameters $\theta = \{\theta_1, \ldots, \theta_{n_p}\}$, $\theta$ is repeatedly updated by a classical optimizer so that the energy (the objective function) $f(\theta) = \langle \Phi(\theta)|H|\Phi(\theta) \rangle$ is reduced to its minimum, where $|\Phi(\theta)\rangle = U(\theta)|\mathbf{0}\rangle$ with an initial state $|\mathbf{0}\rangle$. We study five VQE problems; we begin with a toy problem and then move to four quantum chemistry problems studying $H_2$, LiH, $H_4$, and the transverse field Ising model.

All numerical simulations are performed using the statevector simulator which does not introduce a statistical error due to the measurement process, on Qiskit [35] (version 0.24). As for the classical optimizers, we consider GD, NAG, ADAM (i.e., gradient-based optimizers without normalization), NGD, and the normalized NAG. Note that we use the parameter shift rule [13], [36] to calculate the gradient vector of the objective functions. Also the learning rate of each optimizer is fixed to 0.05. As for the historical NGD, we also set $\eta = \epsilon/\kappa = 0.05$. To solve the quadratic programming problem discussed in Section III C, we use CVXOPT [37], a package for convex optimization; in particular we choose the negative constant $k = -1000$, for the purpose of mitigating the computational instability.

### A. Toy narrow gorge problem

We begin with a toy *narrow gorge* problem, which was studied in [38]. The narrow gorge is a type of the energy landscape, such that the well around the minimum shrinks as the number of qubits increases. As a result, the vanishing gradient issue can be well observed in this problem. That is, the norm of the gradient vector rapidly decreases in all the parameter space except at around the minimum, as the number of qubits increases. Hence we expect to see the informative difference by comparing the convergence speed of the optimizers with and without normalization of the gradient, depending on the number of qubits. We take the same problem setting as [21]. The Hamiltonian is $H = \sum_{k=1}^n \sigma_X^{(k)}$, whose ground state is $|\mathbf{0}\rangle = |0\rangle^{\otimes n}$, where $\sigma_X$ is the Pauli $X$

operator and $n$ is the number of qubits. The PQC is chosen as $U(\theta) = \otimes_{k=1}^n e^{-i\theta_k \sigma_X^{(k)}}$. The goal is to optimize the parameters $\theta = \{\theta_k\}_{k=1}^n$ so that $U(\theta)|\mathbf{0}\rangle = |\mathbf{0}\rangle$.

In this work, we perform the simulation with different number of qubits $n = 2, 4, 8$, to see that NGD indeed resolves the vanishing gradient issue and, at the same time, to evaluate the convergent speed of NGD. Note that the number of qubits is equal to that of parameters, due to the tensor-product structure of the aforementioned PQC. The optimal parameters of this task are all zeros, i.e. $\theta_k = 0$ for all $k$. Hence we set the initial parameters as $\theta_k = \pi/2$ for all $k$, to avoid that the initial point would immediately get close to the optimal point. Also the total number of iterations for the optimization is fixed to 100.

In Fig. 1 we show the energy of the narrow gorge potential, against the iteration steps for each optimizer, where the blue, orange, and green lines in each figure show the case of $n = 2, 4$, and $8$, respectively. The upper three panels show the results for the optimizers without the normalized gradient vector, and the lower three show the optimizers with the normalized gradient vector. From these figures, we observe that the optimizers without normalization crucially suffer from the vanishing gradient issue, except for ADAM. On the other hand, the optimizers with normalization can still decrease the energy even when $n = 8$. Note that the performance of optimizers appears to be degraded as $n$ increases, because the distance between the initial and the optimal points gets larger.

Moreover, Fig. 1 (f) shows that the one-step historical NGD converges faster than the ordinary NGD, as proven in Corollary III.1. Here, we study if the convergence speed can be further improved by increasing the number of gradient information at the previous iteration steps. Fig. 2 shows the energy of narrow gorge potential for the case of $n = 8$, for NGD, NGD2, NGD3, and NGD4. Recall that NGD$m$ means the historical NGD with $m$ past normalized gradient vectors. Clearly, Fig. 2 shows that NGD$m$ with bigger $m$ converges faster. However, this result also shows the issue of computational instability of NGD$m$, as indicated in Section III. In particular, this issue seems more likely to occur, as we utilize more normalized gradient vectors. Thus, there is a room for improvement in our method, to fully exploit the past normalized gradient vectors without suffering from this computational instability issue. Yet, we underscore that NGD4 (as well as the Normalized NAG) is the first to reach within an error of $10^{-2}$ from the minimum energy.

### B. $H_2$ molecule

We next examine the problem of finding the ground state of $H_2$ molecule. The simplified Hamiltonian of $H_2$ is expressed as

$$H = \alpha(\sigma_Z \otimes I + I \otimes \sigma_Z) + \beta(\sigma_X \otimes \sigma_X),$$

where $\sigma_Z$ represents Pauli $Z$ operator and $(\alpha, \beta) = (0.4, 0.2)$. It was reported [25] that the VQE using GD with learning rate $\eta = 0.05$ gets stuck in a plateau for a while and then escapes later, when a single depth $Ry$ ansatz with initial
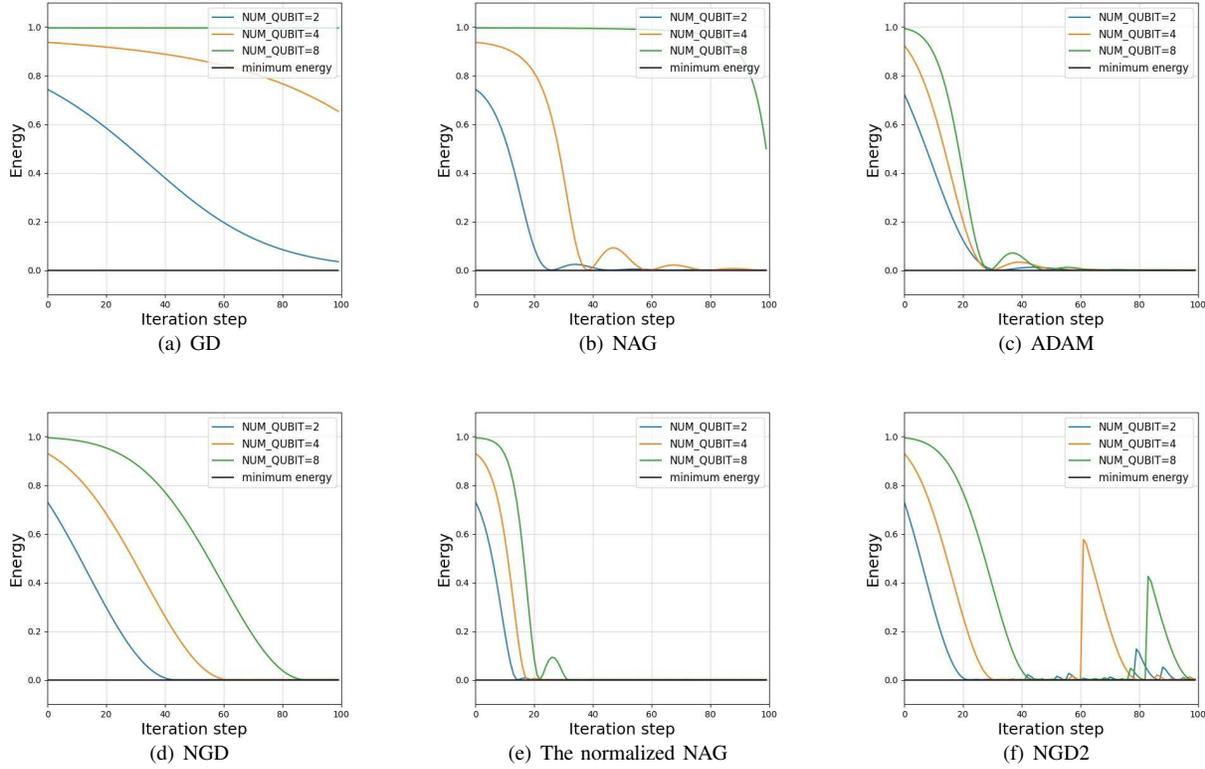
Fig. 1. Energy of the narrow gorge potential in the case of $n = 2$ (blue), 4 (orange), and 8 (green) against the iteration steps for each optimizer. The upper three panels (a, b, c) and lower three panels (d, e, f) show the results for the optimizers without normalized gradient vector and those for the optimizers with normalization, respectively.
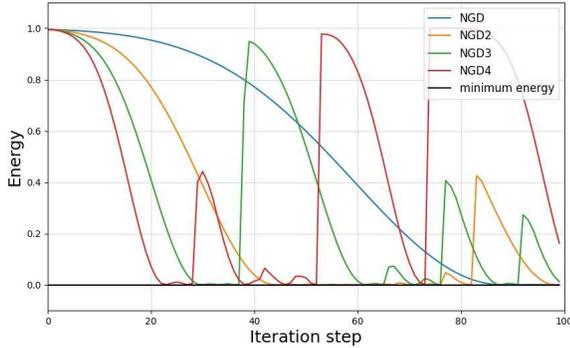


Fig. 2. Energy of the narrow gorge potential in the case of $n = 8$, against the iteration steps for NGD, NGD2, NGD3, and NGD4.
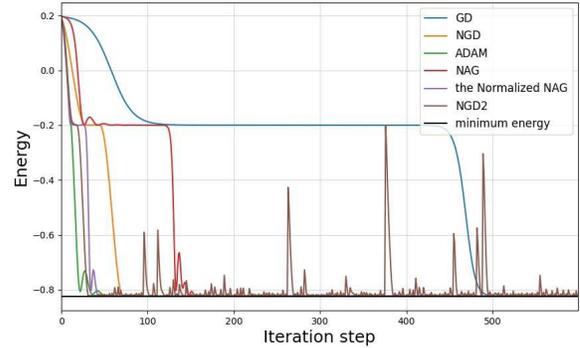


Fig. 3. Energy of $H_2$ molecule against the iteration steps for each optimizer

parameters $(\theta_1, \theta_2, \theta_3, \theta_4) = (7\pi/32, \pi/2, 0, 0)$ are used. This phenomenon arises because GD first arrives in the vicinity of the first excited state, where the gradient vector vanishes. Here we test GD, NAG, ADAM, NGD, the normalized NAG, and NGD2 with the same $R_y$ anzats, to see if they would be trapped in this plateau and, when trapped, how fast they can escape from it.

Figure 3 shows that all optimizers get stuck at the first excited state with energy $-0.2$. But notably, the optimizers with the normalized gradient vector evade the plateau faster than the others without normalization. For example, NGD can get out of the first excited state around 50 iteration steps, while GD requires 450 iterations. Moreover, we can also see the normalized NAG method outperforms the ordinary NAG. Importantly, NGD2 is the first to reach the minimum, while ADAM can get out of the plateau the fastest.
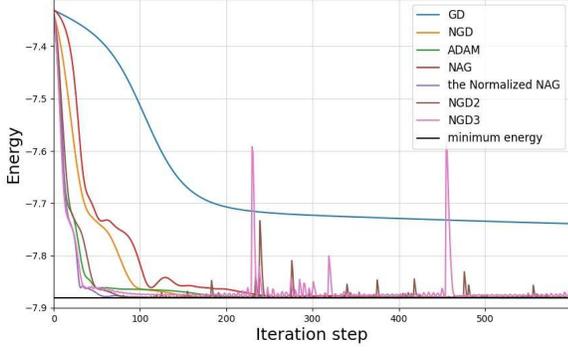
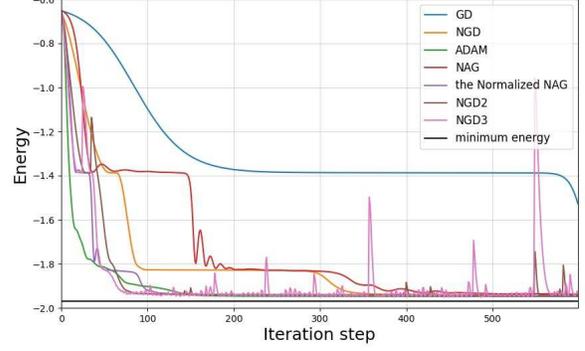Fig. 4. Energy of LiH molecule against the iteration steps for each optimizer



Fig. 5. Energy of $H_4$ molecule against the iteration steps for each optimizer

## C. LiH molecule

The next case-study is on the LiH molecule; we consider the case where the interatomic distance is 1.5 Å, resulting that two unoccupied orbitals are removed and the core is frozen. The Hamiltonian of LiH is constructed in the following way; the fermionic Hamiltonian is first constructed by the Hartree-Fock calculation with STO-3G basis [39] using the PySCF package [40], which is then converted by the parity encoding method [41] to the Hamiltonian. We apply VQE with two-depth $Ry$ ansatz, where the initial parameters are all zeros. The optimizers are GD, NAG, ADAM, NGD, the normalized NAG, NGD2, and NGD3.

The result is shown in Fig. 4. In terms of the convergence speed, NGD and the normalized NAG are superior to GD and NAG, respectively. Also, NGD2 and the Normalized NAG reach the minimum the fastest, while NGD3 and ADAM are competitive with them at the beginning of the optimization process. Notably, NGD3 falls behind NGD2 to converge to the minimum, despite the fact that NGD3 utilizes more gradient information at the past iteration step. The reason may be the computational instability that occurs when the learning rates are computed. In fact, NGD3 shows the fastest convergence at first, but the energy begins to fluctuate at certain timestep.

## D. $H_4$ molecule

We here consider the $H_4$ molecule with square configuration with interatomic distance 1.277 Å. The Hamiltonian of $H_4$ is constructed in the same way to the LiH case. We use the five-depth $Ry$ ansatz as PQC, where the initial parameters are all zeros, and the same set of optimizers.

In this case, all optimizers cannot arrive at the minimum within 600 iterations. However, we still observe the better convergence property of the optimizers that use the normalized gradient vector. Namely, NGD and the normalized NAG method quickly converge to the local minimum about -1.95, in comparison with GD and NAG method, respectively. Moreover, NGD3 reaches the lowest value the fastest. Note that the reason of not achieving the exact minimum may not be attributed to the optimizers, but other factors such as the insufficient expressibility of the PQC ansatz.
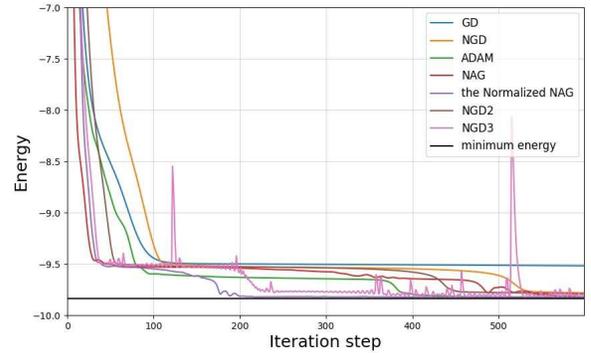


Fig. 6. Energy of the transverse field Ising model, against the iteration steps for each optimizer

## E. Transverse field Ising model

Lastly, we focus on the transverse field Ising model studied in [42], whose Hamiltonian is expressed as

$$H = \sum_{i=1}^{n-1} \sigma_Z^{(i)} \sigma_Z^{(i+1)} + \sum_{i=1}^{n} \sigma_X^{(i)}.$$

Here we set $n = 8$. Again we use the two-depth Ry ansatz with all initial parameters given by $\pi/2$.

Unlike the previous problems, Fig. 6 shows the opposite results at the beginning of the optimization process; GD converges faster than NGD. This is because the chosen initial parameter is so far from the optimal point in the parameter space and the norm of the gradient vector is much bigger than one at around the initial point. This is actually seen in Fig. 7, showing the norm of the gradient vector against the optimization iterations for GD. However, this does not mean that NGD is inferior to GD. In fact, Fig. 6 shows that GD gets stuck in a plateau, while NGD can escape from the plateau after roughly 500 iterations despite the poor convergence in the beginning. This different behaviors can be, again, explained by Fig. 7, showing that the norm of the gradient vector quickly decreases as the optimization proceeds. This reflects
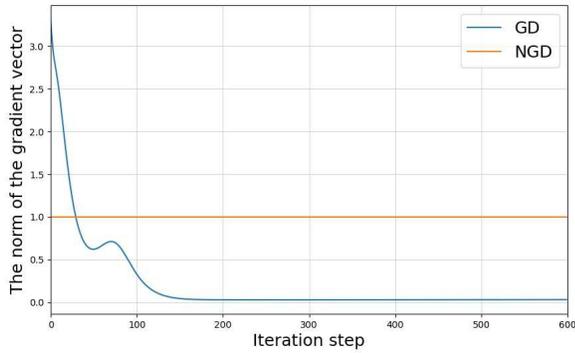
Fig. 7. The norm of the gradient vectors against the iteration steps of GD for the VQE calculation of the transverse field Ising model. For comparison, we also show the norm of the gradient vector for NGD, which is always one.

the effectiveness of the optimizers with normalized gradient vector for the convergence.

## V. Conclusion

In this paper, we apply the NGD method to VQAs, to overcome the vanishing gradient issue. Several numerical simulations actually show that the optimizers with the normalized gradient vector have good convergence property, compared to the optimizers without normalization. Moreover, we proposed a new NGD that uses some normalized gradient vectors computed in the past optimization steps; this historical NGD is guaranteed to have a faster convergence property compared to the ordinary NGD, and actually we have demonstrated that this indeed accelerates the convergence speed of NGD.

We hope this work will pave a new way to deal with the vanishing gradient problem, that often appears in the optimization process of the VQAs. Note also that, since the application of the historical NGD is not limited to VQAs, we hope that our method might be useful in e.g., machine learning.

## Acknowledgment

## References

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

[2] 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," *arXiv preprint arXiv:2012.09265*, 2020.

[3] 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.*, "Noisy intermediate-scale quantum (nisq) algorithms," *arXiv preprint arXiv:2101.08448*, 2021.

[4] S. Endo, Z. Cai, S. C. Benjamin, and X. Yuan, "Hybrid quantum-classical algorithms and quantum error mitigation," *Journal of the Physical Society of Japan*, vol. 90, no. 3, p. 032001, 2021.

[5] 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, p. 4213, 2014.

[6] 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.

[7] A. Kandala, A. Mezzacapo, K. Temme, M. Takita, M. Brink, J. M. Chow, and J. M. Gambetta, "Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets," *Nature*, vol. 549, no. 7671, p. 242, 2017.

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

[9] E. Farhi and A. W. Harrow, "Quantum supremacy through the quantum approximate optimization algorithm," *arXiv preprint arXiv:1602.07674*, 2016.

[10] L. Zhou, S.-T. Wang, S. Choi, H. Pichler, and M. D. Lukin, "Quantum approximate optimization algorithm: Performance, mechanism, and implementation on near-term devices," *Physical Review X*, vol. 10, no. 2, p. 021067, 2020.

[11] V. Havlicek, A. D. Corcoles, K. Temme, A. W. Harrow, A. Kandala, J. M. Chow, and J. M. Gambetta, "Supervised learning with quantum-enhanced feature spaces," *Nature*, vol. 567, no. 7747, p. 209–212, 2019.

[12] M. Schuld, A. Bocharov, K. Svore, and N. Wiebe, "Circuit-centric quantum classifiers," 2018.

[13] K. Mitarai, M. Negoro, M. Kitagawa, and K. Fujii, "Quantum circuit learning," *Physical Review A*, vol. 98, no. 3, 2018.

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

[15] M. R. Hestenes, E. Stiefel *et al.*, *Methods of conjugate gradients for solving linear systems*. NBS Washington, DC, 1952, vol. 49, no. 1.

[16] J. C. Spall *et al.*, "Multivariate stochastic approximation using a simultaneous perturbation gradient approximation," *IEEE transactions on automatic control*, vol. 37, no. 3, pp. 332–341, 1992.

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

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

[19] M. J. Powell, "A direct search optimization method that models the objective and constraint functions by linear interpolation," in *Advances in optimization and numerical analysis*. Springer, 1994, pp. 51–67.

[20] 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.

[21] 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.

[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] G. Verdon, M. Broughton, J. R. McClean, K. J. Sung, R. Babbush, Z. Jiang, H. Neven, and M. Mohseni, "Learning to learn with quantum neural networks via classical neural networks," *arXiv preprint arXiv:1907.05415*, 2019.

[24] S. Hadfield, Z. Wang, B. O'Gorman, E. G. Rieffel, D. Venturelli, and R. Biswas, "From the quantum approximate optimization algorithm to a quantum alternating operator ansatz," *Algorithms*, vol. 12, no. 2, p. 34, 2019.

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

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

[27] E. Hazan, K. Y. Levy, and S. Shalev-Shwartz, "Beyond convexity: Stochastic quasi-convex optimization," *arXiv preprint arXiv:1507.02030*, 2015.

[28] R. Murray, B. Swenson, and S. Kar, "Revisiting normalized gradient descent: Fast evasion of saddle points," *IEEE Transactions on Automatic Control*, vol. 64, no. 11, pp. 4818–4824, 2019.

[29] Y. Nesterov, "A method for unconstrained convex minimization problem with the rate of convergence $o(1/k^2)$," in *Doklady an ussr*, vol. 269, 1983, pp. 543–547.

[30] B. T. Polyak, "Some methods of speeding up the convergence of iteration methods," *Ussr computational mathematics and mathematical physics*, vol. 4, no. 5, pp. 1–17, 1964.

[31] 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.

[32] T. Tieleman and G. Hinton, "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude," *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, pp. 26–31, 2012.

[33] R. Fletcher, *Practical methods of optimization*. John Wiley & Sons, 2013.

[34] D. Kraft, "Algorithm 733: Tomp–fortran modules for optimal control calculations," *ACM Transactions on Mathematical Software (TOMS)*, vol. 20, no. 3, pp. 262–281, 1994.

[35] H. Abraham, I. Y. Akhalwaya, G. Aleksandrowicz, T. Alexander, G. Alexandrowics, E. Arbel, A. Asfaw, C. Azaustre, AzizNgoueya, P. Barkoutsos, G. Barron, L. Bello, Y. Ben-Haim, D. Bevenius, L. S. Bishop, S. Bosch, S. Bravyi, D. Bucher, F. Cabrera, P. Calpin, L. Capelluto, J. Carballo, G. Carrascal, A. Chen, C.-F. Chen, R. Chen, J. M. Chow, C. Claus, C. Clauss, A. J. Cross, A. W. Cross, S. Cross, J. Cruz-Benito, C. Culver, A. D. Córcoles-Gonzales, S. Dague, T. E. Dandachi, M. Dartiailh, DavideFrr, A. R. Davila, D. Ding, J. Doi, E. Drechsler, Drew, E. Dumitrescu, K. Dumon, I. Duran, K. EL-Safty, E. Eastman, P. Eendebak, D. Egger, M. Everitt, P. M. Fernández, A. H. Ferrera, A. Frisch, A. Fuhrer, M. GEORGE, J. Gacon, Gadi, B. G. Gago, J. M. Gambetta, A. Gammanpila, L. Garcia, S. Garion, J. Gomez-Mosquera, S. de la Puente González, I. Gould, D. Greenberg, D. Grinko, W. Guan, J. A. Gunnels, I. Haide, I. Hamamura, V. Havlicek, J. Hellmers, Ł. Herok, S. Hillmich, H. Horii, C. Howington, S. Hu, W. Hu, H. Imai, T. Imamichi, K. Ishizaki, R. Iten, T. Itoko, A. Javadi-Abhari, Jessica, K. Johns, T. Kachmann, N. Kanazawa, Kang-Bae, A. Karazeev, P. Kassebaum, S. King, Knabberjoe, A. Kovyrshin, V. Krishnan, K. Krsulich, G. Kus, R. LaRose, R. Lambert, J. Latone, S. Lawrence, D. Liu, P. Liu, Y. Maeng, A. Malyshev, J. Marecek, M. Marques, D. Mathews, A. Matsuo, D. T. McClure, C. McGarry, D. McKay, S. Meesala, M. Mevissen, A. Mezzacapo, R. Midha, Z. Minev, N. Moll, M. D. Mooring, R. Morales, N. Moran, P. Murali, J. Müggenburg, D. Nadlinger, G. Nannicini, P. Nation, Y. Naveh, P. Neuweiler, P. Niroula, H. Norlen, L. J. O'Riordan, O. Ogunbayo, P. Ollitrault, S. Oud, D. Padilha, H. Paik, S. Perriello, A. Phan, M. Pistoia, A. Pozas-iKerstjens, V. Prutyanov, D. Puzzuoli, J. Pérez, Quintiii, R. Raymond, R. M.-C. Redondo, M. Reuter, J. Rice, D. M. Rodríguez, M. Rossmannek, M. Ryu, T. SAPV, SamFerracin, M. Sandberg, N. Sathaye, B. Schmitt, C. Schnabel, Z. Schoenfeld, T. L. Scholten, E. Schoute, I. F. Sertage, K. Setia, N. Shammah, Y. Shi, A. Silva, A. Simonetto, N. Singstock, Y. Siraichi, I. Sitdikov, S. Sivarajah, M. B. Sletfjerding, J. A. Smolin, M. Soeken, I. O. Sokolov, D. Steenken, M. Stypulkoski, H. Takahashi, I. Tavernelli, C. Taylor, P. Taylour, S. Thomas, M. Tillet, M. Tod, E. de la Torre, K. Trabing, M. Treinish, TrishaPe, W. Turner, Y. Vaknin, C. R. Valcarce, F. Varchon, A. C. Vazquez, D. Vogt-Lee, C. Vuillot, J. Weaver, R. Wieczorek, J. A. Wildstrom, R. Wille, E. Winston, J. J. Woehr, S. Woerner, R. Woo, C. J. Wood, R. Wood, S. Wood, J. Wootton, D. Yeralin, R. Young, J. Yu, C. Zachow, L. Zdanski, C. Zoufal, Zoufalc, azulehner, bcamorrison, brandhsn, chlorophyll zz, dime10, drholmie, elfrocampeador, faisaldebouni, fanizzamarco, gruu, kanejess, klinvill, kurarrr, lerongil, ma5x, merav aharoni, ordmoj, sethmerkel, strickroman, sumitpuri, tigerjack, toural, vvilpas, willhbang, yang.luh, and yotamvakninibm, "Qiskit: An open-source framework for quantum computing," 2019.

[36] M. Schuld, V. Bergholm, C. Gogolin, J. Izaac, and N. Killoran, "Evaluating analytic gradients on quantum hardware," *Physical Review A*, vol. 99, no. 3, p. 032331, 2019.

[37] M. Andersen, J. Dahl, and L. Vandenberghe, "CVXOPT: Convex Optimization," p. ascl:2008.017, Aug. 2020.

[38] 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.

[39] W. J. Hehre, R. F. Stewart, and J. A. Pople, "Self-consistent molecular-orbital methods. i. use of gaussian expansions of slater-type atomic orbitals," *The Journal of Chemical Physics*, vol. 51, no. 6, pp. 2657–2664, 1969.

[40] Q. Sun, T. C. Berkelbach, N. S. Blunt, G. H. Booth, S. Guo, Z. Li, J. Liu, J. D. McClain, E. R. Sayfutyarova, S. Sharma *et al.*, "Pyscf: the python-based simulations of chemistry framework," *Wiley Interdisciplinary Reviews: Computational Molecular Science*, vol. 8, no. 1, p. e1340, 2018.

[41] J. T. Seeley, M. J. Richard, and P. J. Love, "The bravyi-kitaev transformation for quantum computation of electronic structure," *The Journal of chemical physics*, vol. 137, no. 22, p. 224109, 2012.

[42] R. Sweke, F. Wilde, J. J. Meyer, M. Schuld, P. K. Fährmann, B. Meynard-Piganeau, and J. Eisert, "Stochastic gradient descent for hybrid quantum-classical optimization," *Quantum*, vol. 4, p. 314, 2020.