

Generalization Bounds for Quantum Learning via Rényi Divergences

Naqeeb Ahmad Warsi*, Ayanava Dasgupta* and Masahito Hayashi†

Abstract

This work advances the theoretical understanding of quantum learning by establishing a new family of upper bounds on the expected generalization error of quantum learning algorithms, leveraging the framework introduced by Caro et al. (2024) and a new definition for the expected true loss. Our primary contribution is the derivation of these bounds in terms of quantum and classical Rényi divergences, utilizing a variational approach for evaluating quantum Rényi divergences, specifically the Petz and a newly introduced modified sandwich quantum Rényi divergence. Analytically and numerically, we demonstrate the superior performance of the bounds derived using the modified sandwich quantum Rényi divergence compared to those based on the Petz divergence. Furthermore, we provide probabilistic generalization error bounds using two distinct techniques: one based on the modified sandwich quantum Rényi divergence and classical Rényi divergence, and another employing smooth max Rényi divergence.

Index Terms

generalization error, quantum learning algorithms, variational lower-bounds, modified sandwich quantum Rényi divergence, Petz quantum Rényi divergence, smooth max Rényi divergence.

I. INTRODUCTION

Quantum learning theory has emerged as a burgeoning field at the intersection of quantum computation and machine learning, promising enhanced capabilities for data analysis and pattern recognition. Recent years have witnessed significant advancements in developing theoretical frameworks and algorithms for quantum learning, as highlighted in works such as [1]–[7]. Among these, the work by Caro et al. [8] has provided a crucial foundation by introducing a comprehensive quantum learning framework. This framework uniquely addresses the challenges inherent in learning from quantum data, particularly when considering the interplay between classical training data and quantum data states. A key contribution of Caro et al. [8] lies in its analysis of quantum learning algorithms' performance through the lens of quantum information theory. By carefully considering the process of training on classical data in conjunction with quantum systems, and the subsequent generation of both classical and quantum hypotheses, Caro et al. [8] laid the groundwork for a rigorous evaluation of generalization in quantum learning scenarios. Their framework specifically tackles the issue of evaluating the learned hypothesis against a loss operator, a process complicated by the inherent measurements and post-processing that can perturb the quantum data. To overcome this, Caro et al. [8] introduced the concept of a test-train bipartition of the quantum data, allowing for the analysis of potential correlations and entanglement across these partitions and enabling the derivation of upper bounds on the expected generalization error in terms of both classical and quantum information-theoretic quantities. This prior work thus provides a vital starting point for further investigations into the theoretical limits and practical capabilities of quantum learning algorithms.

Building upon these foundational concepts from learning theory, the study of quantum learning aims to understand the capabilities and limitations of learning algorithms that utilize quantum resources. In classical learning theory, a central challenge lies in bridging the gap between the expected empirical loss, directly estimated from training data, and the theoretically critical expected true loss on unseen data. This difference is defined as the generalization error [9]–[13], a key figure of merit for assessing how well a learning algorithm's output generalizes to new, independent data. A substantial body of work in classical learning theory has been dedicated to analyzing this generalization error [14]–[20], [20]–[25]. Notably, an information-theoretic approach, as employed in [15]–[17], [20]–[23], [25], often utilizes the combination of variational lower bounds on divergence [26] and concentration inequalities like Hoeffding's Lemma [27, Lemma 2.2] (or sub-Gaussianity assumptions [27, Section 2.3] for unbounded loss functions).

The generalization error, inherently dependent on the stochastic nature of the training data and the learned hypothesis, is typically studied through two primary ways: its expectation and its probabilistic behavior. Investigations into the expected generalization error in classical learning have successfully employed the aforementioned information-theoretic tools. For instance, the references [15], [16] established bounds on the expected generalization error based on the mutual information between the training data and the resulting hypothesis. Furthermore, Modak et al. [21] derived upper bounds on the expected generalization error by leveraging the variational form of Rényi divergence [28].

Complementarily, a significant amount of research in classical learning theory has also focused on the probabilistic generalization error. Works such as [24], [25] have proven upper bounds on the generalization error that hold with a certain probability. Specifically,

*Indian Statistical Institute, Kolkata 700108, India. Email: naqeebwarsi@isical.ac.in, ayanavadasgupta_r@isical.ac.in

†School of Data Science, The Chinese University of Hong Kong, Shenzhen, Longgang District, Shenzhen, 518172, China

International Quantum Academy, Futian District, Shenzhen 518048, China

Graduate School of Mathematics, Nagoya University, Nagoya, 464-8602, Japan. Email: hmasahito@cuhk.edu.cn

Esposito et al. [25] demonstrated that the generalization error is probabilistically upper-bounded by Rényi divergence (up to additive and multiplicative constants) with high probability, achieving this through a change of measure technique rooted in Hölder’s inequality (refer to Fact 2).

It is worth noting that many existing works in both classical [14]–[20], [20]–[23], [25] and quantum [8] learning theory often assume the loss function (or loss observable in the quantum case) to be sub-Gaussian. However, as observed in the classical setting due to Hoeffding’s inequality (Fact 7), this assumption is less restrictive for bounded loss functions. In this work, we establish a similar result for quantum learning by first proving a quantum analogue of Hoeffding’s inequality, which subsequently allows us to demonstrate that a bounded loss operator is trivially sub-Gaussian in the quantum context.

The main contribution of this paper is to extend classical results on generalization error, specifically those obtained in [21, Theorem 3] and [25, Corollary 2], to the quantum setting within the framework of [8]. We achieve this by proving upper bounds on the expected generalization error for quantum learning algorithms using variational forms involving Rényi-divergence-type quantities. Given that our evaluation is based on measurements, we obtain variational bounds based on measured Rényi divergence. However, these bounds involve optimization over the choice of measurement, making their calculation challenging. To circumvent this, we focus on two specific types of quantum Rényi divergences: Petz quantum Rényi divergence [29] and sandwich quantum Rényi divergence [30], [31]. While sandwich quantum Rényi divergence is generally smaller than Petz quantum Rényi divergence, it only provides an upper bound on measured Rényi divergence for $\alpha \geq 1/2$, whereas Petz quantum Rényi divergence is applicable for $\alpha \in (0, 1) \cup (1, \infty)$ (See Fact 28). To address this limitation, we introduce a modification of the sandwich quantum Rényi divergence using the reverse sandwich quantum Rényi divergence with parameter $\alpha < 1/2$. By combining this modified sandwich quantum Rényi divergence with a quantum Hoeffding bound, we derive tighter variational bounds for the expected generalization error. Furthermore, we investigate the generalization error in probability using various techniques adapted to the quantum setting.

In addition, Caro et al. [8] introduced a definition of true loss (see Definition 16) for a quantum learning algorithm, which we believe is conceptually misleading. Thus, we in this work, propose a new definition of true loss (see Definition 17) and also explain the motivation behind this definition (see Remark 1).

TABLE I: Relationship between results obtained in this work and related studies

Types of upper-bounds		Classical Learning Setting	Quantum Learning Setting
Upper-bounds on the generalization error in expectation	Whole sample (as defined in (53)) based bounds in terms of divergence	[16, Theorem 1] (Proposition 1 in this paper)	[8, Theorem 17], [8, Corollary 23] (Proposition 4 in this paper) and Theorem 1 in this paper.
	Whole sample based bounds in terms of Rényi divergence		Theorem 2 and Corollaries 2 and 3 of this paper.
	Individual sample (as defined in (54)) based bounds in terms of divergence	[17, Proposition 1] (Proposition 2 in this paper)	[8, Corollary 24]
	Individual sample based bounds in terms of Rényi divergence	[21, Theorem 1] (Proposition 3 in this paper)	Corollary 5 of this paper.
Upper-bounds on the generalization error in probability	Using Hölder’s inequality (see Fact 2)	[25, Corollary 2] (Proposition 5 in this paper)	Theorem 4 of this paper.
	Using smooth max Rényi divergence (see Definition 3)	Theorem 3 of this paper.	Theorem 5 of this paper.

Organization of this Paper and our contributions

- In Section II, we introduce the notations, facts, and definitions used in this manuscript. In particular, we present the modified sandwich quantum Rényi divergence. Additionally, we provide a variational lower bound for this modified sandwich quantum Rényi divergence.
- In Section III, we discuss the usefulness of the variational lower bound on KL-divergence. Following this, we extend this discussion to the quantum setting for bounded linear operators. To achieve this, we prove a quantum version of Hoeffding’s lemma in Lemma 1. This lemma is analogous to its classical counterpart. Consequently, it allows us to deduce that every bounded linear operator is sub-Gaussian. Furthermore, we show that the modified sandwich quantum Rényi divergence is asymptotically close to the measured Rényi divergence, even without the i.i.d. assumption.
- In Section IV, we begin by discussing the quantum learning framework introduced by [8]. However, in this section, we also present a new definition for the expected true loss of quantum learning algorithms, offering an alternative to the one found in [8]. Subsequently, we justify and explain the motivation for this new definition.
- In Section V, we discuss a relation between the L_1 distance (between distributions) and the expected generalization error for bounded loss functions. Following this, we discuss a similar relation between the L_1 distance (between quantum states) and the expected quantum generalization error for bounded linear loss observables. Additionally, we review existing works in the literature that study upper bounds on the expected generalization error in both classical and quantum learning scenarios.

- In Section VI, we prove a family of upper bounds for the expected quantum generalization error. These bounds use the modified sandwich quantum Rényi divergence and also the classical Rényi divergence. Although it is possible to replace the modified sandwich quantum Rényi divergence with the Petz quantum Rényi divergence, our bounds using the modified version show better performance. Furthermore, we demonstrate that the bounds obtained in [8] can be derived from our bounds, which are based on quantum divergence. Finally, we numerically compare these three bounds, and the results show that our bounds based on the modified sandwich quantum Rényi divergence are superior.
- Section VII presents two distinct evaluations of the probabilistic behavior of the generalization error for quantum learning algorithms. Specifically, one probabilistic bound uses the modified sandwich quantum Rényi divergence and the classical Rényi divergence. The other probabilistic bound we obtain uses the smooth max Rényi divergence.

In particular, Table I above summarizes the relation between our target problems and existing studies.

II. NOTATIONS, FACTS AND DEFINITIONS

We use \mathcal{H} to denote a finite-dimensional Hilbert space and we denote its dimension with $|\mathcal{H}|$, $\mathcal{D}(\mathcal{H})$ to represent the set of all state density matrix acting on \mathcal{H} , $\mathcal{B}(\mathcal{H})$ to represent the set of all bounded self-adjoint operators acting on \mathcal{H} and $\mathcal{L}(\mathcal{H})$ represents the set of all linear operators over \mathcal{H} , $\mathcal{L}_{\geq 0}(\mathcal{H})$ represents the set of all positive operators over \mathcal{H} ($\mathcal{D}(\mathcal{H}) \subset \mathcal{L}_{\geq 0}(\mathcal{H})$), $T(\mathcal{H})$ represents the set of all trace class operators over \mathcal{H} and $|\mathcal{H}|$ represents the dimension of the Hilbert space \mathcal{H} .

Definition 1. Let P and Q be two probability distributions (measures) over a common metric space \mathcal{X} . Then, for any $p \in [1, \infty)$ the L_p distance between P and Q is defined as follows,

$$\|P - Q\|_p := \left(\sum_{x \in \mathcal{X}} |P(x) - Q(x)|^p \right)^{\frac{1}{p}},$$

where the supremum is taken over all measurable sets. Note that $\|P - Q\|_p$ is a strictly decreasing function of p .

Definition 2 (Rényi Divergence [32, Equation 3.3]). Consider two probability distributions P and Q and $\gamma \in (0, 1) \cup (1, +\infty)$. Then, the Rényi divergence of order γ is defined as follows,

$$D_\gamma^c(P\|Q) := \begin{cases} \frac{1}{\gamma-1} \log \mathbb{E}_Q \left[\left(\frac{dP}{dQ} \right)^\gamma \right], & \text{if } P \ll Q, \\ +\infty, & \text{else.} \end{cases}$$

Definition 3 (Smooth max Rényi divergence [33, Definition 10]). Consider two probability distributions P and Q over a finite set \mathcal{X} . Then, the smooth max Rényi divergence $D_{\max}^{(\varepsilon)}(P\|Q)$ of order ε is defined as follows,

$$D_{\max}^{(\varepsilon)}(P\|Q) := \inf \left\{ a : \Pr_P \left\{ x \in \mathcal{X} : \log \frac{P(x)}{Q(x)} < a \right\} \geq 1 - \varepsilon \right\}. \quad (1)$$

Definition 4 (Convex conjugate of a function [34]). Given a convex function $f : \mathbb{R} \rightarrow \mathbb{R}$, the convex conjugate f^* corresponding to f is given as,

$$f^*(y) := \sup_{x \in \mathbb{R}} (xy - f(x)).$$

Definition 5 (sub-Gaussianity of random variables [27, Section 2.3]). For some $\beta > 0$, a random variable X is defined to be β -sub-Gaussian if, $\forall \lambda \in \mathbb{R}$, the logarithmic moment-generating function of X satisfies the following condition,

$$\log \mathbb{E} \left[e^{\lambda(X - \mathbb{E}[X])} \right] \leq \frac{\lambda^2 \beta^2}{2}. \quad (2)$$

Definition 6 ([35, Equation IV.31]). Given any operator $O \in \mathcal{L}(\mathcal{H})$, for any $p \in [1, \infty)$, the Schatten- p -norm $\|O\|_p$ of O is defined as follows,

$$\|O\|_p := \left(\text{Tr} \left[\left(\sqrt{O^\dagger O} \right)^p \right] \right)^{\frac{1}{p}},$$

where O^\dagger is the conjugate-transpose of O and we let $\|O\|_\infty := \lim_{p \rightarrow \infty} \|O\|_p$, which turns out to be the largest singular value of O , if $O \in \mathcal{B}(\mathcal{H})$.

Definition 7. Consider $\rho, \sigma \in \mathcal{D}(\mathcal{H})$. Quantum Divergence between ρ and σ is defined as follows

$$D(\rho\|\sigma) := \begin{cases} \text{Tr}[\rho(\log \rho - \log \sigma)], & \text{if } \rho \ll \sigma, \\ +\infty, & \text{else.} \end{cases}$$

Then, Measured Quantum Divergence between ρ and σ is defined as follows

$$D^{\mathbb{M}}(\rho\|\sigma) := \sup_{\mathcal{X}, \{\Lambda_x\}_{x \in \mathcal{X}}} \sum_{x \in \mathcal{X}} \text{Tr}[\Lambda_x \rho] (\log \text{Tr}[\Lambda_x \rho] - \log \text{Tr}[\Lambda_x \sigma]),$$

where the supremum is over the choices of finite sets \mathcal{X} and POVMs $\{\Lambda_x\}_{x \in \mathcal{X}}$.

Definition 8 (Measured Rényi Divergence [36, Eqs. 3.116-3.117]). Consider $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\alpha \in (0, 1) \cup (1, +\infty)$. Then, Measured Rényi Divergence of order α between ρ and σ is defined as follows,

$$D_\alpha^{\mathbb{M}}(\rho \parallel \sigma) := \sup_{\mathcal{X}, \{\Lambda_x\}_{x \in \mathcal{X}}} \frac{1}{\alpha - 1} \log \sum_{x \in \mathcal{X}} (\text{Tr}[\Lambda_x \rho])^\alpha (\text{Tr}[\Lambda_x \sigma])^{1-\alpha},$$

where the supremum is over the choices of finite sets \mathcal{X} and POVMs $\{\Lambda_x\}_{x \in \mathcal{X}}$.

Definition 9 (Petz Quantum Rényi Divergence [29]). Consider $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\alpha \in (0, 1) \cup (1, +\infty)$. Then, Petz Quantum Rényi Divergence of order α between ρ and σ is defined as follows,

$$D_\alpha(\rho \parallel \sigma) := \begin{cases} \frac{1}{\alpha-1} \log \text{Tr}[\rho^\alpha \sigma^{1-\alpha}], & \text{if } (\alpha < 1 \cap \rho \not\leq \sigma) \cup (\rho \ll \sigma), \\ +\infty, & \text{else.} \end{cases}$$

Definition 10 (Minimal/ Sandwiched Quantum Rényi Divergence [30], [31, Definition 4]). Consider $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\alpha \in (0, 1) \cup (1, +\infty)$. Then, ‘minimal/sandwiched Quantum Rényi Divergence’ or ‘Quantum Rényi Divergence’ of order α between ρ and σ is defined as follows,

$$\tilde{D}_\alpha(\rho \parallel \sigma) := \begin{cases} \frac{1}{\alpha-1} \log \text{Tr} \left[\left(\sigma^{\frac{1-\alpha}{2\alpha}} \rho \sigma^{\frac{1-\alpha}{2\alpha}} \right)^\alpha \right], & \text{if } (\alpha < 1 \cap \rho \not\leq \sigma) \cup (\rho \ll \sigma), \\ +\infty, & \text{else.} \end{cases}$$

Definition 11 (Reverse Sandwiched Quantum Rényi Divergence). Consider $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\alpha \in (0, 1) \cup (1, +\infty)$. Then, reverse sandwiched Quantum Rényi Divergence of order α between ρ and σ is defined as follows,

$$\tilde{D}_\alpha^R(\rho \parallel \sigma) := \frac{\alpha}{1-\alpha} \tilde{D}_{1-\alpha}(\sigma \parallel \rho).$$

Definition 12 (Modified Sandwiched Quantum Rényi Divergence). Consider $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\alpha \in (0, 1)$. Then, modified sandwiched Quantum Rényi Divergence of order α between ρ and σ is defined as follows,

$$\overline{D}_\alpha(\rho \parallel \sigma) := \begin{cases} \tilde{D}_\alpha^R(\rho \parallel \sigma) & \text{if } (\alpha < 1/2), \\ \tilde{D}_\alpha(\rho \parallel \sigma) & \text{if } (\alpha \geq 1/2). \end{cases}$$

Definition 13 (sub-Gaussianity of observables). For some $\mu > 0$, a self-adjoint operator $O \in \mathcal{B}(\mathcal{H})$ is defined to be μ -sub-Gaussian with respect to a quantum state $\rho \in \mathcal{D}(\mathcal{H})$ if, $\forall \lambda \in \mathbb{R}$, O satisfies the following condition,

$$\log \text{Tr} \left[e^{\lambda(O - \text{Tr}[O\rho] \mathbb{I}_{\mathcal{H}})} \rho \right] \leq \frac{\lambda^2 \mu^2}{2}.$$

Fact 1 (Jensen’s inequality [37]). Given X is a random variable and ψ and ϕ are convex and concave functions, respectively. Then,

$$\psi(\mathbb{E}[X]) \leq \mathbb{E}[\psi(X)], \quad (3)$$

$$\phi(\mathbb{E}[X]) \geq \mathbb{E}[\phi(X)]. \quad (4)$$

Fact 2 (Hölder’s inequality). Given two random variables X and Y and two real numbers $p, q \in [1, \infty)$ such that $\frac{1}{p} + \frac{1}{q} = 1$, we have,

$$\mathbb{E}[|XY|] \leq (\mathbb{E}[|X|^p])^{1/p} (\mathbb{E}[|Y|^q])^{1/q}. \quad (5)$$

Fact 3 (Variational form of L_p distance). Let P and Q be two probability distributions (measures) over a common metric space \mathcal{X} . Then for any $p \in [1, \infty)$ the L_p distance between P and Q has the following variational representation,

$$\|P - Q\|_p := \frac{1}{B} \sup_{\|f\|_q \leq B} (\mathbb{E}_{X \sim P}[f(X)] - \mathbb{E}_{X \sim Q}[f(X)]), \quad (6)$$

where q is the Hölder conjugate of p i.e. q satisfies $\frac{1}{p} + \frac{1}{q} = 1$ and the supremum is taken over all bounded functions f over \mathcal{X} within range $[-B, B]$ i.e. $\|f\|_q \leq B$ (where $\|f\|_q := (\sum_{x \in \mathcal{X}} |f(x)|^q)^{\frac{1}{q}}$). Further, setting $p = 1$, the L_1 distance between P and Q has the following variational representation,

$$\|P - Q\|_1 := \frac{1}{B} \sup_{\|f\|_\infty \leq B} (\mathbb{E}_{X \sim P}[f(X)] - \mathbb{E}_{X \sim Q}[f(X)]), \quad (7)$$

where the supremum is taken over all bounded functions f over \mathcal{X} within range $[-B, B]$ i.e. $\|f\|_\infty \leq B$ (where $\|f\|_\infty := \sup_{x \in \mathcal{X}} |f(x)|$).

Fact 4. D_γ^c is monotonically increasing in $\gamma > 0$ and

$$D_1^c(P\|Q) := \lim_{\gamma \rightarrow 1} D_\gamma^c(P\|Q) = D^c(P\|Q),$$

where,

$$D^c(P\|Q) := \begin{cases} \mathbb{E}_P \left[\log \left(\frac{dP}{dQ} \right) \right], & \text{if } P \ll Q, \\ +\infty, & \text{else.} \end{cases}$$

Fact 5 (Donsker-Varadhan variational form for divergence [26, Corollary 4.15]). *Let P and Q be probability measures on \mathcal{X} . Then, we have the following dual form of $D^c(P\|Q)$,*

$$D^c(P\|Q) = \sup_{G \in \mathcal{M}_b(\mathcal{X})} \left\{ \mathbb{E}_{X \sim P}[G(X)] - \log \mathbb{E}_{X \sim Q}[e^{G(X)}] \right\}.$$

where $\mathcal{M}_b(\mathcal{X})$ denotes the set of bounded measurable real-valued functions on \mathcal{X} . The above dual form of $D^c(P\|Q)$ is also known as a **variational form**, which is a direct consequence of convex duality (see Definition 4).

Fact 6 (Variational form for Rényi divergence [28, Theorem 3.1]). *Let P and Q be probability measures on \mathcal{X} and $\gamma \in \mathbb{R} \setminus \{0, 1\}$. Then, we have,*

$$D_\gamma^c(P\|Q) = \sup_{G \in \Gamma} \left\{ \frac{\gamma}{\gamma - 1} \log \mathbb{E}_{X \sim P}[e^{(\gamma-1)G(X)}] - \log \mathbb{E}_{X \sim Q}[e^{\gamma G(X)}] \right\}. \quad (8)$$

where $\mathcal{M}_b(\mathcal{X}) \subset \Gamma \subset \mathcal{M}(\mathcal{X})$ ($\mathcal{M}(\mathcal{X})$ denotes the set of all real-valued measurable functions \mathcal{X}).

Fact 7 (Hoeffding's Lemma [27, Lemma 2.2]). *Let X_i , for $i \in [n]$, be i.i.d. random variables distributed according to a probability distribution P_X taking values in a bounded interval $[a, b]$ and $\mathbb{E}[X_i] = \mu$. Let $X = \sum_{i=1}^n X_i/n$ denote the average of the $\{X_i\}_{i=1}^n$. Then, $\forall \lambda \in \mathbb{R}$, the following holds,*

$$\log \mathbb{E}[e^{\lambda X}] \leq \lambda \mu + \frac{\lambda^2(b-a)^2}{8n}. \quad (9)$$

Fact 8 ([25], [21]). *Consider a probability distribution P_{AB} over $\mathcal{A} \times \mathcal{B}$ (where $|\mathcal{A}|, |\mathcal{B}| > 0$) and if $P_{AB} \ll P_A \times P_B$ (where P_A and P_B are the corresponding marginal of P_{AB}), then, for any event $E \subseteq \mathcal{A} \times \mathcal{B}$, the following holds,*

$$\Pr_{(A,B) \sim P_{AB}} \{E\} \leq e^{\frac{\gamma-1}{\gamma} (\log(\mathbb{E}_{P_A}[P_{TB \sim P_B}\{E_B\}]) + I_\gamma^c[A; B])}, \quad (10)$$

where $\forall b \in \mathcal{B}$, $E_b := \{a \in \mathcal{A} : (a, b) \in E\}$.

Fact 9 ([38, Proposition 2.5]). *Consider a collection of n τ -sub-Gaussian i.i.d. random variable X_1, \dots, X_n and let $S = \frac{1}{n} \sum_{i=1}^n nX_i$. Then, we have,*

$$\begin{aligned} \Pr\{S - \mathbb{E}[S] > t\} &\leq e^{-\frac{nt^2}{2\tau^2}}, \\ \Pr\{S - \mathbb{E}[S] < -t\} &\leq e^{-\frac{nt^2}{2\tau^2}}. \end{aligned} \quad (11)$$

Fact 10. *Consider a collection of n τ -sub-Gaussian i.i.d. random variable X_1, \dots, X_n and let $S = \frac{1}{n} \sum_{i=1}^n nX_i$. Then, we have,*

$$\Pr\{|S - c| < t\} \leq 2e^{-\frac{n(t - |\mathbb{E}[S] - c|)^2}{2\tau^2}}. \quad (12)$$

Proof. Consider the following series of inequalities,

$$\begin{aligned} \Pr\{|S - c| > t\} &\leq \Pr\{S > c + t\} + \Pr\{S < c - t\} \\ &= \Pr\{S - \mathbb{E}[S] > t - (\mathbb{E}[S] - c)\} + \Pr\{S - \mathbb{E}[S] < (c - \mathbb{E}[S]) - t\} \\ &\stackrel{a}{\leq} e^{-\frac{n(t - (\mathbb{E}[S] - c))^2}{2\tau^2}} + e^{-\frac{n((c - \mathbb{E}[S]) - t)^2}{2\tau^2}} \\ &= e^{-\frac{n(t - (\mathbb{E}[S] - c))^2}{2\tau^2}} + e^{-\frac{n(t - (c - \mathbb{E}[S]))^2}{2\tau^2}} \\ &\leq 2e^{-\frac{n(t - |\mathbb{E}[S] - c|)^2}{2\tau^2}}, \end{aligned}$$

where a follows from Fact 9. This proves Fact 10. ■

Fact 11. (Variational lower-bound for Schatten L_p distance) *Given two $\rho, \sigma \in \mathcal{D}(\mathcal{H})$, for $p \in [1, \infty)$, we have the following variational form for $\|\rho - \sigma\|_p$,*

$$\|\rho - \sigma\|_p \geq \frac{1}{B} \sup_{\substack{H \in \mathcal{B}(\mathcal{H}): \\ \|H\|_q \leq B}} (\text{Tr}[H\rho] - \text{Tr}[H\sigma]), \quad (13)$$

where q is Hölder conjugate of p i.e. q satisfies $\frac{1}{p} + \frac{1}{q} = 1$. Further, setting $p = 1$, we get the following variational form for $\|\rho - \sigma\|_1$,

$$\|\rho - \sigma\|_1 \geq \frac{1}{B} \sup_{\substack{H \in \mathcal{B}(\mathcal{H}): \\ -B\mathbb{I} \leq H \leq B\mathbb{I}}} (\text{Tr}[H\rho] - \text{Tr}[H\sigma]). \quad (14)$$

Proof. For any $H \in \mathcal{B}(\mathcal{H}) : \|H\|_q \leq B$, consider the following series of inequalities,

$$\begin{aligned} \frac{1}{B} \text{Tr}[H(\rho - \sigma)] &\stackrel{a}{\leq} \frac{1}{B} \|H\|_q \|\rho - \sigma\|_p \\ &\stackrel{b}{\leq} \|\rho - \sigma\|_p, \end{aligned} \quad (15)$$

where a follows from Fact 21 and b follows from the fact that $\|H\|_q \leq B$. Then, taking supremum over H on both sides of (15) completes the proof for Fact 11. \blacksquare

Fact 12. D_α is monotonically increasing in $\alpha > 0$ and,

$$D_1(\rho\|\sigma) := \lim_{\alpha \rightarrow 1} D_\alpha(\rho\|\sigma) = D(\rho\|\sigma).$$

Fact 13. \tilde{D}_α is monotonically increasing in $\alpha > 0$ and,

$$\tilde{D}_1(\rho\|\sigma) := \lim_{\alpha \rightarrow 1} \tilde{D}_\alpha(\rho\|\sigma) = D(\rho\|\sigma).$$

Fact 14. Consider two quantum states $\rho := \bigotimes_{i=1}^n \rho^{(i)}$ and $\sigma := \bigotimes_{i=1}^n \sigma^{(i)}$. Then, for any $\alpha \in (0, 1) \cup (1, +\infty)$, the following holds,

$$D_\alpha(\rho\|\sigma) = \sum_{i=1}^n D_\alpha(\rho^{(i)}\|\sigma^{(i)}), \quad (16)$$

$$\tilde{D}_\alpha(\rho\|\sigma) = \sum_{i=1}^n \tilde{D}_\alpha(\rho^{(i)}\|\sigma^{(i)}). \quad (17)$$

Fact 15 (Variational characterization of the quantum divergence [39], [40, Theorem 5.9], [41]). Let $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ be two quantum states. Then, the divergence between ρ and σ can be rewritten as follows,

$$D^\mathbb{M}(\rho\|\sigma) = \sup_{H \in \mathcal{B}(\mathcal{H})} \{ \text{Tr}[H\rho] - \log \text{Tr}[e^H \sigma] \}. \quad (18)$$

Fact 16 ([42, Lemma 4]). Let $\rho, \sigma \in \mathcal{D}(\mathcal{A})$. Then, for $\alpha \in (0, 1) \cup (1, \infty)$, we have,

$$\tilde{D}_\alpha(\rho\|\sigma) = \sup_{\substack{H \in \mathcal{L}(\mathcal{A}): \\ H \succeq 0}} \left\{ \frac{1}{\alpha - 1} \log \left(\alpha \text{Tr}[H\rho] - (\alpha - 1) \text{Tr} \left[\left(H^{\frac{1}{2}} \sigma^{\frac{\alpha-1}{\alpha}} H^{\frac{1}{2}} \right)^{\frac{\alpha}{\alpha-1}} \right] \right) \right\}. \quad (19)$$

Fact 17 ([41, Lemma 3 and Theorem 4]). Let $\rho, \sigma \in \mathcal{D}(\mathcal{A})$. Then, for $\alpha \in (0, 1) \cup (1, \infty)$, we have,

$$D_\alpha^\mathbb{M}(\rho\|\sigma) = \sup_{\substack{H \in \mathcal{L}(\mathcal{A}): \\ H \succ 0}} \left\{ \frac{\alpha}{\alpha - 1} \log \text{Tr} \left[e^{(\alpha-1)H} \rho \right] - \log \text{Tr} \left[e^{\alpha H} \sigma \right] \right\}, \quad \forall \alpha \in (0, 1) \cup (1, \infty).$$

Recently, Fang et al. in [43] obtained a variational expression for measured f -divergences (for operator convex functions).

Fact 18 (Araki-Lieb-Thirring trace inequality [44], [45]). Consider $X, Y \in \mathcal{L}_{\geq 0}(\mathcal{H})$. Then, we have the following inequality:

$$\begin{cases} \text{Tr}[[YXY]^r] &\leq \text{Tr}[Y^r X^r Y^r], & \text{if } r \geq 1, \\ \text{Tr}[[YXY]^r] &\geq \text{Tr}[Y^r X^r Y^r], & \text{if } r \in [0, 1]. \end{cases}$$

Fact 19. From Fact 18, it directly follows that for any $\rho, \sigma \in \mathcal{D}(\mathcal{H}) : \rho \ll \sigma$ and $\forall \alpha \in (0, 1) \cup (1, \infty)$, we have,

$$D_\alpha(\rho\|\sigma) \geq \tilde{D}_\alpha(\rho\|\sigma).$$

Fact 20. Consider a positive operator $A \in \mathcal{L}_{\geq 0}(\mathcal{H})$ and a state-density operator $\rho \in \mathcal{D}(\mathcal{H})$. Then, we have,

$$\log \text{Tr}[\rho A] \geq \text{Tr}[\rho \log A].$$

Proof. Assume the following eigen decomposition of A and ρ .

$$A = \sum_{i=1}^{|\mathcal{H}|} \alpha_i |i\rangle\langle i| \text{ and } \rho = \sum_{j=1}^{|\mathcal{H}|} \beta_j |j\rangle\langle j|, \text{ where } \forall j \in [|\mathcal{H}|], 0 < \beta_j < 1 \text{ and } \sum_{j=1}^{|\mathcal{H}|} \beta_j = 1.$$

Then, consider the following series of inequalities,

$$\begin{aligned}\log \text{Tr}[\rho A] &= \log \left(\sum_{i=1}^{|\mathcal{H}|} \sum_{j=1}^{|\mathcal{H}|} \alpha_i \beta_j |\langle i|j \rangle|^2 \right) \\ &= \log \left(\sum_{i=1}^{|\mathcal{H}|} \sum_{j=1}^{|\mathcal{H}|} \beta_j |\langle i|j \rangle|^2 \alpha_i \right).\end{aligned}\tag{20}$$

For all $i \in [|\mathcal{H}|]$, let $p_i := \sum_{j=1}^{|\mathcal{H}|} \beta_j |\langle i|j \rangle|^2$. It is easy to see that $\forall i \in [|\mathcal{H}|]$, $p_i \geq 0$ and $\sum_{i=1}^{|\mathcal{H}|} p_i = 1$.

Thus, we can now lower-bound (20) as follows,

$$\begin{aligned}\log \text{Tr}[\rho A] &= \log \left(\sum_{i=1}^{|\mathcal{H}|} p_i \alpha_i \right) \\ &\stackrel{a}{\geq} \sum_{i=1}^{|\mathcal{H}|} p_i \log \alpha_i \\ &= \sum_{i=1}^{|\mathcal{H}|} \sum_{j=1}^{|\mathcal{H}|} \beta_j |\langle i|j \rangle|^2 \log \alpha_i \\ &= \text{Tr} \left[\left(\sum_{j=1}^{|\mathcal{H}|} \beta_j |j\rangle \langle j| \right) \left(\sum_{i=1}^{|\mathcal{H}|} \log \alpha_i |i\rangle \langle i| \right) \right] \\ &= \text{Tr}[\rho \log A],\end{aligned}$$

where a follows from Jensen's inequality. ■

Fact 21 (Hölder's Inequality for operators [46, Equation 12.6]). *Given two positive semidefinite operator $A, B \in \mathcal{L}(\mathcal{H})$ and two real numbers $p, q \in [1, \infty)$ such that $\frac{1}{p} + \frac{1}{q} = 1$, we have,*

$$|\text{Tr}[AB]| \leq (\text{Tr}[A^p])^{\frac{1}{p}} (\text{Tr}[B^q])^{\frac{1}{q}}.\tag{21}$$

Fact 22 (Data-processing inequality of sandwiched quantum Rényi divergence [42, Theorem 1]). *For any $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\forall \alpha \in [\frac{1}{2}, 1) \cup (1, \infty)$, $\tilde{D}_\alpha(\rho\|\sigma)$ satisfies the following,*

$$\tilde{D}_\alpha(\rho\|\sigma) \geq \tilde{D}_\alpha(\mathcal{E}(\rho)\|\mathcal{E}(\sigma)),\tag{22}$$

where \mathcal{E} is any completely positive and trace-preserving (CP-TP) map.

Fact 23 ([40, Eq. (3.17)]). *For any $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\forall \alpha \in [\frac{1}{2}, 1) \cup (1, \infty)$, $\tilde{D}_\alpha(\rho\|\sigma)$ satisfies the following,*

$$\tilde{D}_\alpha(\rho\|\sigma) = \lim_{n \rightarrow \infty} \frac{1}{n} D_\alpha^{\mathfrak{M}}(\rho^{\otimes n}\|\sigma^{\otimes n}).\tag{23}$$

Fact 24 (Data-processing inequality of quantum divergence). *For any $\rho, \sigma \in \mathcal{D}(\mathcal{H})$, $D(\rho\|\sigma)$ satisfies the following,*

$$D(\rho\|\sigma) \geq D(\mathcal{E}(\rho)\|\mathcal{E}(\sigma)),\tag{24}$$

where \mathcal{E} is any completely positive and trace-preserving (CP-TP) map.

Fact 25 (Data-processing inequality of modified sandwiched quantum Rényi divergence). *For any $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\forall \alpha \in (0, 1) \cup (1, \infty)$, $\overline{D}_\alpha(\rho\|\sigma)$ satisfies the following,*

$$\overline{D}_\alpha(\rho\|\sigma) \geq \overline{D}_\alpha(\mathcal{E}(\rho)\|\mathcal{E}(\sigma)),\tag{25}$$

where \mathcal{E} is any completely positive and trace-preserving (CP-TP) map.

Fact 26. *For any $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\forall \alpha \in (0, 1) \cup (1, \infty)$, $\overline{D}_\alpha(\rho\|\sigma)$ satisfies the following,*

$$\overline{D}_\alpha(\rho\|\sigma) = \lim_{n \rightarrow \infty} \frac{1}{n} D_\alpha^{\mathfrak{M}}(\rho^{\otimes n}\|\sigma^{\otimes n}).\tag{26}$$

Proof. Consider the following series of inequalities:

$$\begin{aligned}\overline{D}_\alpha(\rho\|\sigma) &\stackrel{a}{=} \lim_{n \rightarrow \infty} \frac{1}{n} \overline{D}_\alpha(\rho^{\otimes n}\|\sigma^{\otimes n}) \\ &\stackrel{b}{=} \lim_{n \rightarrow \infty} \frac{1}{n} D_\alpha^{\mathfrak{M}}(\rho^{\otimes n}\|\sigma^{\otimes n}),\end{aligned}$$

where a follows from (17) of Fact 14 and b follows from Lemma 2 by setting $\rho_n = \rho^{\otimes n}$ and $\sigma_n = \sigma^{\otimes n}$. This proves Fact 26 \blacksquare

Fact 27 (Data-processing inequality of Petz quantum Rényi divergence [29]). *For any $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\forall \alpha \in (0, 1) \cup (1, 2]$, $D_\alpha(\rho \parallel \sigma)$ satisfies the following,*

$$D_\alpha(\rho \parallel \sigma) \geq D_\alpha(\mathcal{E}(\rho) \parallel \mathcal{E}(\sigma)), \quad (27)$$

where \mathcal{E} is any completely positive and trace-preserving (CP-TP) map.

Fact 28 (Measurement-data-processing inequality of Petz quantum Rényi divergence [40, Eq. (3.23)]). *For any $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\forall \alpha \in (0, 1) \cup (1, \infty)$, $D_\alpha(\rho \parallel \sigma)$ satisfies the following,*

$$D_\alpha(\rho \parallel \sigma) \geq D_\alpha^{\mathbb{M}}(\rho \parallel \sigma). \quad (28)$$

The combination of Facts 25, 26, and 28 implies the following fact.

Fact 29. *For any $\rho, \sigma \in \mathcal{D}(\mathcal{H})$ and $\forall \alpha \in (0, 1) \cup (1, \infty)$, $D_\alpha(\rho \parallel \sigma)$ satisfies the following,*

$$D_\alpha(\rho \parallel \sigma) \geq \overline{D}_\alpha(\rho \parallel \sigma) \geq D_\alpha^{\mathbb{M}}(\rho \parallel \sigma). \quad (29)$$

Proof. Fact 28 shows that

$$D_\alpha(\rho \parallel \sigma) = \frac{1}{n} D_\alpha(\rho^{\otimes n} \parallel \sigma^{\otimes n}) \geq \frac{1}{n} D_\alpha^{\mathbb{M}}(\rho^{\otimes n} \parallel \sigma^{\otimes n}). \quad (30)$$

Taking the limit $n \rightarrow \infty$ and using Fact 26, we have

$$D_\alpha(\rho \parallel \sigma) \geq \lim_{n \rightarrow \infty} \frac{1}{n} D_\alpha^{\mathbb{M}}(\rho^{\otimes n} \parallel \sigma^{\otimes n}) = \overline{D}_\alpha(\rho \parallel \sigma). \quad (31)$$

Also, Fact 25 implies

$$\overline{D}_\alpha(\rho \parallel \sigma) \geq D_\alpha^{\mathbb{M}}(\rho \parallel \sigma). \quad (32)$$

\blacksquare

III. DISCUSSION ON VARIATIONAL LOWER-BOUND ON DIVERGENCE AND ITS APPLICATION

In this section, we start with discussing the usefulness of a “change of measure” based variational lower-bounds of divergence, both in classical and quantum settings, by portraying two toy examples. Further, we observe that the upper-bound obtained in the classical example can be further tightened by using a variational lower-bound of Rényi divergence. To perform the same in the quantum scenario, we require a variational lower-bound for Petz quantum Rényi divergence and later in this section, we give a new proof for the variational lower-bound for Petz quantum Rényi divergence.

Variational representation for KL-divergences (as discussed in Fact 5) has several applications in learning theory and related areas [15]–[17], [21], [47]. In particular, consider a scenario where we want to analyze a function $-\beta < f(X) < \beta$, under $X \sim P$. However, the analysis under P might be hard to perform. The variational form discussed in Fact 5 along with Hoeffding’s lemma mentioned in Fact 7 (or the sub-Gaussianity assumption mentioned in Definition 5, if we don’t assume the function to be bounded), are beneficial in such situations. Instead of analyzing $f(X)$ under P , we will analyze it under some distribution Q such that $P \ll Q$. Even though this “change of measure” may allow easier analysis of $f(X)$ under Q it will also force us to incur a penalty in terms of the divergence between P and Q . Formally, consider the following series of inequalities for any $\lambda > 0$,

$$\begin{aligned} \mathbb{E}_P[f(X)] &= \frac{1}{\lambda} \mathbb{E}_P[\lambda f(X)] \\ &\stackrel{a}{\leq} \frac{1}{\lambda} \left(\log \mathbb{E}_Q[e^{\lambda f(X)}] + D^c(P \parallel Q) \right) \\ &\stackrel{b}{\leq} \mathbb{E}_Q[f(X)] + \frac{D^c(P \parallel Q)}{\lambda} + \frac{\lambda \beta^2}{2}, \end{aligned} \quad (33)$$

where a follows from Fact 5 and b follows from Fact 7. Thus, optimizing over the choice of λ , we get,

$$\mathbb{E}_P[f(X)] \leq \mathbb{E}_Q[f(X)] + \beta \sqrt{2D^c(P \parallel Q)}.$$

Similarly, by optimizing over λ (for $\lambda < 0$) we have,

$$\mathbb{E}_P[f(X)] \geq \mathbb{E}_Q[f(X)] - \beta \sqrt{2D^c(P \parallel Q)}. \quad (34)$$

From the derivation of (33) it follows that Fact 5 and Fact 7 (Definition 5, if we don’t assume the function to be bounded) together are equivalent to change of measure. Before mentioning the quantum version of this discussion we note here that to obtain the bound on $\mathbb{E}_P[f(X)]$ we only needed a lower-bound on $D^c(P \parallel Q)$ and this is one of the motivation for our lower-bound discussed in Lemma 3. Further, the proof of the upper-bound on $\mathbb{E}_P[f(X)]$ crucially needs Hoeffding’s lemma (Fact 7). This leads us to prove a quantum version of Hoeffding’s lemma mentioned below.

Lemma 1 (Quantum Hoeffding's lemma). *Given a quantum state $\rho \in \mathcal{D}(\mathcal{H})$ and a self-adjoint operator $L \in \mathcal{B}(\mathcal{H})$ such that $a\mathbb{I} \preceq L \preceq b\mathbb{I}$ (where $a \geq b$ and $a, b \in \mathbb{R}$ and \mathbb{I} denotes the projection over \mathcal{H}). Then, $\forall \lambda \in \mathbb{R}$,*

$$\log \text{Tr} \left[e^{\lambda(L - \text{Tr}[L\rho]\mathbb{I})} \rho \right] \leq \frac{\lambda^2(b-a)^2}{8}, \quad (35)$$

or equivalently,

$$\log \text{Tr} [e^{\lambda L} \rho] \leq \lambda \text{Tr}[L\rho] + \frac{\lambda^2(b-a)^2}{8}. \quad (36)$$

Proof. See Appendix A for the proof. ■

The following corollary directly follows from Lemma 1.

Corollary 1. *Every $L \in \mathcal{B}(\mathcal{H})$ is μ^2 -sub-Gaussian. That is for every $\rho \in \mathcal{D}(\mathcal{H})$,*

$$\log \text{Tr} \left[e^{\lambda(L - \text{Tr}[L\rho]\mathbb{I})} \rho \right] \leq \frac{\lambda^2 \mu^2}{2}, \quad (37)$$

where $\mu := \frac{\|L\|_\infty}{2}$ (where $\|\cdot\|_\infty$ is defined in Definition 6) and $\lambda \in \mathbb{R}$.

We now discuss a quantum version of the problem discussed above in the context of (34). Consider a quantum state $\rho \in \mathcal{D}(\mathcal{H})$ and an observable $L \in \mathcal{L}(\mathcal{H})$ such that $-\mu\mathbb{I} \preceq L \preceq \mu\mathbb{I}$, where $\mu < \infty$. Suppose we want to analyze L under ρ . However, this might be hard to analyze. As discussed in the classical setting, we will perform this analysis by using a ‘change of measure’. For this, we will require a variational form of $D(\rho\|\sigma)$ mentioned in Fact 15, along with a quantum version of Hoeffding's lemma (quantum sub-Gaussianity assumption mentioned in Definition 13, if we don't assume the observable L to be bounded) mentioned above. This change of measure may allow an easier analysis of L under σ , but it will also force us to incur a loss in terms of $D(\rho\|\sigma)$. Formally, consider the following series of inequalities for any $\lambda > 0$,

$$\begin{aligned} \text{Tr}[L\rho] &= \frac{1}{\lambda} \text{Tr}[\lambda L\rho] \\ &\stackrel{a}{\leq} \frac{1}{\lambda} (D(\rho\|\sigma) + \log \text{Tr}[e^{\lambda L}(\sigma)]), \\ &\stackrel{b}{\leq} \frac{1}{\lambda} \left(D(\rho\|\sigma) + \lambda \text{Tr}[L\sigma] + \frac{\lambda^2 \mu^2}{2} \right), \\ &= \frac{D(\rho\|\sigma)}{\lambda} + \frac{\lambda \mu^2}{2} + \text{Tr}[L\sigma], \end{aligned} \quad (38)$$

where a follows from Facts 15 and 24 and b follows from Lemma 1. Thus, in (38), if we minimize the RHS over λ , (for $\lambda > 0$) we get the following bound,

$$\text{Tr}[L\rho] \leq \mu \sqrt{2D(\rho\|\sigma)} + \text{Tr}[L\sigma],$$

Similarly, by optimizing over λ (for $\lambda < 0$) we have,

$$\text{Tr}[L\rho] \geq \text{Tr}[L\sigma] - \mu \sqrt{2D(\rho\|\sigma)}.$$

From the derivation of (38) it follows that Fact 18 and Lemma 1 (Definition 13, if we don't assume the observable L to be bounded) together are equivalent to change of measure.

We can derive a tighter bound in (33) (likewise in (38)), if instead of using a variational form (a lower-bound is sufficient) for the KL divergence (quantum KL divergence), we use a variational form (see Fact 6) for the Rényi divergence (quantum Rényi divergence).

Measured Rényi divergence has a variational lower-bound, which can be used in the context of the problem discussed above. However, it is not so easy to calculate the measured Rényi divergence. Instead of this, we can employ the modified sandwiched quantum Rényi divergence.

Indeed, for $\alpha \in (1, \infty)$ it follows via the data-processing inequality for sandwiched quantum Rényi divergence and Fact 17. But, as mentioned in Fact 22, the sandwiched Rényi divergence satisfies the data processing inequality only for $\alpha \in [\frac{1}{2}, \infty)$. Therefore, we can not use Sandwiched quantum Rényi divergence for $\alpha \in (0, \frac{1}{2})$ to get a variational lower-bound for the variational form for $D_\alpha^{\text{M}}(\cdot\|\cdot)$. Further, as mentioned in Fact 16 (variational form for sandwiched quantum Rényi divergence), the terms involving ρ (the original state) and σ (change in measure state) are sitting inside the log term. Therefore, it is not very clear how to use this variational form for sandwiched quantum Rényi divergence to further tighten this bound. To resolve this problem, we employ modified sandwiched quantum Rényi divergence instead of sandwiched quantum Rényi divergence. In fact, as mentioned in Fact 29, modified sandwiched quantum Rényi divergence is upper bounded by Petz quantum Rényi divergence, which implies that modified sandwiched quantum Rényi divergence gives a better bound than Petz quantum Rényi divergence.

Although Fact 23 holds for the iid case, as a generalization of Fact 26, the following lemma shows that modified sandwiched quantum Rényi divergence gives a good approximation of measured Rényi divergence.

Lemma 2. *Consider sequences of states $\rho_n, \sigma_n \in \mathcal{D}(\mathcal{A})$ with $\rho_n \ll \sigma_n$. We denote the minimum eigenvalue of σ_n by λ_n . When $\log \lambda_n$ behaves as a polynomial order for n , $\forall \alpha \in (0, 1) \cup (1, \infty)$ we have the following,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \bar{D}_\alpha(\rho_n \| \sigma_n) = \lim_{n \rightarrow \infty} \frac{1}{n} D_\alpha^{\mathbb{M}}(\rho_n \| \sigma_n). \quad (39)$$

This lemma suggest the following. When the system size is not so large, we can use measured Rényi divergence for the calculation of the variational form. However, when the system size is too large to calculate measured Rényi divergence, we can use modified sandwiched quantum Rényi divergence as an upper bound of the variational form, which is close to the variational form.

Proof. It is sufficient to show

$$\lim_{n \rightarrow \infty} \frac{1}{n} \tilde{D}_\alpha(\rho_n \| \sigma_n) = \lim_{n \rightarrow \infty} \frac{1}{n} D_\alpha^{\mathbb{M}}(\rho_n \| \sigma_n) \quad (40)$$

for $\forall \alpha \in [1/2, 1) \cup (1, \infty)$ because the part $\alpha \in (0, 1/2)$ follows from the part with $\alpha \in (1/2, 1)$. Since the inequality $\frac{1}{n} \tilde{D}_\alpha(\rho_n \| \sigma_n) \geq \frac{1}{n} D_\alpha^{\mathbb{M}}(\rho_n \| \sigma_n)$ holds, it is sufficient to show

$$\lim_{n \rightarrow \infty} \frac{1}{n} \tilde{D}_\alpha(\rho_n \| \sigma_n) \leq \lim_{n \rightarrow \infty} \frac{1}{n} D_\alpha^{\mathbb{M}}(\rho_n \| \sigma_n). \quad (41)$$

for $\forall \alpha \in [1/2, 1) \cup (1, \infty)$.

We choose a polynomial $p(n)$ such that $-\log \lambda_n \leq p(n)$ and $p(n)$ is an positive integer. We choose the spectral decomposition of σ_n as

$$\sigma_n := \sum_j e^{-s_j} E_j \quad (42)$$

We define the operator $\tilde{\sigma}_n$ as

$$\tilde{\sigma}_n := \sum_j e^{-\lceil \frac{s_j}{p(n)} \rceil p(n)} E_j. \quad (43)$$

Hence, we have

$$e^{\frac{1}{p(n)}} \tilde{\sigma}_n \geq \sigma_n \geq \tilde{\sigma}_n. \quad (44)$$

Also, the number of eigenvalues of $\tilde{\sigma}_n$ is at most $p(n) + 1$. The relation (44) implies

$$\tilde{D}_\alpha(\rho_n \| \tilde{\sigma}_n) \geq \tilde{D}_\alpha(\rho_n \| \sigma_n) \geq \tilde{D}_\alpha(\rho_n \| \tilde{\sigma}_n) - \log p(n). \quad (45)$$

Measured Rényi divergence also satisfies the relation;

$$D_\alpha^{\mathbb{M}}(\rho_n \| \tilde{\sigma}_n) \geq D_\alpha^{\mathbb{M}}(\rho_n \| \sigma_n) \geq D_\alpha^{\mathbb{M}}(\rho_n \| \tilde{\sigma}_n) - \log p(n). \quad (46)$$

In addition, we have [40, Eq. (3.153) and the next equation of Eq. (3.156)]

$$\log(p(n) + 1) + D_\alpha^{\mathbb{M}}(\rho_n \| \tilde{\sigma}_n) \geq \tilde{D}_\alpha(\rho_n \| \tilde{\sigma}_n). \quad (47)$$

The combination of (45), (46), and (47) implies

$$2 \log(p(n) + 1) + D_\alpha^{\mathbb{M}}(\rho_n \| \sigma_n) \geq \log(p(n) + 1) + D_\alpha^{\mathbb{M}}(\rho_n \| \tilde{\sigma}_n) \geq \tilde{D}_\alpha(\rho_n \| \tilde{\sigma}_n) \geq \tilde{D}_\alpha(\rho_n \| \sigma_n). \quad (48)$$

Thus, we have

$$\frac{2 \log(p(n) + 1)}{n} + \frac{1}{n} D_\alpha^{\mathbb{M}}(\rho_n \| \sigma_n) \geq \frac{1}{n} \tilde{D}_\alpha(\rho_n \| \sigma_n). \quad (49)$$

Taking the limit in (49), we obtain (41).

The idea for the discretization $\tilde{\sigma}_n$ of σ_n was essentially used in [48]. This proves Lemma 2. ■

Below *without invoking data-processing inequality and Facts 17 and 28*, we present an alternative simple proof for the following lemma related to a variational lower-bound. The proof of Lemma 3 follows from the Hölder's inequality (Fact 21) and Araki-Lieb-Thirring trace inequality (see Fact 18). The following proof can be considered as another proof for Facts 17 and 28.

Lemma 3. *Consider $\rho, \sigma \in \mathcal{D}(\mathcal{A})$ with $\rho \ll \sigma$. Then, $\forall \alpha \in (0, 1) \cup (1, \infty)$ we have the following,*

$$D_\alpha(\rho \| \sigma) \geq \sup_{\substack{H \in \mathcal{L}(\mathcal{A}): \\ H > 0}} \left\{ \frac{\alpha}{\alpha - 1} \log \text{Tr} \left[e^{(\alpha-1)H} \rho \right] - \log \text{Tr} \left[e^{\alpha H} \sigma \right] \right\} = D_\alpha^{\mathbb{M}}(\rho \| \sigma). \quad (50)$$

Proof. We divide this proof into two cases. In the first case $\alpha \in (0, 1)$ and in the second case $\alpha \in (1, \infty)$.

Case 1 : $\alpha \in (0, 1)$

For any strictly positive operator H , consider the following series of inequalities,

$$\begin{aligned} \text{Tr}[\rho^\alpha \sigma^{1-\alpha}] &\stackrel{a}{=} \text{Tr}\left[e^{\frac{-(1-\alpha)\alpha H}{2}} \rho^\alpha e^{\frac{-(1-\alpha)\alpha H}{2}} e^{\frac{(1-\alpha)\alpha H}{2}} \sigma^{1-\alpha} e^{\frac{(1-\alpha)\alpha H}{2}}\right] \\ &\stackrel{b}{\leq} \left(\text{Tr}\left[\left(e^{\frac{-(1-\alpha)\alpha H}{2}} \rho^\alpha e^{\frac{-(1-\alpha)\alpha H}{2}}\right)^{\frac{1}{\alpha}}\right]\right)^\alpha \left(\text{Tr}\left[\left(e^{\frac{(1-\alpha)\alpha H}{2}} \sigma^{1-\alpha} e^{\frac{(1-\alpha)\alpha H}{2}}\right)^{\frac{1}{1-\alpha}}\right]\right)^{1-\alpha} \\ &\stackrel{c}{\leq} \left(\text{Tr}\left[e^{\frac{-(1-\alpha)H}{2}} \rho e^{\frac{-(1-\alpha)H}{2}}\right]\right)^\alpha \left(\text{Tr}\left[e^{\frac{\alpha H}{2}} \sigma e^{\frac{\alpha H}{2}}\right]\right)^{1-\alpha} \\ &= \left(\text{Tr}\left[e^{(\alpha-1)H} \rho\right]\right)^\alpha \left(\text{Tr}\left[e^{\alpha H} \sigma\right]\right)^{1-\alpha}, \end{aligned}$$

where a follows from cyclicity of trace, b follows from Fact 21, c follows from Fact 18 and $\frac{1}{\alpha} > 1$. Since $(1 - \alpha) > 0$, we have,

$$\frac{1}{1-\alpha} \log \text{Tr}[\rho^\alpha \sigma^{1-\alpha}] \leq \frac{\alpha}{1-\alpha} \log \text{Tr}\left[e^{(\alpha-1)H} \rho\right] + \log \text{Tr}\left[e^{\alpha H} \sigma\right].$$

Thus, from the above inequality, we have,

$$D_\alpha(\rho \| \sigma) = \frac{1}{\alpha-1} \log \text{Tr}[\rho^\alpha \sigma^{1-\alpha}] \geq \frac{\alpha}{\alpha-1} \log \text{Tr}\left[e^{(\alpha-1)H} \rho\right] - \log \text{Tr}\left[e^{\alpha H} \sigma\right],$$

which implies the inequality in (50).

When we fix a basis $\{|u_j\rangle\}$ and restrict the range of H into the Hermitian matrices diagonal under the basis $\{|u_j\rangle\}$, Fact 6 implies

$$D_\alpha(\mathcal{E}_{\{|u_j\rangle\}}(\rho) \| \mathcal{E}_{\{|u_j\rangle\}}(\sigma)) = \sup_{H>0} \left\{ \frac{\alpha}{\alpha-1} \log \text{Tr}\left[e^{(\alpha-1)H} \rho\right] - \log \text{Tr}\left[e^{\alpha H} \sigma\right] \mid H \text{ is diagonal to } \{|u_j\rangle\} \right\},$$

where $\mathcal{E}_{\{|u_j\rangle\}}(\rho) := \sum_j |u_j\rangle \langle u_j| \rho |u_j\rangle \langle u_j|$. Considering the supremum for the choice of the basis $\{|u_j\rangle\}$ in the above equation, we obtain the equality in (50).

Case 2 : $\alpha \in (1, \infty)$

For any strictly positive operator H , consider the following series of inequalities,

$$\begin{aligned} \text{Tr}\left[e^{(\alpha-1)H} \rho\right] &\stackrel{a}{=} \text{Tr}\left[\sigma^{\frac{-(1-\alpha)}{2\alpha}} e^{(\alpha-1)H} \sigma^{\frac{-(1-\alpha)}{2\alpha}} \sigma^{\frac{(1-\alpha)}{2\alpha}} \rho \sigma^{\frac{(1-\alpha)}{2\alpha}}\right] \\ &\stackrel{b}{\leq} \left(\text{Tr}\left[\left(\sigma^{\frac{-(1-\alpha)}{2\alpha}} e^{(\alpha-1)H} \sigma^{\frac{-(1-\alpha)}{2\alpha}}\right)^{\frac{\alpha}{\alpha-1}}\right]\right)^{\frac{\alpha-1}{\alpha}} \left(\text{Tr}\left[\left(\sigma^{\frac{(1-\alpha)}{2\alpha}} \rho \sigma^{\frac{(1-\alpha)}{2\alpha}}\right)^\alpha\right]\right)^{\frac{1}{\alpha}} \\ &\stackrel{c}{\leq} \left(\text{Tr}\left[\sigma^{\frac{1}{2}} e^{\alpha H} \sigma^{\frac{1}{2}}\right]\right)^{\frac{\alpha-1}{\alpha}} \left(\text{Tr}\left[\sigma^{\frac{(1-\alpha)}{2}} \rho^\alpha \sigma^{\frac{(1-\alpha)}{2}}\right]\right)^{\frac{1}{\alpha}} \\ &= \left(\text{Tr}\left[e^{\alpha H} \sigma\right]\right)^{\frac{\alpha-1}{\alpha}} \left(\text{Tr}\left[\rho^\alpha \sigma^{1-\alpha}\right]\right)^{\frac{1}{\alpha}}, \end{aligned}$$

where a follows from cyclicity of trace, b follows from Fact 21, c follows from Fact 18 and because $\frac{\alpha}{\alpha-1} > 1$. Since $\frac{\alpha}{\alpha-1} > 0$, we have,

$$\begin{aligned} \frac{\alpha}{\alpha-1} \log \text{Tr}\left[e^{(\alpha-1)H} \rho\right] &\leq \log \text{Tr}\left[e^{\alpha H} \sigma\right] + \frac{1}{\alpha-1} \log \text{Tr}\left[\rho^\alpha \sigma^{1-\alpha}\right] \\ &= \log \text{Tr}\left[e^{\alpha H} \sigma\right] + D_\alpha(\rho \| \sigma). \end{aligned}$$

Thus, from the above inequality, we have,

$$D_\alpha(\rho \| \sigma) = \frac{1}{\alpha-1} \log \text{Tr}[\rho^\alpha \sigma^{1-\alpha}] \geq \frac{\alpha}{\alpha-1} \log \text{Tr}\left[e^{(\alpha-1)H} \rho\right] - \log \text{Tr}\left[e^{\alpha H} \sigma\right],$$

which implies the inequality in (50). In the same way, we obtain the equality in (50). This proves Lemma 3. ■

Further, we give the following variational lowerbound for modified sandwiched quantum Rényi divergence.

Lemma 4. Let $\rho, \sigma \in \mathcal{D}(\mathcal{A})$. Then, for $\alpha \in (0, 1) \cup (1, \infty)$, we have,

$$\overline{D}_\alpha(\rho \| \sigma) \geq \sup_{\substack{H \in \mathcal{L}(\mathcal{A}): \\ H > 0}} \left\{ \frac{\alpha}{\alpha-1} \log \text{Tr}\left[e^{(\alpha-1)H} \rho\right] - \log \text{Tr}\left[e^{\alpha H} \sigma\right] \right\}, \quad \forall \alpha \in (0, 1) \cup (1, \infty).$$

Proof. The proof follows from the lower-bound of Facts 29 and 17. ■

In the subsequent sections, we will bound functions which are generalization errors of a quantum learning algorithm in terms of the Petz and modified sandwiched quantum Rényi divergence and thus, obtain a family of upper-bounds. As a special case our bounds will recover the bound obtained in [8]. Therefore, in the section below, we first discuss a framework developed by Caro et al. in [8] for quantum learning algorithms.

Before describing the framework discussed in this paper, we first describe a version of the classical learning scenario [16], [21], [25]. This we believe, will help us to get the motivation behind the framework proposed in [8].

A. Generalized Classical Learning Framework

Although generalized classical learning framework covers various settings, we begin with supervised learning as a typical example. Here, we consider an input space \mathcal{X} and an output space \mathcal{Y} , where a training data point is given as a pair (X, Y) with input $X \in \mathcal{X}$ and output $Y \in \mathcal{Y}$. Given n training data points $(X_1, Y_1), \dots, (X_n, Y_n)$, the learner aims to output a hypothesis that explains this data. A common type of hypothesis is a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, such as an affine function $f(x) = ax + b$. In this case, the parameters a and b are often determined by minimizing the prediction error on the training data using the minimum mean square error method:

$$\operatorname{argmin}_{(a,b)} \sum_{j=1}^n (f(x_j) - y_j)^2 = \operatorname{argmin}_{(a,b)} \sum_{j=1}^n (ax_j + b - y_j)^2.$$

This supervised learning setup can be generalized to a broader learning framework.

More generally, in a learning scenario, training data points reside in a sample space \mathcal{Z} . In the supervised learning example, $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, with each data point $z \in \mathcal{Z}$ being an input-output pair (x, y) . We consider a training data sample $S = \{Z_i\}_{i=1}^n$ consisting of n independent and identically distributed (i.i.d.) data points drawn from a distribution P . A learning algorithm \mathcal{A} takes S as input and produces a hypothesis $w \in \mathcal{W}$. In the supervised learning example, \mathcal{W} is the set of functions from \mathcal{X} to \mathcal{Y} . Since the algorithm \mathcal{A} 's output hypothesis w is conditioned on the training data S , \mathcal{A} can generally be represented by a conditional probability distribution $P_{W|S}^{\mathcal{A}}$.

Once the algorithm produces a hypothesis w , its performance is evaluated based on the training data S using a loss function $l : \mathcal{W} \times \mathcal{Z} \rightarrow \mathbb{R}$. In the supervised learning example, $l((a, b), (x, y)) = (ax + b - y)^2$. A primary goal in a learning scenario is to design an algorithm \mathcal{A} that minimizes the empirical loss, defined as

$$\hat{l}_S(w) := \frac{1}{n} \sum_{i=1}^n l(w, Z_i). \quad (51)$$

The expectation of this empirical loss, $\mathbb{E}_{S,W}[\hat{l}_S(W)]$, is the expected empirical loss, which can be readily estimated as it is expected to be close to the observed $\hat{l}_S(w)$.

However, because the algorithm \mathcal{A} determines the hypothesis based on the training data S , it can inadvertently learn patterns specific to this data, leading to a dependency (bias) between W and S that is undesirable in real-world applications. To address this, we need to evaluate the learned hypothesis on independently generated data. For a more rigorous analysis, we consider splitting the available data into a training set S_{tr} and a testing set S_{te} , where $(S_{te}, S_{tr}) = \{(Z_{te,i}, Z_{tr,i})\}_{i=1}^n$ is a sequence of n i.i.d. pairs drawn from a joint distribution $P_{Z_{te}, Z_{tr}}$. The training data S_{tr} is used by the algorithm \mathcal{A} to produce a hypothesis W , while the testing data S_{te} is used to evaluate the loss. Since \mathcal{A} outputs W conditioned on S_{tr} , it is generally represented by the conditional probability distribution $P_{W|S_{tr}}^{\mathcal{A}}$.

After the algorithm produces a hypothesis w , its performance is evaluated on the test data S_{te} using the loss function $l : \mathcal{W} \times \mathcal{Z} \rightarrow \mathbb{R}$. Thus, a key objective is to design an algorithm \mathcal{A} that minimizes the empirical loss on the test data:

$$\hat{l}_{S_{te}}(w) := \frac{1}{n} \sum_{i=1}^n l(w, Z_{te,i}).$$

This is the empirical test loss of \mathcal{A} for hypothesis w on S_{te} . We are also interested in the average performance of \mathcal{A} , given by the expected test empirical loss $\mathbb{E}_{S_{te}, S_{tr}, W}[\hat{l}_{S_{te}}(W)]$. However, the correlation between the test data S_{te} and the training data S_{tr} , along with the fact that \mathcal{A} 's output W depends solely on S_{tr} , can introduce a bias between W and S_{te} . Indeed, W depends on S_{te} through S_{tr} , forming a Markov chain $W - S_{tr} - S_{te}$, which is generally undesirable.

To mitigate the issue of dependency, we consider a scenario with an independent test set \bar{S}_{te} , which is independent of the training set \bar{S}_{tr} and the learned hypothesis \bar{W} . The joint distribution is then $P_{\bar{S}_{te}, \bar{S}_{tr}, \bar{W}} = P_{\bar{S}_{te}} P_{\bar{S}_{tr}, \bar{W}}$. Since the hypothesis \bar{W} is derived from the training data \bar{S}_{tr} , we have $P_{\bar{S}_{tr}, W}(s_{tr}, w) = P_{\bar{S}_{tr}, \bar{W}}(s_{tr}, w)$, implying $\bar{W} - \bar{S}_{tr} \perp \bar{S}_{te}$. The expected loss on this independent test set is calculated as:

$$\begin{aligned} \mathbb{E}_{\bar{S}_{te}, \bar{S}_{tr}, \bar{W}}[\hat{l}_{\bar{S}_{te}}(\bar{W})] &= \mathbb{E}_{\bar{S}_{te}, \bar{W}}[\hat{l}_{\bar{S}_{te}}(\bar{W})] \\ &= \mathbb{E}_{\bar{W}} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Z_{te,i}} l(\bar{W}, \bar{Z}_{te,i}) \right] \\ &= \mathbb{E}_{\bar{W}}[\mathbb{E}_{\bar{Z}_{te}} l(\bar{W}, \bar{Z}_{te})]. \end{aligned} \quad (52)$$

The quantity on the right-hand side of Equation (52) is known as the expected true loss of \mathcal{A} , representing the average performance of the algorithm over the entire data distribution.

Because of the previous discussions, we understand that a good learning algorithm must be unbiased. Specifically, the difference between its expected empirical loss and its expected true loss should be small. As a result, a learning algorithm achieves good generalization in expectation when this difference is small in magnitude. Therefore, in this scenario of classical learning, we define the generalization error as follows:

$$\text{gen}(w, S_{te}) := \hat{l}_{S_{te}}(w) - \mathbb{E}_{\bar{W}}[\mathbb{E}_{\bar{Z}_{te}} l(\bar{W}, \bar{Z}_{te})]. \quad (53)$$

Note that in the above equation, the generalization error $\text{gen}(\cdot, \cdot)$ is defined in terms of the whole sample S_{te} . Thus, it can be called a *whole sample-based generalization error*. However, Equation (53) can be rewritten as:

$$\text{gen}(w, S_{te}) = \frac{1}{n} \sum_{i=1}^n \text{gen}_{\text{ind}}(w, Z_i), \quad (54)$$

where $\text{gen}_{\text{ind}}(w, Z_i) := l(w, Z_i) - \mathbb{E}_{\bar{Z}_{te}}[l(w, \bar{Z}_{te})]$ is denoted as the *individual sample generalization error* corresponding to the i -th sample Z_i .

Then, we define the expected generalization error as follows:

$$\begin{aligned} \bar{\text{gen}} &:= \mathbb{E}_{S_{te}, S_{tr}, W}[\hat{l}_{S_{te}}(W)] - \mathbb{E}_{\bar{W}}[\mathbb{E}_{\bar{Z}_{te}} l(\bar{W}, \bar{Z}_{te})] \\ &= \mathbb{E}_{S_{te}, S_{tr}, W}[\hat{l}_{S_{te}}(W)] - \mathbb{E}_{S_{te}, S_{tr}, \bar{W}}[\hat{l}_{S_{te}}(\bar{W})], \end{aligned}$$

where the first expectation is calculated with respect to the joint distribution $P_{S_{te}, S_{tr}, W} = P_{S_{te}, S_{tr}} P_{W|S_{tr}}$ and the second expectation is calculated with respect to the marginal distribution i.e. $P_{S_{te}} P_{S_{tr}} P_W$.

In a special case of the above learning scenario, the test and training data are perfectly correlated, i.e., S_{te} and S_{tr} can be represented by a single variable S . Then, the learning scenario becomes a weaker version of classical learning, studied in [13], [16], [21], [25] and discussed in the introduction of this manuscript.

B. Quantum Learning Framework

To get a motivation for defining an input to a quantum learning algorithm, let us first define an input to a classical learning algorithm in a quantum formalism. This can be given as the following state,

$$\rho_{\text{class}} := \mathbb{E}_{S_{te}, S_{tr} \sim P_{S_{te}, S_{tr}}^n} \left[|S_{te}\rangle^{\text{Te}} \langle S_{te}| \otimes |S_{tr}\rangle^{\text{Tr}} \langle S_{tr}| \right].$$

In the quantum learning scenario, a quantum learning algorithm \mathcal{A}_Q , apart from the systems Te and Tr, also has access to another quantum system whose state is $\rho(S_{te}, S_{tr}) \in \mathcal{H}^D$ (where \mathcal{H}^D is denoted as data Hilbert space). Thus, the input to \mathcal{A}_Q can be represented by the following classical quantum state,

$$\rho := \mathbb{E}_{S_{te}, S_{tr} \sim P_{S_{te}, S_{tr}}^n} [|S_{te}\rangle \langle S_{te}| \otimes |S_{tr}\rangle \langle S_{tr}| \otimes \rho(S_{te}, S_{tr})]. \quad (55)$$

Typically in the classical learning literature [13], [16], [21], [25], the testing and training data are assumed to be perfectly correlated, i.e. S_{te} and S_{tr} can be represented by a single random variable S . Likewise, in [8], it is assumed that S_{te} and S_{tr} are perfectly correlated. However, in [8], the authors observe that the analysis for the performance of any quantum learning algorithm where S_{te} and S_{tr} are not perfectly correlated, is very similar to the analysis in the case when S_{te} and S_{tr} are perfectly correlated. Therefore, throughout this manuscript, we build on this assumption.

Thus, under the above assumption, the input to a quantum learning algorithm \mathcal{A}_Q can be represented as the following state,

$$\rho := \mathbb{E}_{S \sim P^n} [|S\rangle \langle S| \otimes \rho(S)]. \quad (56)$$

Upon getting the input mentioned in (56), a quantum learning algorithm \mathcal{A}_Q applies certain POVMs and CP-TP maps over the quantum system residing in \mathcal{H}^D for learning. Due to the irreversible perturbation of $\rho(S)$ during learning, the processed state cannot be further used to evaluate the empirical loss. This is because we need the unperturbed quantum state to calculate the empirical loss.

The authors in [8] avoid this issue by bi-partitioning \mathcal{H}^D into \mathcal{H}^{te} and \mathcal{H}^{tr} i.e. $\mathcal{H}^D := \mathcal{H}^{te} \otimes \mathcal{H}^{tr}$, such that $\mathcal{H}^{te} \cong \mathbb{C}^{d_1}$, $\mathcal{H}^{tr} \cong \mathbb{C}^{d_2}$: $1 < d_1, d_2 < \infty$ are the test data Hilbert space and the train data Hilbert space respectively. Thus, $\forall s \in \mathcal{Z}^n$, $\rho(s) \in \mathcal{H}^{te} \otimes \mathcal{H}^{tr}$ now has two components i.e. the testing component residing in \mathcal{H}^{te} and the training component residing in \mathcal{H}^{tr} . The learning algorithm acts only on \mathcal{H}^{tr} while learning, but the loss is calculated over the whole data Hilbert space \mathcal{H}^D . In general, $\forall s \in \mathcal{Z}^n$, $\rho(s)$ might be correlated or even entangled across \mathcal{H}^{te} and \mathcal{H}^{tr} . We now formally illustrate a general quantum learning algorithm in [8].

The quantum learner \mathcal{A}_Q is presumed to have a collection of extractor POVMs $\left\{ \{E_s^{\mathcal{A}_Q}(w)\}_{w \in \mathcal{W}} \right\}_{s \in \mathcal{Z}^n}$ which acts over \mathcal{H}^{tr} , where \mathcal{W} is assumed to be a discrete measurable space which is called the Hypothesis space and each $w \in \mathcal{W}$ is a classical hypothesis given as a map $w : \mathcal{X} \rightarrow \mathcal{Y}$. Thus, \mathcal{W} can be considered as a subspace of $\mathcal{Y}^{\mathcal{X}}$. The quantum learner \mathcal{A}_Q uses this collection of POVMs to extract classical information from training data states.

Conditioned on the classical data $s := (z_1, \dots, z_n) \in \mathcal{Z}^n$, \mathcal{A}_Q performs the measurement using the POVM $\{E_s^{\mathcal{A}_Q}(w)\}_{w \in \mathcal{W}}$ and records the outcome w classically, which results the following post measurement state over \mathcal{H}^D ,

$$\rho^{\mathcal{A}_Q}(w, s) := \frac{\sqrt{\left(\mathbb{I}_{\mathcal{H}^{te}} \otimes E_s^{\mathcal{A}_Q}(w)\right) \rho(s)} \sqrt{\left(\mathbb{I}_{\mathcal{H}^{te}} \otimes E_s^{\mathcal{A}_Q}(w)\right)}}{\text{Tr}[E_s^{\mathcal{A}_Q}(w) \rho_{tr}(s)]}, \quad (57)$$

$\rho_{tr}(s) = \text{Tr}_{te}[\rho(s)]$ and $\forall w \in \mathcal{W}, s \in \mathcal{Z}^n$, we define a conditional probability distribution $P^{\mathcal{A}_Q}(w|s) := \text{Tr}[E_s^{\mathcal{A}_Q}(w) \rho_{tr}(s)]$.

The quantum learner also has an apriori access to a collection of quantum channels (CP-TP) $\{\Lambda_{w,s}^{\mathcal{A}_Q} : \mathcal{T}(\mathcal{H}^{tr}) \rightarrow \mathcal{T}(\mathcal{H}^{hyp})\}_{\substack{w \in \mathcal{W}, \\ s \in \mathcal{Z}^n}}$, where we denote \mathcal{H}^{hyp} to be a quantum hypothesis Hilbert space. Conditioned on both w and s , \mathcal{A}_Q now applies the channel $\Lambda_{w,s}$ on the post-measurement state $\rho^{\mathcal{A}_Q}(w, s)$ and the resultant state is given as follows:

$$\sigma^{\mathcal{A}_Q}(w, s) := (\mathbb{I}_{\mathcal{H}^{te}} \otimes \Lambda_{w,s})(\rho^{\mathcal{A}_Q}(w, s)).$$

Hence, the overall action of \mathcal{A}_Q over the data state ρ leads us to the following CQ state,

$$\begin{aligned} \sigma^{\mathcal{A}_Q} &= \mathbb{E}_{S \sim P^n} \left[|S\rangle\langle S| \otimes \mathbb{E}_{W \sim P^{\mathcal{A}_Q}(\cdot|S)} [|W\rangle\langle W| \otimes \sigma^{\mathcal{A}_Q}(W, S)] \right] \\ &= \mathbb{E}_{(W,S) \sim P_{WS}^{\mathcal{A}_Q}} [|W\rangle\langle W| \otimes |S\rangle\langle S| \otimes \sigma^{\mathcal{A}_Q}(W, S)] \\ &= \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \mathbb{E}_{S \sim P_{S|W}^{\mathcal{A}_Q}(\cdot|W)} [|W\rangle\langle W| \otimes |S\rangle\langle S| \otimes \sigma^{\mathcal{A}_Q}(W, S)], \end{aligned}$$

where $\forall (w, s) \in \mathcal{W} \times \mathcal{Z}^n$, $P_{WS}^{\mathcal{A}_Q}(w, s) := P^{\mathcal{A}_Q}(w|s)P^n(s)$, $P_W^{\mathcal{A}_Q}(w) := \sum_{s \in \mathcal{Z}^n} P_{WS}^{\mathcal{A}_Q}(w, s)$ and $P_{S|W}^{\mathcal{A}_Q}(s|w) := \frac{P_{WS}^{\mathcal{A}_Q}(w, s)}{P_W^{\mathcal{A}_Q}(w)}$ is the posterior distribution of the data given the hypothesis. In the following discussion, we define how to quantize the loss or error induced from the resultant state $\sigma^{\mathcal{A}_Q}$.

In a classical learning scenario, a loss function is generally defined as a map $l : \mathcal{W} \times \mathcal{Z}^n \rightarrow \mathbb{R}$. However, in the quantum learning scenario, since the data and hypothesis induced by the quantum learner are embedded into quantum states, we use observables (operators) and consider the expected value of the observables with respect to the quantum states that represent the data and the hypothesis induced from it. In [8], the authors consider a family of non-negative self-adjoint loss observables $\{\hat{L}(w, s) \in \mathcal{L}(\mathcal{H}^{te} \otimes \mathcal{H}^{hyp})\}_{(w,s) \in \mathcal{W} \times \mathcal{Z}^n}$. Using these loss observables, we now define the two types of loss in terms of average values of the loss observables.

Definition 14 (Empirical (observed) loss). *The expected empirical loss $\hat{l}_\rho(w, s)$ of \mathcal{A}_Q with an input ρ and observable $\hat{L}(w, s)$ is defined as follows,*

$$\hat{l}_\rho(w, s) := \text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)].$$

Definition 15 (Expected empirical (observed) loss [8, Definition 11]). *The expected empirical loss L_ρ of \mathcal{A}_Q with an input ρ and observables $\{\hat{L}(w, s)\}_{(w,s) \in \mathcal{W} \times \mathcal{Z}^n}$ is defined as follows,*

$$\hat{L}_\rho := \mathbb{E}_{(W,S) \sim P_{WS}^{\mathcal{A}_Q}} [\hat{l}_\rho(W, S)].$$

Definition 16 (True loss [8]). *The true loss $l_\rho^{(\text{old})}(w)$ (here we use the phrase (old) since we propose a novel definition for true loss in Definition 17) of \mathcal{A}_Q with an input ρ and observables $\{\hat{L}(w, s)\}_{s \in \mathcal{Z}^n}$ is defined as follows,*

$$l_\rho^{(\text{old})}(w) := \mathbb{E}_{\bar{S} \sim P^n} \left[\text{Tr} \left[\hat{L}(w, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, \bar{S}) \right) \right] \right],$$

Definition 17 (True loss **proposed**). *The true loss $l_\rho(w)$ of \mathcal{A}_Q with an input ρ and observables $\{\hat{L}(w, s)\}_{s \in \mathcal{Z}^n}$ is defined as follows,*

$$l_\rho(w) := \mathbb{E}_{\bar{S} \sim P^n} \left[\text{Tr} \left[\hat{L}(w, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w) \right) \right] \right],$$

where $\rho_{te}(\bar{S}) := \text{Tr}_{tr}[\rho(\bar{S})]$, $\sigma_{hyp}^{\mathcal{A}_Q}(w) := \text{Tr}_{te}[\sigma^{\mathcal{A}_Q}(w)]$, where for any $w \in \mathcal{W}$, $\sigma^{\mathcal{A}_Q}(w) := \mathbb{E}_{S \sim P_{S|W}^{\mathcal{A}_Q}(\cdot|w)} [\sigma^{\mathcal{A}_Q}(w, S)]$.

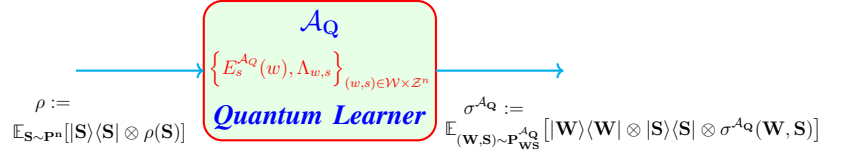


Fig. 1: Quantum learning algorithm structure proposed by [8].

Definition 18 (Expected true loss [8, Definition 12]). The expected true loss $L_\rho^{(\text{old})}$ of \mathcal{A}_Q with an input ρ and observables $\{\hat{L}(w, s)\}_{(w,s) \in \mathcal{W} \times \mathcal{Z}^n}$ is defined as follows,

$$\begin{aligned} L_\rho^{(\text{old})} &:= \mathbb{E}_{\bar{W} \sim P_W^{\mathcal{A}_Q}} [l_\rho^{(\text{old})}(\bar{W})] \\ &= \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^{\mathcal{A}_Q} \times P^n} \left[\text{Tr} \left[\hat{L}(\bar{W}, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(\bar{W}, \bar{S}) \right) \right] \right]. \end{aligned}$$

Definition 19 (Our definition for Expected true loss). The expected true loss L_ρ of \mathcal{A}_Q with an input ρ and observables $\{\hat{L}(w, s)\}_{(w,s) \in \mathcal{W} \times \mathcal{Z}^n}$ is defined as follows,

$$\begin{aligned} L_\rho &:= \mathbb{E}_{\bar{W} \sim P_W^{\mathcal{A}_Q}} [l_\rho(\bar{W})] \\ &= \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^{\mathcal{A}_Q} \times P^n} \left[\text{Tr} \left[\hat{L}(\bar{W}, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(\bar{W}) \right) \right] \right] \\ &= \mathbb{E}_{(\bar{W}, \bar{S}_{tr}, \bar{S}_{te}) \sim P_{W, S}^{\mathcal{A}_Q} \times P^n} \left[\text{Tr} \left[\hat{L}(\bar{W}, \bar{S}_{te}) \left(\rho_{te}(\bar{S}_{te}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(\bar{W}, \bar{S}_{tr}) \right) \right] \right]. \end{aligned}$$

Remark 1.

- Reference [8] adopts Definition 16 as the true loss. However, the quantum testing system is correlated to the quantum hypothesis system in the definition of $l_\rho^{(\text{old})}(w)$ after taking the average with respect to the classical variable \bar{S} . Even though the testing classical data S_{te} is perfectly correlated to the training classical data S_{tr} , the testing classical data \bar{S}_{te} needs to be independent of the training classical data \bar{S}_{tr} and the classical hypothesis \bar{W} in the definition of the true loss. When we have $\bar{W} - \bar{S}_{tr} \perp \bar{S}_{te}$ in the same way as the classical case, it is suitable to define $l_\rho(w)$ in Definition (17) by choosing \bar{S} to be \bar{S}_{te} .
- The expected true loss is defined as the expectation of the true loss. In Definition 18 of the expected true loss, the testing classical data is perfectly correlated with the training classical data. However, in the definition of the expected true loss, the testing classical data \bar{S}_{te} needs to be independent of the training classical data \bar{S}_{tr} and the classical hypothesis \bar{W} in the same way as the classical case. Hence, we adopt Definition 19. The detailed derivation of Definitions 17 and 19 is given in Subsection IV-C.

In practice, after the testing and training systems are well-prepared and the learner applies the quantum operation \mathcal{A}_Q , measuring the observable $L(w, s)$ is expected to yield an outcome close to the expected empirical loss \hat{L}_ρ . However, our primary interest lies in the expected true loss L_ρ of \mathcal{A}_Q . Consequently, we are interested in their difference, which defines the generalization error.

Definition 20 (Generalization Error). For any $(w, s) \in \mathcal{W} \times \mathcal{Z}^n$ we define the generalization error $\text{gen}(w, s)$ with an input ρ and observables $\{\hat{L}(w, s)\}_{(w,s) \in \mathcal{W} \times \mathcal{Z}^n}$ is defined as follows,

$$\text{gen}(w, s) = l_\rho(w) - \hat{l}_\rho(w, s). \quad (58)$$

Definition 21 (Expected generalization error). The expected generalization error $\overline{\text{gen}}$ of \mathcal{A}_Q with an input ρ and observables $\{\hat{L}(w, s)\}_{(w,s) \in \mathcal{W} \times \mathcal{Z}^n}$ is defined as follows,

$$\begin{aligned} \overline{\text{gen}} &:= \mathbb{E}_{(W, S) \sim P_{W, S}^{\mathcal{A}_Q}} [\text{gen}(W, S)] \\ &= \mathbb{E}_{(W, S) \sim P_{W, S}^{\mathcal{A}_Q}} [l_\rho(W) - \hat{l}_\rho(W, S)] \\ &= \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} [l_\rho(W)] - \mathbb{E}_{(W, S) \sim P_{W, S}^{\mathcal{A}_Q}} [\hat{l}_\rho(W, S)] \\ &= L_\rho - \hat{L}_\rho. \end{aligned}$$

Remark 2. See [8, Appendix C] for various applications of the quantum learning framework discussed above.

Remark 3. From Definitions 15, 19 and 21, we can come up with the classical notions of expected empirical loss, expected true loss, and expected generalization error by exploiting two special cases of the quantum learning framework mentioned above.

(a) The first special case is based on a very trivial assumption, i.e. all the involved quantum systems are trivial ($\mathcal{H}^D = \mathcal{H}^{hyp} = \mathbb{C}$). Then, for each $(w, s) \in \mathcal{W} \times \mathcal{Z}^n$, the loss observable $\hat{L}(w, s)$ becomes a scalar quantity (we consider \mathbb{R} instead of \mathbb{C} for simplicity) and can be thought of as the value of a classical loss function when passed (w, s) as an input to it.

(b) Secondly, for each $(w, s) \in \mathcal{W} \times \mathcal{Z}^n$, if we consider the loss observable of a special form $\hat{L}(w, s) := \hat{l}_s(w) \cdot \mathbb{I}^{\mathcal{H}^{te}} \otimes \mathcal{H}^{hyp}$, (where

$\times \mathcal{Z}^n$, if we consider the loss observable of a special form $\hat{L}(w, s)$ is defined in (51)), then, for each $(w, s) \in \mathcal{W} \times \mathcal{Z}^n$, we are left with the following,

$$\begin{aligned}
\text{gen}(w, s) &\stackrel{a}{=} l_\rho(w) - \hat{l}_\rho(w, s) \\
&\stackrel{b}{=} \mathbb{E}_{\bar{S} \sim P^n} \left[\text{Tr} \left[\hat{L}(w, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(w) \right) \right] \right] - \text{Tr}[\hat{L}(w, s) \sigma^{A_Q}(w, s)] \\
&= \mathbb{E}_{\bar{S} \sim P^n} \left[\hat{l}_{\bar{S}}(w) \text{Tr} \left[\left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(w) \right) \right] \right] - \hat{l}_s(w) \text{Tr}[\sigma^{A_Q}(w, s)] \\
&\stackrel{c}{=} \mathbb{E}_{\bar{S} \sim P^n} \left[\hat{l}_{\bar{S}}(w) \right] - \hat{l}_s(w) \\
&= \mathbb{E}_{\bar{Z} \sim P} [l(w, \bar{Z})] - \hat{l}_s(w),
\end{aligned} \tag{59}$$

where a follows from Definition 20, b follows from Definitions 14 and 17, c follows since for each $(w, s) \in \mathcal{W} \times \mathcal{Z}^n$, $\sigma^{A_Q}(w, s)$ and $\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)$ are quantum states. Observe that (59) coincides with the classical generalization error (mentioned in (53)) in absolute value.

C. Justification for the validity of the proposed definition of expected true loss (Definition 19)

To see the validity of Definition 19, we discuss the general one-shot setting. Assume that \mathcal{H}_{te} , \mathcal{H}_{tr} , and \mathcal{H}_{hyp} are the test space, the training space, and the hypothesis space. Initially, we have the initial state ρ on $\mathcal{H}_{te} \otimes \mathcal{H}_{tr}$. The learner applies a CP-TP map Λ from \mathcal{H}_{tr} to \mathcal{H}_{hyp} . Then, we have the final state $(id \otimes \Lambda)(\rho)$ on $\mathcal{H}_{te} \otimes \mathcal{H}_{hyp}$. The loss is determined by a Hermitian matrix L on $\mathcal{H}_{te} \otimes \mathcal{H}_{hyp}$ as follows. The empirical loss is given as

$$\text{Tr}[L(id \otimes \Lambda)(\rho)]. \tag{60}$$

The true loss is given as

$$\text{Tr}[L(\rho_{te} \otimes \Lambda(\rho_{tr}))]. \tag{61}$$

We now consider the case when \mathcal{H}_{te} , \mathcal{H}_{tr} , \mathcal{H}_{hyp} are composed of the classical and quantum parts, \mathcal{H}_{te}^C , \mathcal{H}_{tr}^C , \mathcal{H}_{hyp}^C and \mathcal{H}_{te}^Q , \mathcal{H}_{tr}^Q , \mathcal{H}_{hyp}^Q . Also, we assume that the classical information S of the test space is identical to the classical information of the training space. Hence, the initial state ρ on $\mathcal{H}_{te} \otimes \mathcal{H}_{tr}$ can be written as

$$\rho = \sum_s P_S(s) |s, s\rangle \langle s, s| \otimes \rho(s), \tag{62}$$

where $\rho(s)$ is a state on $\mathcal{H}_{te}^Q \otimes \mathcal{H}_{tr}^Q$. The learner's operation Λ is written as an instrument $\Gamma_s = (\Gamma_{w|s})$ as follows. For a state $\sum_s Q(s) |s\rangle \langle s| \otimes \sigma(s)$, we have,

$$\Lambda \left(\sum_s Q(s) |s\rangle \langle s| \otimes \sigma(s) \right) = \sum_s Q(s) |w\rangle \langle w| \otimes \Gamma_{w|s}(\sigma(s)). \tag{63}$$

Then, the resultant state is as follows,

$$(id \otimes \Lambda)(\rho) = \sum_{s, w} P_S(s) |s\rangle \langle s| \otimes |w\rangle \langle w| \otimes (id \otimes \Gamma_{w|s})(\rho(s)). \tag{64}$$

Also, the operator L can be written as,

$$L = \sum_{s', w'} |s'\rangle \langle s'| \otimes |w'\rangle \langle w'| \otimes L(w', s'). \tag{65}$$

Then, the empirical loss is given as,

$$\begin{aligned}
\text{Tr}[L(id \otimes \Lambda)(\rho)] &= \text{Tr} \left(\sum_{s', w'} |s'\rangle \langle s'| \otimes |w'\rangle \langle w'| \otimes L(w', s') \right) \left(\sum_{s, w} P_S(s) |s\rangle \langle s| \otimes |w\rangle \langle w| \otimes (id \otimes \Gamma_{w|s})(\rho(s)) \right) \\
&= \sum_{s, w} P_S(s) \text{Tr} [L(w, s) (id \otimes \Gamma_{w|s})(\rho(s))].
\end{aligned} \tag{66}$$

Since,

$$\Lambda(\rho_{tr}) = \sum_{s, w} P_S(s) |w\rangle \langle w| \otimes \Gamma_{w|s}(\rho_{tr}(s)), \tag{67}$$

$$\rho_{te} = \sum_{\bar{s}} P_S(\bar{s}) |\bar{s}\rangle \langle \bar{s}| \otimes \rho_{te}(\bar{s}), \tag{68}$$

the true loss is given as,

$$\begin{aligned}
& \text{Tr}[L(\rho_{te} \otimes \Lambda(\rho_{tr}))] \\
&= \text{Tr} \left[\sum_{s', w'} |s'\rangle\langle s'| \otimes |w'\rangle\langle w'| \otimes L(w', s') \right] \left(\left(\sum_{\bar{s}} P_S(\bar{s}) |\bar{s}\rangle\langle \bar{s}| \otimes \rho_{te}(\bar{s}) \right) \otimes \left(\sum_{s, w} P_S(s) |w\rangle\langle w| \otimes \Gamma_{w|s}(\rho_{tr}(s)) \right) \right) \\
&= \sum_{s, \bar{s}, w} P_S(s) P_S(\bar{s}) \text{Tr}[L(w, \bar{s})(\rho_{te, \bar{s}} \otimes \Gamma_{w|s}(\rho_{tr}(s)))] .
\end{aligned} \tag{69}$$

We define the probability $P_{W|S}$ as,

$$P_{W|S}^{\mathcal{A}_Q}(w|s) := \text{Tr}[\Gamma_{w|s}(\rho_{tr}(s))] . \tag{70}$$

We set $\sigma_{hyp}^{\mathcal{A}_Q}(w, S) := \frac{1}{P_{W|S}^{\mathcal{A}_Q}(w|s)} \Gamma_{w|s}(\rho_{tr}(s))$. Then, we can rewrite (69) as follows,

$$\begin{aligned}
\text{Tr}[L(\rho_{te} \otimes \Lambda(\rho_{tr}))] &= \sum_{s, \bar{s}, w} P_{W|S}^{\mathcal{A}_Q}(w, s) P_S(\bar{s}) \text{Tr}[L(w, \bar{s})(\rho_{te, \bar{s}} \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, S))] \\
&= \sum_{\bar{s}, w} P_W^{\mathcal{A}_Q}(w) P_S(\bar{s}) \text{Tr} \left[L(w, \bar{s})(\rho_{te, \bar{s}} \otimes \sum_s P_{S|W=w}^{\mathcal{A}_Q}(s) \sigma_{hyp}^{\mathcal{A}_Q}(w, S)) \right] \\
&= \sum_{\bar{s}, w} P_W^{\mathcal{A}_Q}(w) P_S(\bar{s}) \text{Tr} [L(w, \bar{s})(\rho_{te, \bar{s}} \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))] .
\end{aligned} \tag{71}$$

Therefore, the relation (71) coincides with Definition 19.

In particular, when the state $\rho_S(s)$ on $\mathcal{H}_{te}^Q \otimes \mathcal{H}_{tr}^Q$ is given a product state $\rho_{te}(s) \otimes \rho_{tr}(s)$, we have $\sigma^{\mathcal{A}_Q}(w, s) = \rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s)$. Since

$$\hat{l}_\rho(w, s) = \text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] = \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))], \tag{72}$$

the expected empirical loss is

$$\begin{aligned}
\hat{L}_\rho &= \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}_Q}} [\hat{l}_\rho(W, S)] \\
&= \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}_Q}} \text{Tr}[\hat{L}(W, S)(\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S))],
\end{aligned} \tag{73}$$

Hence, the difference between the expected empirical loss and the expected true loss defined by [8, Definition 12] is characterized by the difference between the joint distribution $P_{WS}^{\mathcal{A}_Q}$ and the product distribution $P_W^{\mathcal{A}_Q} \times P^n$.

V. DISCUSSION ON PREVIOUS WORKS

In this section, we first discuss an upper-bound on the expected generalization error for bounded loss functions/observables. We then use this bound to discuss the classical results obtained in [16], [17], [21] and the quantum results discussed in [8].

A. Classical learning paradigm

In a classical learning scenario, using (52), we can write the following form of expected generalization error $\mathbb{E}_{(W, S) \sim P_{WS}^A} [\text{gen}(W, S)]$ (where $\text{gen}(W, S)$ is defined in (53)),

$$\mathbb{E}_{(W, S) \sim P_{WS}^A} [\text{gen}(W, S)] = \mathbb{E}_{(W, S) \sim P_{WS}^A} [\hat{l}_S(W)] - \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^A \times P_S} [\hat{l}_{\bar{S}}(\bar{W})].$$

If we assume that the loss function is bounded, then, using the variational form for L_1 distance we can obtain an upper-bound on $\mathbb{E}_{(W, S) \sim P_{WS}^A} [\text{gen}(W, S)]$. In particular, assume that $\forall (w, z) \in \mathcal{W} \times \mathcal{Z}$, $l(w, z) < \tau$, where $\tau < \infty$. Then,

$$\begin{aligned}
\mathbb{E}_{(W, S) \sim P_{WS}^A} [\text{gen}(W, S)] &= \mathbb{E}_{(W, S) \sim P_{WS}^A} [\hat{l}_S(W)] - \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^A \times P_S} [\hat{l}_{\bar{S}}(\bar{W})] \\
&= \frac{1}{n} \sum_{i=1} \mathbb{E}_{(W, Z_i) \sim P_{WZ_i}^A} [\hat{l}_{Z_i}(W)] - \mathbb{E}_{(\bar{W}, \bar{Z}_i) \sim P_W^A \times P} [\hat{l}_{\bar{Z}_i}(\bar{W})] \\
&\leq \frac{1}{n} \sum_{i=1} \sup_{f: \|f\|_\infty \leq \tau} \left(\mathbb{E}_{(W, Z_i) \sim P_{WZ_i}^A} [f(W, Z_i)] - \mathbb{E}_{(\bar{W}, \bar{Z}_i) \sim P_W^A \times P} [f(\bar{W}, \bar{Z}_i)] \right) \\
&\stackrel{a}{=} \frac{1}{n} \sum_{i=1} \tau \|P_{WZ_i}^A - P_W^A \times P\|_1,
\end{aligned} \tag{74}$$

where, a follows from (7) of Fact 3. The bound obtained in (74) holds under a strict assumption that the loss function l is bounded. Further, for any $p > 1$, one can obtain the following upper-bound on $\mathbb{E}_{(W,S) \sim P_{WS}^A}[\text{gen}(W, S)]$, under a stricter assumption that $\|l\|_q \leq \tau$, for a $q < \infty$ such that $\frac{1}{p} + \frac{1}{q} = 1$,

$$\begin{aligned} \mathbb{E}_{(W,S) \sim P_{WS}^A}[\text{gen}(W, S)] &\leq \frac{1}{n} \sum_{i=1} \sup_{f: \|f\|_q \leq \tau} \left(\mathbb{E}_{(W, Z_i) \sim P_{WZ_i}^A} [f(W, Z_i)] - \mathbb{E}_{(\bar{W}, \bar{Z}_i) \sim P_{\bar{W}}^A \times P} [f(\bar{W}, \bar{Z}_i)] \right) \\ &\stackrel{a}{=} \frac{1}{n} \sum_{i=1} \tau \|P_{WZ_i}^A - P_W^A \times P\|_p, \end{aligned} \quad (75)$$

where a follows from (6) of Fact 3. Observe that (75) is a comparatively tighter upper-bound than (74), since $\|\cdot\|_p$ is a decreasing function of p .

Xu and Raginsky in [16], relaxed these strict assumptions mentioned above and state an upper-bound on the absolute value of the expected generalization error under the following assumption.

Assumption 1. (classical sub-Gaussianity assumption) For each $w \in \mathcal{W}$, $l(w, Z)$ for some $0 < \tau < \infty$ and any $\lambda \in \mathbb{R}$ under the distribution P , satisfies the following:

$$\log \mathbb{E}_{Z \sim P} \left[e^{\lambda(l(w, Z) - \mathbb{E}_{Z \sim P}[l(w, Z)])} \right] \leq \frac{\lambda^2 \tau^2}{2}.$$

Using Assumption 1, Xu and Raginsky in [16] proved the following.

Proposition 1 ([16, Theorem 1]). Suppose for each $w \in \mathcal{W}$, $l(w, Z)$ satisfies Assumption 1 for some $0 < \tau < \infty$, then,

$$\left| \mathbb{E}_{(W,S) \sim P_{WS}^A}[\text{gen}(W, S)] \right| \leq \sqrt{\frac{2\tau^2}{n} I[S; W]},$$

where $I[S; W]$ is calculated with respect to P_{WS}^A .

Later, Bu et al. in [17] extend Proposition 1 by proposing a tighter individual sample-based upper-bound on the expected generalization error as follows

Proposition 2 ([17, Proposition 1]). Suppose for each $w \in \mathcal{W}$, $l(w, Z)$ satisfies Assumption 1 for some $0 < \tau < \infty$, then,

$$\left| \mathbb{E}_{(W,S) \sim P_{WS}^A}[\text{gen}(W, S)] \right| \leq \frac{1}{n} \sum_{i=1}^n \sqrt{2\tau^2 I[Z_i; W]},$$

where $I[Z_i; W]$ is calculated with respect to $P_{WZ_i}^A$.

Later, Modak et al. [21] extend the techniques of [17] and generalize the individual sample-based upper-bound on the expected generalization error mentioned in Proposition 2 in terms of α -Rényi divergence, given in Proposition 3 below. To do so, they require the following additional sub-Gaussianity assumptions over Assumption 1.

Assumption 2. (classical sub-Gaussianity assumption) For each $w \in \mathcal{W}$, $i \in [n]$, $l(w, Z_i)$ for some $\tau > 0$ and any $\lambda \in \mathbb{R}$ under the distribution P and $P_{Z_i|W=w}$, satisfies the following:

$$\log \mathbb{E}_{Z_i \sim P_{Z_i|W=w}} \left[e^{\lambda(l(w, Z_i) - \mathbb{E}_{Z_i \sim P_{Z_i|W=w}}[l(w, Z_i)])} \right] \leq \frac{\lambda^2 \tau^2}{2}.$$

Proposition 3 ([21, Theorem 1 and Remark 2]). Suppose for each $w \in \mathcal{W}$, $i \in [n]$, $l(w, Z_i)$ satisfies Assumptions 1 and 2 for some $0 < \tau < \infty$, then,

$$\left| \mathbb{E}_{(W,S) \sim P_{WS}^A}[\text{gen}(W, S)] \right| \leq \begin{cases} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W \sim P_W^A} \left[\sqrt{\frac{2\tau^2 D_\alpha^c(P_{Z_i|W} \| P)}{\alpha}} \right], & \text{if } \alpha \in (0, 1), \\ \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W \sim P_W^A} \left[\sqrt{2\tau^2 D_\alpha^c(P_{Z_i|W} \| P)} \right], & \text{if } \alpha \in (1, \infty). \end{cases}$$

It is important to note that for $\alpha \in (0, 1)$, the upper-bound obtained in Proposition 3 can potentially be tighter than that obtained in Proposition 2. Further, since $D_\alpha^c(\cdot \| \cdot)$ is well-defined for $\alpha \rightarrow 1$ (from Fact 4) and thus setting $\alpha \rightarrow 1$ we can recover results mentioned in Proposition 2.

B. Quantum learning paradigm

We now show that one can obtain a bound analogous to the one obtained in (74) for the expected quantum generalization error (defined in Definition 21), if we assume that the loss observables are bounded. Under this strict assumption, one can easily

derive an upper-bound of similar flavour to one of the main results of [8]. Towards this, for each $(w, s) \in \mathcal{W} \times \mathcal{Z}^n$, we assume $-\mu \mathbb{I} \preceq \hat{L}(w, s) \preceq \mu \mathbb{I}$ for some $\mu < \infty$ and we define

$$\text{gen}^{(\text{old})}(w, s) := l_\rho^{(\text{old})}(w) - \hat{l}_\rho(w, s), \quad (76)$$

$$\begin{aligned} \overline{\text{gen}}^{(\text{old})} &:= \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\text{gen}^{(\text{old})}(w, s)] \\ &= L_\rho^{(\text{old})} - \hat{L}_\rho, \end{aligned} \quad (77)$$

where $l_\rho^{(\text{old})}(w)$, $\hat{l}_\rho(w, s)$, $L_\rho^{(\text{old})}$ and \hat{L}_ρ are defined in Definitions 16, 14, 18 and 15). Then, consider the following series of inequalities,

$$\begin{aligned} \overline{\text{gen}}^{(\text{old})} &= \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\text{Tr} [\hat{L}(W, S) \sigma^{A_Q}(W, S)]] - \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^{A_Q} \times P_S} [\text{Tr} [\hat{L}(\bar{W}, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(\bar{W}, \bar{S}))]] \\ &\stackrel{a}{=} \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\text{Tr} [\hat{L}(W, S) \sigma^{A_Q}(W, S)] - \text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))]] \\ &\quad + \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))]] - \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^{A_Q} \times P_S} [\text{Tr} [\hat{L}(\bar{W}, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(\bar{W}, \bar{S}))]] \\ &\stackrel{b}{\leq} \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\mu \|\sigma^{A_Q}(W, S) - \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S)\|_1] \\ &\quad + \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))]] - \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^{A_Q} \times P_S} [\text{Tr} [\hat{L}(\bar{W}, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(\bar{W}, \bar{S}))]] \\ &\stackrel{c}{\leq} \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\mu \|\sigma^{A_Q}(W, S) - \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S)\|_1] + \mu \|P_{WS}^{A_Q} - P_W^{A_Q} \times P_S\|_1, \end{aligned} \quad (78)$$

where in *a* we denote $P_S := P^n$, *b* follows from (14) of Fact 11 and *c* follows from Fact 3 since for each $(w, s) \in \mathcal{W} \times \mathcal{S}$, $-\mu \leq \text{Tr} [\hat{L}(w, s) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(w, s))] \leq \mu$ as $-\mu \mathbb{I} \preceq \hat{L}(w, s) \preceq \mu \mathbb{I}$. Further, using (13) of Fact 11, for any $p > 1$, to obtain an upper-bound analogous to (75), $\overline{\text{gen}}^{(\text{old})}$, one can consider the following stricter assumption.

Assumption 3. The collection of loss observables $\{\hat{L}(w, s)\}_{(w, s) \in \mathcal{W} \times \mathcal{S}}$ for some $0 < \mu, \tau < \infty$ and $q < \infty$ such that $\frac{1}{p} + \frac{1}{q} = 1$, satisfies the following,

$$\|\hat{L}(w, s)\|_q \leq \mu, \quad \forall (w, s) \in \mathcal{W} \times \mathcal{S}, \quad (79)$$

$$\|l_Q\|_q \leq \tau, \quad \text{where } \forall (w, s) \in \mathcal{W} \times \mathcal{S}, \quad l_Q(w, s) := \text{Tr} [\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s))]. \quad (80)$$

Then, under Assumption 3, we can obtain the following upper-bound on $\overline{\text{gen}}^{(\text{old})}$,

$$\begin{aligned} \overline{\text{gen}}^{(\text{old})} &\stackrel{a}{=} \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\text{Tr} [\hat{L}(W, S) \sigma^{A_Q}(W, S)] - \text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))]] \\ &\quad + \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))]] - \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^{A_Q} \times P_S} [\text{Tr} [\hat{L}(\bar{W}, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(\bar{W}, \bar{S}))]] \\ &\stackrel{a}{\leq} \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\mu \|\sigma^{A_Q}(W, S) - \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S)\|_p] \\ &\quad + \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))]] - \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^{A_Q} \times P_S} [\text{Tr} [\hat{L}(\bar{W}, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(\bar{W}, \bar{S}))]] \\ &= \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\mu \|\sigma^{A_Q}(W, S) - \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S)\|_p] + \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [l_Q(W, S)] - \mathbb{E}_{(\bar{W}, \bar{S}) \sim P_W^{A_Q} \times P_S} [l_Q(\bar{W}, \bar{S})] \\ &\stackrel{b}{\leq} \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} [\mu \|\sigma^{A_Q}(W, S) - \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S)\|_p] + \tau \|P_{WS}^{A_Q} - P_W^{A_Q} \times P_S\|_p, \end{aligned} \quad (81)$$

where *a* follows from (79) and (13) of Fact 11, *b* follows from (80) and Fact 3. Observe that when $\mu = \tau$, (81) is a comparatively tighter upper-bound than the obtained in (78). This is because both Schatten- p and L_p norms are a decreasing function of p .

The upper-bounds obtained above in eqs. (78) and (81) require stringent assumptions on loss observables. However, Caro et al. in [8], motivated by the results obtained in [16] (mentioned as Proposition 1), relax the above stringent assumptions by considering the following sub-Gaussianity assumptions.

Assumption 4. (sub-Gaussianity assumptions mentioned in [8]) The collection of loss observables $\{\hat{L}(w, s)\}_{(w, s) \in \mathcal{W} \times \mathcal{S}}$ for some $0 < \mu, \tau < \infty$ and any $\lambda \in \mathbb{R}$ satisfies the following,

$$\log \text{Tr} \left[e^{\lambda (\hat{L}(w, s) - \text{Tr} [\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s))]) (\mathbb{I}_{\mathcal{H}_{te}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}})} (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \leq \frac{\lambda^2 \mu^2}{2}, \quad \forall (w, s) \in \mathcal{W} \times \mathcal{S}, \quad (82)$$

$$\log \mathbb{E}_{S \sim P_S} \left[e^{\lambda (\text{Tr} [\hat{L}(w, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(w, S))]) - \mathbb{E}_{\bar{S} \sim P_S} [\text{Tr} [\hat{L}(w, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(w, \bar{S}))]]} \right] \leq \frac{\lambda^2 \tau^2}{2}, \quad \forall w \in \mathcal{W}. \quad (83)$$

In Assumption 4, (82) is a quantum sub-Gaussianity (see Definition 13) assumption. Similarly, (83) is a classical sub-Gaussianity (see Definition 5) assumption. We note here that if the loss observable has bounded norm, then Assumption 4 directly follows from Corollary 1 and Fact 7. Using Assumption 4, [8] obtained the following upper-bound on the absolute value of $\overline{\text{gen}}^{(\text{old})}$ (mentioned in (77)), which is a quantum version of the result obtained in [16, Theorem 1] (mentioned as Proposition 1).

Proposition 4 ([8, Corollary 23]). *Suppose for each $(w, s) \in \mathcal{W} \times \mathcal{S}$ $\hat{L}(w, s)$ satisfies Assumption 4 for some $0 < \mu, \tau < \infty$. Then, the following holds,*

$$|\overline{\text{gen}}^{(\text{old})}| \leq \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}_Q}} \left[\sqrt{2\mu^2 D(\sigma^{\mathcal{A}_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S))} \right] + \sqrt{2\tau^2 I[S; W]}. \quad (84)$$

Note that under Assumption 4, the upper-bound obtained in Proposition 4 is weaker than the one obtained in (78). Further, if we use the proposed definition of quantum generalization error (see Definition 20), we get a modified version of the result obtained in Proposition 4, which contains an extra quantum information theoretic quantity in terms of Petz quantum Rényi divergence. To prove this result, we require the following sub-Gaussian assumptions,

Assumption 5. *The collection of loss observables $\{\hat{L}(w, s)\}_{(w, s) \in \mathcal{W} \times \mathcal{S}}$ for some $0 < \mu, \tau < \infty$ and any $\lambda \in \mathbb{R}$ satisfies the following,*

$$\log \text{Tr} \left[e^{\lambda(\hat{L}(w, s) - \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))](\mathbb{I}_{\mathcal{H}_{te}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}}))(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))} \right] \leq \frac{\lambda^2 \mu^2}{2}, \quad \forall (w, s) \in \mathcal{W} \times \mathcal{S}, \quad (85)$$

$$\log \text{Tr} \left[e^{\lambda(\hat{L}(w, s) - \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))](\mathbb{I}_{\mathcal{H}_{te}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}}))(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))} \right] \leq \frac{\lambda^2 \mu^2}{2}, \quad \forall (w, s) \in \mathcal{W} \times \mathcal{S}, \quad (86)$$

$$\log \mathbb{E}_{S \sim P_S} \left[e^{\lambda(\text{Tr}[\hat{L}(w, S)(\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))]) - \mathbb{E}_{\bar{S} \sim P_S} [\text{Tr}[\hat{L}(w, \bar{S})(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))]]} \right] \leq \frac{\lambda^2 \tau^2}{2}, \quad \forall w \in \mathcal{W}. \quad (87)$$

Observe that (86) is an additional sub-Gaussian assumption over Assumption 4, which was not required to prove Proposition 4. However, since the first term in the proposed definition of generalization error (see Definition 20) involves a quantum state of the form $(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))$, we require this additional quantum sub-Gaussianity assumption mentioned in (86), which involves a quantum state of the same form. Using Assumption 5, we state the following upper bound on the absolute value of the expected generalization error (see Definition 21).

Theorem 1 (Modified version of Proposition 4). *Suppose for each $(w, s) \in \mathcal{W} \times \mathcal{S}$, $\hat{L}(w, s)$ satisfies Assumption 5 for some $0 < \mu, \tau < \infty$. Then, the following holds,*

$$|\overline{\text{gen}}| \leq \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}_Q}} \left[\sqrt{2\mu^2 D(\sigma^{\mathcal{A}_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S))} + \sqrt{2\mu^2 D(\sigma_{hyp}^{\mathcal{A}_Q}(W, S) \parallel \sigma_{hyp}^{\mathcal{A}_Q}(W))} \right] + \sqrt{2\tau^2 I[S; W]}. \quad (88)$$

Proof. See Appendix B for the proof. ■

Observe that the second term in (88) was not there in (84). This is because of the difference between the definition of quantum generalization error proposed in [8] (see (76)) and our proposed definition (see Definition 20).

Remark 4. As $\alpha \rightarrow 1$, $D_\alpha^{\mathbb{M}}(\cdot \parallel \cdot)$ becomes equal to $D^{\mathbb{M}}(\cdot \parallel \cdot)$ and thus as a consequence of Fact 17 and Fact 24 (data processing inequality for the quantum divergence), the upper-bounds mentioned in (88) of Theorem 1 can be tightened using $D^{\mathbb{M}}(\cdot \parallel \cdot)$.

All the results obtained in this manuscript can be tightened using $D_\alpha^{\mathbb{M}}(\cdot \parallel \cdot)$. However, the definition of $D_\alpha^{\mathbb{M}}(\cdot \parallel \cdot)$ involves optimization over the choice of POVM. Therefore, for simplicity we prove all our results using Petz and modified sandwiched quantum Rényi divergences (Definition 12).

VI. QUANTUM RÉNYI DIVERGENCES BASED BOUNDS ON EXPECTED GENERALIZATION ERROR

In this section, we prove a quantum version of the results obtained in [21] and [25]. In particular, we obtain bounds on generalization error both in expectation and in probability in terms of quantum Rényi divergence and classical Rényi divergence. Thus, generalizing the results of [21] and [25]. Further, we recover the result of [8] for the expected generalization error.

A. Bounds on the expected quantum generalization error

In this subsection, analogous to the classical sub-Gaussianity assumptions of [21] (mentioned as Assumption 2), we require the following sub-Gaussianity assumption.

Assumption 6. *The collection of loss observables $\{\hat{L}(w, s)\}_{(w, s) \in \mathcal{W} \times \mathcal{S}}$ for some $0 < \mu, \tau < \infty$ and any $\lambda \in \mathbb{R}$ satisfies the following,*

$$\log \text{Tr} \left[e^{\lambda(\hat{L}(w, s) - \text{Tr}[\hat{L}(w, s)\sigma^{\mathcal{A}_Q}(w, s)](\mathbb{I}_{\mathcal{H}_{te}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}}))\sigma^{\mathcal{A}_Q}(w, s)} \right] \leq \frac{\lambda^2 \mu^2}{2}, \quad \forall (w, s) \in \mathcal{W} \times \mathcal{S}, \quad (89)$$

$$\log \mathbb{E}_{S \sim P_{S|W}^{\mathcal{A}_Q}(\cdot|w)} \left[e^{\lambda \left(\text{Tr}[\hat{L}(w, S)(\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))] - \mathbb{E}_{\bar{S} \sim P_{\bar{S}|W}^{\mathcal{A}_Q}(\cdot|w)} [\text{Tr}[\hat{L}(w, \bar{S})(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))] \right]} \right] \leq \frac{\lambda^2 \tau^2}{2}, \quad \forall w \in \mathcal{W}. \quad (90)$$

To state a family of upper-bounds on the absolute value of the expected generalization error (mentioned in Definition 21), we require Lemma 5 below as its preparation.

Lemma 5. *Suppose for each $(w, s) \in \mathcal{W} \times \mathcal{S}$, $\hat{L}(w, s)$ satisfies Assumptions 5 and 6 for some $0 < \mu, \tau < \infty$. Then, the following holds,*

$$\begin{aligned} & \left| \text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] - \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] \right| \\ & \leq \begin{cases} \sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma^{\mathcal{A}_Q}(w, s) \|\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s)\|)}{\alpha}}, & \text{if } \alpha \in (0, 1), \\ \sqrt{2\mu^2 \bar{D}_\alpha(\sigma^{\mathcal{A}_Q}(w, s) \|\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s)\|)}, & \text{if } \alpha \in (1, \infty). \end{cases} \end{aligned} \quad (91)$$

$$\begin{aligned} & \left| \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] - \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))] \right| \\ & \leq \begin{cases} \sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma_{hyp}^{\mathcal{A}_Q}(w, s) \|\sigma_{hyp}^{\mathcal{A}_Q}(w)\|)}{\alpha}}, & \text{if } \alpha \in (0, 1), \\ \sqrt{2\mu^2 \bar{D}_\alpha(\sigma_{hyp}^{\mathcal{A}_Q}(w, s) \|\sigma_{hyp}^{\mathcal{A}_Q}(w)\|)}, & \text{if } \alpha \in (1, \infty). \end{cases} \end{aligned} \quad (92)$$

Proof. We first prove (91) in two cases and later we show that the proof of (92) follows similarly.

Case 1 : $\alpha \in (0, 1)$

From Lemma 4, $\forall (w, s) \in \mathcal{W} \times \mathcal{S}$ and $\lambda \in \mathbb{R}$, we have the following,

$$\bar{D}_\alpha(\sigma^{\mathcal{A}_Q}(w, s) \|\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s)\|) \geq \frac{\alpha}{\alpha - 1} \log \text{Tr}[e^{(\alpha-1)\lambda \hat{L}(w, s)} \sigma^{\mathcal{A}_Q}(w, s)] - \log \text{Tr}[e^{\alpha\lambda \hat{L}(w, s)} (\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(s))]. \quad (93)$$

We now bound the first term in the RHS of (93) as follows,

$$\begin{aligned} & \log \text{Tr}[e^{(\alpha-1)\lambda \hat{L}(w, s)} \sigma^{\mathcal{A}_Q}(w, s)] \\ & = \log \text{Tr}[e^{(\alpha-1)\lambda (\hat{L}(w, s) - \text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] (\mathbb{I}_{\mathcal{H}_{te}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}}))} \sigma^{\mathcal{A}_Q}(w, s)] - ((1 - \alpha)\lambda) \text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] \\ & \stackrel{a}{\leq} -((1 - \alpha)\lambda) \text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] + \frac{(\alpha - 1)^2 \lambda^2 \mu^2}{2}, \end{aligned} \quad (94)$$

where a follows from (89). We now bound the second term in the RHS of (93) as follows,

$$\begin{aligned} & \log \text{Tr}[e^{\alpha\lambda \hat{L}(w, s)} (\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] \\ & = \log \text{Tr}[e^{\alpha\lambda (\hat{L}(w, s) - \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s)]) (\mathbb{I}_{\mathcal{H}_{te}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}}))} (\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] \\ & \quad + (\alpha\lambda) \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] \\ & \stackrel{a}{\leq} (\alpha\lambda) \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] + \frac{\alpha^2 \lambda^2 \mu^2}{2}, \end{aligned} \quad (95)$$

where a follows from (85).

From eqs. (94) and (95), for $\alpha \in (0, 1)$, we can rewrite (93) as follows,

$$\begin{aligned} & \bar{D}_\alpha(\sigma^{\mathcal{A}_Q}(w, s) \|\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s)\|) \\ & \geq \alpha\lambda \left(\text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] - \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] \right) + \frac{\alpha(\alpha - 1)\lambda^2 \mu^2}{2} - \frac{\alpha^2 \lambda^2 \mu^2}{2}. \end{aligned}$$

We can rewrite the above inequality as follows,

$$\left(\frac{\alpha\mu^2}{2} \right) \lambda^2 - \alpha \left(\text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] - \text{Tr}[\hat{L}(w, s)(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] \right) \lambda + \bar{D}_\alpha(\sigma^{\mathcal{A}_Q}(w, s) \|\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s)\|) \geq 0. \quad (96)$$

Since the above inequality is a non-negative quadratic equation in λ with the coefficient $\left(\frac{\alpha\mu^2}{2} \right) \geq 0$, therefore its discriminant must be non-positive. Thus, we have the following inequality,

$$\begin{aligned}
& \alpha^2 \left(\text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] - \text{Tr} \left[\hat{L}(w, s) \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right) \right] \right)^2 \leq 4 \left(\frac{\alpha \mu^2}{2} \right) \overline{D}_\alpha \left(\sigma^{\mathcal{A}_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right) \\
& \Rightarrow \left| \text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] - \text{Tr} \left[\hat{L}(w, s) \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right) \right] \right| \leq \sqrt{\frac{2\mu^2 \overline{D}_\alpha \left(\sigma^{\mathcal{A}_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right)}{\alpha}}. \quad (97)
\end{aligned}$$

Case 2 : $\alpha \in (1, \infty)$

For $\alpha \in (1, \infty)$ the term $\log \text{Tr}[e^{(\alpha-1)\lambda \hat{L}(w, s)} \sigma^{\mathcal{A}_Q}(w, s)]$ in LHS of the inequality mentioned in (94) can be upper-bounded as follows,

$$\log \text{Tr} \left[e^{(\alpha-1)\lambda \hat{L}(w, s)} \sigma^{\mathcal{A}_Q}(w, s) \right] \stackrel{a}{\geq} (\alpha-1)\lambda \text{Tr}[\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)],$$

where a follows from Fact 20. The rest of the proof is similar to the Case 1. This proves (91).

We now proceed to prove (92) in the two following cases.

Case 1 : $\alpha \in (0, 1)$

Using eqs. (85) and (86) and a calculation similar to eqs. (94) to (96) we have the following inequality,

$$\begin{aligned}
& \alpha^2 \left(\text{Tr} \left[\hat{L}(w, s) \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right) \right] - \text{Tr} \left[\hat{L}(w, s) \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w) \right) \right] \right)^2 \\
& \leq 4 \left(\frac{\alpha \mu^2}{2} \right) \overline{D}_\alpha \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w) \right) \\
& \Rightarrow \left| \text{Tr} \left[\hat{L}(w, s) \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right) \right] - \text{Tr} \left[\hat{L}(w, s) \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w) \right) \right] \right| \\
& \leq \sqrt{\frac{2\mu^2 \overline{D}_\alpha \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w) \right)}{\alpha}} \\
& \stackrel{a}{=} \sqrt{\frac{2\mu^2 \overline{D}_\alpha \left(\sigma_{hyp}^{\mathcal{A}_Q}(w, s) \parallel \sigma_{hyp}^{\mathcal{A}_Q}(w) \right)}{\alpha}},
\end{aligned}$$

where a follows from (17) of Fact 14.

Case 2 : $\alpha \in (1, \infty)$

For $\alpha \in (1, \infty)$ the term $\log \text{Tr}[e^{(\alpha-1)\lambda \hat{L}(w, s)} \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right)]$ can be upper-bounded as follows,

$$\log \text{Tr} \left[e^{(\alpha-1)\lambda \hat{L}(w, s)} \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right) \right] \stackrel{a}{\geq} (\alpha-1)\lambda \text{Tr} \left[\hat{L}(w, s) \left(\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s) \right) \right],$$

where a follows from Fact 20. The rest of the proof is similar to the Case 1 mentioned above. This proves (92). This completes the proof of Lemma 5. \blacksquare

Theorem 2 (Expected generalization error bound via modified sandwiched Quantum Rényi Divergence). Suppose $\forall (w, s) \in \mathcal{W} \times \mathcal{S}$, $L(w, s)$ satisfies Assumptions 5 and 6. Then, we have the following two upper bounds for $|\overline{\text{gen}}|$,

$$\begin{aligned}
|\overline{\text{gen}}| & \leq \inf_{\alpha \in (0, 1)} \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}_Q}} \left(\sqrt{\frac{2\mu^2 \overline{D}_\alpha \left(\sigma^{\mathcal{A}_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S) \right)}{\alpha}} \right. \\
& \quad \left. + \sqrt{\frac{2\mu^2 \overline{D}_\alpha \left(\sigma_{hyp}^{\mathcal{A}_Q}(W, S) \parallel \sigma_{hyp}^{\mathcal{A}_Q}(W) \right)}{\alpha}} \right) + \inf_{\gamma \in (0, 1)} \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\sqrt{\frac{2\tau^2 D_\gamma^c(P_{S|W}^{\mathcal{A}_Q} \parallel P_S)}{\gamma}} \right], \quad (98)
\end{aligned}$$

and

$$\begin{aligned}
|\overline{\text{gen}}| & \leq \inf_{\alpha \in (1, \infty)} \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}_Q}} \left(\sqrt{2\mu^2 \overline{D}_\alpha \left(\sigma^{\mathcal{A}_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S) \right)} \right. \\
& \quad \left. + \sqrt{2\mu^2 \overline{D}_\alpha \left(\sigma_{hyp}^{\mathcal{A}_Q}(W, S) \parallel \sigma_{hyp}^{\mathcal{A}_Q}(W) \right)} \right) + \inf_{\gamma \in (1, \infty)} \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\sqrt{2\tau^2 D_\gamma^c(P_{S|W}^{\mathcal{A}_Q} \parallel P_S)} \right]. \quad (99)
\end{aligned}$$

Remark 5. From (99), it may appear that the bound for $\alpha, \gamma \in (1, \infty)$ follows trivially from the bound obtained in (98) for $\alpha, \gamma \in (0, 1)$. However, this is not the case since the proof for the case when $\alpha, \gamma \in (1, \infty)$ is different from the proof for the case when $\alpha, \gamma \in (0, 1)$, because of Lemma 5.

Remark 6. Observe that as discussed in Remark 4 the terms involving modified sandwiched quantum Rényi divergence $\bar{D}_\alpha(\cdot|\cdot)$ in the RHS of eqs. (98) and (99) can be replaced by measured Rényi divergence $D_\alpha^{\mathfrak{M}}(\cdot|\cdot)$ because of the lower-bound of Fact 29, which results in tighter upper-bounds.

Theorem 2 above can be viewed as a quantum version of Proposition 3 ([21, Theorem 1]). Proposition 3 is a generalization of the results obtained in Propositions 1 and 2 in terms of Rényi divergence. Similarly, Theorem 2 can be viewed as a generalization of Proposition 4 in terms of modified sandwiched Rényi divergence. Further, observe that, unlike Proposition 3, we get extra quantum information-theoretic quantities as the second terms of (98) and (99). This is because, unlike in the classical learning setting, the generalization error as defined in Definition 20 is asymmetric. Observe that in the classical case, the true loss for a fixed w , is defined as the expectation over the test data of the same loss function with respect to which the empirical loss is defined. However, this is not the case in the quantum setting (see Definitions 14, 16 and 17). This asymmetric nature of the generalization error in the quantum case, as defined in [8] and the new definition of the true loss proposed in Definition 17, is to mitigate the perturbations caused by the measurements and post-processing during the learning from quantum data. That is, when we employ Definition 18, the second terms do not appear in the upper bounds. Since Definition 19 employs $\sigma_{hyp}^{A_Q}(W)$ instead of $\sigma_{hyp}^{A_Q}(W, S)$, our upper bound has the second term, which comes from the difference between the definitions.

Proof. We calculate the absolute value of the expected generalized error as follows,

$$\begin{aligned}
|\text{gen}| &= \left| \mathbb{E}_{W \sim P_W^{A_Q}} \left[\mathbb{E}_{\bar{S} \sim P_S} \left[\text{Tr} \left[\hat{L}(W, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] - \mathbb{E}_{S \sim P_{S|W}^{A_Q}} \left[\text{Tr} \left[\hat{L}(W, S) \sigma_{hyp}^{A_Q}(W, S) \right] \right] \right] \right| \\
&\leq \mathbb{E}_{W \sim P_W^{A_Q}} \left[\left| \text{Tr} \left[\hat{L}(W, S) \sigma_{hyp}^{A_Q}(W, S) \right] - \text{Tr} \left[\hat{L}(W, S) \left(\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right| \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\left| \mathbb{E}_{\bar{S} \sim P_S} \left[\text{Tr} \left[\hat{L}(W, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] - \mathbb{E}_{S \sim P_{S|W}^{A_Q}} \left[\text{Tr} \left[\hat{L}(W, S) \left(\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] \right| \right] \\
&\leq \mathbb{E}_{W \sim P_W^{A_Q}} \left[\left| \text{Tr} \left[\hat{L}(W, S) \sigma_{hyp}^{A_Q}(W, S) \right] - \text{Tr} \left[\hat{L}(W, S) \left(\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right) \right] \right| \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\left| \text{Tr} \left[\hat{L}(W, S) \left(\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right) \right] - \text{Tr} \left[\hat{L}(W, S) \left(\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right| \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\left| \mathbb{E}_{\bar{S} \sim P_S} \left[\text{Tr} \left[\hat{L}(W, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] - \mathbb{E}_{S \sim P_{S|W}^{A_Q}} \left[\text{Tr} \left[\hat{L}(W, S) \left(\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] \right| \right] \\
&\stackrel{a}{\leq} \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right)}{\alpha}} \right] + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)}{\alpha}} \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\left| \mathbb{E}_{\bar{S} \sim P_S} \left[\text{Tr} \left[\hat{L}(W, \bar{S}) \left(\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] - \mathbb{E}_{S \sim P_{S|W}^{A_Q}} \left[\text{Tr} \left[\hat{L}(W, S) \left(\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] \right| \right] \\
&\stackrel{b}{\leq} \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right)}{\alpha}} \right] + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)}{\alpha}} \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{\frac{2\tau^2 D_\gamma^c(P_{S|W}^{A_Q} \| P_S)}{\gamma}} \right] \\
&= \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right)}{\alpha}} \right] + \sqrt{\frac{2\mu^2 \bar{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)}{\alpha}} \\
&\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{\frac{2\tau^2 D_\gamma^c(P_{S|W}^{A_Q} \| P_S)}{\gamma}} \right],
\end{aligned}$$

where a follows from eqs. (91) and (92) and b follows from [21, Lemma 2] under $\gamma \in (0, 1)$ and the classical sub-Gaussian assumptions mentioned in (87),(90). An important observation to make here is that $D(\sigma_{hyp}^{A_Q}(S, W) \| \sigma_{hyp}^{A_Q}(W))$ is well defined, since $\sigma_{hyp}^{A_Q}(S, W) < \sigma_{hyp}^{A_Q}(W)$. Taking the infimum with $\alpha \in (0, 1)$ and $\gamma \in (0, 1)$, we obtain the upper bound (98).

For the case when $\alpha, \gamma \in (1, \infty)$, using Lemma 5 under the choice of $\alpha \in (1, \infty)$, we directly have the following inequality,

$$\begin{aligned} |\overline{\text{gen}}| &\leq \mathbb{E}_{(W,S) \sim P_{WS}^{A_Q}} \left[\sqrt{2\mu^2 \overline{D}_\alpha \left(\sigma^{A_Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right)} + \sqrt{2\mu^2 \overline{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)} \right] \\ &\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\left| \mathbb{E}_{\overline{S} \sim P_S} \left[\text{Tr} \left[\hat{L}(W, \overline{S}) \left(\rho_{te}(\overline{S}) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] - \mathbb{E}_{S \sim P_{S|W}^{A_Q}} \left[\text{Tr} \left[\hat{L}(W, S) \left(\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W) \right) \right] \right] \right| \right] \\ &\stackrel{a}{\leq} \mathbb{E}_{(W,S) \sim P_{WS}^{A_Q}} \left[\sqrt{2\mu^2 \overline{D}_\alpha \left(\sigma^{A_Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right)} + \sqrt{2\mu^2 \overline{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)} \right] \\ &\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{2\tau^2 D_\gamma^c(P_{S|W}^{A_Q} \| P_S)} \right], \end{aligned}$$

where a follows from [21, Remark 2]. Taking the infimum with $\alpha \in (1, \infty)$ and $\gamma \in (1, \infty)$, we obtain the upper bound (99). This completes the proof of Theorem 2. \blacksquare

Discussion on comparison of Theorem 2 with Proposition 4 (obtained in [8])

It is important to note that if there is no correlation between quantum testing and training data i.e. $\rho(S) := \rho_{te}(S) \otimes \rho_{tr}(S)$ the RHS of (84) in Proposition 4 only contains a classical quantity i.e. $O(\sqrt{I[S; W]})$. However, under this assumption we will still have a quantum quantity in terms of modified sandwiched quantum Rényi divergence i.e. $O\left(\sqrt{\overline{D}_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)}\right)$ along with a classical term in terms of Rényi divergence.

Observe that Theorem 1 turns out to be a direct corollary of Theorem 2, when $\alpha, \gamma \rightarrow 1$. This directly follows from Facts 4 and 13. For clarity, we restate Theorem 1 in the form of a corollary below.

Corollary 2. Suppose $\forall (w, s) \in \mathcal{W} \times \mathcal{S}$, $L(w, s)$ satisfies Assumptions 5 and 6 for some $0 < \mu, \tau < \infty$. Then, we have the following upper-bound on $|\overline{\text{gen}}|$,

$$|\overline{\text{gen}}| \leq \mathbb{E}_{(W,S) \sim P_{WS}^{A_Q}} \left[\sqrt{2\mu^2 D \left(\sigma^{A_Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right)} + \sqrt{2\mu^2 D \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)} \right] + \sqrt{2\tau^2 I[S; W]}.$$

Further, as a consequence of Fact 29, we get the following weakened upper-bound obtained in Theorem 2, mentioned as a corollary below.

Corollary 3. Suppose $\forall (w, s) \in \mathcal{W} \times \mathcal{S}$, $L(w, s)$ satisfies Assumptions 5 and 6. Then, we have the following two upper bounds for $|\overline{\text{gen}}|$,

$$\begin{aligned} |\overline{\text{gen}}| &\leq \inf_{\alpha \in (0,1)} \mathbb{E}_{(W,S) \sim P_{WS}^{A_Q}} \left(\sqrt{\frac{2\mu^2 D_\alpha \left(\sigma^{A_Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right)}{\alpha}} \right. \\ &\quad \left. + \sqrt{\frac{2\mu^2 D_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)}{\alpha}} \right) + \inf_{\gamma \in (0,1)} \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{\frac{2\tau^2 D_\gamma^c(P_{S|W}^{A_Q} \| P_S)}{\gamma}} \right], \quad (100) \end{aligned}$$

and

$$\begin{aligned} |\overline{\text{gen}}| &\leq \inf_{\alpha \in (1,\infty)} \mathbb{E}_{(W,S) \sim P_{WS}^{A_Q}} \left(\sqrt{2\mu^2 D_\alpha \left(\sigma^{A_Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S) \right)} \right. \\ &\quad \left. + \sqrt{2\mu^2 D_\alpha \left(\sigma_{hyp}^{A_Q}(W, S) \| \sigma_{hyp}^{A_Q}(W) \right)} \right) + \inf_{\gamma \in (1,\infty)} \mathbb{E}_{W \sim P_W^{A_Q}} \left[\sqrt{2\tau^2 D_\gamma^c(P_{S|W}^{A_Q} \| P_S)} \right]. \quad (101) \end{aligned}$$

Remark 7. In Definition 20, for any $(w, s) \in \mathcal{W} \times \mathcal{S}$, if we assume $\text{gen}(w, s) = \text{gen}^{(\text{old})}(w, s)$, (where $\text{gen}^{(\text{old})}(w, s)$ is defined in (76)) and $\hat{L}(w, s)$ satisfies the sub-Gaussianity assumptions mentioned in eqs. (85), (87), (89) and (90) for some $0 < \mu, \tau < \infty$, then from Theorem 2, we get the following upper-bounds for $|\overline{\text{gen}}|$,

$$|\overline{\text{gen}}| \leq \begin{cases} \mathbb{E}_{(W,S) \sim P_{WS}^{\mathcal{A}_Q}} \left(\sqrt{\frac{2\mu^2 \overline{D}_\alpha(\sigma^{\mathcal{A}_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S))}{\alpha}} \right) \\ \quad + \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\sqrt{\frac{2\tau^2 D_\gamma^c(P_{S|W}^{\mathcal{A}_Q} \parallel P_S)}{\gamma}} \right], & \text{if } \alpha, \gamma \in (0, 1), \\ \mathbb{E}_{(W,S) \sim P_{WS}^{\mathcal{A}_Q}} \left(\sqrt{2\mu^2 \overline{D}_\alpha(\sigma^{\mathcal{A}_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S))} \right) \\ \quad + \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\sqrt{2\tau^2 D_\gamma^c(P_{S|W}^{\mathcal{A}_Q} \parallel P_S)} \right], & \text{if } \alpha, \gamma \in (1, \infty). \end{cases} \quad (102)$$

For $\alpha, \gamma \rightarrow 1$, we obtain the following upper-bound on $|\overline{\text{gen}}|$,

$$|\overline{\text{gen}}| \leq \mathbb{E}_{(W,S) \sim P_{WS}^{\mathcal{A}_Q}} \left[\sqrt{2\mu^2 D(\sigma^{\mathcal{A}_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S))} \right] + \sqrt{2\tau^2 I[S; W]}. \quad (103)$$

The upper-bound obtained in (103) recovers the same upper-bound as mentioned in Proposition 4.

B. Bounds on the expected quantum generalization error under i.i.d assumption of quantum data

The upper-bounds obtained in Theorem 2 and Corollary 2 do not depend on the size of the training data. This happens because in Theorem 2, we have performed all the calculations with respect to the whole data and not each of the classical and quantum data-points residing in the classical and quantum data, since we did not assume any i.i.d. structure of the classical and the quantum data. However, suppose we take an i.i.d. structure of the classical data as well as an i.i.d. structure of the quantum data as mentioned in [8, eqs. (4.60) and (4.61)]. In that case, we can now fully utilize the i.i.d. assumption of classical data along with the quantum data. Formally, we assume that the test and train data Hilbert spaces are factorized as $\mathcal{H}^{te} := (\mathcal{H}^{Z_{te}})^{\otimes n}$ and $\mathcal{H}^{tr} := (\mathcal{H}^{Z_{tr}})^{\otimes n}$ and $\forall s := (z_1, \dots, z_n) \in \mathcal{Z}^n, \rho(s) := \bigotimes_{i=1}^n \rho(z_i)$ (we assume for each $z \in \mathcal{Z}$ $\rho(z) \in \mathcal{H}^{Z_{te}} \otimes \mathcal{H}^{Z_{tr}}$ might be correlated or even entangled across $\mathcal{H}^{Z_{te}}$ and $\mathcal{H}^{Z_{tr}}$).

The overall action of \mathcal{A}_Q over the data state ρ leads us to the following CQ state,

$$\begin{aligned} \sigma^{\mathcal{A}_Q} &= \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[|W\rangle\langle W| \otimes \mathbb{E}_{S \sim P_{S|W}^{\mathcal{A}_Q}(\cdot|W)} [|S\rangle\langle S| \otimes \sigma^{\mathcal{A}_Q}(W, S)] \right] \\ &= \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[|W\rangle\langle W| \otimes \bigotimes_{i=1}^n \mathbb{E}_{Z_i \sim P_{Z_i|W}^{\mathcal{A}_Q}(\cdot|W)} [|Z_i\rangle\langle Z_i| \otimes \sigma^{\mathcal{A}_Q}(W, Z_i)] \right], \end{aligned}$$

where for any $i \in [n]$, $P_{Z_i|W}^{\mathcal{A}_Q}$ is the corresponding marginal of $P_{S|W}^{\mathcal{A}_Q}$ and $\sigma^{\mathcal{A}_Q}(W, Z_i) := (\mathbb{I}_{\mathcal{H}^{Z_{te}}} \otimes \Lambda_{W, Z_i})(\rho^{\mathcal{A}_Q}(W, Z_i))$, where $\{\Lambda_{w, z} : T(\mathcal{H}^{Z_{tr}}) \rightarrow T(\mathcal{H}^{hyp})\}_{(w, z) \in \mathcal{W} \times \mathcal{Z}}$ is a collection of quantum channels.

We consider a family of non-negative self-adjoint loss observables $\{\hat{L}(w, s) \in \mathcal{L}(\mathcal{H}^{te} \otimes \mathcal{H}^{hyp})\}_{(w, s) \in \mathcal{W} \times \mathcal{Z}^n}$ (where $\mathcal{H}^{hyp} := \widehat{\mathcal{H}^{hyp}}^{\otimes n}$), where for each $(w, s) \in \mathcal{W} \times \mathcal{Z}^n$, we $\hat{L}(w, s)$ is of the following local form,

$$\hat{L}(w, s) := \frac{1}{n} \sum_{i=1}^n (\mathbb{I}_{\mathcal{H}^{Z_{te}}} \otimes \mathbb{I}_{\widehat{\mathcal{H}^{hyp}}})^{\otimes(i-1)} \otimes L(w, z_i) \otimes (\mathbb{I}_{\mathcal{H}^{Z_{te}}} \otimes \mathbb{I}_{\widehat{\mathcal{H}^{hyp}}})^{\otimes(n-i)}, \quad (104)$$

where for each $i \in [n]$, $L(w, z_i)$ is a local loss observable acting on the i -th iteration of the test and depends on (w, z_i) .

With respect to the loss observable mentioned in (104), we prove corollaries of Lemma 5 and Theorem 2 in terms of the modified sandwiched Rényi divergence. To prove these corollaries, we further require the following sub-Gaussianity assumption, which are special cases of Assumptions 5 and 6.

Assumption 7. The collection of loss observables $\{L(w, z)\}_{(w, z) \in \mathcal{W} \times \mathcal{Z}}$ for some $0 < \mu, \tau < \infty$ and any $\lambda \in \mathbb{R}$ satisfies the following,

$$\log \text{Tr} \left[e^{\lambda \left(L(w, z) - \text{Tr}[L(w, z) \sigma^{\mathcal{A}Q}(w, z)] (\mathbb{I}_{\mathcal{H}_{Z_{te}}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}}) \right) \sigma^{\mathcal{A}Q}(w, z)} \right] \leq \frac{\lambda^2 \mu^2}{2}, \quad \forall (w, z) \in \mathcal{W} \times \mathcal{Z}, \quad (105)$$

$$\log \text{Tr} \left[e^{\lambda \left(L(w, z) - \text{Tr}[L(w, z) (\rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w, z))] (\mathbb{I}_{\mathcal{H}_{Z_{te}}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}}) \right) (\rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w, z))} \right] \leq \frac{\lambda^2 \mu^2}{2}, \quad \forall (w, z) \in \mathcal{W} \times \mathcal{Z}, \quad (106)$$

$$\log \text{Tr} \left[e^{\lambda \left(L(w, z) - \text{Tr}[L(w, z) (\rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w))] (\mathbb{I}_{\mathcal{H}_{Z_{te}}} \otimes \mathbb{I}_{\mathcal{H}_{hyp}}) \right) (\rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w))} \right] \leq \frac{\lambda^2 \mu^2}{2}, \quad \forall (w, z) \in \mathcal{W} \times \mathcal{Z}, \quad (107)$$

$$\log \mathbb{E}_{Z_i \sim P} \left[e^{\lambda \left(\text{Tr}[L(w, Z_i) (\rho_{Z_{te}}(Z_i) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w))] - \mathbb{E}_{\bar{Z}_i \sim P} [\text{Tr}[L(w, \bar{Z}_i) (\rho_{Z_{te}}(\bar{Z}_i) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w))] \right)} \right] \leq \frac{\lambda^2 \tau^2}{2}, \quad \forall w \in \mathcal{W}, \quad (108)$$

$$\log \mathbb{E}_{Z_i \sim P_{Z_i|W}^{\mathcal{A}Q}(\cdot|w)} \left[e^{\lambda \left(\text{Tr}[L(w, Z_i) (\rho_{Z_{te}}(Z_i) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w))] - \mathbb{E}_{\bar{Z}_i \sim P_{Z_i|W}^{\mathcal{A}Q}(\cdot|w)} [\text{Tr}[L(w, \bar{Z}_i) (\rho_{Z_{te}}(\bar{Z}_i) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w))] \right)} \right] \leq \frac{\lambda^2 \tau^2}{2}, \quad \forall w \in \mathcal{W}, \quad (109)$$

Corollary 4. Suppose $\forall (w, z) \in \mathcal{W} \times \mathcal{Z}$, $L(w, z)$ satisfies Assumption 7 for some $0 < \mu, \tau < \infty$. Then, we have,

$$\begin{aligned} & \left| \text{Tr}[L(w, z) \sigma^{\mathcal{A}Q}(w, z)] - \text{Tr}[L(w, z) (\rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w, z))] \right| \\ & \leq \begin{cases} \sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma^{\mathcal{A}Q}(w, z) \| \rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w, z))}{\alpha}}, & \text{if } \alpha \in (0, 1), \\ \sqrt{2\mu^2 \bar{D}_\alpha(\sigma^{\mathcal{A}Q}(w, z) \| \rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w, z))}, & \text{if } \alpha \in (1, \infty), \end{cases} \end{aligned} \quad (110)$$

$$\begin{aligned} & \left| \text{Tr}[L(w, z) (\rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w, z))] - \text{Tr}[L(w, z) (\rho_{Z_{te}}(z) \otimes \sigma_{hyp}^{\mathcal{A}Q}(w))] \right| \\ & \leq \begin{cases} \sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma_{hyp}^{\mathcal{A}Q}(w, z) \| \sigma_{hyp}^{\mathcal{A}Q}(w))}{\alpha}}, & \text{if } \alpha \in (0, 1), \\ \sqrt{2\mu^2 \bar{D}_\alpha(\sigma_{hyp}^{\mathcal{A}Q}(w, z) \| \sigma_{hyp}^{\mathcal{A}Q}(w))}, & \text{if } \alpha \in (1, \infty). \end{cases} \end{aligned} \quad (111)$$

Corollary 5. Suppose $\forall (w, z) \in \mathcal{W} \times \mathcal{Z}$, $L(w, z)$ satisfies Assumption 7 for some $0 < \mu, \tau < \infty$. Then, we have the following,

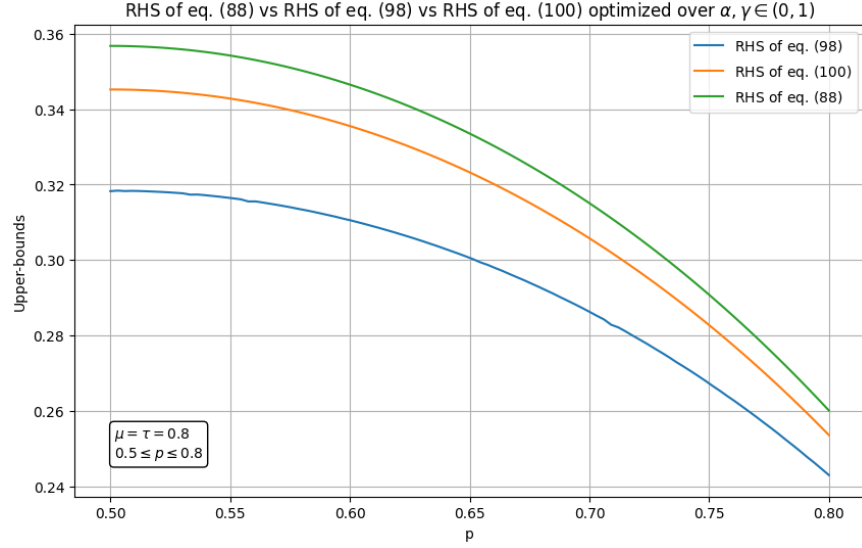
$$|\text{gen}| \leq \begin{cases} \inf_{\alpha \in (0, 1)} \left(\sqrt{\frac{2\mu^2 \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}Q}} [\bar{D}_\alpha(\sigma^{\mathcal{A}Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}Q}(W, S))] }{n\alpha}} \right. \\ \quad \left. + \sqrt{\frac{2\mu^2 \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}Q}} [\bar{D}_\alpha(\sigma_{hyp}^{\mathcal{A}Q}(W, S) \| \sigma_{hyp}^{\mathcal{A}Q}(W))] }{n\alpha}} \right) + \inf_{\gamma \in (0, 1)} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W \sim P_W^A} \left[\sqrt{\frac{2\tau^2 D_\gamma^c(P_{Z_i|W} \| P)}{\gamma}} \right], \\ \inf_{\alpha \in (1, \infty)} \left(\sqrt{\frac{2\mu^2 \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}Q}} [\bar{D}_\alpha(\sigma^{\mathcal{A}Q}(W, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}Q}(W, S))] }{n}} \right. \\ \quad \left. + \sqrt{\frac{2\mu^2 \mathbb{E}_{(W, S) \sim P_{WS}^{\mathcal{A}Q}} [\bar{D}_\alpha(\sigma_{hyp}^{\mathcal{A}Q}(W, S) \| \sigma_{hyp}^{\mathcal{A}Q}(W))] }{n}} \right) + \inf_{\gamma \in (1, \infty)} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W \sim P_W^A} \left[\sqrt{2\tau^2 D_\gamma^c(P_{Z_i|W} \| P)} \right]. \end{cases}$$

Observe that the above result is a quantum version of the individual sample-based upper-bounds mentioned in Proposition 2.

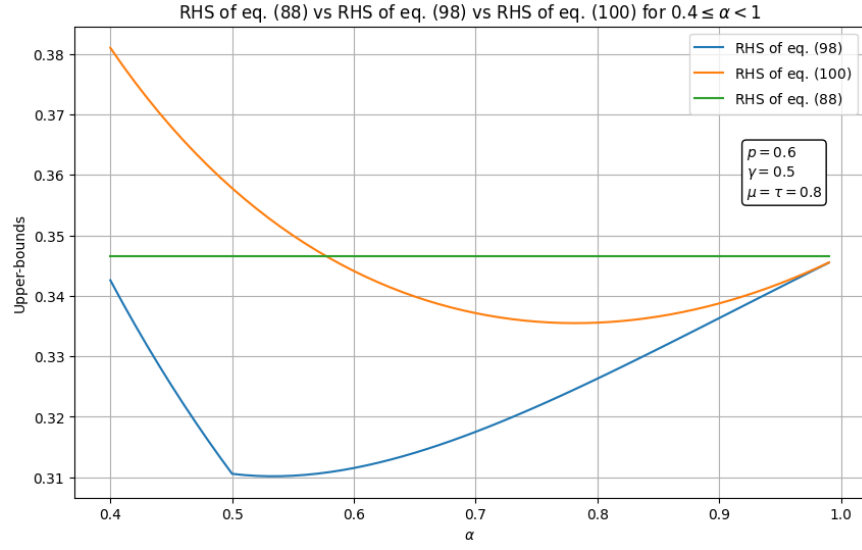
C. Comparison between the results obtained in Theorems 1, 2 and Corollary 3

For the case when $\alpha, \gamma \in (0, 1)$, our bounds are in terms of $\frac{\bar{D}_\alpha(\cdot \| \cdot)}{\alpha}$, $\frac{D_\alpha(\cdot \| \cdot)}{\alpha}$ and $\frac{D_\gamma(\cdot \| \cdot)}{\gamma}$. Because of these extra multiplicative factors of $\frac{1}{\alpha}$ and $\frac{1}{\gamma}$, it is not clear whether the terms $\frac{\bar{D}_\alpha(\cdot \| \cdot)}{\alpha}$, $\frac{D_\alpha(\cdot \| \cdot)}{\alpha}$ and $\frac{D_\gamma^c(\cdot \| \cdot)}{\gamma}$ are smaller than $D(\cdot \| \cdot)$ and $D^c(\cdot \| \cdot)$ respectively. Using simulations for a toy example, in [21], the authors showed cases when the bounds obtained on expected generalization error in terms of $D_\gamma^c(\cdot \| \cdot)$ is smaller than the bounds obtained on expected generalization error in terms of $D(\cdot \| \cdot)$ for $\gamma = 0.5$.

In a similar spirit, to compare the bounds obtained in Theorems 1, 2 and Corollary 3, we consider an example, where for each $z \in \{0, 1\}$ $\rho(z) := \rho_{te}(z) \otimes \rho_{tr}(z)$. Therefore, any measurement and post-processing on the train part will not affect the test part.



(a)



(b)

Fig. 2: RHS of eq. (88) vs RHS of eq. (98) vs RHS of eq. (100).

Hence, the first term of (88), (98) and (100) is equal to zero (however, the second term in the above equations is not zero). For our example, the classical-quantum data state as mentioned in (56) has the following form (parametrized by $p \in (0, 1)$),

$$\rho = p|0\rangle\langle 0| \otimes |\psi_0\rangle\langle\psi_0|_{te} \otimes |\psi_0\rangle\langle\psi_0|_{tr} + (1-p)|1\rangle\langle 1| \otimes |\psi_1\rangle\langle\psi_1|_{te} \otimes |\psi_1\rangle\langle\psi_1|_{tr}. \quad (112)$$

In the above for $z \in \{0, 1\}$,

$$\rho(z) = |\psi_z\rangle\langle\psi_z|_{te} \otimes |\psi_z\rangle\langle\psi_z|_{tr}, \quad (113)$$

where, $|\psi_0\rangle$ and $|\psi_1\rangle$ are defined as follows,

$$\begin{aligned} |\psi_0\rangle &:= \cos\theta|\phi_0\rangle + \sin\theta|\phi_0^\perp\rangle, \\ |\psi_1\rangle &:= \cos\beta|\phi_1\rangle + \sin\beta|\phi_1^\perp\rangle, \end{aligned}$$

where, $\cos^2\theta = 0.45$, $\cos^2\beta = 0.5$, $|\phi_0\rangle = a|0\rangle + b|1\rangle$ and $|\phi_1\rangle = c|0\rangle + d|1\rangle$. Further, $a \approx -0.59 - 0.29i$, $b \approx -0.25 + 0.71i$, $c \approx 0.34 - 0.42i$ and $d \approx -0.83 - 0.12i$ (we list only approximate values of a, b, c and d for simplicity). We now consider the following measurements (dependent on the classical data z), which we will perform on ρ ,

$$\{E_z(w)\}_{(w,z) \in \{0,1\}^2} = \{|\phi_z\rangle\langle\phi_z|, |\phi_z^\perp\rangle\langle\phi_z^\perp|\},$$

where, $\{E_z(w)\}_{(w,z) \in \{0,1\}^2}$ is in the similar spirit as $\{E_s^{A_Q}(w)\}_{(w,s) \in \mathcal{W} \times \mathcal{S}}$, discussed in Subsection IV-B. After performing the measurement, $\forall (w, z) \in \{0, 1\}^2$, we consider $\Lambda_{w,z} = \mathbb{I}$ (where $\Lambda_{w,z}$ is in the similar spirit as $\Lambda_{w,s}$, discussed in Subsection IV-B).

For each $(w, s) \in \mathcal{W} \times \mathcal{S}$, we only assume that the loss observable $\hat{L}(w, s)$, satisfies Assumptions 5 and 6 for $\mu = \tau = 0.8$, because the expressions in the RHS of (88), (98) and (100) depend only on μ and τ for the loss observable $\hat{L}(w, s)$. Therefore, we do not explicitly mention its choice here.

Under these settings, in Figure 2a above, we compare the optimal values of RHS of (88), (98) and (100) respectively by varying p ($0.5 \leq p \leq 0.8$).

Further, under the same settings mentioned above, in Figure 2b above, we compare the RHS of (88), (98) and (100) by varying $\alpha \in [0.4, 1)$ for $p = 0.6$.

For various cases of the example considered above, our simulations show that the modified sandwiched Rényi divergence always gives a better bound compared to the bounds obtained in terms of the Petz Rényi divergence and the quantum relative entropy, respectively.

VII. QUANTUM RÉNYI DIVERGENCES BASED BOUNDS ON GENERALIZATION ERROR IN PROBABILITY

A. Bounds on the generalization error in probability under i.i.d. assumption of quantum data

Expectation bounds derived in the earlier sections only provide average-case guarantees. Therefore, a more relevant metric to study the performance of a learning algorithm would be to obtain bounds on the generalization error in probability. In this section, we will study bounds of the following form,

$$\Pr_{(W,S) \sim P_{WS}} \{|\text{gen}(W, S)| \leq \varepsilon\} \geq 1 - \delta,$$

where W and S are single-drawn according to the distribution $P_{W|S}$ induced by the learning algorithm and the distribution of the data P_S , $\varepsilon > 0$ is the parameter used to denote the error allowed and $\delta > 0$ is the parameter used to denote the confidence $1 - \delta$. This kind of probabilistic upper-bound on the generalization error is called a “single-draw” upper-bound on the generalization error.

In the classical learning scenario, Esposito et al. [25] as a corollary of [25, Theorem 4], give an upper-bound on generalization error in probability mentioned below.

Proposition 5 ([25, Corollary 2]). *Assume that the loss function $l(w, Z)$ satisfies Assumption 5 for some $0 < \tau < \infty$. Furthermore, assume that $P_{WS} \ll P_W P_S$ ($P_S = P_Z^{\otimes n}$). Then, for any $\gamma > 1, \delta \in (0, 1)$,*

$$\mathcal{E} := \left\{ |\text{gen}(W, S)| \leq \sqrt{\frac{2\tau^2}{n} \left(I_\gamma^c[W; S] + \log 2 + \frac{\gamma}{\gamma-1} \log\left(\frac{1}{\delta}\right) \right)} \right\}, \quad (114)$$

satisfies the following,

$$\Pr_{(W,S) \sim P_{WS}} \{\mathcal{E}\} \geq 1 - \delta, \quad (115)$$

where $\forall (w, s) \in \mathcal{W} \times \mathcal{Z}^n$, $\text{gen}(w, s)$ is defined in [25, Definition 8].

From eqs. (114) and (115), we note that if we aim to achieve at most ε ($\varepsilon > 0$) generalization error i.e. $|\text{gen}(W, S)| \leq \varepsilon$ with $(1 - \delta)$ confidence, then, it is necessary to have $n \geq \frac{2\tau^2}{\varepsilon^2} \left(I_\gamma^c[W; S] + \log 2 + \frac{\gamma}{\gamma-1} \log\left(\frac{1}{\delta}\right) \right)$ samples. Thus, for a fixed sample size, there is a trade-off between ε and δ .

The proof of Proposition 5 uses Hölder’s inequality (see Fact 2) non-trivially and is arguably tedious. In comparison, Theorem 3 below proves an alternative strategy to upper-bound the generalization error in probability in terms of smooth max divergence (see Fact 3, just using the definition of smooth max divergence and is simpler than the proof of Proposition 5).

Theorem 3. *Assume that the loss function $l(w, Z)$ satisfies Assumption 5 for some $0 < \tau < \infty$. Furthermore, assume that $P_{WS} \ll P_W P_S$ ($P_S = P_Z^{\otimes n}$). Then, for any $\delta \in (0, 1), \nu < \delta$, the following holds,*

$$\Pr_{(W,S) \sim P_{WS}} \left\{ |\text{gen}(W, S)| \leq \sqrt{\frac{2\tau^2}{n} \left(I_{\max}^{(\nu)}[W; S] + \log 2 + \log\left(\frac{1}{\delta - \nu}\right) \right)} \right\} \geq 1 - \delta, \quad (116)$$

where $I_{\max}^{(\nu)}[W; S] := D_{\max}^{(\nu)}(P_{WS} \| P_W \times P_S)$ (where $D_{\max}^{(\nu)}(\cdot \| \cdot)$ is defined in Fact 3).

Proof. See Appendix D for the proof. ■

From (116), it follows that with probability at most δ , $|\text{gen}(W, S)|$ will go beyond $O(\sqrt{I_{\max}^{(\nu)}[W; S]})$ for any $\nu < \delta$. Further, in earlier sections we observed that in any change of measure based upper-bounds on Generalization error, we try to approximate the joint distribution P_{WS} in terms of the marginal $P_W P_S$ along with some distance measure due to the change of the measure. Moreover, from [49]–[51], it can be realized that whenever we try to make a joint probability distribution (P_{WS}) approximately

close to its marginal ($P_W P_S$), smooth max Rényi divergence (denoted as D_{\max}^ν and defined in Definition 3) naturally comes into the picture. Thus, the upper-bound obtained in (116) in terms of smooth-max Rényi divergence seems to be well-justified.

We now extend Proposition 5 and Theorem 3 in the quantum learning scenario and compare them with the same in the classical scenario. Towards this, we require the following sub-Gaussianity assumption,

$$\log \mathbb{E}_{Z \sim P} \left[e^{\lambda (\text{Tr}[L(w, Z_i) \sigma^{\mathcal{A}_Q}(w, Z)] - \mathbb{E}_{Z \sim P} [\text{Tr}[L(w, \bar{Z}) \sigma^{\mathcal{A}_Q}(w, \bar{Z})])]} \right] \leq \frac{\lambda^2 \tau^2}{2}. \quad (117)$$

where $0 < \tau < \infty$. In the theorem below, we now mention a quantum version of Proposition 5 ([25, Corollary 2]), in the context of above above-discussed “single-draw” probabilistic bounds in a quantum learning scenario.

Theorem 4. *Given the distribution $P_{WS}^{\mathcal{A}_Q}$ induced by a quantum learner \mathcal{A}_Q (such that $P_{WS}^{\mathcal{A}_Q} \ll P_W^{\mathcal{A}_Q} \times P^n$), and for any $\delta \in (0, 1)$, $\alpha \in (0, 1)$, $\gamma > 1$ if the sub-Gaussianity assumptions mentioned in eqs. (105), (106) and (117) holds for some $0 < \mu, \tau < \infty$. Then the event*

$$\mathcal{E} := \left\{ |\text{gen}(W, S)| \leq \sqrt{\frac{2\tau^2}{n} \left(\log 2 + I_\gamma^c[S; W] + \frac{\gamma}{\gamma-1} \log \left(\frac{1}{\delta} \right) \right)} + \inf_{\alpha \in (0, 1)} c_1(\alpha) \right\}, \quad (118)$$

satisfies the following,

$$\Pr_{(W, S) \sim P_{WS}^{\mathcal{A}_Q}} \{ \mathcal{E} \} \geq 1 - \delta,$$

where,

$$c_1(\alpha) := \sup_{w \in \text{supp}(P_W^{\mathcal{A}_Q})} \mathbb{E}_{S \sim P^n} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha \left(\sigma^{\mathcal{A}_Q}(w, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, S) \right)}{n\alpha}} + \sqrt{\frac{2\mu^2 \bar{D}_\alpha \left(\sigma_{hyp}^{\mathcal{A}_Q}(w, S) \| \sigma_{hyp}^{\mathcal{A}_Q}(w) \right)}{n\alpha}} \right]. \quad (119)$$

B. Proof of Theorem 4

Consider the following event

$$E := \{(w, s) \in \mathcal{W} \times \mathcal{Z}^n : |\text{gen}(w, s)| > \varepsilon\},$$

where for any $w \in \mathcal{W}$, we define $E_w := \{s \in \mathcal{Z}^n : (w, s) \in E\}$. Then, from (10) of Fact 8, for some $\gamma > 1$ we can write the following,

$$\Pr_{(W, S) \sim P_{WS}^{\mathcal{A}_Q}} \{E\} \leq \exp \left[\frac{\gamma-1}{\gamma} \left(\log \left(\mathbb{E}_{P_W^{\mathcal{A}_Q}} \left[\Pr_{S \sim P^n} \{E_W\} \right] \right) + I_\gamma^c[S; W] \right) \right]. \quad (120)$$

From (104), it follows that,

$$\begin{aligned} \hat{l}_\rho(w, s) &= \frac{1}{n} \sum_{i=1}^n \text{Tr}[L(w, z_i) \sigma^{\mathcal{A}_Q}(w, z_i)], \\ l_\rho(w) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\bar{Z}_i \sim P} \left[\text{Tr} \left[L(w, \bar{Z}_i) \left(\rho_{Z_{te}}(\bar{Z}_i) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w) \right) \right] \right]. \end{aligned}$$

Further, for any $w \in \mathcal{W}$, $s := (z_1, \dots, z_n) \in \mathcal{Z}^n$, we define $\tilde{l}_\rho(w)$ as follows,

$$\tilde{l}_\rho(w) := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\bar{Z}_i \sim P} \text{Tr}[L(w, \bar{Z}_i) \sigma^{\mathcal{A}_Q}(w, \bar{Z}_i)].$$

Then, for any $w \in \text{supp}(P_W^{\mathcal{A}_Q})$, we have the following,

$$\begin{aligned} \Pr_{S \sim P^n} \{E_w\} &= \Pr_{S \sim P^n} \{|\text{gen}(w, S)| > \varepsilon\} \\ &\stackrel{a}{=} \Pr_{S \sim P^n} \left\{ \left| \hat{l}_\rho(w, S) - l_\rho(w) \right| > \varepsilon \right\} \\ &\stackrel{b}{\leq} 2e^{-\frac{n(\varepsilon - |\hat{l}_\rho(w) - l_\rho(w)|)^2}{2\tau^2}} \\ &\stackrel{c}{\leq} 2 \exp \left[\frac{-n(\varepsilon - c_1(\alpha, w))^2}{2\tau^2} \right], \end{aligned} \quad (121)$$

where a follows from (58), b follows from Fact 10 and in c we define $\tilde{c}(w)$ as follows,

$$c_1(\alpha, w) := \mathbb{E}_{S \sim P^n} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma^{A_Q}(w, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(w, S))}{n\alpha}} + \sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma_{hyp}^{A_Q}(w, S) \| \sigma_{hyp}^{A_Q}(w))}{n\alpha}} \right], \quad (122)$$

and the inequality c follows from the following series of inequalities,

$$\begin{aligned} |\tilde{l}_\rho(w) - l_\rho(w)| &= \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\bar{Z}_i \sim P} \left[\text{Tr}[L(w, \bar{Z}_i) \sigma^{A_Q}(w, \bar{Z}_i)] - \text{Tr}[L(w, \bar{Z}_i) (\rho_{Z_{te}}(\bar{Z}_i) \otimes \sigma_{hyp}^{A_Q}(w)) \right] \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\bar{Z}_i \sim P} \left[\left| \text{Tr}[L(w, \bar{Z}_i) \sigma^{A_Q}(w, \bar{Z}_i)] - \text{Tr}[L(w, \bar{Z}_i) (\rho_{Z_{te}}(\bar{Z}_i) \otimes \sigma_{hyp}^{A_Q}(w, \bar{Z}_i)) \right| \right] \\ &\quad + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\bar{Z}_i \sim P} \left[\left| \text{Tr}[L(w, \bar{Z}_i) (\rho_{Z_{te}}(\bar{Z}_i) \otimes \sigma_{hyp}^{A_Q}(w, \bar{Z}_i)) \right] - \text{Tr}[L(w, \bar{Z}_i) (\rho_{Z_{te}}(\bar{Z}_i) \otimes \sigma_{hyp}^{A_Q}(w)) \right| \right] \\ &\stackrel{a}{\leq} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\bar{Z}_i \sim P} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma^{A_Q}(w, \bar{Z}_i) \| \rho_{Z_{te}}(\bar{Z}_i) \otimes \sigma_{hyp}^{A_Q}(w, \bar{Z}_i))}{\alpha}} + \sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma_{hyp}^{A_Q}(w, \bar{Z}_i) \| \sigma_{hyp}^{A_Q}(w))}{\alpha}} \right] \\ &\leq \mathbb{E}_{S \sim P^n} \left[\sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma^{A_Q}(w, S) \| \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(w, S))}{n\alpha}} + \sqrt{\frac{2\mu^2 \bar{D}_\alpha(\sigma_{hyp}^{A_Q}(w, S) \| \sigma_{hyp}^{A_Q}(w))}{n\alpha}} \right], \end{aligned}$$

where a follows from Corollary 4. Thus, from eqs. (120) and (121) we have,

$$\begin{aligned} \Pr_{(W, S) \sim P_{WS}^{A_Q}} \{E\} &\leq \exp \left[\frac{\gamma-1}{\gamma} \left(\log \left(\mathbb{E}_{P_W^{A_Q}} \left[\Pr_{S \sim P^n} \{E_W\} \right] \right) + I_\gamma^c[S; W] \right) \right] \\ &\leq \exp \left[\frac{\gamma-1}{\gamma} \left(\log \left(\sup_{w \in \text{supp}(P_W^{A_Q})} \left[\Pr_{S \sim P^n} \{E_w\} \right] \right) + I_\gamma^c[S; W] \right) \right] \\ &\leq \exp \left[\frac{\gamma-1}{\gamma} \left(\log \left(\sup_{w \in \text{supp}(P_W^{A_Q})} \left[2 \exp \left[\frac{-n(\varepsilon - c_1(\alpha, w))^2}{2\tau^2} \right] \right] \right) + I_\gamma^c[S; W] \right) \right] \\ &\stackrel{a}{=} \exp \left[\frac{\gamma-1}{\gamma} \left(\log \left(\left[2 \exp \left[\frac{-n(\varepsilon - c_1(\alpha)(w^*))^2}{2\tau^2} \right] \right] \right) + I_\gamma^c[S; W] \right) \right] \\ &= \exp \left[\frac{\gamma-1}{\gamma} \left(\log 2 - \frac{n(\varepsilon - c_1(\alpha))^2}{2\tau^2} + I_\gamma^c[S; W] \right) \right], \quad (123) \end{aligned}$$

where in a , $c_1(\alpha, w^*) := \sup_{w \in \text{supp}(P_W^{A_Q})} c_1(\alpha, w) = c_1(\alpha)$. If we now assume $\delta := e^{\frac{\gamma-1}{\gamma} \left(\log 2 - \frac{n(\varepsilon - c_1(\alpha))^2}{2\tau^2} + I_\gamma^c[S; W] \right)}$. Then, we can write ε as follows,

$$\varepsilon = \sqrt{\frac{2\tau^2}{n} \left(\log 2 + I_\gamma^c[S; W] + \frac{\gamma}{\gamma-1} \log \left(\frac{1}{\delta} \right) \right)} + c_1(\alpha). \quad (124)$$

Hence, from eqs. (123) and (124) we have the following,

$$\Pr_{(W, S) \sim P_{WS}^{A_Q}} \left\{ |\text{gen}(W, S)| \leq \sqrt{\frac{2\tau^2}{n} \left(\log 2 + I_\gamma^c[S; W] + \frac{\gamma}{\gamma-1} \log \left(\frac{1}{\delta} \right) \right)} + c_1(\alpha) \right\} \geq 1 - \delta.$$

Since the above inequality with any $\alpha \in (0, 1)$, we obtain (118). This completes the proof of Theorem 4. \blacksquare

Remark 8. Observe that, unlike the event mentioned in (114), in (118), we have an extra term $c_1(\alpha)$ (defined in eq. (119)). This is because of the asymmetric nature of the generalization error as discussed earlier below the statement of Theorem 2. Further, it is because of this asymmetric nature of the generalization error, we do not get a uniform upper-bound on $\Pr_{S \sim P^n} \{E_w\}$ (mentioned in (121)). However, this is not the case in the classical setting (see [25, eq. (59)]).

Further, we mention a quantum version of Theorem 3 in Theorem 5 below.

Theorem 5. *Given the distribution $P_{WS}^{A_Q}$ induced by a quantum learner \mathcal{A}_Q (such that $P_{WS}^{A_Q} \ll P_W^{A_Q} \times P^n$), then for any $\delta \in (0, 1), \alpha \in (0, 1), \nu < \delta$ if the sub-Gaussianity assumptions mentioned in eqs. (105), (106) and (117) holds for some $0 < \mu, \tau < \infty$, the following holds,*

$$\Pr_{(W,S) \sim P_{WS}^{A_Q}} \left\{ |\text{gen}(W, S)| \leq \sqrt{\frac{2\tau^2}{n} \left(\log 2 + I_{\max}^{(\nu)}[S; W] + \log \left(\frac{1}{\delta - \nu} \right) \right)} + \inf_{\alpha \in (0,1)} c_1(\alpha) \right\} \geq 1 - \delta, \text{ if } \alpha \in (0, 1), \quad (125)$$

where, $c_1(\alpha)$ is defined in (119).

We omit the proof of Theorem 5 for brevity, as it directly follows from the techniques used in the proof of Theorem 3. The results mentioned above in Theorems 4 and 5 give a single-draw upper-bound on the quantum generalization error (defined in Definition 20) in probability and can be thought of as a quantum version of Proposition 5 and Theorem 3 respectively. No results similar to Theorems 4 and 5, have been studied in the literature.

Remark 9. *Here, we remark the comparison of the quantum upper-bounds obtained in Sections VI and VII with their classical counterparts. From the definitions of the generalization error (Definition 20) it follows that if the quantum training and testing data is not entangled i.e. $\forall s \in \mathcal{Z}^n, \rho(s) = \rho_{te}(s) \otimes \rho_{tr}(s)$ and if the quantum hypothesis state ($\sigma_{hyp}^{A_Q}(w, s)$) is classically independent (through S) with the quantum data ($\rho(s)$) state i.e. for each $(w, s) \in \mathcal{W} \times \mathcal{S}, \sigma_{hyp}^{A_Q}(w, s) = \sigma_{hyp}^{A_Q}(w)$, then, all the quantum upper-bounds obtained in Sections VI and VII will boil down to their corresponding classical counterparts as stated earlier in Table I.*

VIII. CONCLUSION

Given the inherent stochasticity of training data and learned hypotheses, we have investigated the generalization error in quantum learning through its expectation and probabilistic behavior as two primary avenues.

TABLE II: Comparison between the result obtained in [8] and our result (mentioned in Sections VI and VII)

Comparison	Upper-bounds obtained in [8]	Upper-bounds obtained in Section VI
Definition of generalization error	Defined in Definitions 20 and 21.	Defined in eqs. (76) and (77).
Assumption on bounded moment generating function	[8, Theorem 17] is proven under a general assumption of bounded moment generating functions of the loss observables (see eqs. (QMGF) and (CMGF) in [8, Theorem 17]). However, [8, Corollaries 23 and 24] assumes sub-Gaussianity of loss observables (see eqs. (82) and (83))	All the results are based on sub-Gaussianity assumptions (see Assumptions 6 and 7).
Quantum terms involved	A single term based on quantum divergence (defined in Fact 12).	Two terms based on modified sandwiched quantum α -Rényi divergence (defined in Definition 12).
Classical terms involved	A single term based on divergence (defined in Fact 4).	A term based on γ -Rényi divergence (defined in Definition 2).
General Result	The bounds are not generalized.	Generalized family of upper-bounds depending on the values of $\alpha, \gamma \in (0, 1) \cap (1, \infty)$.
Recoverability	These results can't recover the results obtained in Section VI.	These results can easily recover results obtained in [8] (see Remark 7 for more details).
Whole sample-based bounds on the expected generalization error	[8, Theorem 17 and Corollary 23] study such upper-bounds.	Theorem 2, Corollaries 2 and 3 in this paper study such upper-bounds.
Individual sample-based bounds on the expected generalization error	[8, Corollary 24] studies such upper-bounds.	Corollary 5 in this paper study such upper-bounds.
Single-draw upper-bounds on the generalization error in probability	Not investigated in [8].	Theorems 4 and 5 in this paper study such upper-bounds.
Techniques involved	Variational lower-bound of classical and quantum relative entropies required (see Facts 8 and 15).	Variational lower-bound of classical Rényi divergence and modified sandwiched quantum Rényi divergence required (see Fact 6 and Lemma 4).

Regarding the expected generalization error, we have built upon the quantum learning framework introduced by [8] and have proposed a novel definition for true loss within the quantum learning context. Leveraging this framework, we have established a family of upper bounds on the expected generalization error, expressed in terms of modified sandwiched, Petz, and classical Rényi divergences. Notably, our bounds have encompassed the bound derived in [8] as a special case. Furthermore, under standard i.i.d. assumptions for both classical and quantum data, we have presented a family of upper bounds on the generalization error by using the modified sandwiched and Petz quantum Rényi divergences and the classical Rényi divergence. A detailed comparison between the initial results of this work and those of [8] has been provided in Table II. For the probabilistic behavior of the generalization error, we have derived a probabilistic bound based on the modified sandwich quantum Rényi divergence and the classical Rényi divergence.

Additionally, we have obtained another probabilistic bound formulated using the smooth max Rényi divergence. Importantly, all the upper bounds derived in this manuscript have held under specific sub-Gaussian assumptions for the loss observables. Within this work, we have demonstrated that these sub-Gaussian assumptions are a direct consequence of the boundedness of the loss observables (though the converse is not necessarily true), achieved by proving a quantum analogue of Hoeffding’s lemma for bounded self-adjoint operators.

Finally, the proofs of all the upper bounds presented herein have necessitated the evaluation of the variational form. To this end, we have newly introduced a variational lower bound for the modified quantum Rényi divergence. Moreover, we have presented an alternative proof for the variational lower bound of the Petz quantum Rényi divergence (Fact 17) and the Measurement-data-processing inequality of Petz quantum Rényi divergence (Fact 28) by employing the operator Hölder’s inequality (Fact 21) and the Araki-Lieb-Thirring inequality (Fact 18), thus circumventing the need for the standard data-processing inequality.

ACKNOWLEDGMENTS

The work of N. A. Warsi was supported in part by MTR/2022/000814, DST/INT/RUS/RSF/P-41/2021 from the Department of Science & Technology, Govt. of India and DCSW grant provided by the Indian Statistical Institute. The work of MH was supported in part by the National Natural Science Foundation of China under Grant 62171212 and the General R&D Projects of 1+1+1 CUHKCUHK(SZ)-GDST Joint Collaboration Fund (Grant No. GRDP2025-022).

REFERENCES

- [1] O. Fawzi, A. Oufkir, and D. S. França, “Lower bounds on learning pauli channels with individual measurements,” *IEEE Transactions on Information Theory*, vol. 71, no. 4, pp. 2642–2661, 2025.
- [2] M. Fanizza, A. Mari, and V. Giovannetti, “Optimal universal learning machines for quantum state discrimination,” *IEEE Transactions on Information Theory*, vol. 65, no. 9, pp. 5931–5944, 2019.
- [3] Y. Du, M.-H. Hsieh, T. Liu, S. You, and D. Tao, “Quantum differentially private sparse regression learning,” *IEEE Transactions on Information Theory*, vol. 68, no. 8, pp. 5217–5233, 2022.
- [4] G. De Palma, M. Marvian, D. Trevisan, and S. Lloyd, “The quantum wasserstein distance of order 1,” *IEEE Transactions on Information Theory*, vol. 67, no. 10, pp. 6627–6643, 2021.
- [5] A. Angrisani and E. Kashefi, “Quantum differential privacy in the local model,” *IEEE Transactions on Information Theory*, vol. 71, no. 5, pp. 3675–3692, 2025.
- [6] S. Sreekumar and M. Berta, “Limit distribution theory for quantum divergences,” *IEEE Transactions on Information Theory*, vol. 71, no. 1, pp. 459–484, 2025.
- [7] C. Hirche, C. Rouzé, and D. S. França, “Quantum differential privacy: An information theory perspective,” *IEEE Transactions on Information Theory*, vol. 69, no. 9, pp. 5771–5787, 2023.
- [8] M. C. Caro, T. Gur, C. Rouzé, D. Stilck França, and S. Subramanian, “Information-theoretic generalization bounds for learning from quantum data,” in *Proceedings of Thirty Seventh Conference on Learning Theory*, ser. Proceedings of Machine Learning Research, S. Agrawal and A. Roth, Eds., vol. 247. PMLR, 30 Jun–03 Jul 2024, pp. 775–839. [Online]. Available: <https://proceedings.mlr.press/v247/caro24a.html>
- [9] V. N. Vapnik and A. Y. Chervonenkis, “On the uniform convergence of relative frequencies of events to their probabilities,” *Theory of Probability & Its Applications*, vol. 16, no. 2, pp. 264–280, 1971. [Online]. Available: <https://doi.org/10.1137/1116025>
- [10] L. G. Valiant, “A theory of the learnable,” *Commun. ACM*, vol. 27, no. 11, p. 1134–1142, Nov. 1984. [Online]. Available: <https://doi.org/10.1145/1968.1972>
- [11] D. A. McAllester, “Some pac-bayesian theorems,” *Machine Learning*, vol. 37, no. 3, pp. 355–363, Dec 1999. [Online]. Available: <https://doi.org/10.1023/A:1007618624809>
- [12] T. Zhang, “Information-theoretic upper and lower bounds for statistical estimation,” *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1307–1321, 2006.
- [13] F. Hellström, G. Durisi, B. Guedj, and M. Raginsky, “Generalization bounds: Perspectives from information theory and pac-bayes,” 2024. [Online]. Available: <https://arxiv.org/abs/2309.04381>
- [14] D. McAllester, “A pac-bayesian tutorial with a dropout bound,” 2013. [Online]. Available: <https://arxiv.org/abs/1307.2118>
- [15] D. Russo and J. Zou, “Controlling bias in adaptive data analysis using information theory,” in *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, A. Gretton and C. C. Robert, Eds., vol. 51. Cadiz, Spain: PMLR, 09–11 May 2016, pp. 1232–1240. [Online]. Available: <https://proceedings.mlr.press/v51/russo16.html>
- [16] A. Xu and M. Raginsky, “Information-theoretic analysis of generalization capability of learning algorithms,” in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2017/file/ad71c82b22f4f65b9398f76d8be4c615-Paper.pdf
- [17] Y. Bu, S. Zou, and V. V. Veeravalli, “Tightening mutual information-based bounds on generalization error,” *IEEE Journal on Selected Areas in Information Theory*, vol. 1, no. 1, p. 121–130, May 2020. [Online]. Available: <http://dx.doi.org/10.1109/JSait.2020.2991139>
- [18] J. Negrea, M. Haghifam, G. K. Dziugaite, A. Khisti, and D. M. Roy, “Information-theoretic generalization bounds for sgld via data-dependent estimates,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [19] M. Haghifam, J. Negrea, A. Khisti, D. M. Roy, and G. K. Dziugaite, “Sharpened generalization bounds based on conditional mutual information and an application to noisy, iterative algorithms,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 9925–9935, 2020.
- [20] F. Hellström and G. Durisi, “A new family of generalization bounds using samplewise evaluated cmi,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 10 108–10 121, 2022.
- [21] E. Modak, H. Asnani, and V. M. Prabhakaran, “Rényi divergence based bounds on generalization error,” in *2021 IEEE Information Theory Workshop (ITW)*, 2021, pp. 1–6.
- [22] M. Sefidgaran, A. Gohari, G. Richard, and U. Simsekli, “Rate-distortion theoretic generalization bounds for stochastic learning algorithms,” in *Proceedings of Thirty Fifth Conference on Learning Theory*, ser. Proceedings of Machine Learning Research, P.-L. Loh and M. Raginsky, Eds., vol. 178. PMLR, 02–05 Jul 2022, pp. 4416–4463. [Online]. Available: <https://proceedings.mlr.press/v178/sefidgaran22a.html>
- [23] M. Sefidgaran, R. Chor, and A. Zaidi, “Rate-distortion theoretic bounds on generalization error for distributed learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 19 687–19 702, 2022.
- [24] O. Catoni, “Pac-bayesian supervised classification: The thermodynamics of statistical learning,” *Lecture Notes-Monograph Series*, vol. 56, pp. i–163, 2007. [Online]. Available: <http://www.jstor.org/stable/20461499>
- [25] A. R. Esposito, M. Gastpar, and I. Issa, “Generalization error bounds via rényi-, f-divergences and maximal leakage,” *IEEE Transactions on Information Theory*, vol. 67, no. 8, pp. 4986–5004, 2021.
- [26] M. D. Donsker and S. R. S. Varadhan, “Asymptotic evaluation of certain markov process expectations for large time, i,” *Communications on Pure and Applied Mathematics*, vol. 28, no. 1, pp. 1–47, 1975. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/cpa.3160280102>

- [27] S. Boucheron, G. Lugosi, and P. Massart, *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 02 2013. [Online]. Available: <https://doi.org/10.1093/acprof:oso/9780199535255.001.0001>
- [28] J. Birrell, P. Dupuis, M. A. Katsoulakis, L. Rey-Bellet, and J. Wang, “Variational representations and neural network estimation of rényi divergences,” *SIAM Journal on Mathematics of Data Science*, vol. 3, no. 4, pp. 1093–1116, 2021. [Online]. Available: <https://doi.org/10.1137/20M1368926>
- [29] D. Petz, “Quasi-entropies for finite quantum systems,” *Reports on Mathematical Physics*, vol. 23, no. 1, pp. 57–65, 1986. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/0034487786900674>
- [30] M. Müller-Lennert, F. Dupuis, O. Szechr, S. Fehr, and M. Tomamichel, “On quantum rényi entropies: A new generalization and some properties,” *Journal of Mathematical Physics*, vol. 54, no. 12, p. 122203, 12 2013. [Online]. Available: <https://doi.org/10.1063/1.4838856>
- [31] M. M. Wilde, A. Winter, and D. Yang, “Strong converse for the classical capacity of entanglement-breaking and hadamard channels via a sandwiched rényi relative entropy,” *Communications in Mathematical Physics*, vol. 331, no. 2, pp. 593–622, Oct 2014. [Online]. Available: <https://doi.org/10.1007/s00220-014-2122-x>
- [32] A. Rényi, “On measures of entropy and information,” *Proceeding of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, pp. 547–561, 1961.
- [33] N. A. Warsi, “Simple one-shot bounds for various source coding problems using smooth rényi quantities,” *Problems of Information Transmission*, vol. 52, no. 1, pp. 39–65, Jan 2016. [Online]. Available: <https://doi.org/10.1134/S0032946016010051>
- [34] W. Fenchel, *On conjugate convex functions*. Springer, 2014.
- [35] R. Bhatia, *Matrix analysis*. Springer Science & Business Media, 2013, vol. 169.
- [36] C. A. Fuchs, “Distinguishability and accessible information in quantum theory,” 1996. [Online]. Available: <https://arxiv.org/abs/quant-ph/9601020>
- [37] J. L. W. V. Jensen, “Sur les fonctions convexes et les inégalités entre les valeurs moyennes,” *Acta Mathematica*, vol. 30, no. 1, pp. 175–193, Dec 1906. [Online]. Available: <https://doi.org/10.1007/BF02418571>
- [38] M. J. Wainwright, *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, ser. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019.
- [39] F. Hiai and D. Petz, “The golden-thompson trace inequality is complemented,” *Linear Algebra and its Applications*, vol. 181, pp. 153–185, 1993. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/002437959390029N>
- [40] M. Hayashi, *Quantum Information Theory*. United States: Springer Cham, 2017.
- [41] M. Berta, O. Fawzi, and M. Tomamichel, “On variational expressions for quantum relative entropies,” *Letters in Mathematical Physics*, vol. 107, no. 12, pp. 2239–2265, Dec 2017. [Online]. Available: <https://doi.org/10.1007/s11005-017-0990-7>
- [42] R. L. Frank and E. H. Lieb, “Monotonicity of a relative rényi entropy,” *Journal of Mathematical Physics*, vol. 54, no. 12, p. 122201, 12 2013. [Online]. Available: <https://doi.org/10.1063/1.4838835>
- [43] K. Fang, H. Fawzi, and O. Fawzi, “Variational expressions and uhlmann theorem for measured divergences,” 2025. [Online]. Available: <https://arxiv.org/abs/2502.07745>
- [44] H. Araki, “On an inequality of lieb and thirring,” *Letters in Mathematical Physics*, vol. 19, no. 2, pp. 167–170, Feb 1990. [Online]. Available: <https://doi.org/10.1007/BF01045887>
- [45] E. Lieb and W. Thirring, *Inequalities for the moments of the eigenvalues of the schrödinger hamiltonian and their relation to sobolev inequalities*. Springer Berlin Heidelberg, 2005, pp. 205–239.
- [46] M. M. Wilde, *Quantum Information Theory*. Cambridge University Press, 2013. [Online]. Available: <https://doi.org/10.1017/CBO9781139525343>
- [47] F. Hellström and G. Durisi, “Generalization bounds via information density and conditional information density,” *IEEE Journal on Selected Areas in Information Theory*, vol. 1, no. 3, pp. 824–839, 2020.
- [48] M. Hayashi and H. Yamasaki, “Generalized quantum stein’s lemma and second law of quantum resource theories,” 2024, arXiv:2408.02722.
- [49] J. Radhakrishnan, P. Sen, and N. A. Warsi, “One-shot private classical capacity of quantum wiretap channel: Based on one-shot quantum covering lemma,” 2017. [Online]. Available: <https://arxiv.org/abs/1703.01932>
- [50] A. Anshu, R. Jain, and N. A. Warsi, “A one-shot achievability result for quantum state redistribution,” *IEEE Transactions on Information Theory*, vol. 64, no. 3, pp. 1425–1435, 2018.
- [51] —, “Building blocks for communication over noisy quantum networks,” *IEEE Transactions on Information Theory*, vol. 65, no. 2, pp. 1287–1306, 2019.

APPENDIX A PROOF OF LEMMA 1

We define $L' := L - \text{Tr}[L\rho]\mathbb{I}$ and note that $\text{Tr}[L'\rho] = 0$. Let L' and ρ have the following eigen decomposition,

$$L' = \sum_{i=1}^{|\mathcal{H}|} \alpha_i |i\rangle\langle i|, \text{ where } \forall i \in [|\mathcal{H}|], a - \text{Tr}[L\rho] < \alpha_i < b - \text{Tr}[L\rho],$$

$$\rho = \sum_{j=1}^{|\mathcal{H}|} \beta_j |j\rangle\langle j|, \text{ where } \forall j \in [|\mathcal{H}|], 0 < \beta_j < 1 \text{ and } \sum_{j=1}^{|\mathcal{H}|} \beta_j = 1.$$

Then, we have,

$$\log \text{Tr} \left[e^{\lambda L'} \rho \right] = \log \left(\sum_{i=1}^{|\mathcal{H}|} \sum_{j=1}^{|\mathcal{H}|} e^{\lambda \alpha_i \beta_j} |\langle i|j\rangle|^2 \right). \quad (126)$$

For all $i \in [|\mathcal{H}|]$, let $p_i := \sum_{j=1}^{|\mathcal{H}|} \beta_j |\langle i|j\rangle|^2$. It is easy to see that $\forall i \in [|\mathcal{H}|]$, $p_i \geq 0$ and $\sum_{i=1}^{|\mathcal{H}|} p_i = 1$. Thus,

$$\mathbb{E}_{\mathbf{a} \sim P}[\mathbf{a}] = \sum_{i=1}^{|\mathcal{H}|} \alpha_i p_i = \text{Tr} \left[\sum_{i=1}^{|\mathcal{H}|} \alpha_i |i\rangle\langle i| \left(\sum_{j=1}^{|\mathcal{H}|} \beta_j |j\rangle\langle j| \right) \right] = \text{Tr}[L'\rho] = 0. \quad (127)$$

We now upper-bound $\log \text{Tr} [e^{\lambda L'} \rho]$ as follows,

$$\begin{aligned}
\log \text{Tr} [e^{\lambda L'} \rho] &= \log \left(\sum_{i=1}^{|\mathcal{H}|} e^{\lambda \alpha_i} p_i \right) \\
&= \log \mathbb{E}_{\mathbf{a} \sim P} [e^{\lambda \mathbf{a}}] \\
&\stackrel{a}{\leq} \mathbb{E}_{\mathbf{a} \sim P} [\lambda \mathbf{a}] + \frac{\lambda^2 (b - \text{Tr}[L\rho] - a + \text{Tr}[L\rho])^2}{8} \\
&\stackrel{b}{=} \frac{\lambda^2 (b - a)^2}{8},
\end{aligned} \tag{128}$$

where a follows from Fact 7 and b follows from (127). This proves (35). We now prove (36) as follows,

$$\begin{aligned}
\log \text{Tr} [e^{\lambda L} \rho] &= \log \text{Tr} [e^{\lambda(L' + \text{Tr}[L\rho]\mathbb{I})} \rho] \\
&\stackrel{a}{=} \log \text{Tr} [e^{\lambda L'} e^{\lambda \text{Tr}[L\rho]\mathbb{I}} \rho] \\
&= \log \left(\text{Tr} [e^{\lambda L'} e^{\lambda \text{Tr}[L\rho]\mathbb{I}} \rho] \right) \\
&= \log \left(e^{\lambda \text{Tr}[L\rho]} \text{Tr} [e^{\lambda L'} \rho] \right) \\
&= \lambda \text{Tr}[L\rho] + \log \text{Tr} [e^{\lambda L'} \rho] \\
&\leq \lambda \text{Tr}[L\rho] + \frac{\lambda^2 (b - a)^2}{8},
\end{aligned}$$

where a follows since L' and \mathbb{I} commute with each other (otherwise, this equality in a is not true).

This completes the proof. ■

APPENDIX B PROOF OF THEOREM 1

Before proceeding to the proof of Theorem 1, we state the following lemma, which will be required to prove Theorem 1.

Lemma 6. *Suppose for each $(w, s) \in \mathcal{W} \times \mathcal{S}$, $\hat{L}(w, s)$ satisfies Assumption 5. Then, the following holds,*

$$\left| \text{Tr} [\hat{L}(w, s) \sigma^{\mathcal{A}_Q}(w, s)] - \text{Tr} [\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] \right| \leq \sqrt{2\mu^2 D(\sigma^{\mathcal{A}_Q}(w, s) \| \rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))}, \tag{129}$$

$$\left| \text{Tr} [\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w, s))] - \text{Tr} [\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(w))] \right| \leq \sqrt{2\mu^2 D(\sigma_{hyp}^{\mathcal{A}_Q}(w, s) \| \sigma_{hyp}^{\mathcal{A}_Q}(w))}. \tag{130}$$

Proof. See Appendix C for the proof. ■

Another way to prove Lemma 6 is via Lemma 5 by taking $\alpha \rightarrow 1$. However, it will require two assumptions, i.e., Assumptions 5 and 6.

We now consider the following series of inequalities,

$$\begin{aligned}
|\text{gen}| &= \left| \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\mathbb{E}_{\bar{S} \sim P_S} \left[\text{Tr} [\hat{L}(W, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W))] \right] \right] - \mathbb{E}_{S \sim P_{S|W}^{\mathcal{A}_Q}} \left[\text{Tr} [\hat{L}(W, S) \sigma^{\mathcal{A}_Q}(W, S)] \right] \right| \\
&\leq \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\left| \text{Tr} [\hat{L}(W, S) \sigma^{\mathcal{A}_Q}(W, S)] - \text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W))] \right| \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\left| \mathbb{E}_{\bar{S} \sim P_S} \left[\text{Tr} [\hat{L}(W, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W))] \right] - \mathbb{E}_{S \sim \hat{P}_{S|W}^{\mathcal{A}_Q}} \left[\text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W))] \right] \right| \right] \\
&\leq \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\left| \text{Tr} [\hat{L}(W, S) \sigma^{\mathcal{A}_Q}(W, S)] - \text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S))] \right| \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\left| \text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W, S))] - \text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W))] \right| \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{\mathcal{A}_Q}} \left[\left| \mathbb{E}_{\bar{S} \sim P_S} \left[\text{Tr} [\hat{L}(W, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W))] \right] - \mathbb{E}_{S \sim \hat{P}_{S|W}^{\mathcal{A}_Q}} \left[\text{Tr} [\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{\mathcal{A}_Q}(W))] \right] \right| \right]
\end{aligned}$$

$$\begin{aligned}
&\stackrel{a}{\leq} \mathbb{E}_{\substack{W \sim P_W^{A_Q} \\ S \sim P_{S|W}^{A_Q}}} \left[\sqrt{2\mu^2 D(\sigma^{A_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))} \right] + \mathbb{E}_{\substack{W \sim P_W^{A_Q} \\ S \sim P_{S|W}^{A_Q}}} \left[\sqrt{2\mu^2 D(\sigma_{hyp}^{A_Q}(W, S) \parallel \sigma_{hyp}^{A_Q}(W))} \right] \\
&\quad + \mathbb{E}_{W \sim P_W^{A_Q}} \left[\mathbb{E}_{\bar{S} \sim P_S} \left[\text{Tr} \left[\hat{L}(W, \bar{S}) (\rho_{te}(\bar{S}) \otimes \sigma_{hyp}^{A_Q}(W)) \right] \right] - \mathbb{E}_{S \sim P_{S|W}^{A_Q}} \left[\text{Tr} \left[\hat{L}(W, S) (\rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W)) \right] \right] \right] \\
&\stackrel{b}{\leq} \mathbb{E}_{\substack{W \sim P_W^{A_Q} \\ S \sim P_{S|W}^{A_Q}}} \left[\sqrt{2\mu^2 D(\sigma^{A_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))} \right] + \mathbb{E}_{\substack{W \sim P_W^{A_Q} \\ S \sim P_{S|W}^{A_Q}}} \left[\sqrt{2\mu^2 D(\sigma_{hyp}^{A_Q}(W, S) \parallel \sigma_{hyp}^{A_Q}(W))} \right] + \sqrt{2\tau^2 I[S; W]} \\
&= \mathbb{E}_{(W, S) \sim P_{WS}^{A_Q}} \left[\sqrt{2\mu^2 D(\sigma^{A_Q}(W, S) \parallel \rho_{te}(S) \otimes \sigma_{hyp}^{A_Q}(W, S))} + \sqrt{2\mu^2 D(\sigma_{hyp}^{A_Q}(W, S) \parallel \sigma_{hyp}^{A_Q}(W))} \right] + \sqrt{2\tau^2 I[S; W]},
\end{aligned}$$

where a follows from eqs. (129) and (130) and b follows from [16, Theorem 1] under the classical sub-Gaussianity assumptions mentioned in (87). This completes the proof of Theorem 1. \blacksquare

APPENDIX C PROOF OF LEMMA 6

We first prove (129) in two cases and later we show that the proof of (130) follows similarly. Towards this, for all $\lambda \in \mathbb{R}$ we consider the following series of inequalities,

$$\begin{aligned}
&D(\sigma^{A_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \\
&\stackrel{a}{\geq} \lambda \text{Tr} \left[\hat{L}(w, s) \sigma^{A_Q}(w, s) \right] - \log \text{Tr} \left[e^{\lambda \hat{L}(w, s)} (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right]
\end{aligned} \tag{131}$$

$$\begin{aligned}
&= \lambda \left(\text{Tr} \left[\hat{L}(w, s) \sigma^{A_Q}(w, s) \right] - \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \right) \\
&\quad + \lambda \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] - \log \text{Tr} \left[e^{\lambda \hat{L}(w, s)} (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \\
&= \lambda \left(\text{Tr} \left[\hat{L}(w, s) \sigma^{A_Q}(w, s) \right] - \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \right) \\
&\quad - \log \text{Tr} \left[e^{\lambda (\hat{L}(w, s) - \text{Tr} [\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s))])} (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \\
&\stackrel{b}{\geq} \lambda \left(\text{Tr} \left[\hat{L}(w, s) \sigma^{A_Q}(w, s) \right] - \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \right) - \frac{\lambda^2 \mu^2}{2},
\end{aligned} \tag{132}$$

where a follows from Fact 18 and 24, and b follows from (85).

We can rewrite the inequality (132) as follows:

$$\frac{\lambda^2 \mu^2}{2} - \lambda \left(\text{Tr} \left[\hat{L}(w, s) \sigma^{A_Q}(w, s) \right] - \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \right) + D(\sigma^{A_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \geq 0,$$

Since the above inequality is a non-negative quadratic equation in λ with the coefficient $\left(\frac{\mu^2}{2}\right) \geq 0$, therefore its discriminant must be non-positive. Thus, we have the following inequality,

$$\begin{aligned}
&\left(\text{Tr} \left[\hat{L}(w, s) \sigma^{A_Q}(w, s) \right] - \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \right)^2 \leq 4 \left(\frac{\mu^2}{2} \right) D(\sigma^{A_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \\
&\Rightarrow \left| \text{Tr} \left[\hat{L}(w, s) \sigma^{A_Q}(w, s) \right] - \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] \right| \leq \sqrt{2\mu^2 D(\sigma^{A_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s))}.
\end{aligned} \tag{133}$$

We now proceed to prove (130). Using eq. (86) and a calculation similar to eqs. (132) and (133) we have the following inequality,

$$\begin{aligned}
&\left| \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s)) \right] - \text{Tr} \left[\hat{L}(w, s) (\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w)) \right] \right| \leq \sqrt{2\mu^2 D(\rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w, s) \parallel \rho_{te}(s) \otimes \sigma_{hyp}^{A_Q}(w))} \\
&\stackrel{a}{=} \sqrt{2\mu^2 D(\sigma_{hyp}^{A_Q}(w, s) \parallel \sigma_{hyp}^{A_Q}(w))},
\end{aligned}$$

where a follows from (16) of Fact 14 and Fact 12. This completes the proof of Lemma 6. \blacksquare

APPENDIX D
PROOF OF THEOREM 3

Consider the following events

$$E := \{(w, s) \in \mathcal{W} \times \mathcal{Z}^n : |\text{gen}(w, s)| > \varepsilon\},$$

$$A(\nu) := \left\{ (w, s) \in \mathcal{W} \times \mathcal{Z}^n : \log \frac{P_{WS}(w, s)}{P_W(w) \times P_S(s)} < I_{\max}^{(\nu)}[W; S] \right\},$$

where $\nu \in (0, 1)$ and for any $w \in \mathcal{W}$, we define $E_w := \{s \in \mathcal{Z}^n : (w, s) \in E\}$. Then we can write the following,

$$\begin{aligned} \Pr_{(W,S) \sim P_{WS}} \{E\} &= \Pr_{(W,S) \sim P_{WS}} \{E \cap A(\nu)\} + \Pr_{(W,S) \sim P_{WS}} \{E \cap A^c(\nu)\} \\ &\leq \Pr_{(W,S) \sim P_{WS}} \{E \cap A(\nu)\} + \Pr_{(W,S) \sim P_{WS}} \{A^c(\nu)\} \\ &\stackrel{a}{\leq} \sum_{(w,s) \in E \cap A(\nu)} P_{WS}(w, s) + \nu \\ &\stackrel{b}{\leq} \sum_{(w,s) \in E \cap A(\nu)} (P_W(w) \times P_S(s)) e^{I_{\max}^{(\nu)}[W; S]} + \nu \\ &\leq e^{I_{\max}^{(\nu)}[W; S]} \Pr_{(W,S) \sim P_W \times P_S} \{E \cap A(\nu)\} + \nu \\ &\leq e^{I_{\max}^{(\nu)}[W; S]} \Pr_{(W,S) \sim P_W \times P_S} \{E\} + \nu \\ &= e^{I_{\max}^{(\nu)}[W; S]} \mathbb{E}_{W \sim P_W} \left[\Pr_{S \sim P_S} \{E_W\} \right] + \nu. \end{aligned} \tag{134}$$

For each $w \in \mathcal{W}$, from eqs. (51) and (53), we now consider the following series of ineqlities,

$$\begin{aligned} \Pr_{S \sim P_S} \{E_w\} &= \Pr_{S \sim P_S} \{|\text{gen}(w, S)| > \varepsilon\} \\ &= \Pr_{S \sim P_S} \left\{ \left| \frac{1}{n} \sum_{i=1}^n l(w, Z_i) - \mathbb{E}_{\bar{Z}}[l(w, \bar{Z})] \right| > \varepsilon \right\} \\ &\stackrel{a}{\leq} 2e^{-\frac{n\varepsilon^2}{2\tau^2}}, \end{aligned} \tag{135}$$

where, a follows from Fact 9 as for each $w \in \mathcal{W}$, $l(w, Z)$ is τ -sub-Gaussian. Thus combining eqs. (134) and (135), we write the following,

$$\begin{aligned} \Pr_{(W,S) \sim P_{WS}} \{E\} &\leq e^{I_{\max}^{(\nu)}[W; S]} 2e^{-\frac{n\varepsilon^2}{2\tau^2}} + \nu \\ &= e^{-\frac{n\varepsilon^2}{2\tau^2} + I_{\max}^{(\nu)}[W; S] + \log 2} + \nu, \end{aligned} \tag{136}$$

If we choose $\delta := e^{-\frac{n\varepsilon^2}{2\tau^2} + I_{\max}^{(\nu)}[W; S] + \log 2} + \nu$, then ε can be written as follows,

$$\varepsilon = \sqrt{\frac{2\tau^2}{n} \left(\log 2 + I_{\max}^{(\nu)}[W; S] + \log \left(\frac{1}{\delta - \nu} \right) \right)}, \tag{137}$$

Hence, from eqs. (136) and (137) we have the following,

$$\Pr_{(W,S) \sim P_{WS}} \left\{ |\text{gen}(W, S)| \leq \sqrt{\frac{2\tau^2}{n} \left(\log 2 + I_{\max}^{(\nu)}[W; S] + \log \left(\frac{1}{\delta - \nu} \right) \right)} \right\} \geq 1 - \delta.$$

This completes the proof of Theorem 3. ■