
Bayesian Parameter Shift Rules in Variational Quantum Eigensolvers

Anonymous Authors¹

Abstract

Parameter shift rules (PSRs) are key techniques for efficient gradient estimation in variational quantum eigensolvers (VQEs). In this paper, we propose its Bayesian variant, where Gaussian processes with appropriate kernels are used to estimate the gradient of the VQE objective. Our *Bayesian PSR* offers flexible gradient estimation from observations at arbitrary locations with uncertainty information, and reduces to the generalized PSR in special cases. In stochastic gradient descent (SGD), the flexibility of Bayesian PSR allows reuse of observations in previous steps, which accelerates the optimization process. Furthermore, the accessibility to the posterior uncertainty, along with our proposed notion of *gradient confident region* (GradCoRe), enables us to minimize the observation costs in each SGD step. Our numerical experiments show that the VQE optimization with Bayesian PSR and GradCoRe significantly accelerates SGD, and outperforms the state-of-the-art methods, including sequential minimal optimization.

1. Introduction

The variational quantum eigensolver (VQE) (Peruzzo et al., 2014; McClean et al., 2016) is a hybrid quantum-classical algorithm for approximating the ground state of the Hamiltonian of a given physical system. The quantum part of VQEs uses parameterized quantum circuits to generate trial quantum states and measures the expectation value of the Hamiltonian, i.e., the energy, while the classical part forms energy minimization with noisy observations from the quantum device. Provided that the parameterized quantum circuits can accurately approximate the ground state, the minimized energy gives a tight upper bound of the ground state energy of the Hamiltonian.

The observation noise in the quantum device comes from multiple sources. One source of noise is *measurement shot noise*, which arises from the statistical nature of quantum measurements—outcomes follow the probabilities specified by the quantum state, and finite sampling introduces fluctuations. Since this noise source is random and inde-

pendent, it can be reduced by increasing the number of measurement shots, to which the variance is inversely proportional. Another source of noise stems from imperfections in the quantum hardware, which have been reduced in recent years by hardware design (Bluvstein et al., 2023), as well as error mitigation (Cai et al., 2023), quantum error correction (Roffe, 2019; Acharya et al., 2024), and machine learning (Nicoli et al., 2025) techniques. In this paper, we do not consider hardware noise, as is common in papers developing optimization methods (Nakanishi et al., 2020; Nicoli et al., 2023b).

Stochastic gradient descent (SGD), sequential minimal optimization (SMO), and Bayesian optimization (BO) have been used to minimize the VQE objective function. Under some mild assumptions (Nakanishi et al., 2020), this objective function is known to have special properties. Based on those properties, SGD methods can use the gradient estimated by so-called *parameter shift rules* (PSRs) (Mitarai et al., 2018), and specifically designed SMO (Platt, 1998) methods, called Nakanishi-Fuji-Todo (NFT) (Nakanishi et al., 2020), perform one-dimensional subspace optimization with only a few observations in each iteration. Iannelli & Jansen (2021) applied BO to solve VQEs as a noisy global optimization problem.

Although Gaussian processes (GPs) have been used in VQEs as a common surrogate function for BO (Frazier, 2018), they have also been used to improve SGD-based and SMO-based methods. Nicoli et al. (2023a) proposed the *VQE kernel*—a physics-informed kernel that fully reflects the properties of VQEs—and combined SMO and BO with the *expected maximum improvement within confident region* (EMICoRe) acquisition function. This allows for the identification of the optimal locations to measure on the quantum computer in each SMO iteration. Tamiya & Yamasaki (2022) combined SGD and BO, and proposed *stochastic gradient line BO* (SGLBO), which uses BO to identify the optimal step size in each SGD iteration. Anders et al. (2024) proposed the *subspace in confident region* (SubsCoRe) approach, where the observation costs are minimized based on the posterior uncertainty estimation in each SMO iteration.

In this paper, we take a different approach by leveraging GPs to introduce a *Bayesian parameter shift rule* (Bayesian PSR), where the gradient of the VQE objective is estimated

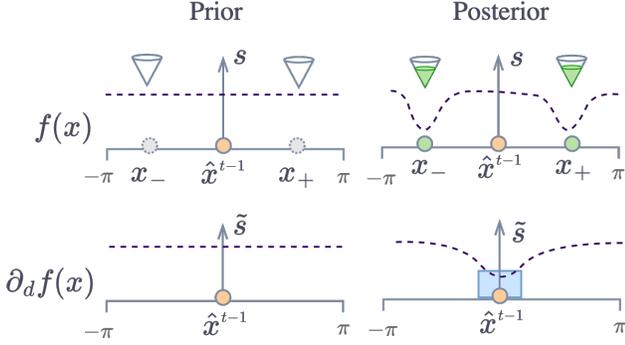


Figure 1. Illustration of our gradient confident region (GradCoRe) approach. Our goal is to minimize the true energy $f^*(\mathbf{x})$ over the set of parameters $\mathbf{x} \in [0, 2\pi)^D$. We use a GP surrogate $f(\mathbf{x})$ for $f^*(\mathbf{x})$. Observing f^* at points \mathbf{x}_- and \mathbf{x}_+ (green circles) along the d -th direction (solid horizontal line) decreases the GP uncertainty (dashed curves) not only at $f(\mathbf{x}_\pm)$, but also at $\partial_d f(\hat{\mathbf{x}}^{t-1})$ which thus falls within the GradCoRe (blue square). Our GradCoRe-based SGD minimizes the total number of measurement shots for optimization.

using a GP with the VQE kernel. The Bayesian PSR translates into a regularized variant of PSRs if the observations are performed at designated locations. However, our approach offers significant advantages—flexibility and direct access to uncertainty—over existing PSRs (Wierichs et al., 2022). More specifically, the Bayesian PSR can use observations at any set of locations, which allows the reuse of observations performed in previous iterations of SGD. Reusing previous observations along with new observations improves the gradient estimation accuracy, and thus accelerates the optimization process. Furthermore, the uncertainty information can be used to adapt the observation cost in each SGD iteration, in a similar spirit to Anders et al. (2024). Adapting the observation cost can significantly reduce the necessary cost of obtaining new observations, while maintaining a required level of accuracy. We implement this adaptive observation cost strategy by introducing a novel notion of *gradient confidence region* (GradCoRe)—the region in which the uncertainty of the gradient estimation is below a specific threshold (see Figure 1). Our empirical evaluations show that our proposed Bayesian PSR improves the gradient estimator, and SGD equipped with our GradCoRe approach outperforms all previous state-of-the-art methods including NFT and its variants.

The main contributions are summarized as follows:

- We propose *Bayesian PSR*, a flexible variant of existing PSRs that provides access to uncertainty information.
- We theoretically establish the relationship between Bayesian PSR and existing PSRs, revealing the op-

timality of the *shift* parameter in first-order PSRs.

- We introduce the notion of *GradCoRe*, and propose an adaptive observation cost strategy for SGD optimization.
- We numerically validate our theory and empirically demonstrate the effectiveness of the proposed Bayesian PSR and GradCoRe.

Related work: Finding the optimal set of parameters for a variational quantum circuit is a challenging problem, prompting the development of various approaches to improve the optimization in VQEs. Gradient-based methods for VQEs often rely on PSRs (Mitarai et al., 2018; Wierichs et al., 2022), which enable reasonably accurate gradient estimation of the output of quantum circuits with respect to their parameters. Nakanishi et al. (2020) proposed an SMO (Platt, 1998) algorithm, known as *NFT*, where, at each step of SMO, one parameter is analytically minimized by performing a few observations. Nicoli et al. (2023a) combined NFT with GP and BO by developing a physics-inspired kernel for GP regression and proposing the EMICoRe acquisition function, relying on the concept of confident regions (CoRe). This method improves upon NFT by leveraging the information from observations in previous steps to identify the optimal locations to perform the next observations. Anders et al. (2024) leveraged the same notion of CoRe, and proposed *SubsCoRe*, where, instead of optimizing the observed locations, the minimal number of measurement shots is identified to achieve the required accuracy defined by the CoRe. The resulting algorithm converges to the same energy as NFT with a smaller quantum computation cost, i.e., the total number of measurement shots on a quantum computer. Tamiya & Yamasaki (2022) combined SGD with BO to tackle the excessive cost of standard SGD approaches and used BO to accelerate the convergence by finding the optimal step size. On a related note, recent works (Jiang et al., 2024) have begun integrating GP with error mitigation techniques, further highlighting the potential of Bayesian approaches for noisy intermediate-scale quantum (NISQ) devices (Preskill, 2018).

The remainder of the paper is structured as follows: in Section 2, we provide the necessary background on GP and VQEs. In Section 3, we propose our Bayesian PSR and provide a theory that relates it to the existing PSRs. In Section 4, we propose our novel SGD-based algorithms based on Bayesian PSR and GradCoRe. In Section 5, we describe the experimental setup and present numerical experiments. Finally, in Section 6, we summarize our findings and provide an outlook for future research.

2. Background

Here we briefly introduce Gaussian process (GP) regression and its derivatives, as well as VQEs with their known properties.

2.1. GP Regression and Derivative GP

Assume we aim to learn an unknown function $f^*(\cdot) : \mathcal{X} \mapsto \mathbb{R}$ from the training data $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathcal{X}^N$, $\mathbf{y} = (y_1, \dots, y_N)^\top \in \mathbb{R}^N$, $\boldsymbol{\sigma} = (\sigma_1^2, \dots, \sigma_N^2) \in \mathbb{R}_{++}^N$ that fulfills

$$y_n = f^*(\mathbf{x}_n) + \varepsilon_n, \quad \varepsilon_n \sim \mathcal{N}_1(y_n; 0, \sigma_n^2), \quad (1)$$

where $\mathcal{N}_D(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the D -dimensional Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$. With the Gaussian process (GP) prior

$$p(f(\cdot)) = \text{GP}(f(\cdot); 0(\cdot), k(\cdot, \cdot)), \quad (2)$$

where $0(\cdot)$ and $k(\cdot, \cdot)$ are the prior zero-mean and the kernel (covariance) functions, respectively, the posterior distribution of the function values $\mathbf{f}' = (f(\mathbf{x}'_1), \dots, f(\mathbf{x}'_M))^\top \in \mathbb{R}^M$ at arbitrary test points $\mathbf{X}' = (\mathbf{x}'_1, \dots, \mathbf{x}'_M) \in \mathcal{X}^M$ is given as

$$p(\mathbf{f}' | \mathbf{X}, \mathbf{y}) = \mathcal{N}_M(\mathbf{f}'; \boldsymbol{\mu}'_{[\mathbf{X}, \mathbf{y}, \boldsymbol{\sigma}]}, \mathbf{S}'_{[\mathbf{X}, \boldsymbol{\sigma}]}, \text{ where } (3)$$

$$\boldsymbol{\mu}'_{[\mathbf{X}, \mathbf{y}, \boldsymbol{\sigma}]} = \mathbf{K}'^\top (\mathbf{K} + \text{Diag}(\boldsymbol{\sigma}))^{-1} \mathbf{y} \quad \text{and} \quad (4)$$

$$\mathbf{S}'_{[\mathbf{X}, \boldsymbol{\sigma}]} = \mathbf{K}'' - \mathbf{K}'^\top (\mathbf{K} + \text{Diag}(\boldsymbol{\sigma}))^{-1} \mathbf{K}' \quad (5)$$

are the posterior mean and covariance, respectively (Rasmussen & Williams, 2006). Here $\text{Diag}(v)$ is the diagonal matrix with v specifying the diagonal entries, and $\mathbf{K} = k(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{N \times N}$, $\mathbf{K}' = k(\mathbf{X}, \mathbf{X}') \in \mathbb{R}^{N \times M}$, and $\mathbf{K}'' = k(\mathbf{X}', \mathbf{X}') \in \mathbb{R}^{M \times M}$ are the train, train-test, and test kernel matrices, respectively, where $k(\mathbf{X}, \mathbf{X}')$ denotes the kernel matrix evaluated at each column of \mathbf{X} and \mathbf{X}' such that $(k(\mathbf{X}, \mathbf{X}'))_{n,m} = k(\mathbf{x}_n, \mathbf{x}'_m)$. We also denote the posterior as $p(f(\cdot) | \mathbf{X}, \mathbf{y}) = \text{GP}(f(\cdot); \mu_{[\mathbf{X}, \mathbf{y}, \boldsymbol{\sigma}]}(\cdot), s_{[\mathbf{X}, \boldsymbol{\sigma}]}(\cdot, \cdot))$ with the posterior mean $\mu_{[\mathbf{X}, \mathbf{y}, \boldsymbol{\sigma}]}(\cdot)$ and covariance $s_{[\mathbf{X}, \boldsymbol{\sigma}]}(\cdot, \cdot)$ functions.

Since the derivative operator is linear, the derivative $\nabla_{\mathbf{x}} f = (\partial_1 f, \dots, \partial_D f)^\top \in \mathbb{R}^D$, where we abbreviate $\partial_d = \frac{\partial}{\partial x_d}$, of GP samples also follows a GP. Therefore, we can straightforwardly handle the derivative outputs at training and test points by modifying the kernel function. Assume that \mathbf{x} is a training or test point with non-derivative output $y = f^*(\mathbf{x}) + \varepsilon$, and \mathbf{x}' and \mathbf{x}'' are training or test points with derivative outputs, $y' = \partial_{d'} f^*(\mathbf{x}') + \varepsilon'$, $y'' = \partial_{d''} f^*(\mathbf{x}'') + \varepsilon''$. Then, the kernel functions should be replaced with

$$\tilde{k}(\mathbf{x}, \mathbf{x}') = \frac{\partial}{\partial x_{d'}} k(\mathbf{x}, \mathbf{x}'), \quad (6)$$

$$\tilde{k}(\mathbf{x}', \mathbf{x}'') = \frac{\partial^2}{\partial x_{d'} \partial x_{d''}} k(\mathbf{x}', \mathbf{x}''). \quad (7)$$

The posterior (3) with appropriately replaced kernel matrix entries gives the posterior distribution of derivatives at test points. We denote the GP posterior of a single component of the derivative as

$$p(\partial_d f(\cdot) | \mathbf{X}, \mathbf{y}) = \text{GP} \left(\partial_d f(\cdot); \tilde{\mu}_{[\mathbf{X}, \mathbf{y}, \boldsymbol{\sigma}]}^{(d)}(\cdot), \tilde{s}_{[\mathbf{X}, \boldsymbol{\sigma}]}^{(d)}(\cdot, \cdot) \right) \quad (8)$$

with the posterior mean $\tilde{\mu}^{(d)}(\cdot)$ and covariance $\tilde{s}^{(d)}(\cdot, \cdot)$ functions for the derivative with respect to x_d . More generally, GP regression can be analytically performed in the case where the training outputs (i.e., observations) and the test outputs (i.e., predictions) contain derivatives with different orders (see Appendix A for more details).

2.2. Variational Quantum Eigensolvers and their Physical Properties

The VQE (Peruzzo et al., 2014; McClean et al., 2016) is a hybrid quantum-classical computing protocol for estimating the ground-state energy of a given quantum Hamiltonian for a Q -qubit system. The quantum computer is used to prepare a parametric quantum state $|\psi_{\mathbf{x}}\rangle$, which depends on D angular parameters $\mathbf{x} \in \mathcal{X} = [0, 2\pi)^D$. This trial state $|\psi_{\mathbf{x}}\rangle$ is generated by applying $D' (\geq D)$ quantum gate operations, $G(\mathbf{x}) = G_{D'} \circ \dots \circ G_1$, to an initial quantum state $|\psi_0\rangle$, i.e., $|\psi_{\mathbf{x}}\rangle = G(\mathbf{x})|\psi_0\rangle$. All gates $\{G_{d'}\}_{d'=1}^{D'}$ are unitary operators, parameterized by at most one variable x_d . Let $d(d') : \{1, \dots, D'\} \mapsto \{1, \dots, D\}$ be the mapping specifying which one of the variables $\{x_d\}$ parameterizes the d' -th gate. We consider parametric gates of the form $G_{d'}(x) = U_{d'}(x_{d(d')}) = \exp(-ix_{d(d')} P_{d'}/2)$, where $P_{d'}$ is an arbitrary sequence of the Pauli operators $\{\mathbf{1}_q, \sigma_q^X, \sigma_q^Y, \sigma_q^Z\}_{q=1}^Q$ acting on each qubit at most once. This general structure covers both single-qubit gates, such as $R_X(x) = \exp(-i\theta\sigma_q^X)$, and entangling gates acting on multiple qubits simultaneously, such as $R_{XX}(x) = \exp(-ix\sigma_{q_1}^X \circ \sigma_{q_2}^X)$ for $q_1 \neq q_2$, commonly realized in trapped-ion quantum hardware setups (Kielpinski et al., 2002; Debnath et al., 2016).

The quantum computer is used to evaluate the energy of the resulting quantum state $|\psi_{\mathbf{x}}\rangle$ by observing

$$y = f^*(\mathbf{x}) + \varepsilon, \quad \text{where}$$

$$f^*(\mathbf{x}) = \langle \psi_{\mathbf{x}} | H | \psi_{\mathbf{x}} \rangle = \langle \psi_0 | G(\mathbf{x})^\dagger H G(\mathbf{x}) | \psi_0 \rangle, \quad (9)$$

and \dagger denotes the Hermitian conjugate. For each observation, repeated measurements, called *shots*, on the quantum computer are performed. Averaging over the number N_{shots} of shots suppresses the variance $\sigma^{*2}(N_{\text{shots}}) \propto N_{\text{shots}}^{-1}$ of the observation noise ε .¹ Since the observation y is the

¹We do not consider the hardware noise, and therefore, the observation noise ε consists only of the *measurement shot* noise.

sum of many random variables, it approximately follows the Gaussian distribution, according to the central limit theorem. The Gaussian likelihood (1) therefore approximates the observation y well if $\sigma_n^2 \approx \sigma^{*2}(N_{\text{shots}})$. Using the noisy estimates of $f^*(\mathbf{x})$ obtained from the quantum computer, a protocol running on a classical computer is used to solve the following minimization problem:

$$\min_{\mathbf{x} \in [0, 2\pi)^D} f^*(\mathbf{x}), \quad (10)$$

thus finding the minimizer $\hat{\mathbf{x}}$, i.e., the optimal parameters for the (rotational) quantum gates. Given the high expense of quantum computing resources, the computation cost is primarily driven by quantum operations. As a result, the optimization cost in VQE is typically measured by the total number of measurement shots required during the optimization process.² We refer to Tilly et al. (2022) for further details about VQEs and their challenges.

Let V_d be the number of gates parameterized by x_d , i.e., $V_d = |\{d' \in \{1, \dots, D\}; d = d(d')\}|$. Mitarai et al. (2018) proved that the VQE objective (9) for $V_d = 1$ satisfies the parameter shift rule (PSR)

$$\partial_d f^*(\mathbf{x}') = \frac{f^*(\mathbf{x}' + \alpha \mathbf{e}_d) - f^*(\mathbf{x}' - \alpha \mathbf{e}_d)}{2 \sin \alpha}, \quad \forall \mathbf{x} \in [0, 2\pi)^D, d = 1, \dots, D, \alpha \in [0, 2\pi), \quad (11)$$

where $\{\mathbf{e}_d\}_{d=1}^D$ are the standard basis, and the *shift* α is typically set to $\frac{\pi}{2}$. Wierichs et al. (2022) generalized the PSR (11) for arbitrary V_d with equidistant observations $\{\mathbf{x}_w = \mathbf{x}' + \frac{2w+1}{2V_d} \pi \mathbf{e}_d\}_{w=0}^{2V_d-1}$:

$$\partial_d f^*(\mathbf{x}') = \frac{1}{2V_d} \sum_{w=0}^{2V_d-1} \frac{(-1)^w f^*(\mathbf{x}_w)}{2 \sin^2 \left(\frac{(2w+1)\pi}{4V_d} \right)}. \quad (12)$$

Most gradient-based approaches rely on those PSRs, which allow reasonably accurate gradient estimation from $\sum_{d=1}^D 2V_d$ observations. Let

$$\psi_\gamma(\theta) = (\gamma, \sqrt{2} \cos \theta, \sqrt{2} \cos 2\theta, \dots, \sqrt{2} \cos V_d \theta, \sqrt{2} \sin \theta, \sqrt{2} \sin 2\theta, \dots, \sqrt{2} \sin V_d \theta)^\top \in \mathbb{R}^{1+2V_d} \quad (13)$$

be the (1-dimensional) V_d -th order Fourier basis for arbitrary $\gamma > 0$. Nakanishi et al. (2020) found that the VQE objective function $f^*(\cdot)$ in Eq. (9) with any³ $G(\cdot)$, H , and $|\psi_0\rangle$ can be expressed exactly as

$$f^*(\mathbf{x}) = \mathbf{b}^\top \mathbf{vec} \left(\otimes_{d=1}^D \psi_\gamma(x_d) \right) \quad (14)$$

²When the Hamiltonian consists of N_{og} groups of non-commuting operators, each of which needs to be measured separately, N_{shots} denotes the number of shots *per operator group*. Therefore, the number of shots *per observation* is $N_{\text{og}} \times N_{\text{shots}}$. In our experiments, we report on the total number of shots per operator group, i.e., the cumulative sum of N_{shots} over all observations, when evaluating the observation cost.

³Any circuit consisting of parametrized rotation gates and non-parametric unitary gates.

for some $\mathbf{b} \in \mathbb{R}^{\prod_{d=1}^D (1+2V_d)}$, where \otimes and $\mathbf{vec}(\cdot)$ denote the tensor product and the vectorization operator for a tensor, respectively. Based on this property, the Nakanishi-Fuji-Todo (NFT) method (Nakanishi et al., 2020) performs SMO (Platt, 1998), where the optimum in a chosen 1D subspace for each iteration is analytically estimated from only $1 + 2V_d$ observations (see Appendix B for the detailed procedure). It was shown that the PSR (11) and the trigonometric polynomial function form (14) are mathematically equivalent (Nicoli et al., 2023a).

Inspired by the function form (14) of the objective, Nicoli et al. (2023a) proposed the VQE kernel

$$k_\gamma(\mathbf{x}, \mathbf{x}') = \sigma_0^2 \prod_{d=1}^D \left(\frac{\gamma^2 + 2 \sum_{v=1}^{V_d} \cos(v(x_d - x'_d))}{\gamma^2 + 2V_d} \right), \quad (15)$$

which is decomposed as $k_\gamma(\mathbf{x}, \mathbf{x}') = \phi_\gamma(\mathbf{x})^\top \phi_\gamma(\mathbf{x}')$ with feature maps $\phi_\gamma(\mathbf{x}) = \frac{\sigma_0}{(\gamma^2 + 2V_d)^{D/2}} \mathbf{vec} \left(\otimes_{d=1}^D \psi_\gamma(x_d) \right)$, for GP regression. The kernel parameter γ^2 controls the smoothness of the function, i.e., suppressing the interaction terms when $\gamma^2 > 1$. When $\gamma^2 = 1$, the Fourier basis (13) is orthonormal, and the VQE kernel (15) is proportional to the product of Dirichlet kernels (Rudin, 1964). The VQE kernel reflects the physical knowledge (14) of VQE, and thus allows us to perform a Bayesian variant of NFT—*Bayesian NFT* or *Bayesian SMO*—where the 1D subspace optimization in each SMO step is performed with GP (see Appendix B for more details and the performance comparison between the original NFT and Bayesian NFT). Nicoli et al. (2023a) furthermore enhanced Bayesian NFT with BO, using the notion of confident region (CoRe),

$$\mathcal{Z}_{[\mathbf{X}, \sigma]}(\kappa^2) = \{ \mathbf{x} \in \mathcal{X}; s_{[\mathbf{X}, \sigma]}(\mathbf{x}, \mathbf{x}) \leq \kappa^2 \}, \quad (16)$$

i.e., the region in which the uncertainty of the GP prediction is lower than a threshold κ . More specifically, they introduced the EMICoRe acquisition function to find the best observation points in each SMO iteration, such that the maximum expected improvement within the CoRe is maximized.

3. Bayesian Parameter Shift Rules

We propose the *Bayesian PSR*, which estimates the gradient of the VQE objective (9) by the GP posterior (8) with the VQE kernel (15) along with its derivatives (6) and (7). The advantages of the Bayesian PSR include the following:

- The gradient estimator has an analytic-form.
- The estimation can be performed using observations at any set of points.
- The estimation is optimal for heteroschedastically noisy observations (from the Bayesian perspective),

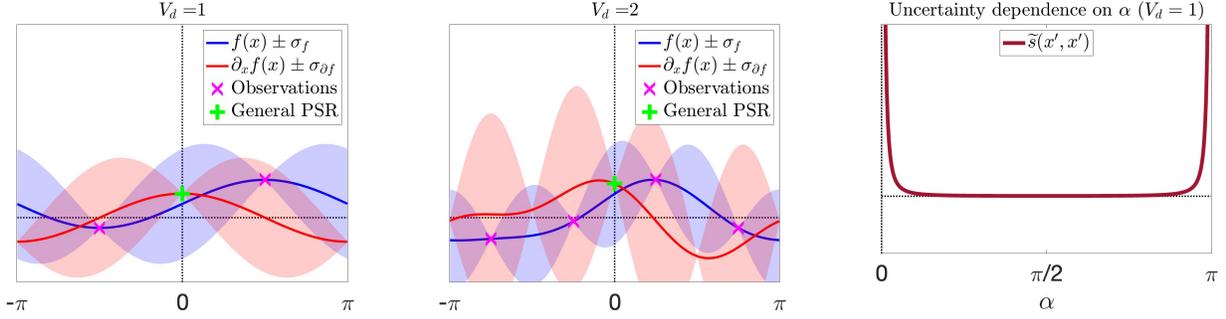


Figure 2. Illustration of the behavior of the Bayesian PSR when $V_d = 1$ (left) and when $V_d = 2$ (middle). Bayesian PSR prediction (red) coincides with general PSR (green cross) for the designed equidistant observations (magenta crosses). The right plot visualizes the variance (20) of derivative GP prediction at x' , as a function of the shift α of observations when $V_d = 1$. Although the optimum is at $\alpha = \frac{\pi}{2}$, the dependence is weak. For all panels, the noise and kernel parameters are set to $\sigma^2 = 0.01, \gamma^2 = 9, \sigma_0^2 = 100$.

as long as the prior with the kernel parameters, γ and σ_0^2 , is appropriately set.

- The posterior uncertainty can be analytically computed before performing the observations.

In Section 4, we propose novel SGD solvers for VQEs that leverage the advantages of the Bayesian PSR.

As naturally expected, our Bayesian PSR is a generalization of existing PSRs, and reduces to the general PSR (12) for noiseless and equidistant observations. Let $\mathbf{1}_D \in \mathbb{R}^D$ be the vector with all entries equal to one.

Theorem 3.1. For any $x' \in [0, 2\pi)^D$ and $d = 1, \dots, D$, the mean and variance of the derivative GP prediction, given observations $\mathbf{y} = (y_0, \dots, y_{2V_d-1})^\top \in \mathbb{R}^{2V_d}$ at $2V_d$ equidistant training points $\mathbf{X} = (\mathbf{x}_0, \dots, \mathbf{x}_{2V_d-1}) \in \mathbb{R}^{D \times 2V_d}$ for $\mathbf{x}_w = \mathbf{x}' + \frac{2w+1}{2V_d}\pi \mathbf{e}_d$ with homoschedastic noise $\boldsymbol{\sigma} = \sigma^2 \cdot \mathbf{1}_{2V_d}$ for $\sigma^2 \ll \sigma_0$, are

$$\tilde{\mu}_{[\mathbf{X}, \mathbf{y}, \boldsymbol{\sigma}]}^{(d)}(\mathbf{x}') = \frac{\sum_{w=0}^{2V_d-1} \frac{(-1)^w y_w}{2 \sin^2\left(\frac{(2w+1)\pi}{4V_d}\right)}}{(\gamma^2 + 2V_d) \frac{\sigma^2}{\sigma_0^2} + 2V_d} + O\left(\frac{\sigma^4}{\sigma_0^4}\right), \quad (17)$$

$$\tilde{s}_{[\mathbf{X}, \boldsymbol{\sigma}]}^{(d)}(\mathbf{x}', \mathbf{x}') = \sigma^2 \frac{2V_d^2 + 1}{6} + O\left(\frac{\sigma^4}{\sigma_0^4}\right). \quad (18)$$

The proof, the non-asymptotic form of the mean and the variance, and a numerical validation are given in Appendix C. Apparently, the mean prediction (17) by Bayesian PSR converges to the general PSR (12) with the uncertainty (18) converging to zero in the noiseless limit, i.e., $\sigma^2 \rightarrow +0$ and hence $y_w = f^*(\mathbf{x}_w)$. In noisy cases, the prior variance $\sigma_0^2 \sim O(\sigma^2)$ suppresses the amplitude of the gradient estimator as a regularizer through the first term in the denominator in Eq. (17).

Figure 2 illustrates the behavior of the Bayesian PSR when $V_d = 1$ (left panel) and when $V_d = 2$ (middle panel). In each panel, given $2V_d$ equidistant observations (magenta

crosses), the blue curve shows the (non-derivative) GP prediction with uncertainty (blue shades), while the red curve shows the derivative GP prediction with uncertainty (red shades). Note the $\frac{\pi}{2V_d}$ shift of the low uncertainty locations between the GP prediction (blue) and the derivative GP prediction (red). The green cross shows the output of the general PSR (12) at $x' = 0$, which almost coincides with the Bayesian PSR prediction (red curve) under this setting. Other examples, including cases where the Bayesian regularization is visible, are given in Appendix C.

In the simplest first-order case, i.e., where $V_d = 1, \forall d = 1, \dots, D$, we can theoretically investigate the optimality of the choice of the shift α in Eq. (11) (the proof is also given in Appendix C).

Theorem 3.2. Assume that $V_d = 1, \forall d = 1, \dots, D$. For any $x' \in [0, 2\pi)^D$ and $d = 1, \dots, D$, the mean and variance of the derivative GP prediction, given observations $\mathbf{y} = (y_1, y_2)^\top \in \mathbb{R}^2$ at two training points $\mathbf{X} = (\mathbf{x}' - \alpha \mathbf{e}_d, \mathbf{x}' + \alpha \mathbf{e}_d) \in \mathbb{R}^{D \times 2}$ with homoschedastic noise $\boldsymbol{\sigma} = (\sigma^2, \sigma^2)^\top$, are

$$\tilde{\mu}_{[\mathbf{X}, \mathbf{y}, \boldsymbol{\sigma}]}^{(d)}(\mathbf{x}') = \frac{(y_2 - y_1) \sin \alpha}{(\gamma^2/2 + 1) \sigma^2 / \sigma_0^2 + 2 \sin^2 \alpha}, \quad (19)$$

$$\tilde{s}_{[\mathbf{X}, \boldsymbol{\sigma}]}^{(d)}(\mathbf{x}', \mathbf{x}') = \frac{\sigma^2}{(\gamma^2/2 + 1) \sigma^2 / \sigma_0^2 + 2 \sin^2 \alpha}. \quad (20)$$

Again, the mean prediction (19) is a regularized version of the PSR (11). The uncertainty prediction (20) implies that $\alpha = \pi/2$ minimizes the uncertainty in the noisy case, regardless of σ^2, σ_0^2 and γ . This supports most of the use cases of the PSR in the literature (Mitarai et al., 2018), and matches the intuition that the maximum span minimizes the uncertainty. However, the right panel in Figure 2, where the variance (20) of the derivative GP prediction at x' is visualized as a function of the shift α of observations for $V_d = 1$, implies that the estimation accuracy is not very sensitive to the choice of α .

4. SGD with Bayesian PSR

In this section, we equip SGD with the Bayesian PSR. In the standard implementation of SGD for VQE, $2V_d$ equidistant points along each direction $d = 1, \dots, D$ are observed for gradient estimation by the general PSR (12) (or by the PSR (11) if $V_d = 1, \forall d$) in each SGD iteration.

Bayesian SGD (Bayes-SGD): A straightforward application of the Bayesian PSR is to replace existing PSRs with the Bayesian PSR for gradient estimation, allowing for the reuse of previous observations. We retain $R \cdot 2V_d \cdot D$ latest observations for a predetermined R in our experiments. We expect that reusing previous observations accumulates the gradient information, and thus improves the gradient estimation accuracy.

4.1. Gradient Confident Region (GradCoRe)

We propose an adaptive observation cost control strategy that leverages the uncertainty information provided by the Bayesian PSR. This strategy adjusts the number of measurement shots for gradient estimation in each SGD iteration so that the variances of the derivative GP prediction at the current optimal point $\hat{\mathbf{x}}$ are below certain thresholds. In a fashion similar to the CoRe (16), we define the *gradient confident region* (GradCoRe)

$$\tilde{\mathcal{Z}}_{[\mathbf{X}, \sigma]}(\boldsymbol{\kappa}) = \left\{ \mathbf{x} \in \mathcal{X}; \tilde{s}_{[\mathbf{X}, \sigma]}^{(d)}(\mathbf{x}, \mathbf{x}) \leq \kappa_d^2, \forall d \right\}, \quad (21)$$

where $\boldsymbol{\kappa} = (\kappa_1^2, \dots, \kappa_D^2)^\top \in \mathbb{R}^D$ are the required accuracy thresholds. Our proposed SGD-based optimizer, named *SGD-GradCoRe*, measures new equidistant points $\tilde{\mathbf{X}} = \left\{ \mathbf{x}_w^{(d)} = \hat{\mathbf{x}} + \frac{2w+1}{2V_d} \pi \mathbf{e}_d \right\}_{w=0}^{2V_d}$ for all directions with the minimum total number of shots such that the current optimal point $\hat{\mathbf{x}}$ is in the GradCoRe (see Figure 1).

Following Anders et al. (2024), we estimate the single-shot observation noise variance $\sigma_1^{*2} = \sigma^{*2}(1)$ before the optimization by collecting measurements at random locations in order to estimate the observation noise variance as a function of the number of shots as

$$\sigma^{*2}(N_{\text{shots}}) = \frac{\sigma_1^{*2}}{N_{\text{shots}}}. \quad (22)$$

Let $(\mathbf{X}^t, \mathbf{y}^t, \boldsymbol{\sigma}^t)$ be the training data (all previous observations) at the t -th SGD iteration step, and let $\tilde{\nu} \in \mathbb{R}^{2V_d D}$ be the vector of the numbers of measurement shots at the new equidistant measurement points $\tilde{\mathbf{X}}$ for all directions. Before measuring at $\tilde{\mathbf{X}}$ in the $(t+1)$ -th SGD iteration, we solve the following:

$$\min_{\tilde{\nu}} \|\tilde{\nu}\|_1 \text{ s.t. } \hat{\mathbf{x}} \in \tilde{\mathcal{Z}}_{[(\mathbf{X}^t, \tilde{\mathbf{X}}), (\boldsymbol{\sigma}^t, \check{\sigma}(\tilde{\nu}))]}(\boldsymbol{\kappa}(t)), \quad (23)$$

where $\check{\sigma}(\tilde{\nu}) = \sigma_1^{*2} \cdot (\tilde{\nu}_1^{-1}, \dots, \tilde{\nu}_{2V_d D}^{-1})^\top$, and $\boldsymbol{\kappa}(t)$ is the required accuracy dependent on the iteration step t . Informally, we minimize the total measurement budget under the

constraint that the posterior gradient variance along each direction d is smaller than the required accuracy threshold. For simplicity, we solve the GradCoRe problem (23) by grid search over $[\kappa_d^2, \sigma_1^{*2}] \forall d$ under the additional constraint that all $2V_d D$ points are measured with an equal number of shots.

We set the required accuracy thresholds $\boldsymbol{\kappa}(t) = \kappa^2(t) \mathbf{1}_D$, where

$$\kappa^2(t) = \max \left(c_0, \frac{c_1}{D} \sum_{d=1}^D \left(\tilde{\mu}_{[\mathbf{X}^t, \mathbf{y}^t, \boldsymbol{\sigma}^t]}^{(d)}(\hat{\mathbf{x}}^t) \right)^2 \right). \quad (24)$$

Namely, $\kappa(t)$ is set proportional to the L2-norm of the estimated gradient at the current optimal point at the t -th SGD iteration, as long as it is larger than a lower bound. The lower bound c_0 and the slope c_1 are hyperparameters to be tuned. This strategy for setting the required accuracy based on the estimated gradient norm was proposed by Tamiya & Yamasaki (2022). Alternatively, one could also set $\kappa_d(t)$ proportional to the absolute value of the estimated gradient separately for *each* direction, i.e., $\kappa_d(t) = \max(c_0, c_1 |\tilde{\mu}_{[\mathbf{X}^t, \mathbf{y}^t, \boldsymbol{\sigma}^t]}^{(d)}(\hat{\mathbf{x}}^t)|)$, and solve the GradCoRe problem (23) direction-wise.

In the experiment plots in Section 5, we will refer to SGD-GradCoRe as *GradCoRe*. Further algorithmic details, including pseudo-code and used hyperparameters, are given in Appendix D.

5. Experiments

5.1. Setup

We demonstrate the performance of our Bayesian PSR and GradCoRe approaches in the same setup used by Nicoli et al. (2023a). For all experiments, we prepared 50 different random initial points from which all methods are initialized. Our Python implementation uses Qiskit (Abraham et al., 2019) for the classical simulation of quantum hardware. The implementation for reproducing our results is attached as supplemental material.

Hamiltonian and Quantum Circuit: We focus on the quantum Heisenberg Hamiltonian with open boundary conditions,

$$H = - \sum_{i \in \{X, Y, Z\}} \left[\sum_{j=1}^{Q-1} (J_i \sigma_j^i \sigma_{j+1}^i) + \sum_{j=1}^Q h_i \sigma_j^i \right], \quad (25)$$

where $\{\sigma_j^i\}_{i \in \{X, Y, Z\}}$ are the Pauli operators acting on the j -th qubit. For the quantum circuit, we use a common ansatz, called the L -layered Efficient SU(2) circuit with open boundary conditions, where $V_d = 1, \forall d$ (see Nicoli et al. (2023a) for more details).

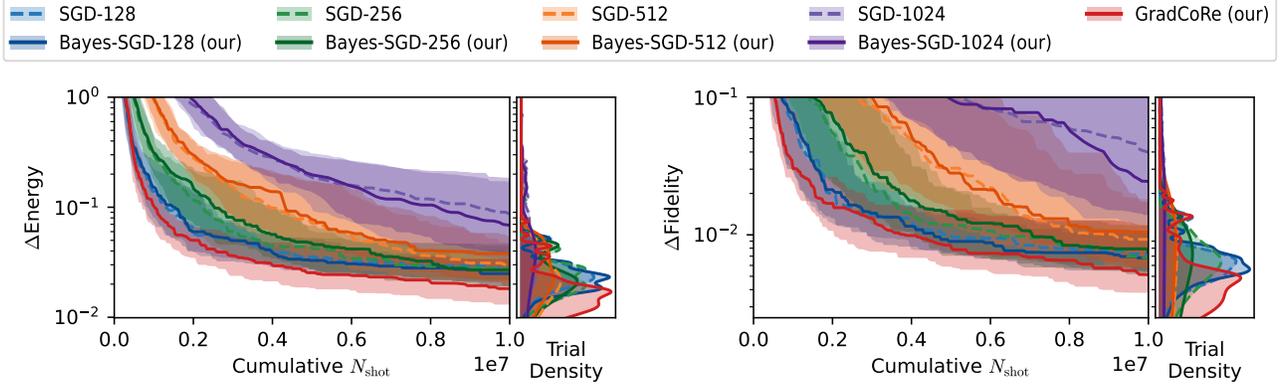


Figure 3. Comparison between SGD with PSR (dashed curves) and SGD with Bayesian PSR (solid curves), as well as GradCoRe (red solid curve), on the Ising Hamiltonian with $(L = 3)$ -layered $(Q = 5)$ -qubits quantum circuit. The energy (left) and fidelity (right) are plotted as function of the cumulative N_{shots} , i.e., the total number of measurement shots. Except GradCoRe equipped with the adaptive shots strategy, the number of shots per observation is set to $N_{\text{shots}} = 128$ (blue), 256 (green), 512 (orange), and 1024 (purple).

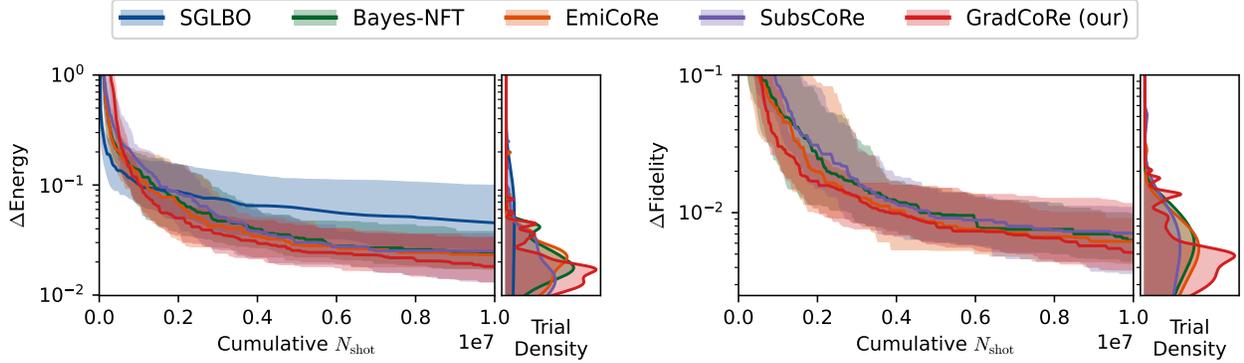


Figure 4. Energy (left) and fidelity (right) achieved within the cumulative number of measurement shots for the Ising Hamiltonian with an $(L = 3)$ -layered $(Q = 5)$ -qubits quantum circuit. The curves corresponds to SGLBO (blue), Bayes-NFT (green), EMICoRe (orange), SubsCoRe (purple), and our proposed GradCoRe (red).

Evaluation Metrics: We compare all methods using two metrics: the cumulatively lowest *true energy* $f^*(\hat{\mathbf{x}})$, for $f^*(\cdot)$ defined in Eq. (9), and the *fidelity* $\langle \psi_{\text{GS}} | \psi_{\hat{\mathbf{x}}} \rangle \in [0, 1]$. The latter is the inner product between the true ground-state wave function $|\psi_{\text{GS}}\rangle$, computed by exact diagonalization of the target Hamiltonian H , and the trial wave function, $|\psi_{\hat{\mathbf{x}}}\rangle$, corresponding to the quantum state generated by the circuit using the optimized parameters $\hat{\mathbf{x}}$. For both metrics, we plot the difference (smaller is better) to the respective target, i.e.,

$$\begin{aligned} \Delta \text{Energy} &= \langle \psi_{\hat{\mathbf{x}}} | H | \psi_{\hat{\mathbf{x}}} \rangle - \langle \psi_{\text{GS}} | H | \psi_{\text{GS}} \rangle \\ &= f^*(\hat{\mathbf{x}}) - \langle \psi_{\text{GS}} | H | \psi_{\text{GS}} \rangle, \end{aligned} \quad (26)$$

$$\begin{aligned} \Delta \text{Fidelity} &= \langle \psi_{\text{GS}} | \psi_{\text{GS}} \rangle - \langle \psi_{\text{GS}} | \psi_{\hat{\mathbf{x}}} \rangle \\ &= 1 - \langle \psi_{\text{GS}} | \psi_{\hat{\mathbf{x}}} \rangle, \end{aligned} \quad (27)$$

in log scale. Here, $|\psi_{\text{GS}}\rangle$ and $\langle \psi_{\text{GS}} | H | \psi_{\text{GS}} \rangle$ are the ground-state wave function and the true energy at the ground-state, respectively, both of which are computed analytically. As a measure of the quantum computation cost, we consider the total number of measurement shots *per operator group*

(see Footnote 2) for all observations over the whole optimization process.

Baseline Methods: We compare our Bayesian SGD and GradCoRe approaches to the baselines, including SGD, NFT (Nakanishi et al., 2020), Bayesian NFT, SGLBO (Tamiya & Yamasaki, 2022), EMICoRe (Nicoli et al., 2023a), and SubsCoRe (Anders et al., 2024). SGD uses the PSR (11) for gradient estimation.

Algorithm Setting: All SGD-based methods use the ADAM optimizer with $l_r = 0.05$, $\beta_s = (0.9, 0.999)$. For the methods not equipped with adaptive cost control (i.e., all methods except SGLBO, SubsCoRe and GradCoRe) we set $N_{\text{shots}} = 1024$ for each observation, the same setting as in Nicoli et al. (2023a), unless specified explicitly. To avoid error accumulation, all SMO-based methods measure the “center”, i.e., the current optimal point without shift, every $D + 1$ iterations (Nakanishi et al., 2020). Bayes-SGD and GradCoRe estimate the gradient from the $R \cdot 2V_d \cdot D$

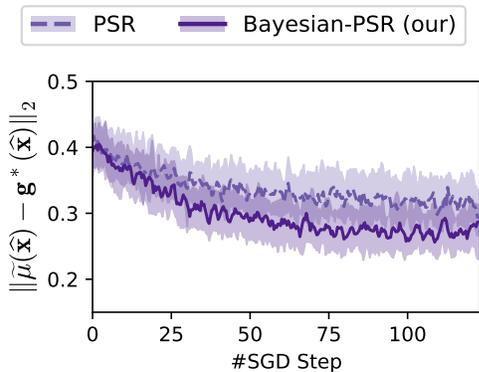


Figure 5. Gradient estimation error by PSR (dashed curve) and Bayesian PSR (solid curve) for $N_{\text{shots}} = 1024$, evaluated by the L2-distance between the estimated gradient $\tilde{\mu}(\hat{\mathbf{x}})$ and the true gradient $\mathbf{g}^*(\hat{\mathbf{x}})$ (computed by PSR with simulated noiseless measurements).

latest observations for $R = 5$. GradCoRe initially uses the fixed threshold $\kappa^2(t) = \sigma_1^2/256$ before starting the cost adaption after D SGD iterations.

Further details on the algorithmic and experimental settings are described in Appendix D and Appendix E, respectively.

5.2. Improvement over SGD with Bayesian PSR and GradCoRe

First, we investigate the potential improvement of our Bayesian PSR and GradCoRe over plain SGD. Figure 3 compares SGD with the standard PSR (SGD) and SGD with Bayesian PSR (Bayes-SGD) on the Ising Hamiltonian, i.e., Eq. (25) with $J_{i \in \{X, Y, Z\}} = (-1, 0, 0)$ and $h_{i \in \{X, Y, Z\}} = (0, 0, -1)$, with $(L = 3)$ -layered $(Q = 5)$ -qubits quantum circuit. Both standard and Bayesian PSR are shown with $N_{\text{shots}} = 128, 256, 512, 1024$ measurement shots. The left and right panels plot the difference to the ground-state in true energy (26) and fidelity (27) achieved by each method as functions of the cumulative N_{shots} , i.e., the total number of measurement shots. The *trial density* to the right of each panel shows a kernel-density estimation of the true energy distribution over the trials after 1×10^7 measurement shots. The median, the 25-th and the 75-th percentiles are shown as a solid curve and shades, respectively. We observe that Bayesian PSR, with a more accurate gradient estimator as shown in Figure 5, is comparable or compares favorably to the original SGD. More importantly, we observe that GradCoRe automatically selects the optimal number of measurement shots in each optimization phase, thus outperforming SGD and Bayes-SGD with fixed N_{shots} through the entire optimization process. The adaptively selected measurement shots and accuracy threshold $\kappa(t)$ for GradCoRe are shown in Appendix F.

5.3. Comparison with State-of-the-art Methods

Figure 4 compares GradCoRe to the baseline methods SGLBO, Bayes-NFT, EMICoRe, and SubsCoRe. Our GradCoRe method, which significantly improves upon SGD as shown in Figure 3, establishes itself as the new state-of-the-art, exhibiting faster convergence and achieving lower overall energy. We excluded the original NFT in this comparison, as it is outperformed by Bayes-NFT in all observed settings (see Figure 6 in Appendix B).

6. Conclusion

The physical properties of variational quantum eigensolvers (VQEs) allow for the use of specialized optimization methods, i.e., stochastic gradient descent (SGD) with parameter shift rules (PSRs) and a specialized sequential minimal optimization (SMO), called NFT (Nakanishi et al., 2020). Contemporary research has shown that those properties can be appropriately captured by the physics-informed VQE kernel, with which NFT has been successfully improved through Bayesian machine learning techniques. For instance, previous observations could be used to determine the optimal measurement points (Nicoli et al., 2023a) and computational cost could be minimized based on the uncertainty prediction (Anders et al., 2024). In this paper, we have shown that a similar approach can also improve SGD-based methods. Specifically, we proposed the Bayesian PSR, where the gradient is estimated by derivative GPs. The Bayesian PSR generalizes existing PSRs to allow for flexible estimation from observations at an arbitrary set of locations. Furthermore, it provides uncertainty information, which enables observation cost adaptation through the novel notion of gradient confident region (GradCoRe). Our theoretical analysis revealed the relation between Bayesian PSR and existing PSRs, while our numerical investigation empirically demonstrated the utility of our approaches. We envisage that Bayesian approaches will facilitate further development of more efficient algorithms for VQEs and, more generally, quantum computing. In future work, we aim to explore the optimal combination of existing methods and strategies for selecting the most suitable approaches for specific tasks, i.e., specific Hamiltonians.

Impact Statement

This paper presents work whose goal is to advance the field of machine learning and quantum computing. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

References

- Abraham, H. et al. Qiskit: An open-source framework for quantum computing. *Zenodo*, 2019. doi: 10.5281/zenodo.2562111.
- Acharya, R., Abanin, D. A., et al. Quantum error correction below the surface code threshold. *Nature*, 2024. doi: 10.1038/s41586-024-08449-y.
- Anders, C. J., Nicoli, K., Wu, B., Elosegui, N., Pedrielli, S., Funcke, L., Jansen, K., Kuhn, S., and Nakajima, S. Adaptive observation cost control for variational quantum eigensolvers. In *Proceedings of 41st International Conference on Machine Learning (ICML2024)*, 2024. doi: 10.5555/3692070.3692133.
- Bluvstein, D., Evered, S. J., Geim, A. A., Li, S. H., Zhou, H., Manovitz, T., Ebadi, S., Cain, M., Kalinowski, M., Hangleiter, D., et al. Logical quantum processor based on reconfigurable atom arrays. *Nature*, pp. 1–3, 2023. doi: 10.1038/s41586-023-06927-3.
- Cai, Z., Babbush, R., Benjamin, S. C., Endo, S., Huggins, W. J., Li, Y., McClean, J. R., and O’Brien, T. E. Quantum error mitigation. *Rev. Mod. Phys.*, 95:045005, Dec 2023. doi: 10.1103/RevModPhys.95.045005.
- Debnath, S., Linke, N. M., Figgatt, C., Landsman, K. A., Wright, K., and Monroe, C. Demonstration of a small programmable quantum computer with atomic qubits. *Nature*, 536(7614):63–66, 2016. doi: 10.1038/nature18648.
- Frazier, P. A tutorial on Bayesian optimization. *ArXiv e-prints*, 2018. doi: 10.48550/arXiv.1807.02811.
- Iannelli, G. and Jansen, K. Noisy Bayesian optimization for variational quantum eigensolvers. *ArXiv e-prints*, 2021. doi: 10.48550/arXiv.2112.00426.
- Jiang, T., Rogers, J., Frank, M. S., Christiansen, O., Yao, Y.-X., and Lanatà, N. Error mitigation in variational quantum eigensolvers using tailored probabilistic machine learning. *Phys. Rev. Res.*, 6:033069, Jul 2024. doi: 10.1103/PhysRevResearch.6.033069.
- Kielpinski, D., Monroe, C., and Wineland, D. J. Architecture for a large-scale ion-trap quantum computer. *Nature*, 417(6890):709–711, 2002. doi: 10.1038/nature00784.
- McClean, J. R., Romero, J., Babbush, R., et al. The theory of variational hybrid quantum-classical algorithms. *New Journal of Physics*, 18(2):023023, 2016. doi: 10.1088/1367-2630/18/2/023023.
- Mitarai, K., Negoro, M., Kitagawa, M., et al. Quantum circuit learning. *Phys. Rev. A*, 98:032309, 2018. doi: 10.1103/PhysRevA.98.032309.
- Nakanishi, K. M., Fujii, K., and Todo, S. Sequential minimal optimization for quantum-classical hybrid algorithms. *Phys. Rev. Res.*, 2:043158, 2020. doi: 10.1103/PhysRevResearch.2.043158.
- Nicoli, K. A., Anders, C. J., Funcke, L., Hartung, T., Jansen, K., Kuhn, S., Müller, K.-R., Stornati, P., Kessel, P., and Nakajima, S. Physics-informed Bayesian optimization of variational quantum circuits. In *Advances in Neural Information Processing Systems (NeurIPS2023)*, 2023a.
- Nicoli, K. A., Anders, C. J., et al. EMICoRe: Expected maximum improvement over confident regions. <https://github.com/emicores/emicores>, 2023b.
- Nicoli, K. A., Wagner, L., and Funcke, L. Machine-learning-enhanced optimization of noise-resilient variational quantum eigensolvers. *ArXiv e-prints*, 2025. doi: 10.48550/arXiv.2501.17689.
- Peruzzo, A., McClean, J., Shadbolt, P., et al. A variational eigenvalue solver on a photonic quantum processor. *Nature Communications*, 5(1):4213, 2014. doi: 10.1038/ncomms5213.
- Platt, J. Sequential minimal optimization : A fast algorithm for training support vector machines. *Microsoft Research Technical Report*, 1998.
- Preskill, J. Quantum computing in the NISQ era and beyond. *Quantum*, 2:79, August 2018. doi: 10.22331/q-2018-08-06-79.
- Rasmussen, C. E. and Williams, C. K. I. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, USA, 2006. doi: 10.7551/mitpress/3206.001.0001.
- Roffe, J. Quantum error correction: An introductory guide. *Contemporary Physics*, 60(3):226–245, 2019. doi: 10.1080/00107514.2019.1667078.
- Rudin, W. *Principles of Mathematical Analysis*. McGraw-Hill, 1964. doi: 10.1017/S0013091500008889.
- Tamiya, S. and Yamasaki, H. Stochastic gradient line Bayesian optimization for efficient noise-robust optimization of parameterized quantum circuits. *npj Quantum Information*, 8(1):90, 2022. doi: 10.1038/s41534-022-00592-6.
- Tilly, J., Chen, H., Cao, S., et al. The variational quantum eigensolver: A review of methods and best practices. *Physics Reports*, 986:1–128, 2022. doi: <https://doi.org/10.1016/j.physrep.2022.08.003>.
- Wierichs, D., Izaac, J., Wang, C., and Lin, C. Y.-Y. General parameter-shift rules for quantum gradients. *Quantum*, 6:677, March 2022. ISSN 2521-327X. doi: 10.22331/q-2022-03-30-677.

A. General Gaussian Processes (GPs) with Derivative Outputs

The derivative GP regression can be straightforwardly extended to the case where both training outputs (i.e., observations), and test outputs (i.e., predictions) contain different orders of derivatives.

Assume that we have a set of input points, and for each input point $\mathbf{x} \in \mathbb{R}^D$, the corresponding output, i.e., observation or prediction, is $f(\mathbf{x})$ or $\partial_{x_d} f(\mathbf{x})$, where $\partial_{x_d} \equiv \frac{\partial}{\partial x_d}$. Let us denote the derivative kernel functions as

$$\tilde{k}^{(d,d')}(\mathbf{x}, \mathbf{x}') = \begin{cases} k(\mathbf{x}, \mathbf{x}') & \text{if } d = 0, d' = 0, \\ \partial_{x_{d'}} k(\mathbf{x}, \mathbf{x}') & \text{if } d = 0, d' = 1, \dots, D, \\ \partial_{x_d} k(\mathbf{x}, \mathbf{x}') & \text{if } d = 1, \dots, D, d' = 0, \\ \partial_{x_d} \partial_{x_{d'}} k(\mathbf{x}, \mathbf{x}') & \text{if } d = 1, \dots, D, d' = 1, \dots, D. \end{cases}$$

For training points $\mathbf{X} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ and test points $\mathbf{X}' = \{\mathbf{x}'^{(m)}\}_{m=1}^M$, we should set the the entries of the train-train $\mathbf{K} \in \mathbb{R}^{N \times N}$, train-test $\mathbf{K}' \in \mathbb{R}^{N \times M}$, and test-test $\mathbf{K}'' \in \mathbb{R}^{M \times M}$ kernels as

$$K_{n,n'} = \tilde{k}^{(d(\mathbf{x}_n), d(\mathbf{x}_{n'}))}(\mathbf{x}_n, \mathbf{x}_{n'}), \quad (28)$$

$$K'_{n,m} = \tilde{k}^{(d(\mathbf{x}_n), d(\mathbf{x}_m))}(\mathbf{x}_n, \mathbf{x}_m), \quad (29)$$

$$K''_{m,m'} = \tilde{k}^{(d(\mathbf{x}_m), d(\mathbf{x}_{m'}))}(\mathbf{x}_m, \mathbf{x}_{m'}), \quad (30)$$

where

$$d(\mathbf{x}) = \begin{cases} 0 & \text{if the corresponding output for the input } \mathbf{x} \text{ is } f(\mathbf{x}), \\ d & \text{if the corresponding output for the input } \mathbf{x} \text{ is } \partial_{x_d} f(\mathbf{x}). \end{cases}$$

Eqs.(3)–(5) with the kernel matrices $\mathbf{K}, \mathbf{K}', \mathbf{K}''$ set as Eqs.(28)–(30) give the posterior GP for the corresponding test outputs.

For higher-order derivative outputs, we can define the kernels in exactly the same way as above, by applying the same derivative operators to the kernels as the ones applied to the outputs, i.e.,

$$\tilde{k}(\mathbf{x}, \mathbf{x}') = \left[\partial_{x_1}^{(r_1)} \dots \partial_{x_D}^{(r_D)} \right] \left[\partial_{x'_1}^{(r'_1)} \dots \partial_{x'_D}^{(r'_D)} \right] k(\mathbf{x}, \mathbf{x}'),$$

if the corresponding outputs at \mathbf{x} and \mathbf{x}' are $\partial_{x_1}^{(r_1)} \dots \partial_{x_D}^{(r_D)} f(\mathbf{x})$ and $\partial_{x'_1}^{(r'_1)} \dots \partial_{x'_D}^{(r'_D)} f(\mathbf{x}')$, respectively, where $\partial_{x_d}^{(r)} \equiv \frac{\partial^r}{\partial x_d^r}$ denotes the r -th order derivative with respect to x_d .

B. Nakanishi-Fuji-Todo (NFT) Algorithm (Nakanishi et al., 2020) and Bayesian NFT

Let $\{e_d\}_{d=1}^D$ be the standard basis. NFT is initialized with a random point $\hat{\mathbf{x}}^0$ with a first observation $\hat{y}^0 = f^*(\hat{\mathbf{x}}^0) + \varepsilon_0$, and iterates the following procedure: for each iteration step t ,

1. Select an axis $d \in \{1, \dots, D\}$ sequentially and observe the objective $\mathbf{y} \in \mathbb{R}^{2V_d}$ at $2V_d$ points $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_{2V_d}) = \{\hat{\mathbf{x}}^{t-1} + \alpha_w e_d\}_{w=1}^{2V_d} \in \mathbb{R}^{D \times 2V_d}$ along the axis d .⁴ Here $\alpha \in [0, 2\pi)^{2V_d}$ is such that $\alpha_w \neq 0$, $\alpha_{w'} \neq \alpha_w$, for all w and $w' \neq w$.
2. Apply the 1D trigonometric polynomial regression $\tilde{f}(\theta) = \tilde{\mathbf{b}}^\top \psi_1(\theta)$ to the $2V_d$ new observations \mathbf{y} , together with the previous best estimated score \hat{y}^{t-1} , and analytically compute the new optimum $\hat{\mathbf{x}}^t = \hat{\mathbf{x}}^{t-1} + \hat{\theta} e_d$, where $\hat{\theta} = \text{argmin}_\theta \tilde{f}(\theta)$.
3. Update the best score by $\hat{y}^t = \tilde{f}(\hat{\theta})$.

Note that if the observation noise is negligible, i.e., $y \approx f^*(\mathbf{x})$, each step of NFT reaches the global optimum in the 1D subspace along the chosen axis d for any choice of α , and thus performs SMO exactly. Otherwise, errors can be accumulated

⁴With slight abuse of notation, we use the set notation to specify the column vectors of a matrix, i.e., $(\mathbf{x}_1, \dots, \mathbf{x}_N) = \{\mathbf{x}_n\}_{n=1}^N$.

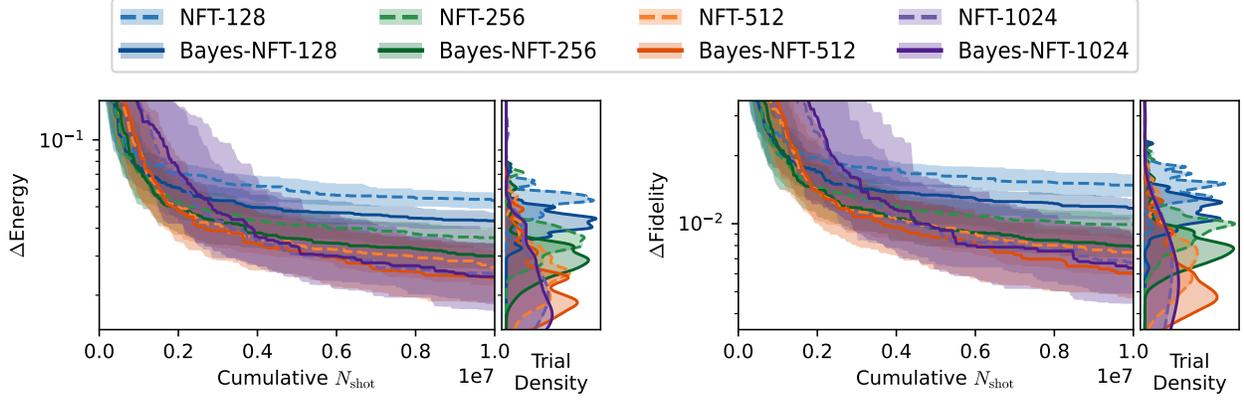


Figure 6. Comparison between NFT (Nakanishi et al., 2020) and Bayes-NFT for the Ising Hamiltonian with an ($L = 3$)-layered ($Q = 5$)-qubits quantum circuit. The energy (left) and fidelity (right), in the forms of Eqs.(26) and (27), respectively, are plotted as functions of the cumulative N_{shots} , i.e., the total number of measurement shots. The number of shots per observation is set to $N_{\text{shots}} = 128$ (blue), 256 (green), 512 (orange), and 1024 (purple).

in the best score \hat{y}^t , and therefore an additional measurement may need to be performed at \hat{x}^t after a certain iteration interval.

Bayesian NFT (Bayes-NFT) performs the 1D trigonometric polynomial regression and optimization in Step 2 with GP with the VQE kernel (15), where all previous observations are used for training. Using previous observations allows prediction with smaller uncertainty and thus more accurate subspace optimization. Figure 6 compares the original NFT and Bayesian NFT on the Ising Hamiltonian with an ($L = 3$)-layered ($Q = 5$)-qubits quantum circuit with different number of shots per observation. We observe that using GP generally accelerates the optimization process.

C. Proofs

Here, we give proofs of theorems in Section 3, and numerically validate them.

C.1. Proof of Theorem 3.1

We start from a more general theorem than Theorem 3.1, which is proven in Appendix C.3.

Theorem C.1. Assume that, for any given point $\hat{x} \in [0, 2\pi)^D$, we have observations $\mathbf{y} = (y_0, \dots, y_{2V_d-1})^T \in \mathbb{R}^{2V_d}$ at $2V_d$ equidistant training points $\mathbf{X} = (\mathbf{x}_0, \dots, \mathbf{x}_{2V_d-1}) \in \mathbb{R}^{D \times 2V_d}$ for $\mathbf{x}_w = \hat{x} + \frac{2w+1}{2V_d}\pi \mathbf{e}_d$ with homoschedastic noise $\sigma = \sigma^2 \cdot \mathbf{1}_{2V_d} \in \mathbb{R}^{2V_d}$. Then, the mean and variance of the derivative $\partial_d f(\mathbf{x}')$ prediction at $\mathbf{x}' = \hat{x} + \alpha' \mathbf{e}_d$ for any $d = 1, \dots, D$ and $\alpha' \in [0, 2\pi)$ are given as

$$\tilde{\mu}_{[\mathbf{X}, \mathbf{y}, \sigma]}^{(d)}(\mathbf{x}') = \frac{\sum_{w=0}^{2V_d-1} (-1)^w y_w \left(\frac{\cos(V_d \alpha')}{2 \sin^2\left(\frac{(2w+1)\pi}{4V_d} - \alpha'/2\right)} + \frac{V_d \sin\left(\frac{(2w+1)\pi}{4V_d} - (V_d+1/2)\alpha'\right)}{\sin\left(\frac{(2w+1)\pi}{4V_d} - \alpha'/2\right)} - \frac{4V_d^2 \cos V_d \alpha'}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d} \right)}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d}, \quad (31)$$

$$\tilde{s}_{[\mathbf{X}, \sigma]}^{(d)}(\mathbf{x}', \mathbf{x}') = \sigma^2 \left(\frac{V_d(V_d+1)(2V_d+1)}{3((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d)} - \frac{4V_d^3 \cos(2V_d \alpha')}{((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d)((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d)} \right) - \sigma_0^2 \frac{8V_d^4 (\cos(2V_d \alpha') - 1)}{(\gamma^2 + 2V_d)((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d)((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d)}. \quad (32)$$

Regardless of the observations, the predictive uncertainty (32) is periodic with respect to α' with the period of π/V_d . We can easily get the following corollaries.

Corollary C.2. For the test point at $\mathbf{x}' = \hat{x}$, i.e., $\alpha' = 0$, the mean of the derivative GP prediction is

$$\tilde{\mu}_{[\mathbf{X}, \mathbf{y}, \sigma]}^{(d)}(\mathbf{x}') = \frac{\sum_{w=0}^{2V_d-1} (-1)^w y_w \left(\frac{1}{2 \sin^2\left(\frac{(2w+1)\pi}{4V_d}\right)} + \frac{V_d(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d} \right)}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d}, \quad (33)$$

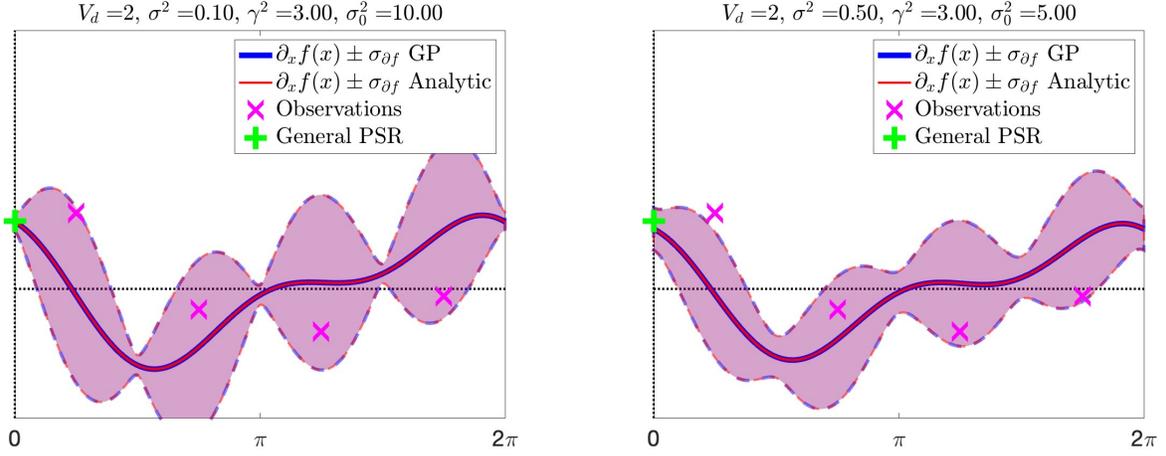


Figure 7. Numerical validation of Theorem C.1 under two parameter settings (see above each panel). Given the $2V_d$ equidistant observations (magenta crosses), the derivative GP prediction (blue curve) with uncertainty (blue shades) is compared to their analytic forms (31) and (32), i.e., the mean function (red curve) and the variance function (red shades), respectively. We observe that our theory perfectly matches the numerical computation. The green cross shows the prediction by the general PSR (12), which almost coincides with Bayesian PSR prediction when $\sigma^2/\sigma_0^2 = 0.01$ (left panel), while a significant difference is observed when $\sigma^2/\sigma_0^2 = 0.1$ (right panel).

Corollary C.3. For the test point at $\mathbf{x}' = \hat{\mathbf{x}} + \alpha' \mathbf{e}_d, \forall \alpha' = 0, \pi/V_d, 2\pi/V_d, \dots, (2V_d - 1)\pi/V_d$, the variance of the derivative GP prediction is

$$\tilde{s}_{[\mathbf{x}, \sigma]}^{(d)}(\mathbf{x}', \mathbf{x}') = \sigma^2 \left(\frac{V_d(V_d + 1)(2V_d + 1)}{3((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d)} - \frac{4V_d^3}{((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d)((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d)} \right). \quad (34)$$

Ignoring high order terms with respect to σ^2/σ_0^2 in Eqs.(33) and (34) gives Theorem 3.1. \square

Figure 7 shows numerical validation of Theorem C.1, where the derivative GP prediction (blue curve) with uncertainty (blue shades) is compared to their analytic forms, i.e., the mean function (31) (red curve) and the variance function (32) (red shades), respectively, under two settings of noise and kernel parameters. We observe that our theory perfectly matches the numerical computation. When $\sigma^2/\sigma_0^2 = 0.01$ (left panel), the regularization is small enough and the Bayesian PSR prediction (red curve) almost coincides with the general PSR prediction (green cross). On the other hand, when $\sigma^2/\sigma_0^2 = 0.1$ (right panel), the Bayesian PSR prediction (red) does not match the general PSR prediction (green cross), because of the regularization.

C.2. Mathematical Preparations

Before proving Theorem C.1, we give some mathematical identities on the trigonometric functions.

C.2.1. ROOT OF UNITY

For a natural number $N \in \{1, 2, \dots\}$, let us define a root of unity $\rho_N = e^{2\pi i/N}$ such that $\rho_N^N = 1$. Then, the following hold:

$$\sum_{n=0}^{N-1} \rho_N^{nk} = \frac{1 - \rho_N^{kN}}{1 - \rho_N^k} = 0 \quad \text{for} \quad k = 1, \dots, N-1, \quad (35)$$

$$\sum_{n=0}^{N-1} \rho_N^{(n+\phi)k} = \rho_N^{k\phi} \sum_{n=0}^{N-1} \rho_N^{nk} = \rho_N^{k\phi} \frac{1 - \rho_N^{kN}}{1 - \rho_N^k} = 0 \quad \text{for} \quad k = 1, \dots, N-1, \quad (36)$$

It also holds for even N that

$$\sum_{n=0}^{N-1} \rho_N^{(n+1/2)k+nN/2} = \rho_N^{k/2} \sum_{n=0}^{N-1} \rho_N^{n(k+N/2)} = \rho_N^{k/2} \frac{1 - \rho_N^{(k+N/2)N}}{1 - \rho_N^{(k+N/2)}} = 0 \quad \text{for} \quad k = 1, \dots, N/2 - 1. \quad (37)$$

C.2.2. PROPERTIES OF DIRICHLET KERNEL

The summation in the Dirichlet kernel can be analytically performed as

$$\begin{aligned}
 1 + 2 \sum_{n=1}^N \cos(nx) &= 1 + 2 \sum_{n=1}^N \frac{e^{inx} + e^{-inx}}{2} = \sum_{n=-N}^N e^{inx} \\
 &= e^{-iNx} \frac{1 - e^{i(2N+1)x}}{1 - e^{ix}} \\
 &= \frac{e^{-i(N+1/2)x} - e^{i(N+1/2)x}}{e^{-ix/2} - e^{ix/2}} \\
 &= \frac{\sin((N+1/2)x)}{\sin(x/2)}. \tag{38}
 \end{aligned}$$

Therefore, it also holds that

$$\begin{aligned}
 2 \sum_{n=1}^N n \sin(nx) &= - \sum_{v=1}^{V_d} \frac{\partial}{\partial x} (1/V_d + 2 \cos(nx)) \\
 &= - \frac{\partial}{\partial x} \left(1 + 2 \sum_{v=1}^{V_d} \cos(vx) \right) \\
 &= - \frac{(N+1/2) \cos((N+1/2)x) \sin(x/2) - \frac{1}{2} \sin((N+1/2)x) \cos(x/2)}{\sin^2(x/2)} \\
 &= - \frac{N \cos((N+1/2)x) \sin(x/2) - \frac{1}{2} \sin(Nx)}{\sin^2(x/2)} \\
 &= \frac{\sin(Nx)}{2 \sin^2(x/2)} - \frac{N \cos((N+1/2)x)}{\sin(x/2)}. \tag{39}
 \end{aligned}$$

C.3. Proof of Theorem C.1

For derivative predictions $\partial_d f(\mathbf{x}')$, $\partial_d f(\mathbf{x}'')$, the test kernels should be modified as Eqs.(6) and (7). For the VQE kernel (15), they are

$$\tilde{k}(\mathbf{x}, \mathbf{x}') = \partial_{x'_d} k(\mathbf{x}, \mathbf{x}') = \sigma_0^2 \left(\frac{2 \sum_{v=1}^{V_d} v \sin(v(x_d - x'_d))}{\gamma^2 + 2V_d} \right) \prod_{d' \neq d} \left(\frac{\gamma^2 + 2 \sum_{v=1}^{V_{d'}} \cos(v(x_{d'} - x'_{d'}))}{\gamma^2 + 2V_{d'}} \right), \tag{40}$$

$$\tilde{k}(\mathbf{x}', \mathbf{x}'') = \partial_{x'_d} \partial_{x''_d} k(\mathbf{x}', \mathbf{x}'') = \sigma_0^2 \left(\frac{2 \sum_{v=1}^{V_d} v^2 \cos(v(x'_d - x''_d))}{\gamma^2 + 2V_d} \right) \prod_{d' \neq d} \left(\frac{\gamma^2 + 2 \sum_{v=1}^{V_{d'}} \cos(v(x'_{d'} - x''_{d'}))}{\gamma^2 + 2V_{d'}} \right). \tag{41}$$

The training kernel matrix for $\{\mathbf{x}_w = \hat{\mathbf{x}} + \frac{2w+1}{2V_d} \pi \mathbf{e}_d\}_{w=0}^{2V_d-1}$ is Toeplitz as

$$\mathbf{K} = \sigma_0^2 \begin{pmatrix} \tau_0 & \tau_1 & \tau_2 & \cdots & \tau_{2V_d-2} & \tau_{2V_d-1} \\ \tau_1 & \tau_0 & \tau_1 & & & \\ \tau_2 & \tau_1 & \tau_0 & & & \vdots \\ \vdots & & & \ddots & & \\ \tau_{2V_d-2} & & & & \tau_0 & \tau_1 \\ \tau_{2V_d-1} & \cdots & & & \tau_1 & \tau_0 \end{pmatrix} \in \mathbb{R}^{2V_d \times 2V_d},$$

where

$$\tau_w = \frac{\gamma^2 + 2 \sum_{v=1}^{V_d} \cos\left(\frac{vw}{2V_d} 2\pi\right)}{\gamma^2 + 2V_d}. \tag{42}$$

For a test point at $\mathbf{x}' = \hat{\mathbf{x}} + \alpha' \mathbf{e}_d$, the test kernel components are

$$\tilde{\mathbf{k}}' = \sigma_0^2 \begin{pmatrix} \kappa_0 \\ \kappa_1 \\ \vdots \\ \kappa_{2V_d-1} \end{pmatrix},$$

$$\tilde{k}'' = \sigma_0^2,$$

where

$$\kappa_w = \frac{2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{2w+1}{2V_d} \pi - \alpha' \right) \right)}{\gamma^2 + 2V_d}. \quad (43)$$

The first identity (35) for the root of unity implies that

$$\sum_{v=0}^{2V_d-1} e^{vw \frac{2\pi i}{2V_d}} = 0 \quad \text{for} \quad w = 1, \dots, 2V_d - 1,$$

and therefore

$$\sum_{v=0}^{2V_d-1} \cos \left(vw \frac{2\pi}{2V_d} \right) = \begin{cases} 2V_d & \text{for } w = 0, 2V_d, \\ 0 & \text{for } w = 1, \dots, 2V_d - 1, \end{cases} \quad (44)$$

$$\sum_{v=0}^{2V_d-1} \sin \left(vw \frac{2\pi}{2V_d} \right) = 0 \quad \text{for } w = 0, \dots, 2V_d. \quad (45)$$

The second identity (36) for the root of unity gives

$$\sum_{v=0}^{2V_d-1} e^{(v+1/2)w \frac{2\pi i}{2V_d}} = \sum_{v=0}^{2V_d-1} e^{(2v+1)w \frac{\pi i}{2V_d}} = 0 \quad \text{for } w = 1, \dots, 2V_d - 1,$$

and therefore

$$\sum_{v=0}^{2V_d-1} \cos \left((2v+1)w \frac{\pi}{2V_d} \right) = \begin{cases} 2V_d & \text{for } w = 0, \\ -2V_d & \text{for } w = 2V_d, \\ 0 & \text{for } w = 1, \dots, 2V_d - 1, \end{cases} \quad (46)$$

$$\sum_{v=0}^{2V_d-1} \sin \left((2v+1)w \frac{\pi}{2V_d} \right) = 0 \quad \text{for } w = 0, \dots, 2V_d. \quad (47)$$

The third identity (37) for the root of unity gives

$$\sum_{v=0}^{2V_d-1} e^{((v+1/2)w+vV_d) \frac{2\pi i}{2V_d}} = \sum_{v=0}^{2V_d-1} e^{v\pi i} e^{(2v+1)w \frac{\pi i}{2V_d}} = \sum_{v=0}^{2V_d-1} (-1)^v e^{(2v+1)w \frac{\pi i}{2V_d}} = 0 \quad \text{for } w = 1, \dots, V_d - 1,$$

and therefore

$$\sum_{v=0}^{2V_d-1} (-1)^v \cos \left((2v+1)w \frac{\pi}{2V_d} \right) = 0 \quad \text{for } w = 0, \dots, V_d, \quad (48)$$

$$\sum_{v=0}^{2V_d-1} (-1)^v \sin \left((2v+1)w \frac{\pi}{2V_d} \right) = \begin{cases} 2V_d & \text{for } w = V_d, \\ 0 & \text{for } w = 0, \dots, V_d - 1. \end{cases} \quad (49)$$

From the symmetry of the trigonometric functions, it holds that

$$\sum_{v=1}^{V_d} \cos\left(vw \frac{2\pi}{2V_d}\right) = \begin{cases} -1 & \text{for } w = 1, 3, 5, \dots, 2V_d - 1, \\ 0 & \text{for } w = 2, 4, 6, \dots, 2V_d, \end{cases} \quad (50)$$

$$\sum_{v=1}^{V_d} \sin\left(vw \frac{2\pi}{2V_d}\right) = - \sum_{v=V_d+1}^{2V_d} \sin\left(vw \frac{2\pi}{2V_d}\right). \quad (51)$$

Note that the factor -1 in the odd w case in Eq. (50) is because the summand is -1 for $v = V_d$, while the summands for the other v are canceled each other.

By using Eq. (50), Eq. (42) can be written as

$$\tau_w = \frac{\gamma^2 + 2 \sum_{v=1}^{V_d} \cos\left(\frac{vw}{2V_d} 2\pi\right)}{\gamma^2 + 2V_d} = \begin{cases} 1 & \text{for } w = 0, \\ \frac{\gamma^2 - 2}{\gamma^2 + 2V_d} & \text{for } w = 1, 3, 5, \dots, 2V_d - 1, \\ \frac{\gamma^2}{\gamma^2 + 2V_d} & \text{for } w = 2, 4, 6, \dots, 2V_d - 2, \end{cases}$$

and therefore

$$\begin{aligned} \mathbf{K} &= \frac{\sigma_0^2}{\gamma^2 + 2V_d} (2V_d \mathbf{I}_{2V_d} + (\gamma^2 - 1) \mathbf{1}\mathbf{1}^\top + \mathbf{c}\mathbf{c}^\top) \\ &= \frac{\sigma_0^2}{\gamma^2 + 2V_d} \left(2V_d \mathbf{I}_{2V_d} + (\mathbf{1} \ \mathbf{c}) \begin{pmatrix} \gamma^2 - 1 & 0 \\ 0 & 1 \end{pmatrix} (\mathbf{1} \ \mathbf{c})^\top \right), \end{aligned} \quad (52)$$

where

$$\mathbf{c} = \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \\ \vdots \\ 1 \\ -1 \end{pmatrix} \in \mathbb{R}^{2V_d}.$$

With the training kernel expression (52), the matrix inversion lemma gives

$$\begin{aligned} (\mathbf{K} + \sigma^2 \mathbf{I}_{2V_d})^{-1} &= \frac{\gamma^2 + 2V_d}{\sigma_0^2} \left((\gamma^2 + 2V_d)(\sigma^2/\sigma_0^2 + 2V_d) \mathbf{I}_{2V_d} + (\mathbf{1} \ \mathbf{c}) \begin{pmatrix} \gamma^2 - 1 & 0 \\ 0 & 1 \end{pmatrix} (\mathbf{1} \ \mathbf{c})^\top \right)^{-1} \\ &= \frac{\gamma^2 + 2V_d}{\sigma_0^2} \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} \\ &\quad \left(\mathbf{I}_{2V_d} + \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} (\mathbf{1} \ \mathbf{c}) \begin{pmatrix} \gamma^2 - 1 & 0 \\ 0 & 1 \end{pmatrix} (\mathbf{1} \ \mathbf{c})^\top \right)^{-1} \\ &= \frac{\gamma^2 + 2V_d}{\sigma_0^2} \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} \\ &\quad \left\{ \mathbf{I}_{2V_d} - \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} (\mathbf{1} \ \mathbf{c}) \begin{pmatrix} \gamma^2 - 1 & 0 \\ 0 & 1 \end{pmatrix} \right. \\ &\quad \left. \left(\mathbf{I}_2 + (\mathbf{1} \ \mathbf{c})^\top \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} (\mathbf{1} \ \mathbf{c}) \begin{pmatrix} \gamma^2 - 1 & 0 \\ 0 & 1 \end{pmatrix} \right)^{-1} (\mathbf{1} \ \mathbf{c})^\top \right\} \\ &= \frac{\gamma^2 + 2V_d}{\sigma_0^2} \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} \\ &\quad \left\{ \mathbf{I}_{2V_d} - \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} (\mathbf{1} \ \mathbf{c}) \begin{pmatrix} \gamma^2 - 1 & 0 \\ 0 & 1 \end{pmatrix} \right. \end{aligned}$$

$$\begin{aligned}
 & \left(\mathbf{I}_2 + \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} \begin{pmatrix} 2V_d(\gamma^2 - 1) & 0 \\ 0 & 2V_d \end{pmatrix} \right)^{-1} (\mathbf{1} \quad \mathbf{c})^\top \Big\} \\
 &= \frac{\gamma^2 + 2V_d}{\sigma_0^2} \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} \\
 & \quad \left\{ \mathbf{I}_{2V_d} - (\mathbf{1} \quad \mathbf{c}) \begin{pmatrix} \gamma^2 - 1 & 0 \\ 0 & 1 \end{pmatrix} \right. \\
 & \quad \left. \left(\begin{pmatrix} (\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d\gamma^2 & 0 \\ 0 & (\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d \end{pmatrix} \right)^{-1} (\mathbf{1} \quad \mathbf{c})^\top \right\} \\
 &= \frac{1}{\sigma_0^2} a(\mathbf{I}_{2V_d} + b\mathbf{1}\mathbf{1}^\top + c\mathbf{c}\mathbf{c}^\top),
 \end{aligned}$$

where

$$\begin{aligned}
 a &= \frac{\gamma^2 + 2V_d}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d}, \\
 b &= -\frac{\gamma^2 - 1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d\gamma^2}, \\
 c &= -\frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d}.
 \end{aligned}$$

For the test kernels

$$\begin{aligned}
 \tilde{\mathbf{k}}' &= \sigma_0^2 \begin{pmatrix} \kappa_0 \\ \kappa_1 \\ \vdots \\ \kappa_{2V_d-1} \end{pmatrix}, \\
 \tilde{\mathbf{k}}'' &= \sigma_0^2 \left(\frac{2 \sum_{v=1}^{V_d} v^2}{\gamma^2 + 2V_d} \right) = \frac{\sigma_0^2 V_d (V_d + 1) (2V_d + 1)}{3(\gamma^2 + 2V_d)},
 \end{aligned}$$

with

$$\kappa_w = \frac{2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\gamma^2 + 2V_d},$$

(53)

we have

$$\begin{aligned}
 \|\tilde{\mathbf{k}}'\|^2 &= \sigma_0^4 \sum_{w=0}^{2V_d-1} \left(\frac{2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\gamma^2 + 2V_d} \right)^2 \\
 &= \frac{\sigma_0^4}{(\gamma^2 + 2V_d)^2} \sum_{w=0}^{2V_d-1} \left\{ 4 \sum_{v=1}^{V_d} \sum_{v'=1}^{V_d} v v' \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) \sin \left(v' \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) \right\} \\
 &= \frac{\sigma_0^4}{(\gamma^2 + 2V_d)^2} \sum_{w=0}^{2V_d-1} \left\{ 2 \sum_{v=1}^{V_d} \sum_{v'=1}^{V_d} v v' \left(\cos \left((v-v') \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) - \cos \left((v+v') \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) \right) \right\} \\
 &= \frac{\sigma_0^4}{(\gamma^2 + 2V_d)^2} \left\{ 2 \sum_{v=1}^{V_d} \sum_{v'=1}^{V_d} v v' \sum_{w=0}^{2V_d-1} \left(\cos \left((v-v') \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) - \cos \left((v+v') \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) \right) \right\} \\
 &= \frac{\sigma_0^4}{(\gamma^2 + 2V_d)^2} \left\{ \right.
 \end{aligned}$$

$$\begin{aligned}
 & 2 \sum_{v=1}^{V_d} \sum_{v'=1}^{V_d} v v' \sum_{w=0}^{2V_d-1} \left(\cos \frac{(2w+1)(v-v')\pi}{2V_d} \cos((v-v')\alpha') + \sin \frac{(2w+1)(v-v')\pi}{2V_d} \sin((v-v')\alpha') \right. \\
 & \left. - \cos \frac{(2w+1)(v+v')\pi}{2V_d} \cos((v+v')\alpha') - \sin \frac{(2w+1)(v+v')\pi}{2V_d} \sin((v+v')\alpha') \right) \Bigg\} \quad (54)
 \end{aligned}$$

$$\begin{aligned}
 & = \frac{\sigma_0^4}{(\gamma^2 + 2V_d)^2} \left\{ 2 \sum_{v=1}^{V_d} \sum_{v'=1}^{V_d} v v' \sum_{w=0}^{2V_d-1} \left(\cos \frac{(2w+1)(v-v')\pi}{2V_d} \cos((v-v')\alpha') \right. \right. \\
 & \left. \left. - \cos \frac{(2w+1)(v+v')\pi}{2V_d} \cos((v+v')\alpha') \right) \right\} \quad (55)
 \end{aligned}$$

$$\begin{aligned}
 & = \frac{\sigma_0^4}{(\gamma^2 + 2V_d)^2} 2(2V_d) \left(\left(\sum_{v=1}^{V_d} v^2 \right) + V_d^2 \cos(2V_d\alpha') \right) \quad (56)
 \end{aligned}$$

$$\begin{aligned}
 & = \frac{\sigma_0^4}{(\gamma^2 + 2V_d)^2} 2(2V_d) \left(\frac{V_d(V_d+1)(2V_d+1)}{6} + V_d^2 \cos(2V_d\alpha') \right)
 \end{aligned}$$

$$\begin{aligned}
 & = \sigma_0^4 \frac{4V_d^2}{(\gamma^2 + 2V_d)^2} \left(\frac{(V_d+1)(2V_d+1)}{6} + V_d \cos(2V_d\alpha') \right).
 \end{aligned}$$

Here we used Eqs.(46) and (47) to obtain Eqs.(55) and (56) from Eq. (54).

We also have

$$\begin{aligned}
 \|\tilde{\mathbf{k}}'\|_1 & = \tilde{\mathbf{k}}'^{\top} \mathbf{1}_{2V_d} = \sigma_0^2 \sum_{w=0}^{2V_d-1} \frac{2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\gamma^2 + 2V_d} \\
 & = \sigma_0^2 \frac{2 \sum_{v=1}^{V_d} v \sum_{w=0}^{2V_d-1} \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\gamma^2 + 2V_d} \\
 & = \sigma_0^2 \frac{2 \sum_{v=1}^{V_d} v \sum_{w=0}^{2V_d-1} \left(\sin \frac{(2w+1)v\pi}{2V_d} \cos v\alpha' - \cos \frac{(2w+1)v\pi}{2V_d} \sin v\alpha' \right)}{\gamma^2 + 2V_d} \\
 & = 0,
 \end{aligned}$$

and

$$\begin{aligned}
 \tilde{\mathbf{k}}'^{\top} \mathbf{c} & = \sigma_0^2 \sum_{w=0}^{2V_d-1} (-1)^w \frac{2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\gamma^2 + 2V_d} \\
 & = \sigma_0^2 \frac{2 \sum_{v=1}^{V_d} v \sum_{w=0}^{2V_d-1} (-1)^w \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\gamma^2 + 2V_d} \\
 & = \sigma_0^2 \frac{2 \sum_{v=1}^{V_d} v \sum_{w=0}^{2V_d-1} (-1)^w \left(\sin \frac{(2w+1)v\pi}{2V_d} \cos v\alpha' - \cos \frac{(2w+1)v\pi}{2V_d} \sin v\alpha' \right)}{\gamma^2 + 2V_d} \\
 & = \sigma_0^2 \frac{2V_d 2V_d \cos V_d\alpha'}{\gamma^2 + 2V_d} \\
 & = \sigma_0^2 \frac{4V_d^2 \cos V_d\alpha'}{\gamma^2 + 2V_d}.
 \end{aligned}$$

Here, we used Eqs.(48) and (49) in the second last equation. Therefore, the mean of the derivative is

$$\begin{aligned}
 \tilde{\mu}_{[\mathbf{X}, \mathbf{y}, \sigma]}^{(d)}(\mathbf{x}') & = \tilde{\mathbf{k}}'^{\top} (\mathbf{K} + \sigma^2 \mathbf{I}_{2V_d})^{-1} \mathbf{y} \\
 & = \tilde{\mathbf{k}}'^{\top} \frac{a}{\sigma_0^2} (\mathbf{I}_{2V_d} + b \mathbf{1}_{2V_d} \mathbf{1}_{2V_d}^{\top} + \mathbf{c} \mathbf{c}^{\top}) \mathbf{y}
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{a}{\sigma_0^2} \left(\tilde{\mathbf{k}}'^\top \mathbf{y} + b \tilde{\mathbf{k}}'^\top \mathbf{1}_{2V_d} \mathbf{1}_{2V_d}^\top \mathbf{y} + c \tilde{\mathbf{k}}'^\top \mathbf{c} \mathbf{c}^\top \mathbf{y} \right) \\
 &= a \left(\sum_{w=0}^{2V_d-1} y_w \frac{2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\gamma^2 + 2V_d} + c \frac{4V_d^2 \cos V_d \alpha'}{\gamma^2 + 2V_d} \sum_{w=0}^{2V_d-1} (-1)^w y_w \right). \\
 &= a \left(\sum_{w=0}^{2V_d-1} y_w \left(\frac{2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\gamma^2 + 2V_d} + c \frac{4V_d^2 (-1)^w \cos V_d \alpha'}{\gamma^2 + 2V_d} \right) \right) \\
 &= \frac{a}{\gamma^2 + 2V_d} \left(\sum_{w=0}^{2V_d-1} y_w \left(\left\{ 2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) \right\} + 4cV_d^2 (-1)^w \cos V_d \alpha' \right) \right) \\
 &= \frac{\sum_{w=0}^{2V_d-1} y_w \left(\left\{ 2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) \right\} - \frac{4V_d^2 (-1)^w \cos V_d \alpha'}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d} \right)}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d}. \tag{57}
 \end{aligned}$$

Eq. (39) implies that, for $w = 0, 1, \dots, 2V_d - 1$, it holds that

$$\begin{aligned}
 2 \sum_{v=1}^{V_d} v \sin \left(v \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right) &= \frac{\sin(V_d \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right))}{2 \sin^2 \left(\left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) / 2 \right)} - \frac{V_d \cos \left((V_d + 1/2) \left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) \right)}{\sin \left(\left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) / 2 \right)} \\
 &= \frac{\sin \left(\frac{(2w+1)\pi}{2} - V_d \alpha' \right)}{2 \sin^2 \left(\left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) / 2 \right)} - \frac{V_d \cos \left(\frac{(2w+1)\pi}{2} + \frac{(2w+1)\pi}{4V_d} - (V_d + 1/2) \alpha' \right)}{\sin \left(\left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) / 2 \right)} \\
 &= \frac{\sin \left((-1)^w \frac{\pi}{2} - V_d \alpha' \right)}{2 \sin^2 \left(\left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) / 2 \right)} - \frac{V_d \cos \left((-1)^w \frac{\pi}{2} + \frac{(2w+1)\pi}{4V_d} - (V_d + 1/2) \alpha' \right)}{\sin \left(\left(\frac{(2w+1)\pi}{2V_d} - \alpha' \right) / 2 \right)} \\
 &= (-1)^w \left(\frac{\cos(V_d \alpha')}{2 \sin^2 \left(\frac{(2w+1)\pi}{4V_d} - \alpha' / 2 \right)} + \frac{V_d \sin \left(\frac{(2w+1)\pi}{4V_d} - (V_d + 1/2) \alpha' \right)}{\sin \left(\frac{(2w+1)\pi}{4V_d} - \alpha' / 2 \right)} \right). \tag{58}
 \end{aligned}$$

Substituting Eq. (58) into Eq. (57) gives Eq. (31).

The posterior variance can be computed as

$$\begin{aligned}
 \hat{s}_{[\mathbf{x}, \sigma]}^{(d)}(\mathbf{x}', \mathbf{x}') &= \tilde{\mathbf{k}}'' - \tilde{\mathbf{k}}'^\top (\mathbf{K} + \sigma^2 \mathbf{I}_{2V_d})^{-1} \tilde{\mathbf{k}}' \\
 &= \tilde{\mathbf{k}}'' - \tilde{\mathbf{k}}'^\top \frac{1}{\sigma_0^2} a (\mathbf{I}_{2V_d} + b \mathbf{1}_{2V_d} \mathbf{1}_{2V_d}^\top + c c c^\top) \tilde{\mathbf{k}}' \\
 &= \tilde{\mathbf{k}}'' - \frac{1}{\sigma_0^2} a \left(\|\tilde{\mathbf{k}}'\|^2 + b (\tilde{\mathbf{k}}'^\top \mathbf{1}_{2V_d})^2 + c (\tilde{\mathbf{k}}'^\top \mathbf{c})^2 \right) \\
 &= \frac{\sigma_0^2 V_d (V_d + 1) (2V_d + 1)}{3(\gamma^2 + 2V_d)} \\
 &\quad - \frac{1}{\sigma_0^2} a \left\{ \sigma_0^4 \frac{4V_d^2}{(\gamma^2 + 2V_d)^2} \left(\frac{(V_d + 1)(2V_d + 1)}{6} + V_d \cos(2V_d \alpha') \right) + c \sigma_0^4 \left(\frac{4V_d^2 \cos V_d \alpha'}{\gamma^2 + 2V_d} \right)^2 \right\} \\
 &= \frac{\sigma_0^2 V_d (V_d + 1) (2V_d + 1)}{3(\gamma^2 + 2V_d)} - \sigma_0^2 a \frac{4V_d^2}{(\gamma^2 + 2V_d)^2} \frac{(V_d + 1)(2V_d + 1)}{6} \\
 &\quad - \sigma_0^2 a \left\{ \frac{4V_d^3 \cos(2V_d \alpha')}{(\gamma^2 + 2V_d)^2} + c \frac{16V_d^4 \cos^2 V_d \alpha'}{(\gamma^2 + 2V_d)^2} \right\} \\
 &= \frac{\sigma_0^2 V_d (V_d + 1) (2V_d + 1)}{3(\gamma^2 + 2V_d)} - \sigma_0^2 \frac{\gamma^2 + 2V_d}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} \frac{4V_d^2}{(\gamma^2 + 2V_d)^2} \frac{(V_d + 1)(2V_d + 1)}{6}
 \end{aligned}$$

$$\begin{aligned}
 & -\sigma_0^2 \frac{\gamma^2 + 2V_d}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d} \left\{ \frac{4V_d^3 \cos(2V_d\alpha')}{(\gamma^2 + 2V_d)^2} - \frac{1}{(\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d} \frac{8V_d^4(1 + \cos 2V_d\alpha')}{(\gamma^2 + 2V_d)^2} \right\} \\
 & = \sigma^2 \frac{V_d(V_d + 1)(2V_d + 1)}{3((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d)} - \sigma^2 \frac{4V_d^3 \cos(2V_d\alpha')}{((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d)((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d)} \\
 & \quad - \sigma_0^2 \frac{8V_d^4(\cos(2V_d\alpha') - 1)}{(\gamma^2 + 2V_d)((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 2V_d)((\gamma^2 + 2V_d)\sigma^2/\sigma_0^2 + 4V_d)}, \tag{59}
 \end{aligned}$$

which gives Eq. (32). \square

C.4. Proof of Theorem 3.2

In the first order case with $V_d = 1, \forall d = 1, \dots, D$, the test VQE kernels for predicting derivatives $\partial_d f(\mathbf{x}')$, $\partial_d f(\mathbf{x}'')$ are

$$\begin{aligned}
 \tilde{k}(\mathbf{x}, \mathbf{x}') &= \frac{\partial}{\partial x'_d} k(\mathbf{x}, \mathbf{x}') = \sigma_0^2 \left(\frac{2 \sin(x_d - x'_d)}{\gamma^2 + 2} \right) \prod_{d' \neq d} \left(\frac{\gamma^2 + 2 \cos(x_{d'} - x'_{d'})}{\gamma^2 + 2} \right), \\
 \tilde{k}(\mathbf{x}', \mathbf{x}'') &= \frac{\partial^2}{\partial x'_d \partial x''_d} k(\mathbf{x}', \mathbf{x}'') = \sigma_0^2 \left(\frac{2 \cos(x'_d - x''_d)}{\gamma^2 + 2} \right) \prod_{d' \neq d} \left(\frac{\gamma^2 + 2 \cos(x'_{d'} - x''_{d'})}{\gamma^2 + 2} \right).
 \end{aligned}$$

Then, the kernels with the two training points $\mathbf{X} = (\mathbf{x}' - \alpha \mathbf{e}_d, \mathbf{x}' + \alpha \mathbf{e}_d)$ and the one test point \mathbf{x}' are

$$\mathbf{K} = \sigma_0^2 \begin{pmatrix} 1 & \frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2} \\ \frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2} & 1 \end{pmatrix}, \quad \tilde{\mathbf{k}}' = \frac{2\sigma_0^2 \sin \alpha}{\gamma^2 + 2} \begin{pmatrix} -1 \\ 1 \end{pmatrix}, \quad \tilde{k}'' = \frac{2\sigma_0^2}{\gamma^2 + 2}.$$

With these kernels, the posterior mean is

$$\begin{aligned}
 \tilde{\mu}_{[\mathbf{X}, \mathbf{y}, \sigma]}^{(d)}(\mathbf{x}') &= \tilde{\mathbf{k}}'^{\top} (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y} \\
 &= \frac{2 \sin \alpha}{\gamma^2 + 2} \begin{pmatrix} -1 & 1 \end{pmatrix} \begin{pmatrix} 1 + \sigma^2/\sigma_0^2 & \frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2} \\ \frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2} & 1 + \sigma^2/\sigma_0^2 \end{pmatrix}^{-1} \mathbf{y} \\
 &= \frac{2 \sin \alpha}{\gamma^2 + 2} \begin{pmatrix} -1 & 1 \end{pmatrix} \frac{1}{(1 + \sigma^2/\sigma_0^2)^2 - \left(\frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2}\right)^2} \begin{pmatrix} 1 + \sigma^2/\sigma_0^2 & -\frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2} \\ -\frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2} & 1 + \sigma^2/\sigma_0^2 \end{pmatrix} \mathbf{y} \\
 &= \frac{2 \sin \alpha}{\gamma^2 + 2} \frac{1}{(1 + \sigma^2/\sigma_0^2) - \left(\frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2}\right)} \begin{pmatrix} -1 & 1 \end{pmatrix} \mathbf{y} \\
 &= 2 \sin \alpha \frac{y_2 - y_1}{(1 + \sigma^2/\sigma_0^2)(\gamma^2 + 2) - (\gamma^2 + 2 \cos 2\alpha)} \\
 &= \frac{(y_2 - y_1) \sin \alpha}{(\gamma^2/2 + 1)\sigma^2/\sigma_0^2 + 2 \sin^2 \alpha}.
 \end{aligned}$$

The posterior variance is

$$\begin{aligned}
 \tilde{s}_{[\mathbf{X}, \sigma]}^{(d)}(\mathbf{x}', \mathbf{x}') &= \tilde{k}'' - \tilde{\mathbf{k}}'^{\top} (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \tilde{\mathbf{k}}' \\
 &= \frac{2\sigma_0^2}{\gamma^2 + 2} - \frac{4\sigma_0^2 \sin^2 \alpha}{(\gamma^2 + 2)^2} \begin{pmatrix} -1 & 1 \end{pmatrix} \begin{pmatrix} 1 + \sigma^2/\sigma_0^2 & \frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2} \\ \frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2} & 1 + \sigma^2/\sigma_0^2 \end{pmatrix}^{-1} \begin{pmatrix} -1 \\ 1 \end{pmatrix} \\
 &= \frac{2\sigma_0^2}{\gamma^2 + 2} - \frac{4\sigma_0^2 \sin^2 \alpha}{(\gamma^2 + 2)^2} \frac{1}{(1 + \sigma^2/\sigma_0^2) - \left(\frac{\gamma^2 + 2 \cos 2\alpha}{\gamma^2 + 2}\right)} \begin{pmatrix} -1 & 1 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \end{pmatrix} \\
 &= \frac{2\sigma_0^2}{\gamma^2 + 2} \left(1 - \frac{4 \sin^2 \alpha}{(\gamma^2 + 2)\sigma^2/\sigma_0^2 + 2 - 2 \cos 2\alpha} \right)
 \end{aligned}$$

1045
$$= \frac{2\sigma_0^2}{\gamma^2 + 2} \left(\frac{(\gamma^2 + 2)\sigma^2/\sigma_0^2}{(\gamma^2 + 2)\sigma^2/\sigma_0^2 + 4\sin^2 \alpha} \right)$$

1046

1047

1048
$$= \frac{\sigma^2}{(\gamma^2/2 + 1)\sigma^2/\sigma_0^2 + 2\sin^2 \alpha}.$$

1049

1050 Thus we obtained Eqs.(19) and (20).

□

1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099

D. Algorithm Details

D.1. GradCoRe Pseudo-Code

Algorithm 1 describes SGD-GradCoRe in detail. SGD-GradCoRe uses the GradCoRe measurement selection subroutine described in Algorithm 2, which selects measurement points and respective minimum required number of shots to estimate the quantum circuit parameter derivative required for the SGD.

Algorithm 1 (SGD-GradCoRe) Improved SGD algorithm using a VQE-derivative kernel GP with the GradCoRe measurement selection subroutine, as described in Algorithm 2. The algorithm finds the minimum number of shots required to estimate the gradient wrt. parameter configurations $\hat{\mathbf{x}}$ of the quantum circuit to optimize with SGD. The optimization stops when the total number of measurement shots reaches the maximum number of observation shots allowed, i.e., $N_{\text{tot-shots}}$. To avoid cluttering notation, the algorithm is restricted to the case where $V_d = 1$. Generalization to an arbitrary V_d is straightforward.

Input :

- $\hat{\mathbf{x}}^0$: initial starting point (best point)

Parameters :

- $V_d = 1$
- D : number of parameters to optimize, i.e., $\hat{\mathbf{x}} \in \mathbb{R}^D$.
- $N_{\text{tot-shots}}$: Total # of shots, i.e., maximum allowed quantum computing budget.
- σ_1^{*2} : measurement variance using a single shot.
- κ_0 : Initial GradCoRe threshold at step $t = 0$.
- T_{initial} : Number of steps in beginning to use initial GradCoRe threshold κ_0 .
- c_0 : smallest allowed GradCoRe threshold
- c_1 : GradCoRe threshold scaling parameter

Output :

- $\hat{\mathbf{x}}^*$: optimal choice of parameters for the quantum circuit.

```

1131 1  $n \leftarrow 0$  /* initialize consumed shot budget */
1132 2  $t \leftarrow 0$  /* initialize optimization step */
1133 3  $\kappa^0 \leftarrow \mathbf{1}_D \kappa_0$  /* initial  $\kappa_0$  to use for  $T_{\text{initial}}$  steps */
1134 4  $\mathbf{X}^0, \mathbf{y}^0, \boldsymbol{\sigma}^0 \leftarrow (), (), ()$  /* initialize empty Gaussian process */
1135 5 while  $n < N_{\text{tot-shots}}$  do
1136 6     /* choose measurement points & number of shots s.t.  $\hat{\mathbf{x}}^t$  is in the GradCoRe of  $\boldsymbol{\kappa}^t$  */
1137 7      $\tilde{\mathbf{X}}, \tilde{\nu} \leftarrow \text{gradcore\_measurements}(\mathbf{X}^t, \mathbf{y}^t, \boldsymbol{\sigma}^t, \hat{\mathbf{x}}^t, \boldsymbol{\kappa}^t)$  /* (Algorithm 2) */
1138 8     for  $i \in \{1, \dots, |\tilde{\mathbf{X}}|\}$  do
1139 9          $\check{y}_i \leftarrow \text{quantum\_circuit}(\text{parameters}=\tilde{\mathbf{X}}_i, \text{shots}=\tilde{\nu}_i)$  /* measure chosen points */
1140 10         $\check{\sigma}_i \leftarrow \frac{\sigma_1^{*2}}{\tilde{\nu}_i}$ 
1141 11     end
1142 12      $\check{\mathbf{y}}, \check{\boldsymbol{\sigma}} \leftarrow (\check{y}_1, \dots, \check{y}_{|\tilde{\mathbf{X}}|}), (\check{\sigma}_1, \dots, \check{\sigma}_{|\tilde{\mathbf{X}}|})$  /* concatenate observed values & noise */
1143 13      $\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \boldsymbol{\sigma}^{t+1} \leftarrow (\mathbf{X}^t, \tilde{\mathbf{X}}), (\mathbf{y}^t, \check{\mathbf{y}}), (\boldsymbol{\sigma}^t, \check{\boldsymbol{\sigma}})$  /* add new observations to Gaussian process */
1144 14      $\hat{\mathbf{x}}^{t+1} \leftarrow \hat{\mathbf{x}}^t - \rho \tilde{\mu}_{[\mathbf{X}^{t+1}, \boldsymbol{\sigma}^{t+1}, \mathbf{y}^{t+1}]}(\hat{\mathbf{x}}^t)$  /* SGD (or variant) step using GP derivative */
1145 15     if  $t \geq T_{\text{initial}}$  then
1146 16          $\boldsymbol{\kappa}^{t+1} \leftarrow \mathbf{1}_D \max \left[ c_0, \frac{c_1}{D} \sum_{d=1}^D \left( \tilde{\mu}_{[\mathbf{X}^{t+1}, \boldsymbol{\sigma}^{t+1}, \mathbf{y}^{t+1}]}^{(d)}(\hat{\mathbf{x}}^t) \right)^2 \right]$  /* adapt GradCoRe threshold */
1147 17     end
1148 18      $t \leftarrow t + 1$  /* update the step */
1149 19      $n \leftarrow n + \sum_d \tilde{\nu}_d$  /* update the consumed shot budget */
1150 20 end
1151 21 return  $\hat{\mathbf{x}}^*$ 

```

Algorithm 2 (GradCoRe measurement selection subroutine) Select the points to measure and respective minimum number of required shots such that when updating the GP with these new measurements, the GP’s derivative uncertainty at the current best point is smaller than the threshold κ , i.e., the current point is within the GradCoRe.

Input :

- $\mathbf{X}, \mathbf{y}, \sigma$: Gaussian process at current step
- $\hat{\mathbf{x}}$: current best point
- $\kappa = (\kappa_1^2, \dots, \kappa_D^2)$: GradCoRe thresholds at current step

Parameters :

- $V_d = 1$
- σ_1^{*2} : measurement variance using a single shot.
- $\hat{\alpha}$: shift from best point at the previous step (default to $\hat{\alpha} = \frac{\pi}{2}$)

Output :

- $\check{\mathbf{X}}$: points which should be measured and added to the GP to compute the derivative.
- $\check{\nu}$: number of shots for the measured points.

1 **begin**

2 **for** $d \in \{1, \dots, D\}$ **do**

3 $\check{\mathbf{X}}_d \leftarrow (\hat{\mathbf{x}} - \hat{\alpha} \cdot \mathbf{e}_d, \hat{\mathbf{x}} + \hat{\alpha} \cdot \mathbf{e}_d)$ /* choose points to measure along d */

4 $\check{\sigma}_\pm \leftarrow \kappa_d$ /* initialize measurement noise to minimum (most expensive, $\kappa_d \ll \sigma_1^{*2}$) */

5 **for** $\check{\sigma} \in [\sigma_1^*, \kappa_d]$ **do**

6 /* create temporary GP copies, add points with $\check{\sigma}$ observation noise */

7 $\mathbf{X}', \mathbf{y}', \sigma' \leftarrow (\mathbf{X}, \check{\mathbf{X}}_d), (\mathbf{y}, 0, 0), (\sigma, \check{\sigma}, \check{\sigma})$

8 /* find largest observation noise for which $\hat{\mathbf{x}}$ is in the GradCoRe */

9 **if** $(\tilde{s}_{[\mathbf{X}', \sigma']}(\hat{\mathbf{x}}, \hat{\mathbf{x}}) \leq \kappa_d^2) \wedge (\check{\sigma}_\pm > \check{\sigma})$ **then**

10 | $\check{\sigma}_\pm \leftarrow \check{\sigma}$

11 **end**

12 **end**

13 $\check{\nu}_d \leftarrow \left(\frac{\sigma_1^{*2}}{\check{\sigma}_\pm}, \frac{\sigma_1^{*2}}{\check{\sigma}_\pm} \right)$ /* compute shots from variance through single shot variance σ_1^{*2} */

14 **end**

15 $\check{\mathbf{X}} \leftarrow (\check{\mathbf{X}}_1, \dots, \check{\mathbf{X}}_D)$ /* concatenate points to measure */

16 $\check{\nu} \leftarrow (\check{\nu}_1, \dots, \check{\nu}_D)$ /* concatenate shots to measure per point */

17 **return** $\check{\mathbf{X}}, \check{\nu}^{t+1}$

18 **end**

D.2. Parameter Setting

Every algorithm used in our benchmarking analysis has several hyperparameters to be set. For transparency and to allow the reproduction of our experiments, we detail the choice of parameters for EMICoRe, SubsCoRe and GradCoRe in Table 1. The SGLBO results were obtained using the original code from Tamiya & Yamasaki (2022) and we used the default setting from the original paper. For NFT, Bayes-NFT and Bayes-SGD runs, we used the default parameters specified in Table 2. For algorithmic efficiency, all Bayesian-SMO methods use the inducer option introduced in Nicoli et al. (2023a), retaining only the last $R \cdot 2V_d \cdot D - 1 = 399$ measured points once more than $R \cdot 2V_d \cdot D - 1 + D = 439$ points were stored in the GP, where we chose $R = 5$. Since the discarded points are replaced with a single point predicted from them, the number of the training points for the GP is kept constant at $R \cdot 2V_d \cdot D = 400$. On the other hand, Bayesian-SGD methods measure (at most, in the SGD-GradCoRe case) $2V_d D = 80$ points per SGD step, and we retain $R \cdot 2V_d \cdot D = 400$ points after more than $(R + 1) \cdot 2V_d \cdot D = 480$ points are measured. Unlike the Bayesian-SMO methods, we do not add additional inducer based on the prediction from the discarded points, and therefore the number of the training points for the GP is kept constant at $R \cdot 2V_d \cdot D = 400$.

⁵a.k.a., “readout” in Nicoli et al. (2023a).

⁶All hyperparameters not specified in the table are set to the default in Nicoli et al. (2023a).

Table 1. Algorithm specific parameter choice for EMICoRe, SubsCoRe and GradCoRe for all experiments (unless specified otherwise).

	Algorithmic specific parameters	
--acq-params	EMICoRe params	as in Nicoli et al. (2023a)
func	ei	Base acq. func. type
optim	emicore	Optimizer type
pairsize (J_{SG})	20	# of candidate points
gridsize (J_{OG})	100	# of evaluation points
corethresh-strategy	grad	Gradient strategy for κ
pnorm	2	Order of gradient norm
corethresh (κ)	256	CoRe threshold κ
corethresh_width (T_{initial})	40	# initial steps with fixed κ
coremin_scale (C_0)	2048	Coefficient C_0 for updating κ
corethresh_scale (C_1)	1.0	Coefficient C_1 for updating κ
stabilize_interval	41	Stabilization interval in SMO steps
samplesize (N_{MC})	100	# of MC samples
smo-steps (T_{NFT})	0	# of initial NFT steps
smo-axis	True	Sequential direction choice
--acq-params	SubsCoRe params	as in Anders et al. (2024)
optim	subscore ⁵	Optimizer type
readout-strategy	center	Alg type SubsCoRe
corethresh-strategy	grad	Gradient strategy for κ
pnorm	2	Order of gradient norm
corethresh (κ)	256	Initial N_{shots} for CoRe
corethresh_width (T_{initial})	40	# initial steps with fixed κ
coremin_scale (C_0)	2048	Coefficient C_0 for updating κ
corethresh_scale (C_1)	1.0	Coefficient C_1 for updating κ
stabilize_interval	41	Stabilization interval in SMO steps
coremetric	readout	Metric to set CoRe
--acq-params	GradCoRe params	this paper ⁶
optim	gradcore	Optimizer type
corethresh-strategy	grad	Gradient strategy for κ
pnorm	2	Order of gradient norm
corethresh (κ)	256	Initial N_{shots} for CoRe
corethresh_width (T_{initial})	40	# initial steps with fixed κ
coremin_scale (C_0)	2048	Coefficient C_0 for updating κ
corethresh_scale (C_1)	1.2	Coefficient C_1 for updating κ
coremetric	readout	Metric to set CoRe
lr	0.05	learning rate for SGD
gdoptim	adam	Optimizer for SGD

Table 2. Default choice of circuit parameters and kernel hyperparameters for all experiments (unless specified otherwise).

	Default params	
--n-qbits	5	# of qubits
--n-layers	3	# of circuit layers
--circuit	esu2	Circuit name
--pbc	False	Open Boundary Conditions
--n-iter	1*10**7	# max number of readouts
--kernel	vqe	Name of the kernel

--kernel-params	Bayes-NFT	EmiCoRe	SubsCoRe	GradCore	Bayes-SGD
gamma	3	3	3	3	1
sigma_0	10	10	10	10	10

E. Experimental Details

As discussed in the main text, our experiments focus on the same experimental setup as in [Nicoli et al. \(2023a\)](#) and [Anders et al. \(2024\)](#). Specifically, starting from the quantum Heisenberg Hamiltonian, we reduce it to the special case of the Ising Hamiltonian at the critical point by choosing the suitable couplings, namely

$$\text{Ising Hamiltonian at criticality: } J = (-1.0, 0.0, 0.0); h = (0.0, 0.0, -1.0).$$

It is important to note that due to the finite size of the system at hand, this choice of parameters does not imply criticality but already represents a challenging setup, as discussed in Sec. I.2 in [Nicoli et al. \(2023a\)](#). We stop the optimization when a maximum number of cumulative shots (total measurement budget on the quantum computer) is reached; unless specified otherwise, we set this cutoff to $N_{\text{shots}}^{\text{max}} = 1 \cdot 10^7$.

Our implementation of GradCoRe can be found in the supplementary zip file and will be made available on Github upon acceptance. In our experiments, the kernel parameters σ_0 and γ are fixed to the values in Table 2. Furthermore, NFT, Bayes-NFT, Bayes-SGD, SubsCoRe and GradCoRe require fixed shifts for the points to measure at each iteration. In our experiments, we always used $\alpha = \frac{2\pi}{3}$ for SMO based methods (as this makes the uncertainty uniform in the 1D-subspace, as discussed in [Anders et al. \(2024\)](#)), and $\alpha = \frac{\pi}{2}$ for SGD based methods (as this minimizes the uncertainty in the noisy case, as discussed in Section 3), unless explicitly stated otherwise.

Each experiment shown in the paper was repeated 50 times (trials) with differently seeded starting points. We aggregated the statistics from these independent trials and presented them in our plots. We used the same starting point for every algorithm in each trial to ensure a fair comparison between all approaches. Note that SGD-based methods do not require measurements at the starting point, but SMO-based methods do. Therefore, each starting point is further paired with a fixed initial measurement.

All experiments were conducted on Intel Xeon Silver 4316 @ 2.30GHz CPUs.

F. Detailed behavior of GradCoRe

Figure 8 shows the behavior of the GradCoRe threshold $\kappa(t)$ (left), and the number $\nu(t)$ of measurement shots (right) that GradCoRe used in each SGD iteration.

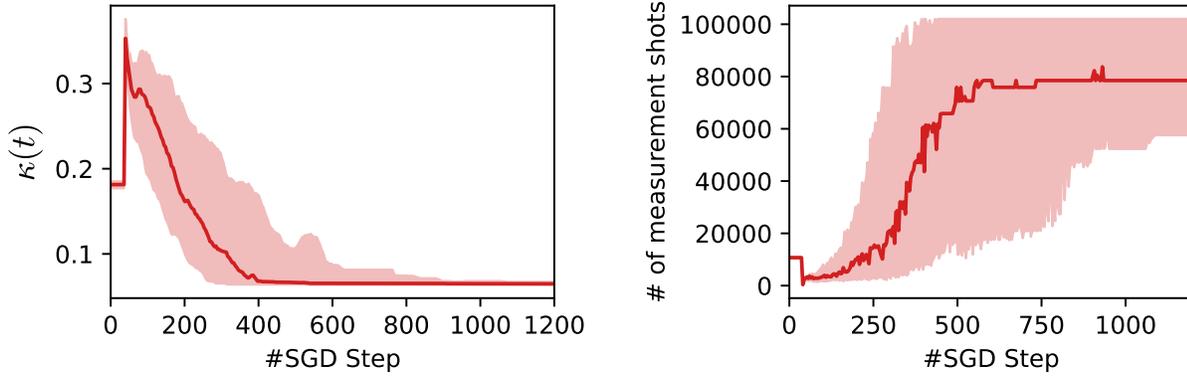


Figure 8. The GradCoRe threshold $\kappa(t)$ (left), set according to Eq. (24), and the number of measurement shots (right) per SGD iteration used by GradCoRe. As expected, the number of shots gradually increases as the GradCoRe threshold decreases, reflecting the flatness of the objective function via the gradient norm estimation.