

Stability and Generalization of Quantum Neural Networks

Jiaqi Yang, Wei Xie*, Xiaohua Xu*

University of Science and Technology of China, Hefei 230027, China
yangjiaqi@mail.ustc.edu.cn, xxieww@ustc.edu.cn, xiaohuaxu@ustc.edu.cn

Abstract

Quantum neural networks (QNNs) play an important role as an emerging technology in the rapidly growing field of quantum machine learning. While their empirical success is evident, the theoretical explorations of QNNs, particularly their generalization properties, are less developed and primarily focus on the uniform convergence approach. In this paper, we exploit an advanced tool in classical learning theory, i.e., algorithmic stability, to study the generalization of QNNs. We first establish high-probability generalization bounds for QNNs via uniform stability. Our bounds shed light on the key factors influencing the generalization performance of QNNs and provide practical insights into both the design and training processes. We next explore the generalization of QNNs on near-term noisy intermediate-scale quantum (NISQ) devices, highlighting the potential benefits of quantum noise. Moreover, we argue that our previous analysis characterizes worst-case generalization guarantees, and we establish a refined optimization-dependent generalization bound for QNNs via on-average stability. Numerical experiments on various real-world datasets support our theoretical findings.

1 Introduction

Quantum machine learning (QML) is a rapidly growing field that has generated great excitement [1]. With the aim of solving complex problems beyond the reach of classical computers, firm and steady progress has been achieved over the past decade [2–4]. Quantum neural network (QNN), or equivalently, the parameterized quantum circuit (PQC) with a classical optimizer, has received great attention thanks to its potential to achieve quantum advantages on near-term noisy intermediate scale quantum (NISQ) devices [5–7]. Typically, stochastic gradient descent (SGD) is employed as the classical optimizer in QNNs due to its simplicity and efficiency. Driven by the significance of understanding the power of QNNs, a growing body of literature has been conducted to investigate their expressivity [8–11], trainability [12–15], and generalization [16–21]. Investigating the generalization of QNNs is crucial for understanding their underlying working principles and capabilities from a theoretical perspective. Nevertheless, to date, the theoretical establishment in this area is still in its infancy.

In classical learning theory, a variety of techniques for generalization analysis are known. One of the most popular approach is via the uniform convergence analysis, which studies the uniform generalization gaps in a hypothesis space. This approach typically uses complexity measures such as VC-dimension [22], covering numbers [23], Rademacher complexity [24] to develop capacity-dependent bounds. To the best of our knowledge, essentially all generalization bounds derived for

*Corresponding author.

QNNs so far are of the uniform kind, which, however, are not sufficient to explain the generalization behavior of current-scale QNNs [25]. Impressive alternatives have been proposed, including sample compression [26], PAC-Bayes [27], and algorithmic stability [28]. In particular, the algorithmic stability tools have advantages in some aspects, such as dimension-independent and adaptability for broad learning paradigms. Therefore, it is natural to explore the generalization of QNNs through the lens of algorithmic stability.

In this paper, we study the generalization of QNNs trained using the SGD algorithm based on algorithmic stability. Our contributions are summarized as follows.

- We analyze the uniform stability of QNNs trained using the SGD algorithm. Our results reveal that the negative effects arising from complex QNN models on stability can be mitigated by setting appropriate step sizes. We investigate the simplified stability results for two commonly used step sizes.
- We establish high-probability generalization bounds for QNNs via uniform stability. Our bounds shed light on the key factors influencing the generalization performance of QNNs. Furthermore, we provide practical insights into both the design and training processes for developing powerful QNNs. Notably, our bounds help explain why over-parameterized QNNs trained by SGD exhibit excellent generalization, which was not explained by existing bounds for QNNs with the same setting [19, 20].
- We further extend our analysis to the noise scenario. Considering the standard depolarizing noise model (easily extended to other noise models), we similarly establish high-probability generalization bounds for QNNs. Our bounds highlight the potential benefits of quantum noise. Specifically, we provide new insight that the noise naturally occurring in quantum devices can be effectively tuned as a form of regularization. Moreover, our results theoretically explain the effectiveness of the method proposed in [29], which enhances QNN performance by controlling the noise level in quantum hardware as hyperparameters during the training process.
- We argue that previous analysis characterizes worst-case generalization guarantees. As a complement to our previous results, we establish a refined optimization-dependent generalization bound via on-average stability. This new bound reveals that the worst-case bounds can be improved in certain regimes. In addition, we corroborate the connection of our bound to the generalization performance of the recent experiments in [25]. Our bound captures the effect of randomized labels on generalization in terms of the on-average variance of SGD. As a corollary, our results validates the intuition that if we are good at the initialization point, the QNN model is more stable and thus generalizes better.
- We conduct numerical experiments on real-world datasets, and the empirical results indeed support our theoretical findings.

2 Related Work

Generalization analysis of QNNs. The theory of generalization for QNNs is less developed and primarily focuses on the uniform convergence approach. Recent studies on the generalization of QNNs typically employ complexity measures such as pseudo-dimension [16], effective dimension [30],

VC-dimension [21], Rademacher complexity [18, 31, 32], covering number [19, 20], and all their uniform relatives. Of the most interest to ours is [19], where the authors provide generalization bounds for QNNs based on covering number, all derived with the same setting as ours. Later, [20] uses quantum channels to derive more general results. Their generalization bounds exhibit a sublinear dependence on the number of trainable quantum gates. Unfortunately, this limits their results in providing meaningful guarantees for over-parameterized QNNs. More recently, it has been argued in [25] that the traditional measures of model complexity are not sufficient to explain the generalization behavior of current-scale QNNs. Their empirical findings highlight the need to shift perspective toward non-uniform generalization measures in QML. Therefore, it is natural to leverage the concept of algorithmic stability [28], which additionally enjoys desirable properties on flexibility (dimension-independence) and adaptivity (suited for diverse learning scenarios). While completing our work, we became aware of an anonymous work on OpenReview that derives a generalization bound for data re-uploading QNNs via uniform stability. We develop our approach independently of them, which differs in multiple key aspects and thus obtains more general results. Unlike their focus on constant step sizes, we also analyze the *practically relevant* scenarios of decaying step sizes. We further explore the generalization of QNNs in the noise scenario, highlighting the potential benefits of quantum noise. Moreover, we introduce another concept of on-average stability to provide the refined optimization-generalization bound.

Stability and Generalization. The classical framework of quantifying generalization via stability was established in an influential paper [28], where the celebrated concept of uniform stability was introduced. Subsequently, the uniform stability measure was extended to study stochastic algorithms [33]. The seminal work [34] pioneered the generalization analysis of SGD via uniform stability, which inspired several follow-up studies to understand stochastic optimization algorithms based on different algorithmic stability measures, e.g., on-average stability [35, 36], locally elastic stability [37], and argument stability [36, 38]. Stability-based generalization analysis was also developed for transfer learning [35], pairwise learning [39, 40], and minimax problems [41, 42]. The power of stability analysis is especially reflected by its ability to derive optimal generalization bounds in expectation [43]. Recent studies show that uniform stability can yield almost optimal high-probability bounds [44, 45]. While significant progress has been achieved, there is a lack of analysis on generalization from the perspective of stability in the context of QNNs.

3 Problem Setup

3.1 Quantum Computation Basics and Notations

We briefly introduce some basic concepts of quantum computation that are necessary for this work. Interested readers are recommended to the celebrated textbook by Nielsen and Chuang [46]. Qubit is the fundamental unit of quantum computation and quantum information. An N -qubit quantum state is represented by the density matrix, which is a Hermitian, positive semi-definite matrix $\rho \in \mathbb{C}^{2^N \times 2^N}$ with $\text{Tr}(\rho) = 1$. Quantum gates are unitary matrices used to transform quantum states. Common single-qubit gates include Pauli rotations $\left\{ R_P(\theta) = e^{-i\frac{\theta}{2}P} \mid P \in \{X, Y, Z\} \right\}$, which are in the exponential form of Pauli matrices,

$$X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad Y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, \quad Z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$

Common two-qubit gates include controlled-X gate, CNOT = $|0\rangle\langle 0| \otimes I + |1\rangle\langle 1| \otimes X$, is used to generate quantum entanglement among qubits. The evolution of a quantum state ρ can be mathematically described by employing a quantum circuit: $\rho' = U(\boldsymbol{\theta})\rho U^\dagger(\boldsymbol{\theta})$, where the unitary $U(\boldsymbol{\theta})$ is usually parameterized by a series of single-qubit rotation gate angles $\boldsymbol{\theta}$ and basic two-qubit gates, and U^\dagger denotes the conjugate transpose of U . Quantum measurement is a means to extract classical (observable) information from the quantum states. An observable is represented by a Hermitian matrix $O \in \mathbb{C}^{2^N \times 2^N}$. Since measurement is an irreversible process, it is typically introduced at the end of the quantum circuit. The expectation of the output is calculated as $\text{Tr}(O\rho')$.

Quantum Neural Network. The quantum neural network typically contains three parts, i.e., an N -qubit quantum circuit $U(\boldsymbol{\theta})$, an observable $O \in \mathbb{C}^{2^N \times 2^N}$, and a classical optimizer that update trainable parameters $\boldsymbol{\theta}$ to minimize the predefined objective function. For classical data \mathbf{x} , the input is first encoded into a quantum state via a map $\mathbf{x} \mapsto \rho(\mathbf{x}) \in \mathbb{C}^{2^N \times 2^N}$. Define $U(\boldsymbol{\theta}) = \prod_{k=1}^{K_g} U_k(\boldsymbol{\theta})$, where U_k refers to the k -th quantum gate. In general, $U(\boldsymbol{\theta})$ is formed by K trainable gates, $K_g - K$ fixed gates, and $\boldsymbol{\theta} \in \mathbb{R}^K$. Under the above definitions, the output function of the QNN can be written as follows

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \text{Tr}\left(O U(\boldsymbol{\theta})\rho(\mathbf{x})U^\dagger(\boldsymbol{\theta})\right).$$

In addition, the quantum gates in NISQ chips are prone to errors [47]. The noise can be simulated by applying certain quantum channels to each quantum gate, specifically the depolarization channel. Given a quantum state $\rho \in \mathbb{C}^{2^N \times 2^N}$, the depolarization channel \mathcal{N}_p acts on a 2^N -dimensional Hilbert space follows $\mathcal{N}_p(\rho) = (1-p)\rho + p\mathbb{I}/2^N$, where p is the noise level and $\mathbb{I}/2^N$ is the maximally mixed state.

3.2 Stability and Generalization Analysis of Randomized Algorithm

Let \mathcal{D} be a probability distribution defined on a sample space $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, where $\mathcal{X} \subseteq \mathbb{R}^d$ is an input space and $\mathcal{Y} \subseteq \mathbb{R}$ is an output space. We quantify the loss of $\boldsymbol{\theta}$ on a single example $\mathbf{z} = (\mathbf{x}, y)$ by $\ell(\boldsymbol{\theta}; \mathbf{z})$. The objective is to learn a model $\boldsymbol{\theta} \in \mathbb{R}^K$ minimizing the *population risk* defined by

$$R_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}}[\ell(\boldsymbol{\theta}; \mathbf{z})].$$

In practice, we do not know the distribution \mathcal{D} but instead have access to a dataset $S = \{\mathbf{z}_i = (\mathbf{x}_i, y_i) : i = 1, \dots, m\}$ independently drawn from \mathcal{D} . Then, we approximate $R_{\mathcal{D}}$ by *empirical risk*

$$R_S(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \ell(\boldsymbol{\theta}; \mathbf{z}_i).$$

For a randomized algorithm A , denote by $A(S)$ its output model based on the training dataset S . Let $R(\boldsymbol{\theta}_*) = \inf_{\boldsymbol{\theta}} R_{\mathcal{D}}(\boldsymbol{\theta})$ and $R(\boldsymbol{\theta}_*^S) = \inf_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta})$, then the *excess risk* is $\mathbb{E}_A[R_{\mathcal{D}}(A(S)) - R(\boldsymbol{\theta}_*)]$. Here, the expectation $\mathbb{E}_A[\cdot]$ is taken only over the internal randomness of A . It can be decomposed as

$$\mathbb{E}_A[R_{\mathcal{D}}(A(S)) - R(\boldsymbol{\theta}_*)] \leq \mathbb{E}_A[R_{\mathcal{D}}(A(S)) - R_S(A(S))] + \mathbb{E}_A[R_S(A(S)) - R_S(\boldsymbol{\theta}_*^S)],$$

where we have used the fact that $R_S(\boldsymbol{\theta}_*^S) \leq R_S(\boldsymbol{\theta}_*)$ by the definition of $\boldsymbol{\theta}_*^S$. The first term is called the *generalization (error) gap*, as it quantifies the generalization shift from training to testing

behavior. The second term is called the *optimization error*, as it measures how effectively the algorithm minimizes empirical risk. This paper focuses on bounding the generalization gap, for which a popular approach is based on the stability analysis of the algorithm.

Algorithmic Stability. Algorithmic stability plays an important role in classical learning theory, which measures the sensitivity of an algorithm to the perturbation of training sets. Our analysis relies on two widely used stability measures, namely *uniform stability* [28, 34] and *on-average stability* [36]. Below, we recall these notions.

Definition 3.1 (Uniform Stability). We say a randomized algorithm A is ϵ -uniformly stable if for any datasets $S, S' \in \mathcal{Z}^m$ that differ by at most one example, we have

$$\sup_{\mathbf{z}} |\mathbb{E}_A [\ell(A(S); \mathbf{z}) - \ell(A(S'); \mathbf{z})]| \leq \epsilon.$$

Definition 3.2 (On-Average Stability). Let $S = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ and $S' = \{\mathbf{z}'_1, \dots, \mathbf{z}'_m\}$ be drawn independently from \mathcal{D} . For any $i \in [m]$, define $S^{(i)} = \{\mathbf{z}_1, \dots, \mathbf{z}_{i-1}, \mathbf{z}'_i, \mathbf{z}_{i+1}, \dots, \mathbf{z}_m\}$ as the set formed from S by replacing the i -th element with \mathbf{z}'_i . We say a randomized algorithm A is on-average ϵ -stable if

$$\mathbb{E}_{S, S', A} \left[\frac{1}{m} \sum_{i=1}^m \|A(S) - A(S^{(i)})\|_2 \right] \leq \epsilon.$$

The celebrated relationship between algorithmic stability and generalization was established in the following lemma.

Lemma 3.3 (Stability and Generalization). *Let A be a randomized algorithm, $\epsilon > 0$ and $\delta \in (0, 1)$.*

(a) *If A is ϵ -uniformly stable, then the expected generalization gap satisfies*

$$|\mathbb{E}_{S, A} [R_{\mathcal{D}}(A(S)) - R_S(A(S))]| \leq \epsilon.$$

(b) *Assume $\ell(\boldsymbol{\theta}; \mathbf{z}) \in [0, M]$ for all $\boldsymbol{\theta} \in \mathbb{R}^K$ and $\mathbf{z} \in \mathcal{Z}$. If A is ϵ -uniformly stable, then with probability at least $1 - \delta$, the generalization gap satisfies*

$$\mathbb{E}_A [R_{\mathcal{D}}(A(S)) - R_S(A(S))] = \mathcal{O} \left(\epsilon \log m \log(1/\delta) + M m^{-\frac{1}{2}} \sqrt{\log(1/\delta)} \right).$$

(c) *Assume $|\ell(\boldsymbol{\theta}_1; \mathbf{z}) - \ell(\boldsymbol{\theta}_2; \mathbf{z})| \leq L \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2$ for all $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^K$ and $\mathbf{z} \in \mathcal{Z}$. If A is on-average ϵ -stable, then the expected generalization gap satisfies*

$$|\mathbb{E}_{S, A} [R_{\mathcal{D}}(A(S)) - R_S(A(S))]| \leq L\epsilon.$$

Remark 3.4. Part (a) and Part (b) establish the connection between uniform stability and generalization [28, 45]. Part (c) establishes the connection between on-average stability and generalization [36]. Both Part (a) and Part (c) provide generalization bounds in expectation. Under the assumption $|\ell(\boldsymbol{\theta}_1; \mathbf{z}) - \ell(\boldsymbol{\theta}_2; \mathbf{z})| \leq L \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2$ for all $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^K$ and $\mathbf{z} \in \mathcal{Z}$, we notice that ϵ -uniformly stable implies at least (ϵ/L) -on-average stable. Thus, we can recover the worst-case generalization bound as in Part (a). Part (b) provides an almost optimal high-probability generalization bound.

Typically, SGD is employed as the classical optimizer in QNNs. Given a training dataset $S \in \mathcal{Z}^m$, the QNN minimizes the following objective function,

$$\min_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i); y_i).$$

Definition 3.5 (Stochastic Gradient Descent). Let $\boldsymbol{\theta}_0 \in \mathbb{R}^K$ be an initial point. SGD updates $\{\boldsymbol{\theta}_t\}$ as follows

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}),$$

where η_t is the step size, and $\mathbf{z}_{i_t} = (\mathbf{x}_{i_t}, y_{i_t})$ is the sample chosen in iteration t . There are two popular schemes for choosing the example indices i_t . One is to choose uniformly from $\{1, \dots, m\}$ at each step. The other is to choose a random permutation of $\{1, \dots, m\}$ and then process the examples in order. Our results hold for both variants.

3.3 Main Assumptions

We aim to analyze the stability and generalization of QNNs trained using the SGD algorithm. To achieve this, we first introduce the necessary assumptions.

Assumption 3.6 (α_ℓ -Lipschitzness). We assume that the loss function ℓ is α_ℓ -Lipschitz, if for all $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^K$, and $\mathbf{z} \in \mathcal{Z}$, we have

$$|\ell(f_{\boldsymbol{\theta}_1}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y)| \leq \alpha_\ell |f_{\boldsymbol{\theta}_1}(\mathbf{x}) - f_{\boldsymbol{\theta}_2}(\mathbf{x})|.$$

Assumption 3.7 (ν_ℓ -smoothness). We assume that the loss function ℓ is ν_ℓ -smooth, if for all $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^K$, and $\mathbf{z} \in \mathcal{Z}$, we have

$$\left| \frac{\partial}{\partial f} \ell(f_{\boldsymbol{\theta}_1}(\mathbf{x}); y) - \frac{\partial}{\partial f} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y) \right| \leq \nu_\ell |f_{\boldsymbol{\theta}_1}(\mathbf{x}) - f_{\boldsymbol{\theta}_2}(\mathbf{x})|.$$

Remark 3.8. Unlike the classic setup of [34], we define Lipschitzness and smoothness with respect to the output function $f_{\boldsymbol{\theta}}(\cdot)$ rather than the parameters $\boldsymbol{\theta}$. By relaxing their assumptions, we provide a fine-grained analysis of stability and generalization in the context of QNNs.

4 Main Results

In this section, we present our main results on the stability and generalization bounds for QNNs trained using SGD algorithm. We note that our results are highly general and cover diverse ansatz of QNNs, as long as it is composed of the parameterized single-qubit gates and two-qubit gates. Due to space limitations, please refer to the Appendix C for detailed theoretical proofs.

4.1 Uniform Stability-Based Generalization Bounds

In this subsection, we analyze the uniform stability of QNNs and subsequently derive the high-probability generalization bounds.

Theorem 4.1 (Uniform Stability Bound). *Suppose that Assumption 3.6 and 3.7 hold. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t for T iterations, then $A(S)$ is ϵ -uniformly stable with*

$$\epsilon \leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\sqrt{2}\eta_t \alpha_\ell^2 K \|O\|^2}{m},$$

where

$$\kappa := \alpha_\ell K \|O\| + \sqrt{2\nu_\ell} K \|O\|^2. \quad (1)$$

Remark 4.2. The stability bound expands with coefficient $1 + \eta_t \kappa$, where κ depends on the spectral norm of the observable $\|O\|$ and the number of trainable quantum gates K . In practice, $\|O\|$ is typically bounded such that $0 \leq \|O\| \leq 1$. Moreover, the negative effects of complex QNN models with large K on stability can be mitigated by setting appropriate step sizes.

Notice that the step size is a vital parameter in Theorem 4.1. We further investigate the simplified stability results for two commonly used step sizes in the following corollary.

Corollary 4.3. *Suppose that Assumption 3.6 and 3.7 hold, and $\ell(\cdot, \cdot) \in [0, M]$. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t for T iterations.*

(a) *If we choose the constant step sizes $\eta_t = \eta$, then $A(S)$ is ϵ -uniformly stable with*

$$\epsilon \leq \frac{2\sqrt{2}\alpha_\ell^2 K \|O\|^2}{\kappa m} (1 + \eta \kappa)^T.$$

(b) *If we choose the monotonically non-increasing step sizes $\eta_t \leq c/(t+1)$, $c > 0$, then $A(S)$ is ϵ -uniformly stable with*

$$\epsilon \leq \frac{1 + 1/c\kappa}{m} M^{\frac{c\kappa}{c\kappa+1}} \left(2\sqrt{2}c\alpha_\ell^2 K \|O\|^2 \right)^{\frac{1}{c\kappa+1}} T^{\frac{c\kappa}{c\kappa+1}},$$

where κ is defined by equation (1).

Remark 4.4. Part (a) and Part (b) consider constant and decaying step sizes, respectively. The constant step size setting is widely used in prior works on the stability and generalization analysis of classical graph convolutional networks (GCNs) [48–50]. However, the resulting stability bound exhibits an exponential dependence on T . The exponential dependence can be eliminated by using decaying step sizes, with two main approaches for setting the decay. One approach is the classic analysis from [34], which considers the timing of encountering different examples in datasets. It employs a polynomially decaying step sizes $\eta_t = \mathcal{O}(t^{-1})$. The other approach controls the smallest eigenvalues of the loss’s Hessian matrix, which allows larger step sizes $\eta_t = \mathcal{O}(t^{-\beta})$, $\beta \in (0, 1)$ [51–53]. However, their analysis is limited to shallow neural networks and requires an additional restrictive condition that the network must be sufficiently wide. To develop a highly general stability bound for QNNs, we adopt the same strategy as [34], setting the step sizes $\eta_t \leq c/(t+1)$, $c > 0$ in Part (b).

By combining Lemma 3.3, Part (b), we can now easily derive the generalization bound. The high-probability generalization bounds for QNNs are subsequently established in the following theorem.

Theorem 4.5 (Generalization Bound). *Suppose that Assumption 3.6 and 3.7 hold, and $\ell(\cdot, \cdot) \in [0, M]$. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t for T iterations.*

- (a) *if we choose the constant step sizes $\eta_t = \eta$, then the following generalization bound of $A(S)$ holds with probability at least $1 - \delta$ for $\delta \in (0, 1)$,*

$$\mathbb{E}_A [R_{\mathcal{D}}(A(S)) - R_S(A(S))] \leq \mathcal{O} \left(\frac{(1 + \eta\kappa)^T}{m} \log m \log\left(\frac{1}{\delta}\right) + M \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{m}} \right),$$

- (b) *if we choose the monotonically non-increasing step sizes $\eta_t \leq c/(t + 1)$, $c > 0$, then the following generalization bound of $A(S)$ holds with probability at least $1 - \delta$ for $\delta \in (0, 1)$,*

$$\mathbb{E}_A [R_{\mathcal{D}}(A(S)) - R_S(A(S))] \leq \mathcal{O} \left(\frac{T^{\frac{c\kappa}{c\kappa+1}}}{m} \log m \log\left(\frac{1}{\delta}\right) + M \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{m}} \right),$$

where κ is defined by equation (1).

Remark 4.6. Theorem 4.5 provides practical insights into both the design and training processes for developing powerful QNNs. Specifically, our bounds shed light on the key factors influencing the generalization performance of QNNs, which can be divided into two categories. The first is associated with the design of QNNs, including the number of trainable quantum gates K , and the spectral norm of the observable $\|O\|$. The second is associated with the training of QNNs, including the step sizes η_t , the number of iterations T , and the number of training examples m .

- *Design.* In practice, $\|O\|$ is bounded such that $0 \leq \|O\| \leq 1$, and it is normally chosen as Pauli Strings. The dependence on K reflects an Occam's razor principle in the quantum version [54]. It is evident that, a larger K leads to an increase in the upper bound of the generalization gap. This provides guidance for designing well-performing QNNs with a proper number of trainable quantum gates. Moreover, increasing the number of trainable quantum gates K directly enhances the expressivity of QNNs [19, 55]. This leads to a trade-off between expressivity and generalization in designing powerful QNNs, analogous to the bias-variance trade-off in classical machine learning.
- *Training.* First, increasing the number of training examples m directly enhances the generalization performance of QNNs. Second, it is crucial to carefully choose the step sizes. Clearly, the decaying step size setting is a better choice than the constant step size setting. In the constant step size setting, our bound is primarily controlled by the term $\mathcal{O}((1 + \eta\kappa)^T/m)$. This suggests setting the step size as $\eta = \mathcal{O}(1/K)$ to ensure $\eta\kappa = \mathcal{O}(1)$, thereby preventing the generalization bound from becoming vacuous due to the inherent complexity of QNNs. Lastly, we underscore the importance of reducing training time. Decreasing the number of iterations T is an effective way to reduce the generalization gap. Therefore, in practice, *early stopping* is an important training technique for QNNs, where stop training early after reach a low training error.

Remark 4.7. We compare Theorem 4.5 with related works [19, 20]. The state-of-the-art generalization bound for QNNs in the same setting is of the order $\mathcal{O}\left(\sqrt{K/m}\right)$. However, the *sublinear* dependence on K makes their results trivially loose for the *over-parameterized* QNNs, where $K \gg m$. In contrast, our bounds help explain why over-parameterized QNNs exhibit excellent generalization. Specifically, our bounds highlight that the negative effects of large K on generalization can be mitigated by setting appropriate step sizes, thereby providing meaningful generalization guarantees as well. As mentioned earlier, in the constant step size setting, our bound suggest setting the step sizes as $\eta = \mathcal{O}(1/K)$ to balance the negative effect of large K . In the decaying step size setting, our bound is primarily controlled by the term $\mathcal{O}\left(T^{\frac{c\kappa}{c\kappa+1}}/m\right)$. It is worth noting that the negative effect of K is significantly reduced. No matter how large K becomes, the generalization gap converges to zero as $m \rightarrow \infty$, provided that T is not too large. We show T can grow as m^a for a small $a > 1$.

Our previous analysis can be easily extended to the noise scenario. Considering the standard depolarizing noise model, we similarly establish the high-probability generalization bounds for QNNs in the following corollary. We remark that while our results are presented assuming the depolarization noise, they can be easily extended to other noise models.

Corollary 4.8 (Generalization Bound Under Depolarizing Noise). *Suppose that Assumption 3.6 and 3.7 hold, and $\ell(\cdot, \cdot) \in [0, M]$. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t under depolarizing noise level $p \in [0, 1]$ for T iterations.*

- (a) *if we choose the constant step sizes $\eta_t = \eta$, then the following generalization bound of $A(S)$ holds with probability at least $1 - \delta$ for $\delta \in (0, 1)$,*

$$\mathbb{E}_A [R_{\mathcal{D}}(A(S)) - R_S(A(S))] \leq \mathcal{O} \left(\frac{\left(1 + (1-p)^{K_g} \eta \kappa\right)^T}{m} \log m \log\left(\frac{1}{\delta}\right) + M \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{m}} \right),$$

- (b) *if we choose the monotonically non-increasing step sizes $\eta_t \leq c/(t+1)$, $c > 0$, then the following generalization bound of $A(S)$ holds with probability at least $1 - \delta$ for $\delta \in (0, 1)$,*

$$\mathbb{E}_A [R_{\mathcal{D}}(A(S)) - R_S(A(S))] \leq \mathcal{O} \left(\frac{T^{\frac{c\kappa}{c\kappa+1/(1-p)^{K_g}}}}{m} \log m \log\left(\frac{1}{\delta}\right) + M \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{m}} \right),$$

where κ is defined by equation (1).

Remark 4.9. Corollary 4.8 reveals the potential benefits of the quantum noise. Our bounds shows that an increase in the noise level p enhances the generalization performance of QNNs. Moreover, previous work has indicated that larger noise leads to poorer expressivity of QNNs [19, 55]. Therefore, our bounds provide new insight that the noise naturally occurring in quantum devices can be effectively tuned as a form of ‘*quantum regularization*’, with the ability balancing expressivity and generalization of QNNs. It has been demonstrated that adding noise to the initial data or the weights of QNNs can have an effect analogous to the technique of *regularization* employed in classical machine learning [56, 57]. Unlike these existing techniques, [29] proposed an approach to control the noise level in quantum hardware as hyperparameters during the training process. They numerically investigate this method for a regression task by using several modeled noise channels and demonstrate an improvement in model performance. Our results theoretically explain the potential reasons for its effectiveness, acting akin to regularization in classical neural networks.

4.2 Towards Optimization-Dependent Generalization Bounds

In the previous subsection, we focused on uniform stability and derived the key results of this paper. However, uniform stability does not depend on the data, but captures only intrinsic characteristics of the learning algorithm and global properties of the objective function. Consequently, previous analysis characterizes worst-case generalization guarantees. As a complement to our previous results, we further investigate a refined optimization-dependent generalization bound, leveraging the less restrictive on-average stability [35, 36].

To capture the impact of the variance of the stochastic gradients, we adopt the following standard assumption from stochastic optimization theory [58, 59].

Assumption 4.10 (Bounded Empirical Variance). For any dataset $S \in \mathcal{Z}^m$ and $\theta \in \mathbb{R}^K$, their exist $\sigma^2 > 0$ such that

$$\frac{1}{m} \sum_{i=1}^m \|\nabla_{\theta} \ell(f_{\theta}(\mathbf{x}_i); \mathbf{y}_i) - \nabla_{\theta} R_S(\theta)\|_2^2 \leq \sigma^2. \quad (2)$$

Remark 4.11. Assumption 4.10 essentially bounds the variance of the stochastic gradients for the particular dataset S . Notably, it is always satisfied if $\|\nabla_{\theta} \ell(f_{\theta}(\mathbf{x}); \mathbf{y})\|_2 \leq G$ for any $\theta \in \mathbb{R}^K$ and $z \in \mathcal{Z}$, with $\sigma^2 = G^2$.

The following theorem establishes the generalization bound in expectation for QNNs and first links generalization gap to optimization error.

Theorem 4.12 (Optimization-Dependent Generalization Bound). *Suppose that Assumption 3.6, 3.7, and 4.10 hold. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t for T iterations, then the expected generalization gap satisfies*

$$|\mathbb{E}_{S,A} [R_{\mathcal{D}}(A(S)) - R_S(A(S))]| \leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\eta_t \alpha_{\ell} \sqrt{K} \|O\| (\mathbb{E}_S [\|\nabla_{\theta} R_S(\theta_t)\|_2] + \sigma)}{m},$$

where κ is defined by equation (1).

Remark 4.13. From Theorem 4.12, we notice that we can bound σ and the gradient norms $\|\nabla_{\theta} R_S(\theta_t)\|_2$ by G , where $G = \sqrt{2} \alpha_{\ell} \sqrt{K} \|O\|$. Recall from Lemma 3.3, Part (a), uniform stability directly implies generalization in expectation. Thus, we can recover (up to a factor 2), the worst-case generalization bound from Theorem 4.1. This illustrates well the worst-case notion mentioned earlier, but also the fact that the generalization bound in expectation of Theorem 4.12 can be better than the one of Theorem 4.1. This will notably be the case in “low noise” regimes, when $\sigma \ll G$ and the expected gradient norm $\mathbb{E}_S \|\nabla_{\theta} R_S(\theta_t)\|_2$ reach small values. In particular, the expected gradient norm $\mathbb{E}_S \|\nabla_{\theta} R_S(\theta_t)\|_2$ is related to the optimization error, which decreases as the parameter θ_t minimizes the empirical risk $R_S(\theta)$.

Remark 4.14. Theorem 4.12 helps explain the observations in classification experiments with randomized labels [25]. They conduct randomization experiments by training QNNs on a set of quantum states with varying levels of label corruption. Without changing the QNN ansatz, the number of training examples, or the optimization algorithm, they observe a steady increase in the generalization gap as the random label probability increases. Our bound properly captures how the generalization gap changes with the fraction of random labels via the on-average variance of SGD. Specifically, as the random label probability increases, the on-average variance σ keeps increasing and the generalization gap also increases.

Note that Theorem 4.12 can be similarly extended to various step size settings and noise scenario. Additionally, we show that the expected gradient norm $\mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2]$ are influenced by the choice of the initialization point.

Lemma 4.15 (Link with Initialization Point). *Suppose that Assumption 3.6, 3.7, and 4.10 hold. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes $\eta_t \leq 1/\kappa$ for T iterations, then the following bound holds*

$$\sum_{t=0}^{T-1} \eta_t \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2] \leq 2 \sqrt{\left(\sum_{t=0}^{T-1} \eta_t \right) \left(R_S(\boldsymbol{\theta}_0) - R_S(\boldsymbol{\theta}^*) + \frac{\kappa \sigma^2}{2} \sum_{t=0}^{T-1} \eta_t^2 \right)},$$

where $\boldsymbol{\theta}^*$ is the empirical risk minimizer of $R_S(\boldsymbol{\theta})$, κ is defined by equation (1).

Remark 4.16. Lemma 4.15 validates the intuition that if we are good at the initialization point $\boldsymbol{\theta}_0$, the QNN model is more stable and thus generalizes better. It suggests choosing a initialization point with low empirical risk in practice.

5 Experimental Evaluation

In this section, we conduct numerical simulations to validate our theoretical results. Note that we do not aim to optimize the test accuracy, but rather a simple, interpretable experimental setup.

Experiments Setup. We consider two real-world datasets: MNIST [60] and Fashion MNIST [61], which have also been explored in the anonymous work as well as in other QML studies [62–64]. For these datasets, we conduct binary classification tasks on the digit 0/1 and the T-shirt/Trouser. The implementation of QNN is as follows. The classical information \boldsymbol{x} is encoded into a quantum state $\rho(\boldsymbol{x})$ via angle encoding. We adopt the most widely used hardware-efficient ansatz such that the construction of $U(\boldsymbol{\theta})$ follows a layerwise structure using single-qubit Pauli rotation gates and two-qubit CNOT gates. We empirically estimate the generalization gap by calculating the absolute difference between the training and test errors. Each experiment setting is repeated 50 times to obtain statistical results.

Number of Trainable Quantum Gates. We investigate the impact of the number of trainable quantum gates K on generalization gap by varying the QNN layer numbers $L = \{4, 6, 8, 10, 12\}$. The step size is set to $\eta = 0.01$. In Figure 1, we observe that as L increases, i.e., K increases, the generalization gap also increases. This observation is consistent with our generalization bounds in Theorem 4.5 regarding the impact of K on the generalization gap.

Step Size. To avoid introducing additional hyperparameters, we mainly focus on constant step sizes to investigate the impact of the step size on generalization gap. We try various step sizes $\eta = \{0.001, 0.005, 0.01, 0.05, 0.1\}$. The QNN layer number is set to $L = 8$. It is clear from Figure 2 that the generalization gap increases with larger step sizes. This observation aligns well with our generalization bounds in Theorem 4.5. Additionally, since the number of trainable quantum gates K is fixed throughout the experiment, this further validating our theoretical explanation presented after Theorem 4.5, that is, the negative effects arising from the inherent complexity of the QNN models on generalization can be mitigated by setting appropriate step sizes.

Quantum Noise. We investigate the impact of the quantum noise on generalization gap by varying the noise level $p = \{0.001, 0.005, 0.01, 0.05, 0.1\}$. The QNN layer number is set to $L = 4$ and the step size is set to $\eta = 0.01$. It is depicted in Figure 3 that as the noise level p increases,

the generalization gap decreases. This result echoes with Corollary 4.8 and reinforces the insights that quantum noise can be effectively tuned as a form of regularization.

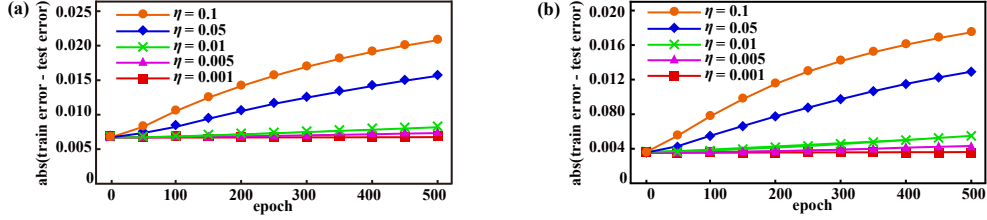


Figure 1: Generalization gap for varying numbers of trainable quantum gates: (a) MNIST, (b) Fashion MNIST.

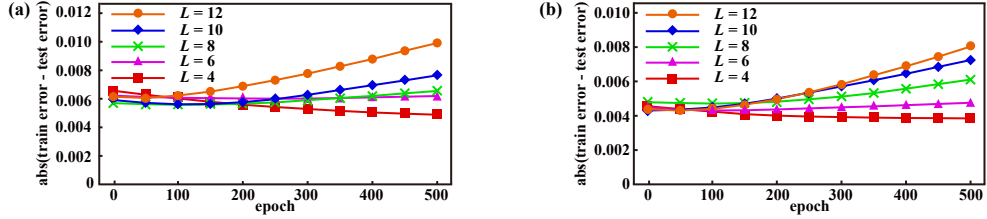


Figure 2: Generalization gap for varying step sizes: (a) MNIST, (b) Fashion MNIST.

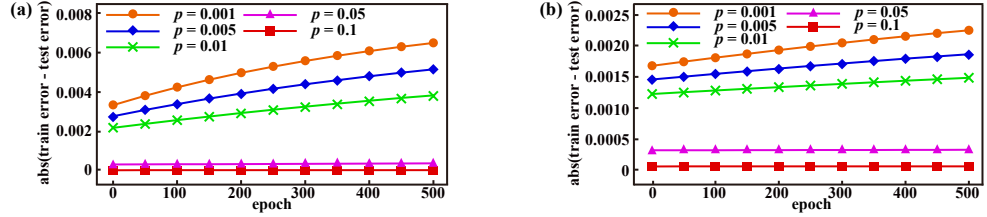


Figure 3: Generalization gap with varying noise levels: (a) MNIST, (b) Fashion MNIST.

Random Label. We conduct experiments to validate our explanation for the observations in [25] that a classification dataset with randomized labels can substantially degrade the generalization performance of QNNs. For all the data labels in each dataset, we replace their underlying true labels with random labels with probability r . The QNN layer number is set to $L = 8$, the step size is set to $\eta = 0.01$, and the number of iterations is set to $T = 500$. In Figure 4, we present the results under the random label probability $r = \{0.1, 0.2, 0.3, 0.4, 0.5\}$. It can be seen from these results that the on-average variance σ consistently increases as the random label probability r increases. At the same time, the generalization gap also increases. This validate our theoretical explanation for the observations in [25] presented after Theorem 4.12, that is, the on-average variance σ can capture how the generalization gap changes with the fraction of random labels.

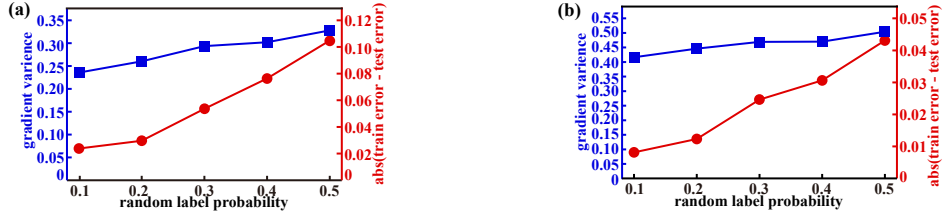


Figure 4: On-average variance and generalization gap for varying random label probabilities: (a) MNIST, (b) Fashion MNIST.

6 Conclusion

In this paper, we study the generalization of QNNs through algorithmic stability. We first establish high-probability generalization bounds for QNNs via uniform stability. Our bounds provide practical insights into both the design and training processes for developing powerful QNNs. We next extend our analysis to the noise scenario, highlighting the potential benefits of quantum noise. We finally argue that previous results are coming from worst-case analysis and propose a refined optimization-dependent generalization bound. While our generalization bounds hold for arbitrary data distributions, an interesting direction is to explore the generalization of QNNs with a specific data distribution.

7 Acknowledgements

This work was partially supported by the National Natural Science Foundation of China (Grant No. 62102388), the Innovation Program for Quantum Science and Technology (Grant No. 2021ZD0302901), and the USTC Kunpeng & Ascend Center of Excellence.

References

- [1] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd. Quantum machine learning. *Nature*, 549(7671):195–202, 2017.
- [2] Aram W Harrow and Ashley Montanaro. Quantum computational supremacy. *Nature*, 549(7671):203–209, 2017.
- [3] Vedran Dunjko and Hans J Briegel. Machine learning & artificial intelligence in the quantum domain: a review of recent progress. *Reports on Progress in Physics*, 81(7):074001, 2018.
- [4] Carlo Ciliberto, Mark Herbster, Alessandro Davide Ialongo, Massimiliano Pontil, Andrea Rocchetto, Simone Severini, and Leonard Wossnig. Quantum machine learning: a classical perspective. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 474(2209):20170551, 2018.
- [5] Edward Farhi and Hartmut Neven. Classification with quantum neural networks on near term processors. *arXiv preprint arXiv:1802.06002*, 2018.

- [6] Vojtěch Havlíček, Antonio D Córcoles, Kristan Temme, Aram W Harrow, Abhinav Kandala, Jerry M Chow, and Jay M Gambetta. Supervised learning with quantum-enhanced feature spaces. *Nature*, 567(7747):209–212, 2019.
- [7] Kerstin Beer, Dmytro Bondarenko, Terry Farrelly, Tobias J Osborne, Robert Salzmann, Daniel Scheiermann, and Ramona Wolf. Training deep quantum neural networks. *Nature communications*, 11(1):808, 2020.
- [8] Sukin Sim, Peter D Johnson, and Alán Aspuru-Guzik. Expressibility and entangling capability of parameterized quantum circuits for hybrid quantum-classical algorithms. *Advanced Quantum Technologies*, 2(12):1900070, 2019.
- [9] Carlos Bravo-Prieto, Josep Lumbrecas-Zarapico, Luca Tagliacozzo, and José I Latorre. Scaling of variational quantum circuit depth for condensed matter systems. *Quantum*, 4:272, 2020.
- [10] Yadong Wu, Juan Yao, Pengfei Zhang, and Hui Zhai. Expressivity of quantum neural networks. *Physical Review Research*, 3(3):L032049, 2021.
- [11] Dylan Herman, Rudy Raymond, Muyuan Li, Nicolas Robles, Antonio Mezzacapo, and Marco Pistoia. Expressivity of variational quantum machine learning on the boolean cube. *IEEE Transactions on Quantum Engineering*, 4:1–18, 2023.
- [12] Jarrod R McClean, Sergio Boixo, Vadim N Smelyanskiy, Ryan Babbush, and Hartmut Neven. Barren plateaus in quantum neural network training landscapes. *Nature communications*, 9(1):4812, 2018.
- [13] Marco Cerezo, Akira Sone, Tyler Volkoff, Lukasz Cincio, and Patrick J Coles. Cost function dependent barren plateaus in shallow parametrized quantum circuits. *Nature communications*, 12(1):1791, 2021.
- [14] Kunal Sharma, Marco Cerezo, Lukasz Cincio, and Patrick J Coles. Trainability of dissipative perceptron-based quantum neural networks. *Physical Review Letters*, 128(18):180505, 2022.
- [15] Martin Larocca, Nathan Ju, Diego García-Martín, Patrick J Coles, and Marco Cerezo. Theory of overparametrization in quantum neural networks. *Nature Computational Science*, 3(6):542–551, 2023.
- [16] Matthias C Caro and Ishaun Datta. Pseudo-dimension of quantum circuits. *Quantum Machine Intelligence*, 2(2):14, 2020.
- [17] Leonardo Bianchi, Jason Pereira, and Stefano Pirandola. Generalization in quantum machine learning: A quantum information standpoint. *PRX Quantum*, 2(4):040321, 2021.
- [18] Kaifeng Bu, Dax Enshan Koh, Lu Li, Qingxian Luo, and Yaobo Zhang. Statistical complexity of quantum circuits. *Physical Review A*, 105(6):062431, 2022.
- [19] Yuxuan Du, Zhuozhuo Tu, Xiao Yuan, and Dacheng Tao. Efficient measure for the expressivity of variational quantum algorithms. *Physical Review Letters*, 128(8):080506, 2022.

- [20] Matthias C Caro, Hsin-Yuan Huang, Marco Cerezo, Kunal Sharma, Andrew Sornborger, Lukasz Cincio, and Patrick J Coles. Generalization in quantum machine learning from few training data. *Nature communications*, 13(1):4919, 2022.
- [21] Casper Gyurik, Vedran Dunjko, et al. Structural risk minimization for quantum linear classifiers. *Quantum*, 7:893, 2023.
- [22] Vladimir Vapnik. Statistical learning theory. *John Wiley & Sons google schola*, 2:831–842, 1998.
- [23] Felipe Cucker and Ding Xuan Zhou. *Learning theory: an approximation theory viewpoint*, volume 24. Cambridge University Press, 2007.
- [24] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [25] Elies Gil-Fuster, Jens Eisert, and Carlos Bravo-Prieto. Understanding quantum machine learning also requires rethinking generalization. *Nature Communications*, 15(1):2277, 2024.
- [26] Nick Littlestone and Manfred Warmuth. Relating data compression and learnability. Technical report, University of California Santa Cruz, 1986.
- [27] David A McAllester. Some pac-bayesian theorems. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 230–234, 1998.
- [28] Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2:499–526, 2002.
- [29] Wilfrid Somogyi, Ekaterina Pankovets, Viacheslav Kuzmin, and Alexey Melnikov. Method for noise-induced regularization in quantum neural networks. *arXiv preprint arXiv:2410.19921*, 2024.
- [30] Amira Abbas, David Sutter, Christa Zoufal, Aurélien Lucchi, Alessio Figalli, and Stefan Woerner. The power of quantum neural networks. *Nature Computational Science*, 1(6):403–409, 2021.
- [31] Kaifeng Bu, Dax Enshan Koh, Lu Li, Qingxian Luo, and Yaobo Zhang. Rademacher complexity of noisy quantum circuits. *arXiv preprint arXiv:2103.03139*, 2021.
- [32] Kaifeng Bu, Dax Enshan Koh, Lu Li, Qingxian Luo, and Yaobo Zhang. Effects of quantum resources and noise on the statistical complexity of quantum circuits. *Quantum Science and Technology*, 8(2):025013, 2023.
- [33] Andre Elisseeff, Theodoros Evgeniou, Massimiliano Pontil, and Leslie Pack Kaelbling. Stability of randomized learning algorithms. *Journal of Machine Learning Research*, 6(1), 2005.
- [34] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR, 2016.

- [35] Ilja Kuzborskij and Christoph Lampert. Data-dependent stability of stochastic gradient descent. In *International Conference on Machine Learning*, pages 2815–2824. PMLR, 2018.
- [36] Yunwen Lei and Yiming Ying. Fine-grained analysis of stability and generalization for stochastic gradient descent. In *International Conference on Machine Learning*, pages 5809–5819. PMLR, 2020.
- [37] Zhun Deng, Hangfeng He, and Weijie Su. Toward better generalization bounds with locally elastic stability. In *International Conference on Machine Learning*, pages 2590–2600. PMLR, 2021.
- [38] Raef Bassily, Vitaly Feldman, Cristóbal Guzmán, and Kunal Talwar. Stability of stochastic gradient descent on nonsmooth convex losses. *Advances in Neural Information Processing Systems*, 33:4381–4391, 2020.
- [39] Yunwen Lei, Mingrui Liu, and Yiming Ying. Generalization guarantee of sgd for pairwise learning. *Advances in neural information processing systems*, 34:21216–21228, 2021.
- [40] Zhenhuan Yang, Yunwen Lei, Puyu Wang, Tianbao Yang, and Yiming Ying. Simple stochastic and online gradient descent algorithms for pairwise learning. *Advances in Neural Information Processing Systems*, 34:20160–20171, 2021.
- [41] Yunwen Lei, Zhenhuan Yang, Tianbao Yang, and Yiming Ying. Stability and generalization of stochastic gradient methods for minimax problems. In *International Conference on Machine Learning*, pages 6175–6186. PMLR, 2021.
- [42] Yue Xing, Qifan Song, and Guang Cheng. On the algorithmic stability of adversarial training. *Advances in neural information processing systems*, 34:26523–26535, 2021.
- [43] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and uniform convergence. *The Journal of Machine Learning Research*, 11:2635–2670, 2010.
- [44] Vitaly Feldman and Jan Vondrak. Generalization bounds for uniformly stable algorithms. *Advances in Neural Information Processing Systems*, 31, 2018.
- [45] Olivier Bousquet, Yegor Klochkov, and Nikita Zhivotovskiy. Sharper bounds for uniformly stable algorithms. In *Conference on Learning Theory*, pages 610–626. PMLR, 2020.
- [46] Michael A Nielsen and Isaac L Chuang. *Quantum computation and quantum information*, volume 2. Cambridge university press Cambridge, 2001.
- [47] John Preskill. Quantum computing in the nisq era and beyond. *Quantum*, 2:79, 2018.
- [48] Saurabh Verma and Zhi-Li Zhang. Stability and generalization of graph convolutional neural networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1539–1548, 2019.
- [49] Michael K Ng, Hanrui Wu, and Andy Yip. Stability and generalization of hypergraph collaborative networks. *Machine Intelligence Research*, 21(1):184–196, 2024.

- [50] Guangrui Yang, Ming Li, Han Feng, and Xiaosheng Zhuang. Deeper insights into deep graph convolutional networks: Stability and generalization. *arXiv preprint arXiv:2410.08473*, 2024.
- [51] Dominic Richards and Mike Rabbat. Learning with gradient descent and weakly convex losses. In *International Conference on Artificial Intelligence and Statistics*, pages 1990–1998. PMLR, 2021.
- [52] Dominic Richards and Ilja Kuzborskij. Stability & generalisation of gradient descent for shallow neural networks without the neural tangent kernel. *Advances in neural information processing systems*, 34:8609–8621, 2021.
- [53] Yunwen Lei, Rong Jin, and Yiming Ying. Stability and generalization analysis of gradient methods for shallow neural networks. *Advances in Neural Information Processing Systems*, 35:38557–38570, 2022.
- [54] Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Occam’s razor. *Information processing letters*, 24(6):377–380, 1987.
- [55] Yuxuan Du, Min-Hsiu Hsieh, Tongliang Liu, Shan You, and Dacheng Tao. Learnability of quantum neural networks. *PRX quantum*, 2(4):040337, 2021.
- [56] Nam H Nguyen, Elizabeth C Behrman, and James E Steck. Quantum learning with noise and decoherence: a robust quantum neural network. *Quantum Machine Intelligence*, 2(1):1, 2020.
- [57] Valentin Heyraud, Zejian Li, Zakari Denis, Alexandre Le Boité, and Cristiano Ciuti. Noisy quantum kernel machines. *Physical Review A*, 106(5):052421, 2022.
- [58] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
- [59] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT’2010: 19th International Conference on Computational Statistics Paris France, August 22-27, 2010 Keynote, Invited and Contributed Papers*, pages 177–186. Springer, 2010.
- [60] Yann LeCun, Corinna Cortes, Chris Burges, et al. Mnist handwritten digit database, 2010.
- [61] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [62] Hsin-Yuan Huang, Michael Broughton, Masoud Mohseni, Ryan Babbush, Sergio Boixo, Hartmut Neven, and Jarrod R McClean. Power of data in quantum machine learning. *Nature communications*, 12(1):2631, 2021.
- [63] Samuel Yen-Chi Chen, Chih-Min Huang, Chia-Wei Hsing, and Ying-Jer Kao. An end-to-end trainable hybrid classical-quantum classifier. *Machine Learning: Science and Technology*, 2(4):045021, 2021.
- [64] Tak Hur, Leeseok Kim, and Daniel K Park. Quantum convolutional neural network for classical data classification. *Quantum Machine Intelligence*, 4(1):3, 2022.

- [65] Kosuke Mitarai, Makoto Negoro, Masahiro Kitagawa, and Keisuke Fujii. Quantum circuit learning. *Physical Review A*, 98(3):032309, 2018.

A Notions

The main notations of this paper are summarized in Table 1.

Table 1: Summary of main notations involved in this paper.

Notion	Description
\mathcal{Z}	the sample space associated with input space \mathcal{X} and output space \mathcal{Y}
\mathcal{D}	the probability distribution defined on the sample space \mathcal{Z}
$z = (\mathbf{x}, y)$	the random example sampling from \mathcal{Z}
$\boldsymbol{\theta}$	the parameter of QNN
K_g	the number of quantum gates in QNN
K	the number of trainable quantum gates in QNN
O	the observable operator
S	the training dataset defined as $S = \{z_i = (\mathbf{x}_i, y_i) : i = 1, \dots, m\}$ independently drawn from \mathcal{D}
m	the number of training examples
$f_{\boldsymbol{\theta}}(\cdot)$	the output function of QNN
$\ell(\cdot)$	the loss function of QNN
$\nabla_{\boldsymbol{\theta}} \ell$	the gradient of $\ell(\cdot)$ to the argument $\boldsymbol{\theta}$
$R_{\mathcal{D}}, R_S$	the population risk and empirical risk based on training dataset S , respectively
T	the number of iterations for SGD
$\boldsymbol{\theta}_t$	the parameter of QNN learned using the SGD algorithm for t iterations
η_t	the step size in iteration t
$A, A(S)$	the learning algorithm and its output model based on training dataset S , respectively
$\alpha_{\ell}, \nu_{\ell}$	the parameters of Lipschitz continuity and smoothness, respectively

B Auxiliary Lemmas

In this section, we provide some auxiliary lemmas from quantum information theory that are essential for our main proofs.

To translate between the spectral norm of unitaries and the diamond norm of the corresponding channels, we employ the following lemma from [20].

Lemma B.1 (Spectral norm and diamond norm of unitary channels; see [20], Lemma 5). *Let $\mathcal{U}(\rho) = U\rho U^\dagger$ and $\mathcal{V}(\rho) = V\rho V^\dagger$ be unitary channels. Then, $\frac{1}{2}\|\mathcal{U}(|\psi\rangle\langle\psi|) - \mathcal{V}(|\psi\rangle\langle\psi|)\|_1 \leq \|(U - V)|\psi\rangle\|_2$ for any pure state $|\psi\rangle$. Therefore,*

$$\frac{1}{2}\|\mathcal{U} - \mathcal{V}\|_{\diamond} \leq \|U - V\|.$$

Moreover, we recall the following lemma from [46].

Lemma B.2 (see [46], section 4.5.3). Define the distance of two unitary matrices U_1, U_2 as the spectral norm of the matrix $U_1 - U_2$, i.e., $E(U_1, U_2) = \|U_1 - U_2\|$. Then,

$$E(U_K U_{K-1} \dots U_1, V_K V_{K-1} \dots V_1) \leq \sum_{j=1}^K E(U_j, V_j),$$

where $U_1, U_2, \dots, U_K, V_1, V_2, \dots, V_K$ are unitary matrices.

Next, we state and prove the following key lemma that provides a bound on the distance between two $U(\boldsymbol{\theta})$ with different parameter settings.

Lemma B.3 (Bound in Parameters Change). Suppose a parameterized unitary $U(\boldsymbol{\theta}) = \prod_{k=1}^K V_k e^{-i\boldsymbol{\theta}_k P_k / 2} V_{K+1}$, we have the upper bound on two different parameter sets $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^K$

$$\|U(\boldsymbol{\theta}_1) - U(\boldsymbol{\theta}_2)\| \leq \frac{\sqrt{K}}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2,$$

where V_k are fixed quantum gates, $P_k \in \{X, Y, Z\}$ denotes a single-qubit Pauli gate. For ease of readability, the tensor factors of identities accompanying the parametrized quantum gates $e^{-i\boldsymbol{\theta}_k P_k / 2}$ are omitted.

Proof. Note that we can write $U(\boldsymbol{\theta}_1) = V_1 e^{-i\boldsymbol{\theta}_{1,1} P_1 / 2} V_2 e^{-i\boldsymbol{\theta}_{1,2} P_2 / 2} \dots V_K e^{-i\boldsymbol{\theta}_{1,K} P_K / 2} V_{K+1}$, $U(\boldsymbol{\theta}_2) = V_1 e^{-i\boldsymbol{\theta}_{2,1} P_1 / 2} V_2 e^{-i\boldsymbol{\theta}_{2,2} P_2 / 2} \dots V_K e^{-i\boldsymbol{\theta}_{2,K} P_K / 2} V_{K+1}$. By Lemma B.2, $\|U(\boldsymbol{\theta}_1) - U(\boldsymbol{\theta}_2)\|$ can be bounded as

$$\begin{aligned} \|U(\boldsymbol{\theta}_1) - U(\boldsymbol{\theta}_2)\| &\leq \sum_{k=1}^{K+1} \|V_k - V_k\| + \sum_{k=1}^K \left\| e^{-i\frac{\boldsymbol{\theta}_{1,k} P_k}{2}} - e^{-i\frac{\boldsymbol{\theta}_{2,k} P_k}{2}} \right\| \\ &= \sum_{k=1}^K \left\| e^{-i\frac{\boldsymbol{\theta}_{1,k} P_k}{2}} - e^{-i\frac{\boldsymbol{\theta}_{2,k} P_k}{2}} \right\|. \end{aligned} \quad (3)$$

Additionally, the sub-term $\|e^{-i\frac{\boldsymbol{\theta}_{1,k} P_k}{2}} - e^{-i\frac{\boldsymbol{\theta}_{2,k} P_k}{2}}\|$ can be written as

$$\begin{aligned} \|e^{-i\frac{\boldsymbol{\theta}_{1,k} P_k}{2}} - e^{-i\frac{\boldsymbol{\theta}_{2,k} P_k}{2}}\| &= \|I - e^{i\frac{(\boldsymbol{\theta}_{1,k} - \boldsymbol{\theta}_{2,k}) P_k}{2}}\| \\ &= \|I - \cos\left(\frac{\boldsymbol{\theta}_{1,k} - \boldsymbol{\theta}_{2,k}}{2}\right) I + i \sin\left(\frac{\boldsymbol{\theta}_{1,k} - \boldsymbol{\theta}_{2,k}}{2}\right) P_k\| \\ &= \sqrt{\left(1 - \cos\left(\frac{\boldsymbol{\theta}_{1,k} - \boldsymbol{\theta}_{2,k}}{2}\right)\right)^2 + \left(\sin\left(\frac{\boldsymbol{\theta}_{1,k} - \boldsymbol{\theta}_{2,k}}{2}\right)\right)^2} \\ &= \left| 2 \sin\left(\frac{\boldsymbol{\theta}_{1,k} - \boldsymbol{\theta}_{2,k}}{4}\right) \right|. \end{aligned}$$

Plugging it into equation (3), we further get

$$\begin{aligned} \|U(\boldsymbol{\theta}_1) - U(\boldsymbol{\theta}_2)\| &\leq \sum_{k=1}^K \left| 2 \sin\left(\frac{\boldsymbol{\theta}_{1,k} - \boldsymbol{\theta}_{2,k}}{4}\right) \right| \\ &\leq \sum_{k=1}^K \left| \frac{\boldsymbol{\theta}_{1,k} - \boldsymbol{\theta}_{2,k}}{2} \right| \\ &\leq \frac{\sqrt{K}}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|. \end{aligned}$$

This completes the proof of Lemma B.3. \square

To examine the effects of depolarizing noise from the perspective of QNN generalization, we recall the following lemma from [55].

Lemma B.4 (see [55], Lemma 6). *Let \mathcal{E}_p be the depolarization channel. There always exists a depolarization channel $\mathcal{E}_{\tilde{p}}$ with $\tilde{p} = 1 - (1 - p)^{K_g}$ that satisfies*

$$\mathcal{E}_p \left\{ U_{K_g}(\boldsymbol{\theta}) \dots U_2(\boldsymbol{\theta}) \mathcal{E}_p \left[U_1(\boldsymbol{\theta}) \rho U_1^\dagger(\boldsymbol{\theta}) \right] U_2^\dagger(\boldsymbol{\theta}) \dots U_{K_g}^\dagger(\boldsymbol{\theta}) \right\} = \mathcal{E}_{\tilde{p}} \left[U(\boldsymbol{\theta}) \rho U^\dagger(\boldsymbol{\theta}) \right],$$

where ρ is the input quantum state.

C Proofs of Main Results

In this section, we provide the proofs of main results in our paper. We require several useful lemmas to prove the main results.

Lemma C.1 (From Loss Stability to Parameter Stability). *Let $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}'_t$ be the parameters of QNNs learned using the SGD algorithm for t iterations on training datasets S and S' , respectively. Then, the output difference of the QNNs is bounded by,*

$$\left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \leq \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2.$$

Proof. The difference between the two output functions of the QNNs can be represented as follows

$$\begin{aligned} \left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| &= \left| \text{Tr} \left(O U(\boldsymbol{\theta}_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}_t) \right) - \text{Tr} \left(O U(\boldsymbol{\theta}'_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}'_t) \right) \right| \\ &= \left| \text{Tr} \left(O \left(U(\boldsymbol{\theta}_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}_t) - U(\boldsymbol{\theta}'_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}'_t) \right) \right) \right| \\ &\leq \|O\| \left\| U(\boldsymbol{\theta}_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}_t) - U(\boldsymbol{\theta}'_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}'_t) \right\|_1, \end{aligned}$$

where the last inequality uses the Cauchy-Schwartz inequality.

Let $\mathcal{E}(\rho) = U(\boldsymbol{\theta}_t) \rho U^\dagger(\boldsymbol{\theta}_t)$, $\mathcal{E}'(\rho) = U(\boldsymbol{\theta}'_t) \rho U^\dagger(\boldsymbol{\theta}'_t)$ be unitary channels. By Lemma B.1, we have

$$\begin{aligned} \left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| &\leq \|O\| \|\mathcal{E}(\rho(\mathbf{x})) - \mathcal{E}'(\rho(\mathbf{x}))\|_1 \\ &\leq \|O\| \|\mathcal{E} - \mathcal{E}'\|_\diamond \\ &\leq 2\|O\| \|U(\boldsymbol{\theta}_t) - U(\boldsymbol{\theta}'_t)\|. \end{aligned}$$

Note that any $U(\boldsymbol{\theta})$ is of the form $U(\boldsymbol{\theta}) = V_1 U_1 V_2 U_2 V_3 \dots V_K U_K V_{K+1}$, where $U_k, 1 \leq k \leq K$, are a particular choice of the trainable single-qubit Pauli rotations and $V_k, 1 \leq k \leq K + 1$, are the non-trainable n -qubit unitaries. (For ease of readability, we have not written out the tensor factors of identities accompanying the U_k .) We can write $U(\boldsymbol{\theta}_t) = V_1 e^{-i\boldsymbol{\theta}_{t,1} P_1/2} V_2 e^{-i\boldsymbol{\theta}_{t,2} P_2/2} \dots V_K e^{-i\boldsymbol{\theta}_{t,K} P_K/2} V_{K+1}$,

$U(\boldsymbol{\theta}'_t) = V_1 e^{-i\theta'_{t,1} P_1/2} V_2 e^{-i\theta'_{t,2} P_2/2} \dots V_K e^{-i\theta'_{t,K} P_K/2} V_{K+1}$, where $P_k \in \{X, Y, Z\}$ denotes a single-qubit Pauli gate. By Lemma B.3, we further get

$$\begin{aligned} \left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| &\leq 2\|O\| \sum_{k=1}^K \left| 2 \sin \left(\frac{\theta_{t,k} - \theta'_{t,k}}{4} \right) \right| \\ &\leq 2\|O\| \sum_{k=1}^K \left| \frac{\theta_{t,k} - \theta'_{t,k}}{2} \right| \\ &\leq \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2. \end{aligned} \quad (4)$$

This completes the proof of Lemma C.1. \square

Lemma C.2 (QNN Same Sample Loss Stability Bound). *Suppose that Assumption 3.6 and 3.7 hold. Let $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}'_t$ be the parameters of two QNNs learned using the SGD algorithm for t iterations on two training datasets S and S' , respectively. Then, the loss derivative difference of the QNNs with respect to the same sample is bounded by,*

$$\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y)\|_2 \leq \kappa \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2,$$

where $\kappa = \alpha_\ell K \|O\| + \sqrt{2} \nu_\ell K \|O\|^2$.

Proof. Using the Assumption 3.6 and 3.7 that the loss function is Lipschitz continuous and smoothness, we have

$$\begin{aligned} &\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \right\|_2 \\ &= \left\| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2 \\ &\leq \left\| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) (\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x})) \right\|_2 + \left\| \left(\frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \right) \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2 \\ &= \left| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) \right| \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2 + \left| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \right| \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2 \\ &\leq \alpha_\ell \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2 + \nu_\ell \left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2. \end{aligned} \quad (5)$$

The term $\left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2$ can be bounded by the individual terms $\left| \nabla_{\theta_j} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|$, which can be computed using parameter-shift rules [65],

$$\begin{aligned} &\left| \nabla_{\theta_j} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \\ &= \frac{1}{2} \left| \left(f_{\boldsymbol{\theta}_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right) - \left(f_{\boldsymbol{\theta}'_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right) \right| \\ &= \frac{1}{2} \left| \left(f_{\boldsymbol{\theta}_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right) - \left(f_{\boldsymbol{\theta}_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right) \right| \\ &\leq \frac{1}{2} \left[\left| f_{\boldsymbol{\theta}_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right| + \left| f_{\boldsymbol{\theta}_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right| \right], \end{aligned} \quad (6)$$

where \mathbf{e}_j is the unit vector along the $\boldsymbol{\theta}_j$ axis.

According to Lemma C.1, we have

$$\left| f_{\boldsymbol{\theta}_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right| \leq \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2,$$

and

$$\left| f_{\boldsymbol{\theta}_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right| \leq \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2.$$

Plugging it into equation (6), we further get

$$\left| \nabla_{\theta_j} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \leq \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2.$$

Therefore, we can bound the term $\|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x})\|_2$ as follows

$$\begin{aligned} \|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x})\|_2 &= \sqrt{\sum_{j=1}^K \left| \nabla_{\theta_j} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|^2} \\ &\leq \sqrt{K \left(\sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2 \right)^2} \\ &= K \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2. \end{aligned} \tag{7}$$

Similarly, the term $\left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2$ can be bounded by the individual terms $\left\| \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2$. Additionally, by equation (4) in Lemma C.1, we have

$$\begin{aligned} \left| \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| &= \frac{1}{2} \left| f_{\boldsymbol{\theta}'_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right| \\ &\leq \|O\| \left(\sum_{k \neq j} \left| 2 \sin \left(\frac{\theta'_{t,k} - \theta'_{t,k}}{4} \right) \right| + \left| 2 \sin \left(\frac{(\theta'_{t,j} + \frac{\pi}{2}) - (\theta'_{t,j} - \frac{\pi}{2})}{4} \right) \right| \right) \\ &= \sqrt{2} \|O\|. \end{aligned} \tag{8}$$

Therefore, we can bound the term $\left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2$ as follows

$$\begin{aligned} \|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x})\|_2 &= \sqrt{\sum_{j=1}^K \left| \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|^2} \\ &\leq \sqrt{K \left(\sqrt{2} \|O\| \right)^2} \\ &= \sqrt{2K} \|O\|. \end{aligned} \tag{9}$$

Finally, according to Lemma C.1, the term $\left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|$ can be bounded as

$$\left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \leq \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2 \tag{10}$$

Plugging equation (7), (9), and (10) back into equation (5) completes the proof of Lemma C.2. \square

Lemma C.3 (QNN Differnet Sample Loss Stability Bound). *Suppose that Assumption 3.6 hold. Let $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}'_t$ be the parameters of two QNNs learned using the SGD algorithm for t iterations on two training datasets S and S' , respectively. Then, the loss derivative difference of the QNNs with respect to the different sample is bounded by,*

$$\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}'); y')\|_2 \leq 2\sqrt{2}\alpha_\ell \sqrt{K} \|O\|.$$

Proof. Using the Assumption 3.6 that the loss function is Lipschitz continuous, we have

$$\begin{aligned} & \left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}'); y') \right\|_2 \\ &= \left\| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}'_t}(\mathbf{x}'); y') \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}') \right\|_2 \\ &\leq \left\| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) \right\|_2 + \left\| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}'_t}(\mathbf{x}'); y') \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}') \right\|_2 \\ &= \left| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) \right| \|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x})\|_2 + \left| \frac{\partial \ell}{\partial f}(f_{\boldsymbol{\theta}'_t}(\mathbf{x}'); y') \right| \|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}')\|_2 \\ &\leq \alpha_\ell \left(\|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x})\|_2 + \|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}')\|_2 \right). \end{aligned} \tag{11}$$

By equation (9) in Lemma C.2, we similarly bound $\|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x})\|_2 \leq \sqrt{2K} \|O\|$, $\|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}')\|_2 \leq \sqrt{2K} \|O\|$. Plugging this into equation (11) completes the proof of Lemma C.3. \square

C.1 Proof of Theorem 4.1

In this subsection, we prove Theorem 4.1 on the uniform stability of QNNs. We begin by quoting the result which we set to prove.

Theorem 4.1 (Uniform Stability Bound). *Suppose that Assumption 3.6 and 3.7 hold. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t for T iterations, then $A(S)$ is ϵ -uniformly stable with*

$$\epsilon \leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\sqrt{2}\eta_t \alpha_\ell^2 K \|O\|^2}{m},$$

where

$$\kappa := \alpha_\ell K \|O\| + \sqrt{2}\nu_\ell K \|O\|^2.$$

Proof. Let S and S' be two datasets of size m differing in only a single sample. Consider two sequences of the parameters, $\{\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T\}$ and $\{\boldsymbol{\theta}'_0, \boldsymbol{\theta}'_1, \dots, \boldsymbol{\theta}'_T\}$, learned by the QNN running SGD on S and S' , respectively. Let $\delta_t = \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2$.

Using the Assumption 3.6 that the loss function is Lipschitz continuous, the linearity of expectation and Lemma C.1, we have

$$\begin{aligned} \left| \mathbb{E}_A \left[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y) \right] \right| &\leq \alpha_\ell \mathbb{E}_A \left[\left\| f_{\boldsymbol{\theta}_T}(\mathbf{x}) - f_{\boldsymbol{\theta}'_T}(\mathbf{x}) \right\| \right] \\ &\leq \alpha_\ell \sqrt{K} \|O\| \mathbb{E}_A \left[\|\boldsymbol{\theta}_T - \boldsymbol{\theta}'_T\|_2 \right] \\ &\leq \alpha_\ell \sqrt{K} \|O\| \mathbb{E}_A [\delta_T]. \end{aligned} \tag{12}$$

Then, we focus on the term $\mathbb{E}_A[\delta_T]$. Observe that at iteration t , with probability $1 - 1/m$, the example selected by SGD is the same in both S and S' . With probability $1/m$ the selected example is different. Therefore, we have

$$\begin{aligned}
\mathbb{E}_A[\delta_{t+1}] &\leq \left(1 - \frac{1}{m}\right) \mathbb{E}_A \left[\left\| \left(\boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) \right) - \left(\boldsymbol{\theta}'_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \right) \right\|_2 \right] \\
&\quad + \frac{1}{m} \mathbb{E}_A \left[\left\| \left(\boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}'); y') \right) - \left(\boldsymbol{\theta}'_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}''); y'') \right) \right\|_2 \right] \\
&= \mathbb{E}_A[\delta_t] + \left(1 - \frac{1}{m}\right) \eta_t \mathbb{E}_A \left[\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \right\|_2 \right] \\
&\quad + \frac{1}{m} \eta_t \mathbb{E}_A \left[\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}'); y') - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}''); y'') \right\|_2 \right].
\end{aligned} \tag{13}$$

According to Lemma C.2 and Lemma C.3, we further get

$$\begin{aligned}
\mathbb{E}_A[\delta_{t+1}] &\leq \mathbb{E}_A[\delta_t] + \left(1 - \frac{1}{m}\right) \eta_t \cdot \mathbb{E}_A[\kappa \delta_t] + \frac{1}{m} \eta_t \cdot \mathbb{E}_A[2\sqrt{2}\alpha_\ell \sqrt{K} \|O\|] \\
&\leq (1 + \eta_t \kappa) \mathbb{E}_A[\delta_t] + \frac{2\sqrt{2}\eta_t \alpha_\ell \sqrt{K} \|O\|}{m}.
\end{aligned}$$

Unraveling the recursion gives

$$\mathbb{E}_A[\delta_T] \leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\sqrt{2}\eta_t \alpha_\ell \sqrt{K} \|O\|}{m}.$$

Plugging it into equation (12), we obtain

$$\left| \mathbb{E}_A \left[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y) \right] \right| \leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\sqrt{2}\eta_t \alpha_\ell^2 K \|O\|^2}{m}.$$

By the definition of uniform stability as shown in Definition 3.1, we obtain the desired bound on the uniform stability of QNNs. This completes the proof of Theorem 4.1. \square

C.2 Proof of Corollary 4.3

In this subsection, we prove Corollary 4.3, which provides simplified uniform stability results for two commonly used step sizes of QNNs. We first introduce an extension of Lemma 3.11 in [34], specifically adapted to the unique context of quantum neural networks. The following lemma is motivated by the fact that SGD typically runs several iterations before encountering the different example between S and S' .

Lemma C.4. *Suppose that Assumption 3.6 and 3.7 hold, and $\ell(\cdot, \cdot) \in [0, M]$. Let S and S' of size m differing in only a single example. Consider two sequences of parameters, $\{\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T\}$ and $\{\boldsymbol{\theta}'_0, \boldsymbol{\theta}'_1, \dots, \boldsymbol{\theta}'_T\}$, learned by the QNN running SGD on S and S' , respectively. Let $\delta_t = \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2$. Then, for any $\mathbf{z} \in \mathcal{Z}$ and $t_0 \in \{0, 1, \dots, m\}$, we have*

$$\left| \mathbb{E}_A[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y)] \right| \leq \alpha_\ell \sqrt{K} \|O\| \mathbb{E}_A[\delta_T \mid \delta_{t_0} = 0] + \frac{t_0 M}{m}.$$

Proof. Let \mathcal{E} denote the event that $\delta_{t_0} = 0$. Then we have

$$\begin{aligned}
& |\mathbb{E}[\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)]| \\
& \leq \mathbb{E}[|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)|] \\
& = \mathbb{P}[\mathcal{E}] \mathbb{E}\left[|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)| \mid \mathcal{E}\right] \\
& \quad + \mathbb{P}[\mathcal{E}^c] \cdot \mathbb{E}\left[|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)| \mid \mathcal{E}^c\right] \\
& \leq \mathbb{E}\left[|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)| \mid \mathcal{E}\right] \\
& \quad + \mathbb{P}[\mathcal{E}^c] \cdot \sup\left|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)\right|.
\end{aligned} \tag{14}$$

Using the Assumption 3.6 that the loss function is Lipschitz continuous and Lemma C.1, the first term $\mathbb{E}\left[|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)| \mid \mathcal{E}\right]$ can be bounded as

$$\mathbb{E}\left[|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)| \mid \mathcal{E}\right] \leq \alpha_\ell \sqrt{K} \|O\| \mathbb{E}[\delta_T \mid \mathcal{E}]. \tag{15}$$

It remains to bound the second term $\mathbb{P}[\mathcal{E}^c] \cdot \sup\left|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)\right|$. Let i^* be the position where S and S' are different and denote the first time SGD uses the example \mathbf{z}_{i^*} by the random variable I . Note that when $I > t_0$, then we must have that $\delta_{t_0} = 0$, since the execution on S and S' is identical until iteration t_0 . We have that

$$\mathbb{P}[\mathcal{E}^c] = \mathbb{P}[\delta_{t_0} \neq 0] \leq \mathbb{P}[I \leq t_0] \leq \frac{t_0}{m}$$

According the condition $\ell(\cdot, \cdot) \in [0, M]$, the second term can be bounded as

$$\mathbb{P}[\mathcal{E}^c] \cdot \sup\left|\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)\right| \leq \frac{t_0 M}{m}. \tag{16}$$

Plugging equation (15) and (16) back into equation (14) completes the proof of Lemma C.4 \square

We are now ready to prove Corollary 4.3.

Corollary 4.3. *Suppose that Assumption 3.6 and 3.7 hold, and $\ell(\cdot, \cdot) \in [0, M]$. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t for T iterations.*

(a) *If we choose the constant step sizes $\eta_t = \eta$, then $A(S)$ is ϵ -uniformly stable with*

$$\epsilon \leq \frac{2\sqrt{2}\alpha_\ell^2 K \|O\|^2}{\kappa m} (1 + \eta\kappa)^T.$$

(b) *If we choose the monotonically non-increasing step sizes $\eta_t \leq c/(t+1)$, $c > 0$, then $A(S)$ is ϵ -uniformly stable with*

$$\epsilon \leq \frac{1 + 1/c\kappa}{m} M^{\frac{c\kappa}{c\kappa+1}} \left(2\sqrt{2}c\alpha_\ell^2 K \|O\|^2\right)^{\frac{1}{c\kappa+1}} T^{\frac{c\kappa}{c\kappa+1}},$$

where κ is defined by equation (1).

Proof. Let S and S' be two datasets of size m differing in only a single sample. Consider two sequences of the parameters, $\{\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T\}$ and $\{\boldsymbol{\theta}'_0, \boldsymbol{\theta}'_1, \dots, \boldsymbol{\theta}'_T\}$, learned by the QNN running SGD on S and S' , respectively. Let $\delta_t = \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2$.

For the constant step sizes $\eta_t = \eta$, by Theorem 4.1, we have

$$\begin{aligned} \epsilon &\leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta\kappa) \right] \frac{2\sqrt{2}\eta\alpha_\ell^2 K \|O\|^2}{m} \\ &\leq \frac{2\sqrt{2}\eta\alpha_\ell^2 K \|O\|^2}{m} \sum_{t=0}^{T-1} (1 + \eta\kappa)^t \\ &\leq \frac{2\sqrt{2}\eta\alpha_\ell^2 K \|O\|^2}{m} \cdot \frac{(1 + \eta\kappa)^T - 1}{\eta\kappa} \\ &\leq \frac{2\sqrt{2}\alpha_\ell^2 K \|O\|^2}{\kappa m} (1 + \eta\kappa)^T. \end{aligned}$$

We immediately get the claimed upper bound on the uniform stability under the constant step size setting. This completes the proof of Corollary 4.3, Part(a).

For the decaying step sizes $\eta_t \leq c/(t+1)$, by Lemma C.4, we have for every $t_0 \in \{0, 1, \dots, m\}$,

$$\left| \mathbb{E}_A[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y)] \right| \leq \alpha_\ell \sqrt{K} \|O\| \mathbb{E}_A[\delta_T \mid \delta_{t_0} = 0] + \frac{t_0 M}{m}. \quad (17)$$

Let $\Delta_t = \mathbb{E}[\delta_t \mid \delta_{t_0} = 0]$. Observe that at iteration t , with probability $1 - 1/m$, the example selected by SGD is the same in both S and S' . With probability $1/m$ the selected example is different. Therefore, we have

$$\begin{aligned} \Delta_{t+1} &\leq \left(1 - \frac{1}{m}\right) \mathbb{E}_A \left[\left\| (\boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y)) - (\boldsymbol{\theta}'_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y)) \right\|_2 \mid \delta_{t_0} = 0 \right] \\ &\quad + \frac{1}{m} \mathbb{E}_A \left[\left\| (\boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}'); y')) - (\boldsymbol{\theta}'_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}''; y'')) \right\|_2 \mid \delta_{t_0} = 0 \right] \\ &= \Delta_t + \left(1 - \frac{1}{m}\right) \eta_t \mathbb{E}_A \left[\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \right\|_2 \mid \delta_{t_0} = 0 \right] \\ &\quad + \frac{1}{m} \eta_t \mathbb{E}_A \left[\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}'); y') - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}''; y'')) \right\|_2 \mid \delta_{t_0} = 0 \right]. \end{aligned} \quad (18)$$

According to Lemma C.2 and Lemma C.3, we further get

$$\begin{aligned} \Delta_{t+1} &\leq \Delta_t + \left(1 - \frac{1}{m}\right) \eta_t \cdot \mathbb{E}_A[\kappa \delta_t \mid \delta_{t_0} = 0] + \frac{1}{m} \eta_t \cdot \mathbb{E}_A[2\sqrt{2}\alpha_\ell \sqrt{K} \|O\| \mid \delta_{t_0} = 0] \\ &\leq (1 + \eta_t \kappa) \Delta_t + \frac{2\sqrt{2}\eta_t \alpha_\ell \sqrt{K} \|O\|}{m}. \end{aligned}$$

Using the fact that $\Delta_{t_0} = 0$, we can unwind this recurrence relation from T down to $t_0 + 1$. This gives

$$\Delta_{t+1} = \sum_{t=t_0+1}^T \left[\prod_{j=t+1}^T (1 + \eta_j \kappa) \right] \frac{2\sqrt{2}\eta_t \alpha_\ell \sqrt{K} \|O\|}{m}.$$

By the elementary inequality $1 + a \leq \exp(a)$ and $\eta_t \leq c/(t + 1)$, we further derive

$$\begin{aligned}
\Delta_{t+1} &\leq \sum_{t=t_0+1}^T \left[\prod_{j=t+1}^T \exp\left(\frac{c\kappa}{j}\right) \right] \frac{2\sqrt{2}c\alpha_\ell\sqrt{K}\|O\|}{tm} \\
&\leq \sum_{t=t_0+1}^T \exp\left(\sum_{j=t+1}^T \frac{c\kappa}{j}\right) \frac{2\sqrt{2}c\alpha_\ell\sqrt{K}\|O\|}{tm} \\
&\leq \sum_{t=t_0+1}^T \exp\left(c\kappa \log\left(\frac{T}{t}\right)\right) \frac{2\sqrt{2}c\alpha_\ell\sqrt{K}\|O\|}{tm} \\
&\leq \frac{2\sqrt{2}c\alpha_\ell\sqrt{K}\|O\|}{m} T^{c\kappa} \sum_{t=t_0+1}^T \frac{1}{t^{c\kappa+1}} \\
&\leq \frac{2\sqrt{2}c\alpha_\ell\sqrt{K}\|O\|}{\kappa m} \left(\frac{T}{t_0}\right)^{c\kappa}.
\end{aligned}$$

Plugging this bound into equation (17), we get

$$\left| \mathbb{E}_A[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y)] \right| \leq \frac{2\sqrt{2}c\alpha_\ell^2 K \|O\|^2}{\kappa m} \left(\frac{T}{t_0}\right)^{c\kappa} + \frac{t_0 M}{m}. \quad (19)$$

The right hand side is approximately minimized when

$$t_0 = \left(\frac{2\sqrt{2}c\alpha_\ell^2 K \|O\|^2}{M} \right)^{\frac{1}{c\kappa+1}} T^{\frac{c\kappa}{c\kappa+1}}.$$

Plugging it into equation (19) we have (for simplicity we assume the above t_0 is an integer)

$$\left| \mathbb{E}_A[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y)] \right| \leq \frac{1 + \frac{1}{c\kappa}}{m} M^{\frac{c\kappa}{c\kappa+1}} \left(2\sqrt{2}c\alpha_\ell^2 K \|O\|^2\right)^{\frac{1}{c\kappa+1}} T^{\frac{c\kappa}{c\kappa+1}}.$$

By the definition of uniform stability as shown in Definition 3.1, we obtain the desired bound on the uniform stability under the decaying step size setting. This completes the proof of Corollary 4.3, Part(b). \square

C.3 Proof of Corollary 4.8

In this subsection, we prove Corollary 4.8, which extends previous analysis to noise scenario. We first introduce several useful lemmas, as extensions of Lemma C.1, C.2, C.3, and C.4, under the depolarizing noise setting.

Lemma C.5 (From Loss Stability to Parameter Stability, Depolarizing Noise). *Let $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}'_t$ be the parameters of QNNs learned using the SGD algorithm for t iterations on training datasets S and S' , under depolarizing noise level $p \in [0, 1]$, respectively. Then, the output difference of the QNNs is bounded by,*

$$\left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \leq (1 - p)^{K_g} \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2.$$

Proof. We consider the noisy quantum channel is simulated by the local depolarization noise, i.e., the depolarization channel $\mathcal{E}_p(\cdot)$ is applied to each quantum gate in $U(\boldsymbol{\theta})$. By Lemma B.4, the difference between the two output functions of the QNNs can be represented as follows,

$$\begin{aligned} \left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| &= \left| \text{Tr} \left(O \mathcal{E}_p \left(U(\boldsymbol{\theta}_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}_t) \right) \right) - \text{Tr} \left(O \mathcal{E}_p \left(U(\boldsymbol{\theta}'_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}'_t) \right) \right) \right| \\ &= (1-p)^{K_g} \left| \text{Tr} \left(O U(\boldsymbol{\theta}_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}_t) \right) - \text{Tr} \left(O U(\boldsymbol{\theta}'_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}'_t) \right) \right| \\ &= (1-p)^{K_g} \left| \text{Tr} \left(O \left(U(\boldsymbol{\theta}_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}_t) - U(\boldsymbol{\theta}'_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}'_t) \right) \right) \right| \\ &\leq (1-p)^{K_g} \|O\| \left\| U(\boldsymbol{\theta}_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}_t) - U(\boldsymbol{\theta}'_t) \rho(\mathbf{x}) U^\dagger(\boldsymbol{\theta}'_t) \right\|_1, \end{aligned}$$

where the last inequality uses the Cauchy-Schwartz inequality.

Next, by combining the proof technique used in Lemma C.1, we have

$$\begin{aligned} \left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| &\leq 2(1-p)^{K_g} \|O\| \sum_{k=1}^K \left| 2 \sin \left(\frac{\theta_{t,k} - \theta'_{t,k}}{4} \right) \right| \\ &\leq 2(1-p)^{K_g} \|O\| \sum_{k=1}^K \left| \frac{\theta_{t,k} - \theta'_{t,k}}{2} \right| \\ &\leq (1-p)^{K_g} \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2. \end{aligned}$$

This completes the proof of Lemma C.5. \square

Lemma C.6 (QNN Same Sample Loss Stability Bound, Depolarizing Noise). *Suppose that Assumption 3.6 and 3.7 hold. Let $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}'_t$ be the parameters of two QNNs learned using the SGD algorithm for t iterations on two training datasets S and S' , under depolarizing noise level $p \in [0, 1]$, respectively. Then, the loss derivative difference of the QNNs with respect to the same sample is bounded by,*

$$\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y)\|_2 \leq (1-p)^{K_g} \kappa \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2,$$

where $\kappa = \alpha_\ell K \|O\| + \sqrt{2} \nu_\ell K \|O\|^2$.

Proof. According to equation (5) in Lemma C.2, we have

$$\begin{aligned} &\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \right\|_2 \\ &\leq \alpha_\ell \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2 + \nu_\ell \left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2. \end{aligned} \tag{20}$$

Combing the equation (6) in Lemma C.2 and Lemma C.5, we can bound the individual terms $\left| \nabla_{\theta_j} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|$ as follows

$$\begin{aligned} &\left| \nabla_{\theta_j} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \\ &\leq \frac{1}{2} \left[\left| f_{\boldsymbol{\theta}_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t + \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right| + \left| f_{\boldsymbol{\theta}_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t - \frac{\pi}{2} \mathbf{e}_j}(\mathbf{x}) \right| \right] \\ &\leq (1-p)^{K_g} \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2. \end{aligned} \tag{21}$$

Therefore, the term $\left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2$ can be bounded as

$$\begin{aligned} \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2 &= \sqrt{\sum_{j=1}^K \left| \nabla_{\theta_j} f_{\boldsymbol{\theta}_t}(\mathbf{x}) - \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|^2} \\ &\leq \sqrt{K \left((1-p)^{K_g} \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2 \right)^2} \\ &= (1-p)^{K_g} K \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2. \end{aligned} \quad (22)$$

Combing the equation (8) in Lemma C.2 and Lemma C.5, we can similarly bound the individual terms $\left| \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|$ as follows

$$\begin{aligned} \left| \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| &\leq (1-p)^{K_g} \|O\| \left(\sum_{k \neq j} \left| 2 \sin \left(\frac{\theta'_{t,k} - \theta'_{t,k}}{4} \right) \right| + \left| 2 \sin \left(\frac{(\theta'_{t,j} + \frac{\pi}{2}) - (\theta'_{t,j} - \frac{\pi}{2})}{4} \right) \right| \right) \\ &= \sqrt{2} (1-p)^{K_g} \|O\|. \end{aligned} \quad (23)$$

Therefore, the term $\left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2$ can be bounded as

$$\begin{aligned} \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right\|_2 &= \sqrt{\sum_{j=1}^K \left| \nabla_{\theta_j} f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|^2} \\ &\leq \sqrt{K \left(\sqrt{2} (1-p)^{K_g} \|O\| \right)^2} \\ &= \sqrt{2K} (1-p)^{K_g} \|O\|. \end{aligned} \quad (24)$$

Finally, by Lemma C.5, the term $\left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right|$ can be bounded as

$$\left| f_{\boldsymbol{\theta}_t}(\mathbf{x}) - f_{\boldsymbol{\theta}'_t}(\mathbf{x}) \right| \leq (1-p)^{K_g} \sqrt{K} \|O\| \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2 \quad (25)$$

Plugging equation (22), (24), and (25) back into equation (5), we further get

$$\begin{aligned} \left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}); y) \right\|_2 &\leq (1-p)^{K_g} \alpha_{\ell} K \|O\| + (1-p)^{2K_g} \sqrt{2} \nu_{\ell} K \|O\|^2 \\ &\leq (1-p)^{K_g} \kappa, \end{aligned}$$

where $\kappa = \alpha_{\ell} K \|O\| + \sqrt{2} \nu_{\ell} K \|O\|^2$.

This completes the proof of Lemma C.6. \square

Lemma C.7 (QNN Diffnet Sample Loss Stability Bound, Depolarizing Noise). *Suppose that Assumption 3.6 hold. Let $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}'_t$ be the parameters of two QNNs learned using the SGD algorithm for t iterations on two training datasets S and S' , under depolarizing noise level $p \in [0, 1]$, respectively. Then, the loss derivative difference of the QNNs with respect to the different sample is bounded by,*

$$\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}'); y') \right\|_2 \leq 2\sqrt{2} (1-p)^{K_g} \alpha_{\ell} \sqrt{K} \|O\|.$$

Proof. According to equation (11) in Lemma C.3, we have

$$\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}'_t}(\mathbf{x}'); y') \right\|_2 \leq \alpha_\ell \left(\|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x})\|_2 + \left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}') \right\|_2 \right). \quad (26)$$

By equation (24) in Lemma C.6, we similarly bound $\|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}(\mathbf{x})\|_2 \leq \sqrt{2K}(1-p)^{K_g} \|O\|$, $\left\| \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}'_t}(\mathbf{x}') \right\|_2 \leq \sqrt{2K}(1-p)^{K_g} \|O\|$. Plugging this into equation (26) completes the proof of Lemma C.3. \square

Lemma C.8. *Suppose that Assumption 3.6 and 3.7 hold, and $\ell(\cdot, \cdot) \in [0, M]$. Let S and S' of size m differing in only a single example. Consider two sequences of parameters, $\{\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T\}$ and $\{\boldsymbol{\theta}'_0, \boldsymbol{\theta}'_1, \dots, \boldsymbol{\theta}'_T\}$, learned by the QNN running SGD on S and S' , under depolarizing noise level $p \in [0, 1]$, respectively. Let $\delta_t = \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2$. Then, for any $\mathbf{z} \in \mathcal{Z}$ and $t_0 \in \{0, 1, \dots, m\}$, we have*

$$\left| \mathbb{E}_A[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y)] \right| \leq (1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\| \mathbb{E}_A[\delta_T \mid \delta_{t_0} = 0] + \frac{t_0 M}{m}.$$

Proof. Let \mathcal{E} denote the event that $\delta_{t_0} = 0$. By equation (14) and (16) in Lemma C.4, we have

$$\left| \mathbb{E}[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y)] \right| \leq \mathbb{E} \left[\left| \ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y) \right| \mid \mathcal{E} \right] + \frac{t_0 M}{m}. \quad (27)$$

Using the Assumption 3.6 that the loss function is Lipschitz continuous and Lemma C.5, the first term $\mathbb{E} \left[\left| \ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y) \right| \mid \mathcal{E} \right]$ can be bounded as

$$\mathbb{E} \left[\left| \ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y) \right| \mid \mathcal{E} \right] \leq (1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\| \mathbb{E}[\delta_T \mid \mathcal{E}]. \quad (28)$$

Plugging it into equation (27) completes the proof of Lemma C.8. \square

We now are ready to prove Corollary 4.8.

Corollary 4.8 (Generalization Bound Under Depolarizing Noise). *Suppose that Assumption 3.6 and 3.7 hold, and $\ell(\cdot, \cdot) \in [0, M]$. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t under depolarizing noise level $p \in [0, 1]$ for T iterations.*

- (a) *if we choose the constant step sizes $\eta_t = \eta$, then the following generalization bound of $A(S)$ holds with probability at least $1 - \delta$ for $\delta \in (0, 1)$,*

$$\begin{aligned} & \mathbb{E}_A [R_{\mathcal{D}}(A(S)) - R_S(A(S))] \\ & \leq \mathcal{O} \left(\frac{\left(1 + (1-p)^{K_g} \eta \kappa\right)^T}{m} \log m \log\left(\frac{1}{\delta}\right) + M \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{m}} \right), \end{aligned}$$

- (b) *if we choose the monotonically non-increasing step sizes $\eta_t \leq c/(t+1)$, $c > 0$, then the following generalization bound of $A(S)$ holds with probability at least $1 - \delta$ for $\delta \in (0, 1)$,*

$$\begin{aligned} & \mathbb{E}_A [R_{\mathcal{D}}(A(S)) - R_S(A(S))] \\ & \leq \mathcal{O} \left(\frac{T^{\frac{c\kappa}{c\kappa+1/(1-p)^{K_g}}}}{m} \log m \log\left(\frac{1}{\delta}\right) + M \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{m}} \right), \end{aligned}$$

where κ is defined by equation (1).

Proof. Let S and S' be two datasets of size m differing in only a single sample. Consider two sequences of the parameters, $\{\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T\}$ and $\{\boldsymbol{\theta}'_0, \boldsymbol{\theta}'_1, \dots, \boldsymbol{\theta}'_T\}$, learned by the QNN running SGD on S and S' , respectively. Let $\delta_t = \|\boldsymbol{\theta}_t - \boldsymbol{\theta}'_t\|_2$.

For the constant step sizes $\eta_t = \eta$, using the Assumption 3.6 that the loss function is Lipschitz continuous, the linearity of expectation and Lemma C.5, we have

$$\begin{aligned} \left| \mathbb{E}_A \left[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y) \right] \right| &\leq \alpha_\ell \mathbb{E}_A \left[\|f_{\boldsymbol{\theta}_T}(\mathbf{x}) - f_{\boldsymbol{\theta}'_T}(\mathbf{x})\| \right] \\ &\leq (1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\| \mathbb{E}_A [\|\boldsymbol{\theta}_T - \boldsymbol{\theta}'_T\|_2] \\ &\leq (1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\| \mathbb{E}_A [\delta_T]. \end{aligned} \quad (29)$$

Combining equation (13), Lemma C.6, and Lemma C.7, we further get

$$\begin{aligned} \mathbb{E}_A[\delta_{t+1}] &\leq \mathbb{E}_A[\delta_t] + \left(1 - \frac{1}{m}\right) \eta \cdot \mathbb{E}_A[(1-p)^{K_g} \kappa \delta_t] + \frac{1}{m} \eta \cdot \mathbb{E}_A[2\sqrt{2}(1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\|] \\ &\leq (1 + (1-p)^{K_g} \eta \kappa) \mathbb{E}_A[\delta_t] + \frac{2\sqrt{2}(1-p)^{K_g} \eta \alpha_\ell \sqrt{K} \|O\|}{m}. \end{aligned}$$

Unraveling the recursion gives

$$\begin{aligned} \mathbb{E}_A[\delta_T] &\leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + (1-p)^{K_g} \eta \kappa) \right] \frac{2\sqrt{2}(1-p)^{K_g} \eta \alpha_\ell \sqrt{K} \|O\|}{m} \\ &\leq \frac{2\sqrt{2}(1-p)^{K_g} \eta \alpha_\ell \sqrt{K} \|O\|}{m} \sum_{t=0}^{T-1} (1 + (1-p)^{K_g} \eta \kappa)^t \\ &\leq \frac{2\sqrt{2}(1-p)^{K_g} \eta \alpha_\ell \sqrt{K} \|O\|}{m} \cdot \frac{(1 + (1-p)^{K_g} \eta \kappa)^T - 1}{(1-p)^{K_g} \eta \kappa} \\ &\leq \frac{2\sqrt{2} \alpha_\ell \sqrt{K} \|O\|}{\kappa m} (1 + (1-p)^{K_g} \eta \kappa)^T. \end{aligned}$$

Plugging it into equation (29), we obtain

$$\left| \mathbb{E}_A \left[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y) \right] \right| \leq \frac{2\sqrt{2} \alpha_\ell^2 K \|O\|^2}{\kappa m} (1 + (1-p)^{K_g} \eta \kappa)^T.$$

By the definition of uniform stability as shown in Definition 3.1, we obtain the desired bound on the uniform stability of QNNs under the constant step size and depolarizing noise setting. Combining Lemma 3.3, Part (b) completes the proof of Corollary 4.8, Part (a).

For the decaying step sizes $\eta_t \leq c/(t+1)$, by Lemma C.8, we have for every $t_0 \in \{0, 1, \dots, m\}$,

$$\left| \mathbb{E}_A[\ell(f_{\boldsymbol{\theta}_T}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}'_T}(\mathbf{x}); y)] \right| \leq (1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\| \mathbb{E}_A[\delta_T \mid \delta_{t_0} = 0] + \frac{t_0 M}{m}. \quad (30)$$

Let $\Delta_t = \mathbb{E}[\delta_t \mid \delta_{t_0} = 0]$. Combining equation (18), Lemma C.6, and Lemma C.7, we further get

$$\begin{aligned} \Delta_{t+1} &\leq \Delta_t + \left(1 - \frac{1}{m}\right) \eta_t \cdot \mathbb{E}_A[(1-p)^{K_g} \kappa \delta_t \mid \delta_{t_0} = 0] + \frac{1}{m} \eta_t \cdot \mathbb{E}_A[2\sqrt{2}(1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\| \mid \delta_{t_0} = 0] \\ &\leq (1 + (1-p)^{K_g} \eta_t \kappa) \Delta_t + \frac{2\sqrt{2}(1-p)^{K_g} \eta_t \alpha_\ell \sqrt{K} \|O\|}{m}. \end{aligned}$$

Using the fact that $\Delta_{t_0} = 0$, we can unwind this recurrence relation from T down to $t_0 + 1$. This gives

$$\Delta_{t+1} = \sum_{t=t_0+1}^T \left[\prod_{j=t+1}^T (1 + (1-p)^{K_g} \eta_j \kappa) \right] \frac{2\sqrt{2}(1-p)^{K_g} \eta_t \alpha_\ell \sqrt{K} \|O\|}{m}.$$

By the elementary inequality $1 + a \leq \exp(a)$ and $\eta_t \leq c/(t+1)$, we further derive

$$\begin{aligned} \Delta_{t+1} &\leq \sum_{t=t_0+1}^T \left[\prod_{j=t+1}^T \exp\left(\frac{c(1-p)^{K_g} \kappa}{j}\right) \right] \frac{2\sqrt{2}c(1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\|}{tm} \\ &\leq \sum_{t=t_0+1}^T \exp\left(\sum_{j=t+1}^T \frac{c(1-p)^{K_g} \kappa}{j}\right) \frac{2\sqrt{2}c(1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\|}{tm} \\ &\leq \sum_{t=t_0+1}^T \exp\left(c(1-p)^{K_g} \kappa \log\left(\frac{T}{t}\right)\right) \frac{2\sqrt{2}c(1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\|}{tm} \\ &\leq \frac{2\sqrt{2}c(1-p)^{K_g} \alpha_\ell \sqrt{K} \|O\|}{m} T^{c(1-p)^{K_g} \kappa} \sum_{t=t_0+1}^T \frac{1}{t^{c(1-p)^{K_g} \kappa+1}} \\ &\leq \frac{2\sqrt{2}\alpha_\ell \sqrt{K} \|O\|}{\kappa m} \left(\frac{T}{t_0}\right)^{c(1-p)^{K_g} \kappa}. \end{aligned}$$

Plugging this bound into equation (30), we get

$$\left| \mathbb{E}_A[\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)] \right| \leq \frac{2\sqrt{2}(1-p)^{K_g} \alpha_\ell^2 K \|O\|^2}{\kappa m} \left(\frac{T}{t_0}\right)^{c(1-p)^{K_g} \kappa} + \frac{t_0 M}{m}. \quad (31)$$

The right hand side is approximately minimized when

$$t_0 = \left(\frac{2\sqrt{2}c(1-p)^{K_g} \alpha_\ell^2 K \|O\|^2}{M} \right)^{\frac{1}{c(1-p)^{K_g} \kappa+1}} T^{\frac{c\kappa}{c\kappa+1/(1-p)^{K_g}}}.$$

Plugging it into the equation (30) we have (for simplicity we assume the above t_0 is an integer)

$$\begin{aligned} &\left| \mathbb{E}_A[\ell(f_{\theta_T}(\mathbf{x}); y) - \ell(f_{\theta'_T}(\mathbf{x}); y)] \right| \\ &\leq \frac{1 + \frac{1}{c(1-p)^{K_g} \kappa}}{m} M^{\frac{c\kappa}{c\kappa+1/(1-p)^{K_g}}} \left(2\sqrt{2}c(1-p)^{K_g} \alpha_\ell^2 K \|O\|^2\right)^{\frac{1}{c(1-p)^{K_g} \kappa+1}} T^{\frac{c\kappa}{c\kappa+1/(1-p)^{K_g}}}. \end{aligned}$$

By the definition of uniform stability as shown in Definition 3.1, we obtain the desired bound on the uniform stability of QNNs under the decaying step size and depolarizing noise setting. Combining Lemma 3.3, Part (b) completes the proof of Corollary 4.8, Part (b). \square

C.4 Proof of Theorem 4.12

In this subsection, we prove Theorem 4.12, which establishes the generalization bound in expectation for QNNs and first links generalization gap to optimization error with the help of less restrictive on-average stability. We begin by quoting the result which we set to prove.

Theorem 4.12 (Optimization-Dependent Generalization Bound). *Suppose that Assumption 3.6, 3.7, and 4.10 hold. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes η_t for T iterations, then the expected generalization gap satisfies*

$$\begin{aligned} & |\mathbb{E}_{S,A} [R_{\mathcal{D}}(A(S)) - R_S(A(S))]| \\ & \leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\eta_t \alpha \ell \sqrt{K} \|O\| (\mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2] + \sigma)}{m}, \end{aligned}$$

where κ is defined by equation (1).

Proof. Let $S = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ and $S' = \{\mathbf{z}'_1, \dots, \mathbf{z}'_m\}$ be drawn independently from \mathcal{D} . For any $i \in [m]$, define $S^{(i)} = \{\mathbf{z}_1, \dots, \mathbf{z}_{i-1}, \mathbf{z}'_i, \mathbf{z}_{i+1}, \dots, \mathbf{z}_m\}$ as the set formed from S by replacing the i -th element with \mathbf{z}'_i . Consider two sequences of the parameters, $\{\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T\}$ and $\{\boldsymbol{\theta}_0^{(i)}, \boldsymbol{\theta}_1^{(i)}, \dots, \boldsymbol{\theta}_T^{(i)}\}$, learned by the QNN running SGD on S and $S^{(i)}$, respectively. Let the example index selected by SGD at iteration t denoted by i_t .

Observe that at iteration t , with probability $1 - 1/m$, $i_t \neq i$, the example selected by SGD is the same in both S and $S^{(i)}$. With probability $1/m$, $i_t = i$, the selected example is different. Therefore, we have

$$\begin{aligned} \mathbb{E}_A [\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t+1}^{(i)}\|_2] & \leq (1 - \frac{1}{m}) \mathbb{E}_A \left[\left\| (\boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t})) - (\boldsymbol{\theta}_t^{(i)} - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t^{(i)}}(\mathbf{x}_{i_t}); y_{i_t})) \right\|_2 \right] \\ & \quad + \frac{1}{m} \mathbb{E}_A \left[\left\| (\boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i)) - (\boldsymbol{\theta}_t^{(i)} - \eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t^{(i)}}(\mathbf{x}'_i); y'_i)) \right\|_2 \right] \\ & = \mathbb{E}_A [\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_t^{(i)}\|_2] + (1 - \frac{1}{m}) \eta_t \mathbb{E}_A \left[\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t^{(i)}}(\mathbf{x}_{i_t}); y_{i_t}) \right\|_2 \right] \\ & \quad + \frac{1}{m} \eta_t \mathbb{E}_A \left[\left\| \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t^{(i)}}(\mathbf{x}'_i); y'_i) \right\|_2 \right]. \end{aligned}$$

According to Lemma C.2 and Lemma C.3, we further get

$$\begin{aligned} \mathbb{E}_A [\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t+1}^{(i)}\|_2] & \leq \mathbb{E}_A [\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_t^{(i)}\|_2] + (1 - \frac{1}{m}) \eta_t \cdot \mathbb{E}_A [\kappa \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_t^{(i)}\|_2] \\ & \quad + \frac{1}{m} \eta_t \cdot \mathbb{E}_A [\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i)\|_2 + \|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t^{(i)}}(\mathbf{x}'_i); y'_i)\|_2] \\ & \leq (1 + \eta_t \kappa) \mathbb{E}_A [\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_t^{(i)}\|_2] + \frac{\eta_t (\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i)\|_2 + \|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t^{(i)}}(\mathbf{x}'_i); y'_i)\|_2)}{m}. \end{aligned}$$

Unraveling the recursion gives

$$\begin{aligned} \mathbb{E}_A [\|\boldsymbol{\theta}_T - \boldsymbol{\theta}_T^{(i)}\|_2] & \leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{\eta_t (\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i)\|_2 + \|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t^{(i)}}(\mathbf{x}'_i); y'_i)\|_2)}{m} \\ & \leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\eta_t \|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i)\|_2}{m}. \end{aligned}$$

Using the Assumption 4.10 and take an average over $i \in [m]$, we have

$$\begin{aligned}
\frac{1}{m} \sum_{i=1}^m \mathbb{E}_A [\|\boldsymbol{\theta}_T - \boldsymbol{\theta}_T^{(i)}\|_2] &\leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\eta_t \sum_{i=1}^m \|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i)\|_2}{m^2} \\
&\leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\eta_t \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i)\|_2]}{m} \\
&\leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\eta_t (\mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2] + \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_i); y_i) - \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2])}{m} \\
&\leq \sum_{t=0}^{T-1} \left[\prod_{j=t+1}^{T-1} (1 + \eta_j \kappa) \right] \frac{2\eta_t (\mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2] + \sigma)}{m}.
\end{aligned}$$

By the definition of on-average stability as shown in Definition 3.2, we immediately get the claimed upper bound on the on-average stability of QNNs.

Using the Assumption 3.6 that the loss function is Lipschitz continuous and Lemma C.1, we have for $\forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^K$

$$\begin{aligned}
|\ell(f_{\boldsymbol{\theta}_1}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y)| &\leq \alpha_\ell \|f_{\boldsymbol{\theta}_1}(\mathbf{x}) - f_{\boldsymbol{\theta}_2}(\mathbf{x})\| \\
&\leq \alpha_\ell \sqrt{K} \|O\| \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2.
\end{aligned} \tag{32}$$

Combining equation (32) and Lemma 3.3, Part (c) completes the proof of Theorem 4.12. \square

C.5 Proof of Lemma 4.15

In this subsection, we prove Lemma 4.15, which shows that the expected gradient norms $\mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2]$ are influence by the choice of the initialization point. We first state and prove the following key lemma.

Lemma C.9 (Descent Lemma). *Suppose that Assumption 3.6 and 3.7 hold. Then, for $\forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^K$, and $\mathbf{z} \in \mathcal{Z}$, we have*

$$\ell(f_{\boldsymbol{\theta}_1}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y) \leq \langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \frac{\kappa}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2,$$

where $\kappa = \alpha_\ell K \|O\| + \sqrt{2} \nu_\ell K \|O\|^2$.

Proof. Let $\tilde{\boldsymbol{\theta}}$ be a point in the line segment of $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$, $\tilde{\boldsymbol{\theta}}(u) = \boldsymbol{\theta}_2 + u(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)$, then

$$\begin{aligned}
\ell(f_{\boldsymbol{\theta}_1}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y) &= \int_0^1 \langle \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \nabla_{\boldsymbol{\theta}} \ell(f_{\tilde{\boldsymbol{\theta}}(u)}(\mathbf{x}); y) \rangle du \\
&= \int_0^1 \langle \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y) + \nabla_{\boldsymbol{\theta}} \ell(f_{\tilde{\boldsymbol{\theta}}(u)}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y) \rangle du \\
&= \langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \int_0^1 \langle \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \nabla_{\boldsymbol{\theta}} \ell(f_{\tilde{\boldsymbol{\theta}}(u)}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y) \rangle du \\
&\leq \langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \int_0^1 \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| \cdot \|\nabla_{\boldsymbol{\theta}} \ell(f_{\tilde{\boldsymbol{\theta}}(u)}(\mathbf{x}); y) - \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y)\| du.
\end{aligned}$$

According to Lemma C.2, we further get

$$\begin{aligned}
\ell(f_{\boldsymbol{\theta}_1}(\mathbf{x}); y) - \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y) &\leq \langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \int_0^1 \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| \cdot \kappa \|\tilde{\boldsymbol{\theta}}(u) - \boldsymbol{\theta}_2\| du \\
&= \langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \int_0^1 \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| \cdot \kappa u \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| du \\
&= \langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \kappa \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|^2 \int_0^1 u du \\
&= \langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_2}(\mathbf{x}); y), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \frac{\kappa}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|^2.
\end{aligned}$$

This completes the proof of Lemma C.9. \square

We are now ready to prove Lemma 4.15.

Lemma 4.15 (Link with Initialization Point). *Suppose that Assumption 3.6, 3.7, and 4.10 hold. Let $A(S)$ be the QNN model trained on the dataset $S \in \mathcal{Z}^m$ using the SGD algorithm with step sizes $\eta_t \leq 1/\kappa$ for T iterations, then the following bound holds*

$$\begin{aligned}
&\sum_{t=0}^{T-1} \eta_t \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2] \\
&\leq 2 \sqrt{\left(\sum_{t=0}^{T-1} \eta_t \right) \left(R_S(\boldsymbol{\theta}_0) - R_S(\boldsymbol{\theta}^*) + \frac{\kappa \sigma^2}{2} \sum_{t=0}^{T-1} \eta_t^2 \right)},
\end{aligned}$$

where $\boldsymbol{\theta}^*$ is the empirical risk minimizer of $R_S(\boldsymbol{\theta})$, κ is defined by equation (1).

Proof. We can bound $\sum_{t=0}^{T-1} \eta_t \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2]$ as follows

$$\begin{aligned}
\sum_{t=0}^{T-1} \eta_t \mathbb{E}_S \left[\sqrt{\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2} \right] &\leq \sum_{t=0}^{T-1} \frac{\left(1 - \frac{\eta_t \kappa}{2}\right)}{\left(1 - \frac{\eta_t \kappa}{2}\right)} \cdot \eta_t \sqrt{\mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2]} \\
&\leq 2 \sum_{t=0}^{T-1} \left(\eta_t - \frac{\eta_t^2 \kappa}{2} \right) \sqrt{\mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2]} \\
&= \frac{2 \sum_{t=0}^{T-1} \left(\eta_t - \frac{\eta_t^2 \kappa}{2} \right)}{\sum_{t=0}^{T-1} \left(\eta_t - \frac{\eta_t^2 \kappa}{2} \right)} \sum_{t=0}^{T-1} \left(\eta_t - \frac{\eta_t^2 \kappa}{2} \right) \sqrt{\mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2]} \tag{33} \\
&\leq 2 \sqrt{\sum_{t=0}^{T-1} \eta_t} \sqrt{\sum_{t=0}^{T-1} \left(\eta_t - \frac{\eta_t^2 \kappa}{2} \right) \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2]},
\end{aligned}$$

where the first and last inequality use the Jensen's inequality, and the second inequality uses the condition $\eta_t \leq 1/\kappa$.

Then, we focus on the term $\sqrt{\sum_{t=0}^{T-1} \left(\eta_t - \frac{\eta_t^2 \kappa}{2}\right) \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2]}$. By Lemma C.9 and the update rule of the SGD, we obtain

$$\begin{aligned}
R_S(\boldsymbol{\theta}_{t+1}) - R_S(\boldsymbol{\theta}_t) &\leq \langle \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t, \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t) \rangle + \frac{\kappa}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|_2^2 \\
&= \langle -\eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}), \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t) \rangle + \frac{\kappa \eta_t^2}{2} \|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t})\|_2^2 \\
&= \langle -\eta_t \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}), \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t) \rangle \\
&\quad + \frac{\kappa \eta_t^2}{2} (\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}) - \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2 + \|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2 \\
&\quad - 2\langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}) - \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t), \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t) \rangle) \\
&= -(\eta_t + \eta_t^2 \kappa) \langle \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}), \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t) \rangle + \frac{3\eta_t^2 \kappa}{2} \|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}) - \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2
\end{aligned}$$

Using the Assumption 4.10 and take an average over $i_t \in [m]$, we have

$$\begin{aligned}
\left(\eta_t - \frac{\eta_t^2 \kappa}{2}\right) \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2] &\leq R_S(\boldsymbol{\theta}_t) - R_S(\boldsymbol{\theta}_{t+1}) + \frac{\eta_t^2 \kappa}{2} \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}_t}(\mathbf{x}_{i_t}); y_{i_t}) - \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2] \\
&\leq R_S(\boldsymbol{\theta}_t) - R_S(\boldsymbol{\theta}_{t+1}) + \frac{\eta_t^2 \kappa}{2} \sigma^2.
\end{aligned}$$

We can apply the above inequality recursively and derive

$$\sum_{t=0}^{T-1} \left(\eta_t - \frac{\eta_t^2 \kappa}{2}\right) \mathbb{E}_S [\|\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_t)\|_2^2] \leq R_S(\boldsymbol{\theta}_0) - R_S(\boldsymbol{\theta}^*) + \frac{\kappa \sigma^2}{2} \sum_{t=0}^{T-1} \eta_t^2.$$

Plugging it into equation (33) completes the proof of Lemma 4.15. \square