

Quantum secure learning with classical samples

Wooyeong Song,^{1,*} Youngrong Lim,^{2,*} Hyukjoon Kwon,³ Gerardo Adesso,⁴
Marcin Wieśniak,^{5,6} Marcin Pawłowski,⁶ Jaewan Kim,² and Jeongho Bang^{7,†}

¹*Department of Physics, Hanyang University, Seoul 04763, Korea*

²*School of Computational Sciences, Korea Institute for Advanced Study, Seoul 02455, Korea*

³*QOLS, Blackett Laboratory, Imperial College London, London SW7 2AZ, United Kingdom*

⁴*School of Mathematical Sciences and Centre for the Mathematics
and Theoretical Physics of Quantum Non-Equilibrium Systems,
University of Nottingham, University Park, Nottingham NG7 2RD, United Kingdom*

⁵*Institute of Theoretical Physics and Astrophysics, Faculty of Mathematics,
Physics and Informatics, University of Gdańsk, 80-308 Gdańsk, Poland*

⁶*International Centre for Theory of Quantum Technologies, University of Gdańsk, 80-308 Gdańsk, Poland*

⁷*Electronics and Telecommunications Research Institute, Daejeon 34129, Korea*

(Received April 12, 2021)

Studies addressing the question “Can a learner complete the learning securely?” have recently been spurred from the standpoints of fundamental theory and potential applications. In the relevant context of this question, we present a classical-quantum hybrid sampling protocol and define a security condition that allows only legitimate learners to prepare a finite set of samples that guarantees the success of the learning; the security condition excludes intruders. We do this by combining our security concept with the bound of the so-called probably approximately correct (PAC) learning. We show that while the lower bound on the learning samples guarantees PAC learning, an upper bound can be derived to rule out adversarial learners. Such a secure learning condition is appealing, because it is defined only by the size of samples required for the successful learning and is independent of the algorithm employed. Notably, the security stems from the fundamental quantum no-broadcasting principle. No such condition can thus occur in any classical regime, where learning samples can be copied. Owing to the hybrid architecture, our scheme also offers a practical advantage for implementation in noisy intermediate-scale quantum devices.

I. INTRODUCTION

The hybridization of machine learning and quantum theory has been intensively studied, especially to explore the possibility of exploiting quantum learning speedups. Very recently, the incorporation of useful quantum-algorithm-kernel (e.g., quantum linear solvers [1]) into data processing tasks in machine learning has yielded encouraging results [2–5]. Within a span of a few years, such approaches have become increasingly important in quantum computation, leading to the advent of quantum machine learning [6, 7].

In parallel, the issue of security has been of considerable interest to the machine learning community. The term “secure learning” is usually used to indicate that the learning is allowed only for the legitimate learner, who wants to rule out adversarial learners. The main objective of these adversaries is to acquire ability to become equals of the legitimate learner or to render the learning of the legitimate learner counterproductive. In this context, one of the open issues is how to define a secure learning condition for detecting and preventing these adversaries. While this problem has been widely studied in classical learning [8, 9], only a few quantum mechanical studies have been conducted so far [10–12].

We indicate that the legitimate learning mates can communicate a (classically) encrypted dataset after generating a secret key via a well-established quantum-key-distribution (QKD) scheme. In that case, it would be impractical for the adversarial learner(s) to extract critical learning information once the QKD is completed. However, the adversarial learner(s) may want to spoil the learning by disrupting the communication. Such a purpose can be achieved simply by disrupting the encrypted data after the key is distributed. This is actually one of the distinctive aspects of the learning security [8]. Thus, the learning security can neither be fully achieved nor defined by the QKD alone.

Having the above in mind, we in this paper construct a secure learning condition with favorable quantum properties. To this end, we first design a protocol for secure sampling that runs between two legitimate learning parties. We cast a classical-quantum hybrid oracle that allows large-size classical inputs with a small-scale quantum system [13]. As the main result, we derive a secure learning condition such that only the original legitimate learner is guaranteed success for learning; we designate the condition as the *secure probably-approximately-correct (PAC) learning* condition. The beauty of this condition is that the security is derived only from the size of learning samples the legitimate learner requires and it stems from the quantum no-broadcasting principle [14, 15]. Therefore, such condition cannot be defined in any classical regime. Our paper also leads to an intriguing classical-quantum interplay, namely, in which

* The first two authors contributed equally to this work

† jbang@etri.re.kr

the (large) input data remain classical while the useful quantum properties are explored for a small quantum system [16, 17]. Such architecture helps avoid the use of a largely superposed sample and is well suited to noisy intermediate-scale quantum (NISQ) technologies [18].

II. PROBLEM

Given a (Boolean) function $c \in \mathcal{C}$ that maps the input $\mathbf{x} = x_0x_1 \cdots x_{n-1}$ to a binary value $c(\mathbf{x}) \in \{0, 1\}$, learning is defined as the process of identifying a hypothesis $h \in \mathcal{H}$ close to c . The binary number $x_j \in \{0, 1\}$ ($j = 0, 1, \dots, n-1$) can be considered as the “feature” and the size of the hypothesis set $|\mathcal{H}|$, called “model complexity,” is assumed to be finite. Such a problem covers a wide variety of learning tasks. In particular, this binary setting of the problem can, in principle, be extended to a more general situation such as multi-class tasks [19]. For this reason, the binary classification framework has generally been used in computational learning theory [20, 21].

In such a problem, the learner, say Alice (\mathcal{A}), should first sample a set T of input-target pairs, where $T = \{(\mathbf{x}, c(\mathbf{x}))\}$. To accomplish this sampling, \mathcal{A} employs a black box, called the oracle. The oracle is responsible for accessing critical information, namely, $c(\mathbf{x})$ for a given \mathbf{x} . Here, we assume that the oracle is owned by \mathcal{A} 's distant partner, say Bob (\mathcal{B}). Such an assumption, namely, of the two learning parties being located far apart, is commonly invoked in secure learning [8]. The issue is then how \mathcal{A} can sample a clean dataset T with \mathcal{B} in a manner that is secure against any malicious attack; in other words, how can \mathcal{A} learn c securely?

III. SECURE SAMPLING PROTOCOL

We introduce a classical-quantum hybrid oracle $O(c)$, which consists of input and output channels for n -bit classical data \mathbf{x} and for a single qubit, denoted by $\mathcal{C}^{\mathcal{A}\mathcal{B}}$ and $\mathcal{Q}^{\mathcal{A}\mathcal{B}}$, respectively. This oracle $O(c)$ implements $(\mathbf{x}, |\alpha\rangle) \rightarrow (\mathbf{x}, |c(\mathbf{x}) \oplus \alpha\rangle)$ for $\alpha \in \{0, 1\}$ and $(\mathbf{r}, |\alpha\rangle) \rightarrow (\mathbf{r}, |\alpha\rangle)$ for $\alpha \in \{+, -\}$, where $|c(\mathbf{x})\rangle$ is the oracle-answer for a given \mathbf{x} . Here, \mathbf{r} is a *random* input which is casted for the purpose of testing the existence of any malicious intruder who disturbs the communication. Thus, \mathbf{r} is chosen such that $(\mathbf{r}, y) \notin T$ (for any $y \in \{0, 1\}$). The construction of such an operation is fairly common, e.g., in QKD or quantum secure direct communication schemes [22, 23]. Note that it is not permissible to extract any information by looking into $O(c)$. A useful hybrid oracle architecture is presented in Appendix A.

We now present the secure sampling protocol, which proceeds as follows. First, \mathcal{A} prepares the state $|\alpha\rangle$ as an eigenstate of $\hat{\sigma}_z$ or $\hat{\sigma}_x$ (i.e., $|\alpha\rangle \in \{|0\rangle, |1\rangle, |\pm\rangle = \frac{1}{\sqrt{2}}(|0\rangle \pm |1\rangle)\}$) at random. The prepared state $|\alpha\rangle$ is transferred to \mathcal{B} through $\mathcal{Q}^{\mathcal{A}\mathcal{B}}$. If $|\alpha\rangle = |0\rangle$ or $|1\rangle$, \mathcal{A}

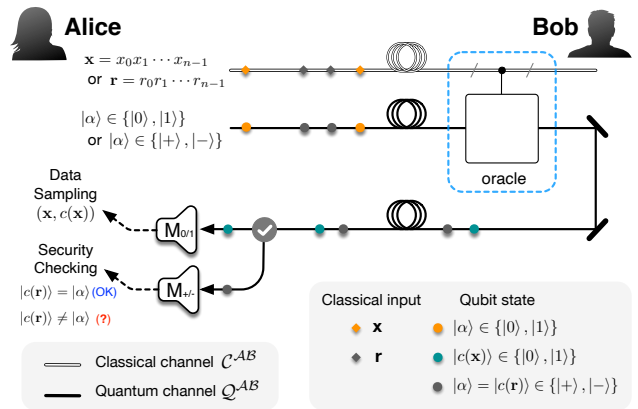


FIG. 1. Schematic of our sampling protocol. Alice (\mathcal{A}) has facilities for the preparation of inputs, $(\mathbf{x}, |\alpha\rangle \in \{|0\rangle, |1\rangle\})$ or $(\mathbf{r}, |\alpha\rangle \in \{|+\rangle, |-\rangle\})$. \mathcal{A} can also perform a single-qubit measurement to identify the returning qubit. Bob (\mathcal{B}) owns the oracle. Here, we consider a classical-quantum hybrid architecture (blue dashed and solid boxes) with a classical input (\mathbf{x} or \mathbf{r}) and an ancillary qubit state ($|\alpha\rangle$). The oracle does not reveal its structure. \mathcal{A} and \mathcal{B} communicate via classical and quantum channels, denoted by $\mathcal{C}^{\mathcal{A}\mathcal{B}}$ and $\mathcal{Q}^{\mathcal{A}\mathcal{B}}$, respectively.

sends the input \mathbf{x} through $\mathcal{C}^{\mathcal{A}\mathcal{B}}$ together with $|\alpha\rangle$, and if $|\alpha\rangle = |\pm\rangle$, \mathcal{A} draws a random input \mathbf{r} . Subsequently, $(\mathbf{x}, |\alpha\rangle \in \{|0\rangle, |1\rangle\})$ or $(\mathbf{r}, |\alpha\rangle \in \{|+\rangle, |-\rangle\})$ are passed through the oracle $O(c)$, and the output states $|c(\mathbf{x})\rangle$ or $|\pm\rangle$ of the qubit are returned to \mathcal{A} . For $|\alpha\rangle \in \{|0\rangle, |1\rangle\}$, \mathcal{A} obtains a sample pair $(\mathbf{x}, c(\mathbf{x}))$ by performing $\hat{\sigma}_z$ measurement, and for $|\alpha\rangle \in \{|+\rangle, |-\rangle\}$, \mathcal{A} should receive $|\alpha\rangle = |\pm\rangle$ from \mathcal{B} . Therefore, by checking the returned state $|\pm\rangle$ with the $\hat{\sigma}_x$ measurement, \mathcal{A} can sense any adversarial learner, often referred to as Eve (\mathcal{E}), who alters the qubits moving $\mathcal{A} \rightarrow \mathcal{B}$ or $\mathcal{B} \rightarrow \mathcal{A}$ (see Fig. 1). Note that $(\mathbf{r}, y) \notin T$ for any $y \in \{0, 1\}$ obtained by $\hat{\sigma}_z$ measurement, and it cannot be a valid sample.

IV. NO-BROADCASTING OF LEARNING SAMPLES

With the protocol described above, we present our first result:

Theorem 1. *In our protocol, for any given $c \in \mathcal{C}$, \mathcal{B} cannot distribute the full set of learning samples, namely, $T = \{(\mathbf{x}, c(\mathbf{x}))\}$, to \mathcal{A} and other (external) learners. Therefore, the condition*

$$T = T^{(k)} \quad (\forall k \in [1, L]), \quad (1)$$

where $T^{(k)}$ is the set of samples that the k learner (i.e., \mathcal{A} or \mathcal{E}) finally gets for strategy \mathcal{E} , cannot be satisfied.

For proving this theorem, we let $\hat{\rho}_0 = |c(\mathbf{x})\rangle \langle c(\mathbf{x})|$ and $\hat{\rho}_1 = |\alpha\rangle \langle \alpha|$, each of which is defined in terms of a state of the ideal oracle output in a trial for a given input (\mathbf{x} or \mathbf{r}). Here, $\hat{\rho}_0 \in \{|0\rangle \langle 0|, |1\rangle \langle 1|\}$ and $\hat{\rho}_1 \in \{|+\rangle \langle +|, |-\rangle \langle -|\}$.

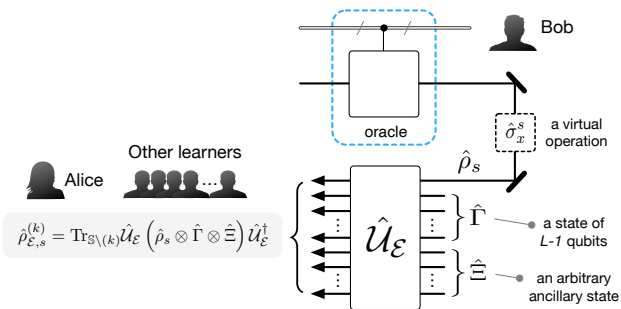


FIG. 2. General attack by adversarial learners. Here, we consider $L-1$ adversarial learners who can freely access $\mathcal{C}^{\mathcal{A}, \mathcal{B}}$ and $\mathcal{Q}^{\mathcal{A}, \mathcal{B}}$. Each adversarial learner has his or her own (in principle, infinite size) ancillary system and is assumed to be an expert in quantum theory. We further assume that the adversarial learners can team up to process an optimal strategy \mathcal{E} for their own or for the group's benefit.

Suppose \mathcal{B} adopts a strategy \mathcal{E} to distribute the samples in T among learners (including \mathcal{A}), with $L \geq 2$. In general, \mathcal{E} can be represented as a completely positive and trace-preserving map with an overall unitary $\hat{U}_{\mathcal{E}}$ and an arbitrary ancilla state $\hat{\Xi}$ (see Fig. 2). The distributed states can be written such that

$$\hat{\rho}_{\mathcal{E},s}^{(k)} = \text{Tr}_{\mathbb{S} \setminus \{k\}} \hat{U}_{\mathcal{E}} \left(\hat{\rho}_s \otimes \hat{\Gamma} \otimes \hat{\Xi} \right) \hat{U}_{\mathcal{E}}^\dagger, \quad (2)$$

where $\text{Tr}_{\mathbb{S} \setminus \{k\}}$ denotes the partial trace with respect to all systems \mathbb{S} except the one labeled with the k th learner, and $\hat{\Gamma}$ represents a state of $L-1$ qubits, each of which is distributed to the corresponding learner, except \mathcal{A} (here, $k=1$ denotes \mathcal{A}). Then, it is true that \mathcal{B} cannot *broadcast* the states $\hat{\rho}_s$ ($s=0,1$) to a (k -indexed) learner. This is confirmed by the principle that the states $\hat{\rho}_0$ and $\hat{\rho}_1$ are not distinguishable [14, 24, 25]. Therefore, a sample pair $(\mathbf{x}, c(\mathbf{x}))$ cannot be shared for a given \mathbf{x} . Thus, the full set of T cannot be distributed in the complete form and Theorem 1 holds.

V. SECURE PROBABLY-APPROXIMATELY-CORRECT LEARNING

Suppose that \mathcal{A} is the only legitimate learner, and the other $L-1$ learners are malicious intruders. Without loss of generality, we let $k \in \{\mathcal{A}, \mathcal{E}\}$ with $L=2$, or equivalently, by assuming that all $L-1$ intruders team up together as one \mathcal{E} . In this setting, we can assume that \mathcal{E} is a general attack strategy adopted by \mathcal{E} . Then, Theorem 1 describes the following situation: if \mathcal{E} disturbs the protocol, the samples prepared by \mathcal{A} (and also \mathcal{E}) must be *noisy*; specifically, a portion $\eta^{\mathcal{A}}$ (and $\eta^{\mathcal{E}}$) of the contaminated samples, for example, $(\mathbf{x}, c(\mathbf{x}) \oplus 1)$, would be included in \mathcal{A} 's (and \mathcal{E} 's) samples. Note that \mathcal{A} and \mathcal{E} cannot identify these contaminations. Here, $\eta^{(k)} \leq \frac{1}{2}$ ($k \in \{\mathcal{A}, \mathcal{E}\}$) and is determined by \mathcal{E} 's strategy \mathcal{E} . It

can be written as (for $T \gg 1$)

$$\eta^{(k)} = 1 - \frac{|T_S^{(k)}|}{|T^{(k)}|} \geq 1 - \min_s F(\hat{\rho}_s, \hat{\rho}_{\mathcal{E},s}^{(k)}), \quad (3)$$

where $T_S^{(k)}$ denotes the set of uncontaminated samples in $T^{(k)}$; thus, $T_S^{(k)} \subseteq T^{(k)}$ and $T_S^{(k)} \subseteq T$. $F(\hat{\rho}, \hat{\sigma})$ is the fidelity between the states $\hat{\rho}$ and $\hat{\sigma}$ [26]. Here, the inequality in the rightmost side is introduced because \mathcal{E} would make the contaminated samples even in cases where $\hat{\rho}_s$ are correctly cloned [?]. The equality is always saturated for \mathcal{A} . We then assume that our protocol forbids any strategy \mathcal{E} that allows the condition

$$(\eta_c > \eta^{\mathcal{A}}) \wedge (\eta_c > \eta^{\mathcal{E}}) \quad (4)$$

with a critical factor η_c . This assumption is true when η_c is chosen such that $\eta_c = 1 - F_{\text{opt}}$ where F_{opt} is the optimal fidelity achievable by a $(1 \rightarrow 2)$ $\hat{\rho}_s$ cloner [?]. Then, Eq. (4) can be rewritten by using Eq. (3) as

$$(F_{\text{opt}} < F(\hat{\rho}_s, \hat{\rho}_{\mathcal{E},s}^{\mathcal{A}})) \wedge (F_{\text{opt}} < F(\hat{\rho}_s, \hat{\rho}_{\mathcal{E},s}^{\mathcal{E}})), \quad (5)$$

which immediately contradicts the quantum no-cloning principle [15]. We note that if Alice could acquire information about Eve's attack scenario (if any), it might be possible to consider a more useful η_c setting. If $\eta_c = 0$, Eq. (4) becomes equivalent to the condition Eq. (1) and we encounter Theorem 1.

We now discuss secure learning in the framework of the so-called PAC learning [21, 27]. In a PAC learning, the concept class C is said to be (ϵ, δ) -PAC learnable [we call the learner a (ϵ, δ) -PAC learner] if an ϵ -approximated correct solution (i.e., hypothesis) $h \in H$ can be found with a probability $1 - \delta$; in other words, C is said to be (ϵ, δ) -PAC learnable if $P[E(h, c) \leq \epsilon] \geq 1 - \delta$ is satisfied for any $c \in C$, where $E(h, c)$ is an error function that indicates how h and c differ [21]. Such a theorem of PAC learning indicates that if a learner is allowed to use a certain size, say $M_b(\epsilon, \delta)$, of contaminated samples with η , he or she is guaranteed to be a (ϵ, δ) -PAC learner. In this case, η is defined as the percentage of contaminated samples in the entire set of samples [refer to Eq. (3)]. Usually, $M_b(\epsilon, \delta)$ is referred to as "sample complexity" [20, 28]. Here, $M_b(\epsilon, \delta)$ is divided into two categories depending on whether the samples are ideal (i.e., $\eta = 0$) or noisy (i.e., $\eta \in (0, \frac{1}{2})$) (For more details, see Appendix B, Refs. [21, 27], and the informative summary in Chap. 5 of Ref. [29]). The latter, namely, the noisy PAC learning model, provides a useful framework and is suitable for our paper because contaminations, either from \mathcal{E} or from imperfection intrinsic to the channels, can be included in the expression for $\eta^{(k)}$.

It is noteworthy that the (full) quantum model of the PAC learning, namely, quantum PAC learning, was also developed by using a quantum oracle that allows the (large) superposition of the inputs \mathbf{x} [29]. However, no study has been conducted on secure learning in a classical or a quantum PAC learning framework.

We now present our second result:

Theorem 2. For any given $c \in C$, let $M_b^{\mathcal{A}}(\epsilon, \delta)$ and $M_b^{\mathcal{E}}(\epsilon, \delta)$ denote the “optimal” sample complexities of \mathcal{A} and \mathcal{E} , respectively [?]. Then, during the running of our protocol, if \mathcal{A} becomes a (ϵ, δ) -PAC learner by identifying the samples smaller than $M_b^{\mathcal{E}}(\epsilon, \delta)$, \mathcal{E} cannot become a (ϵ, δ) -PAC learner for the same ϵ and δ .

The proof of this theorem is as follows. First, consider the case $\eta^{\mathcal{A}} \geq \eta^{\mathcal{E}}$, which will lead to $M_b^{\mathcal{A}}(\epsilon, \delta) \geq M_b^{\mathcal{E}}(\epsilon, \delta)$. In this case, it is impossible for \mathcal{A} to be a (ϵ, δ) -PAC learner with M samples smaller than $M_b^{\mathcal{E}}(\epsilon, \delta)$. Second, in the case of $\eta^{\mathcal{A}} < \eta^{\mathcal{E}}$, if \mathcal{A} completes the learning with M samples and becomes a (ϵ, δ) -PAC learner satisfying $M_b^{\mathcal{E}}(\epsilon, \delta) > M \geq M_b^{\mathcal{A}}(\epsilon, \delta)$, then \mathcal{E} cannot simultaneously be a (ϵ, δ) -PAC learner because the protocol will be terminated before \mathcal{E} obtains a sufficient number of samples (i.e., larger than $M_b^{\mathcal{E}}(\epsilon, \delta)$) to be a (ϵ, δ) -PAC learner. This proves Theorem 2.

On the basis of the above analysis, we present a definition for a secure learner:

Definition 1. For any $c \in C$, suppose \mathcal{A} identifies h with M samples, with

$$M_c(\epsilon, \delta) \geq M \geq M_b(\epsilon, \delta). \quad (6)$$

Here, $M_b(\epsilon, \delta)$ and $M_c(\epsilon, \delta)$ are defined as $M_b^{(k)}(\epsilon, \delta)$ when $\eta^{(k)} \rightarrow 0$ and $\eta^{(k)} \rightarrow \eta_c$, respectively, where k is either \mathcal{A} or \mathcal{E} . Then, we call \mathcal{A} a quantum secure (ϵ, δ) -PAC learner.

In this definition, the lower bound of the sample size [i.e., $M \geq M_b(\epsilon, \delta)$] is necessary for \mathcal{A} to be a (ϵ, δ) -PAC learner. The upper bound [i.e., $M_c(\epsilon, \delta) \geq M$] is adopted for security, and it follows from Theorem 2 and Eq. (4).

For wide applicability of Theorem 1, 2 and Definition 1, we apply two additional rules: **(R.1)** When the number of trials for $(\mathbf{r}, |\alpha\rangle)$ reaches $M_b(\epsilon, \delta) - \Gamma$, then \mathcal{A} tests whether $\frac{M_{c(\mathbf{r}) \neq \alpha}}{M_b(\epsilon, \delta) - \Gamma}$ is larger than $\eta_c - \Delta$, where $M_{c(\mathbf{r}) \neq \alpha}$ is the number of inconsistent results [i.e., $c(\mathbf{r}) \neq \alpha$] in \mathcal{A} 's $\hat{\sigma}_x$ measurement. If $\frac{M_{c(\mathbf{r}) \neq \alpha}}{M_b(\epsilon, \delta) - \Gamma} \geq \eta_c - \Delta$, \mathcal{A} suspends the process by confirming that the state change, namely, $|\pm\rangle \rightarrow |\mp\rangle$, occurs by \mathcal{E} ; otherwise, \mathcal{A} continues the process. Here, we approximate

$$\eta^{\mathcal{A}} \simeq \frac{M_{c(\mathbf{r}) \neq \alpha}}{M_b(\epsilon, \delta)} \quad (7)$$

by assuming $M_{c(\mathbf{x}) \rightarrow c(\mathbf{x}) \oplus 1} = M_{c(\mathbf{r}) \neq \alpha}$, where $M_{c(\mathbf{x}) \rightarrow c(\mathbf{x}) \oplus 1}$ denotes the number of contaminated pairs in \mathcal{A} 's sample set after a certain number of trials. This assumption is reasonable because \mathcal{A} generates $(\mathbf{r}, |\alpha\rangle \in \{|+\rangle, |-\rangle\})$ or $(\mathbf{x}, |\alpha\rangle \in \{|0\rangle, |1\rangle\})$ with probability $\frac{1}{2}$, which cannot be discriminated by \mathcal{E} . **(R.2)** If the learning is not completed until the number of trials for $(\mathbf{x}, |\alpha\rangle)$ reaches $M_c(\epsilon, \delta)$, \mathcal{A} quits the process. It is to be noted that the factors Γ and Δ in **(R.1)** are introduced to limit the quality of \mathcal{E} 's learning.

We can now analyze the possible situations. First, let us consider the case (i) $\eta^{\mathcal{A}} \geq \eta^{\mathcal{E}}$. Then, the following

two subcases can be considered:

$$(i-a) \eta^{\mathcal{A}} \geq \eta_c - \Delta \geq \eta^{\mathcal{E}} \text{ and } (i-b) \eta^{\mathcal{A}} \geq \eta^{\mathcal{E}} \geq \eta_c - \Delta.$$

However, cases (i-a) and (i-b) do not actually happen because **(R.1)** will halt the process when $\eta^{\mathcal{A}} \geq \eta_c - \Delta$; hence \mathcal{E} is not allowed to become a (ϵ, δ) -PAC learner. Second, for the case (ii) $\eta^{\mathcal{A}} < \eta^{\mathcal{E}}$, we can also consider the following two subcases:

$$(ii-a) \eta^{\mathcal{E}} > \eta_c - \Delta \geq \eta^{\mathcal{A}} \text{ and } (ii-b) \eta^{\mathcal{E}} > \eta^{\mathcal{A}} \geq \eta_c - \Delta.$$

In case (ii-a), if \mathcal{A} can learn $h \simeq c$ (for any given ϵ and δ) with M samples, with M satisfying Eq. (6), \mathcal{A} becomes a secure (ϵ, δ) -PAC learner according to Definition 1, while \mathcal{E} cannot. However, *at least in theory*, it is not impossible for \mathcal{E} to obtain the samples with a size identical to \mathcal{A} 's after the completion of \mathcal{A} 's learning. Nevertheless, \mathcal{E} cannot be a (ϵ, δ) -PAC learner at the same level as \mathcal{A} since $\eta^{\mathcal{E}}$ cannot be smaller than $\eta^{\mathcal{A}} + \Delta$. The condition $\eta^{\mathcal{A}} \geq \eta_c - \Delta$ in (ii-b) will also halt the protocol because of rule **(R.1)**. Thus, our results (i.e., Theorem 1 and 2 and Definition 1) can be practically applied to the protocol against any \mathcal{E} . Further, by using Γ and Δ , we can set the minimum gap between the level of \mathcal{A} 's and \mathcal{E} 's PAC learning in the worst case, and it would prevent \mathcal{E} from becoming a slightly weaker PAC learner than \mathcal{A} . The subcases $\eta_c - \Delta \geq \eta^{\mathcal{A}} \geq \eta^{\mathcal{E}}$ and $\eta_c - \Delta \geq \eta^{\mathcal{E}} > \eta^{\mathcal{A}}$ are not expected to occur since they contradict Eq. (4).

VI. MULTI-CLASS CLASSIFICATION

We also consider the multi-class problem by assuming that the input \mathbf{x} belongs to 2^m different classes ($m \geq 2$). Here, we briefly sketch two strategies:

(i) First, the multi-class classification problem is commonly solved by decomposing it into several binary problems. For instance, the “one-vs-all (OVA)” reduction is often used [19], where the problem is decomposed into 2^m decisions of h_i ($i \in \{0, 1, \dots, 2^m - 1\}$) that separates the learning data of the i th class from the other ones. Then, datum \mathbf{x} is classified with $\arg \max_i h_i(\mathbf{x})$. Here, the condition for secure PAC learning in Eq. (6) can be applied to each decision of h_i . However, a long learning time is required because the condition in Eq. (6) should be satisfied for every 2^m decisions.

(ii) In another way, we can consider a single-machine approach, where the oracle can answer for all 2^m labels, that is, $\mathbf{y} \in \{0, 1\}^m$, by allowing m qubits conditioned by the same \mathbf{x} -input channels. In such generalization, our theorems and the condition in Eq. (6) can also be applied consistently for the states of an arbitrary number of qubits. However, in this case, the region that satisfies the secure PAC learning, i.e., $|M_c(\epsilon, \delta) - M_b(\epsilon, \delta)|$, narrows. In other words, the security condition becomes more stringent. For detailed analysis, see Appendix C.

VII. REMARKS

We have presented a concept of secure learning that safeguards against any malicious manipulation of learning samples. In contrast to other studies on secure learning, we constructed an analytic framework based on a computational model of learning theory, called PAC learning. This allowed us to establish the link between sample complexity and the condition for learning security. Our approach is appealing because the security condition is defined solely by the sample size; in particular, it is independent of \mathcal{A} 's (or \mathcal{E} 's) learning algorithms.

Our derivations of Theorem 1 and 2 were based on the quantum principle of no-broadcasting of states, and using these theorems, we introduced the concept of secure PAC learning. Such a security condition cannot exist in the classical regime where \mathcal{E} can create as many copies of the learning samples as he or she wishes.

It is noteworthy that our protocol was designed based on a classical-quantum hybridization, where the input data remain classical but only a single-qubit system is employed. Such a hybridization differs considerably from those of other hybrid models. This architecture renders our protocol suitable for NISQ implementation, without the requirement of an excessively large superposition of samples and/or without accessing a novel quantum gadget, called quantum random-access memory [30, 31].

We finally point out that determining a more practical form of $M_c(\epsilon, \delta)$ in Eq. (6) continues to be an open problem, and it will be considered in a follow-up study. Notably, it is related to the determination of the optimal sample complexity, which has been a long-standing interest in computational learning theory, especially in the case where the samples are noisy. We believe that our paper will contribute to expanding the frontiers for quantum secure machine learning.

ACKNOWLEDGEMENTS

W.S., Y.L. and J.B. are grateful to Nana Liu for valuable discussions. This work was partly supported by National Research Foundation of Korea (NRF) grants (2019R1A2C2005504, NRF-2019M3E4A1079666, and 2020M3E4A1079939), funded by the MSIP (Ministry of Science, ICT and Future Planning) and Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. 2020-0-00890, "Development of trusted node core and interfaces for the interoperability among QKD protocols"). W.S., Y.L., and J.B. acknowledge the research project on developing quantum machine learning and quantum algorithm (No. 2019-100) by the ETRI affiliated research institute. Y.L. acknowledges National Research Foundation of Korea a grant funded by the Ministry of Science and ICT (NRF-2020M3E4A1077861) and KIAS Individual Grant (CG073301) at Korea Institute for Advanced Study. H.

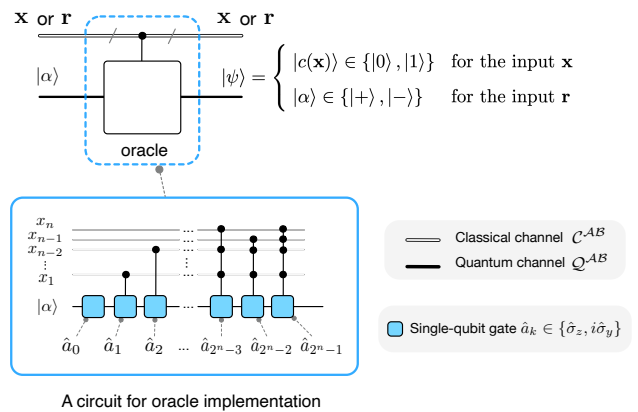


FIG. 3. **Schematic of a hybrid oracle.** The oracle consists of two different input and output channels: classical input data $\mathbf{x} = x_1 x_2 \cdots x_n$ ($x_j \in \{0, 1\} \forall j = 1, \dots, n$) and a single qubit to produce the oracle output states. This oracle applies 2^n unitary gates $\hat{a}_k \in \{\hat{\sigma}_z, i\hat{\sigma}_y\}$ ($k = 0, 1, \dots, 2^n - 1$) conditioned on the values of the classical bits x_j in \mathbf{x} to the qubit channel. In a purely classical case, these gates are either identity or logical-NOT gates.

K. is supported by the KIST Open Research Program. M.P. and M.W. acknowledge the ICTQT IRAP project of FNP, financed by structural funds of EU. G.A. acknowledges financial support from the European Research Council (ERC) under the Starting Grant GQCOP (Grant No. 637352) and the Foundational Questions Institute (FQXi) under the Intelligence in the Physical World Programme (Grant No. FQXiRFP-IPW-1907). M.W. was supported by NCN grants 2015/19/B/ST2/01999 and 2017/26/E/ST2/01008. M.P. was supported under FNP grant First Team/2016-1/5. J.K. was supported in part by KIAS Advanced Research Program (No. CG014604). J.B. was also supported by a KIAS Individual Grant (No. CG061003).

Appendix A: Useful Classical-Quantum Hybrid Oracle Architecture

Here, we present an example of a classical-quantum hybrid oracle, which can be applied to our study of secure learning. This oracle allows the classical inputs \mathbf{x} and a single qubit $|\alpha\rangle$. It performs the mapping

$$(\mathbf{x}, |\alpha\rangle) \rightarrow (\mathbf{x}, |c(\mathbf{x}) \oplus \alpha\rangle) \text{ for } \alpha \in \{0, 1\}, \quad (\text{A1})$$

and

$$(\mathbf{r}, |\alpha\rangle) \rightarrow (\mathbf{r}, |\alpha\rangle) \text{ for } \alpha \in \{+, -\}, \quad (\text{A2})$$

where \mathbf{r} is a random datum that is to be used for performing a security check. Note that \mathbf{x} remains unaltered during and after the sampling process.

This hybrid oracle can be implemented by a circuit having a specific architecture, such as that shown in

Fig. 3. This circuit contains 2^n gates acting on the ancilla qubit: the single-qubit gate \hat{a}_0 and $2^n - 1$ gates \hat{a}_k ($k = 1, 2, \dots, 2^n - 1$) conditioned on the classical-bit values x_1, x_2, \dots, x_n in \mathbf{x} . The gates \hat{a}_k are given by

$$\hat{a}_k \in \{\hat{\sigma}_z, i\hat{\sigma}_y\}, \text{ for all } k = 0, 1, \dots, 2^n - 1, \quad (\text{A3})$$

where $\hat{\sigma}_x$, $\hat{\sigma}_y$, and $\hat{\sigma}_z$ are the Pauli operators. This architecture of the oracle is inspired by the general expression of a Boolean function [32], which is given by

$$h^*(\mathbf{x}) = a_0 \oplus a_1 x_1 \oplus a_2 x_2 \oplus a_3 x_1 x_2 \oplus \dots \oplus a_{2^n - 1} x_1 x_2 \dots x_n, \quad (\text{A4})$$

where $a_k \in \{0, 1\}$ ($k = 0, 1, \dots, 2^n - 1$) are known as the Reed–Muller coefficients. Here, each coefficient has a corresponding gate operation \hat{a}_k . More specifically, $a_k = 0$ implies that \hat{a}_k leaves the bit signal unchanged (identity), while $a_k = 1$ indicates that \hat{a}_k flips the bit signal (logical-not) [33]. The oracle is thus characterized by a fixed set of gates \hat{a}_k for a given c . Information on the gates \hat{a}_k and how they run is not provided, and it should be learned. Such an oracle architecture indeed differs from other hybrid schemes. It has been argued that such hybridization can offer the advantage of being NISQ implementable and of achieving speedups [16, 17].

Appendix B: Probably-Approximately-Correct (PAC) Learning Model

In the PAC learning model [27], a learner samples a finite set of training data $\{(\mathbf{x}_i, c(\mathbf{x}_i))\}$ ($i = 1, 2, \dots, M$) by accessing an oracle. Here, \mathbf{x}_i is typically assumed to be drawn uniformly. For any $c \in \mathcal{C}$, a learning algorithm is a (ϵ, δ) -PAC learner (under uniform distribution) if it can obtain an ϵ -approximated correct $h \in \mathcal{H}$ with probability $1 - \delta$. More specifically, a learning algorithm is a (ϵ, δ) -PAC learner if it satisfies the condition

$$\text{Prob}(E(h, c) \leq \epsilon) \geq 1 - \delta, \quad (\text{B1})$$

where $E(h, c)$ denotes the error, for example, the distance between h and c . If the obtained h agrees with

$$M \geq \frac{1}{\epsilon} \ln \frac{|\mathcal{H}|}{\delta} \quad (\text{B2})$$

of samples constructed from the oracle, then Eq. (B1) holds. Here, $|\mathcal{H}|$ denotes the cardinality of \mathcal{H} , often-called model complexity. In the standard context, Eq. (B2) is known as the ‘‘sample complexity’’ [21, 27]. In other words, it yields the minimum number of training samples required to successfully learn $h \in \mathcal{H}$, satisfying Eq. (B1). Such a sample complexity derived from previous *classical* studies can be directly used in our scenario. In our classical-quantum hybrid query scheme, the same sample complexity exists since \mathbf{x}_i and $c(\mathbf{x}_i)$ identified by the measurement performed by Alice are classical. The beauty of this theorem is that the condition for being a PAC

learner depends only on the number of samples, not on any specific learning algorithm.

In the case where the oracle outputs are contaminated, the sample complexity in Eq. (B2) is modified as follows: First, we draw a sequence of training data,

$$\{(\mathbf{x}_1, m_1), (\mathbf{x}_2, m_2), \dots, (\mathbf{x}_M, m_M)\}, \quad (\text{B3})$$

where $m_i \in \{c(\mathbf{x}_i), c(\mathbf{x}_i) \oplus 1\}$ denotes the outcome of the measurement performed by Alice. Subsequently, if sampling is performed with

$$M \geq \frac{2\xi}{\epsilon^2} \ln \left(\frac{2|\mathcal{H}|}{\delta} \right), \quad (\text{B4})$$

we can verify that Eq. (B1) holds for the algorithm that obtains $h \in \mathcal{H}$. It has been proven that the additional factor ξ is given by [28]

$$\xi = \frac{1}{(1 - 2\eta)^2}. \quad (\text{B5})$$

Such a *noisy* PAC learning model provides a useful framework for our study of secure learning. It is noteworthy that in our scenario, the contamination of the output because of an attack by Eve and that resulting from imperfections intrinsic to the oracle can be incorporated together into the factor η .

Appendix C: Extension to multi-class classification

Each training datum can be considered to belong to one of 2^m different classes ($m \geq 2$), and the goal is to learn a hypothesis that, given a (new) data point, can correctly decide the class to which the data point belongs. This problem is called the multi-class classification problem.

1. One-vs-All (OVA) reduction

The conventional approach used to solve the multi-class classification problem is to decompose the problem into several binary classification problems. The most simple, but powerful, method is the so-called OVA reduction [19], where each binary classifier (e.g., RLSC, SVM) is trained to distinguish the examples in a single class from those in all remaining classes. More specifically, in such strategy, the problem is decomposed to 2^m decisions of h_i , ($i \in \{0, 1, \dots, 2^m - 1\}$) that separates the training data of the i th class from those of the other classes (see Fig. 4), and (new) data are classified using

$$h(\mathbf{x}) = \arg \max_i h_i(\mathbf{x}), \quad (\text{C1})$$

where $h_i(\mathbf{x})$ is a hypothesis identified in each trial and $h(\mathbf{x})$ is a decision for the classification of the input \mathbf{x} . Here, $h_i(\mathbf{x})$ is interpreted as the probability of a given input being included in the i th class, which is very suitable for our PAC learning framework. To achieve OVA

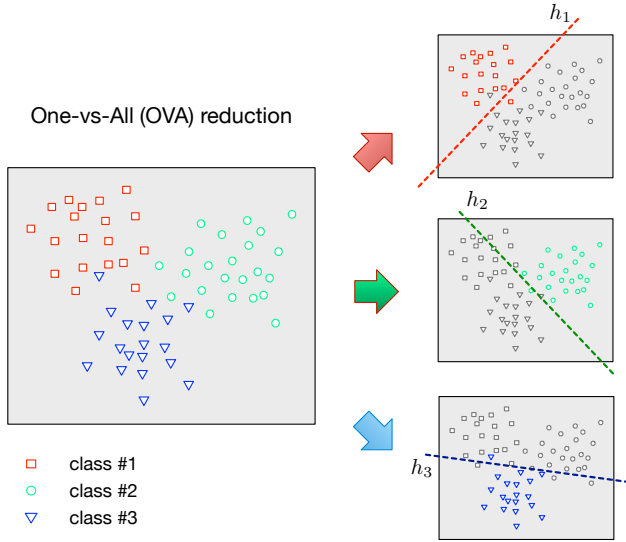


FIG. 4. Schematic of OVA reduction for three classes.

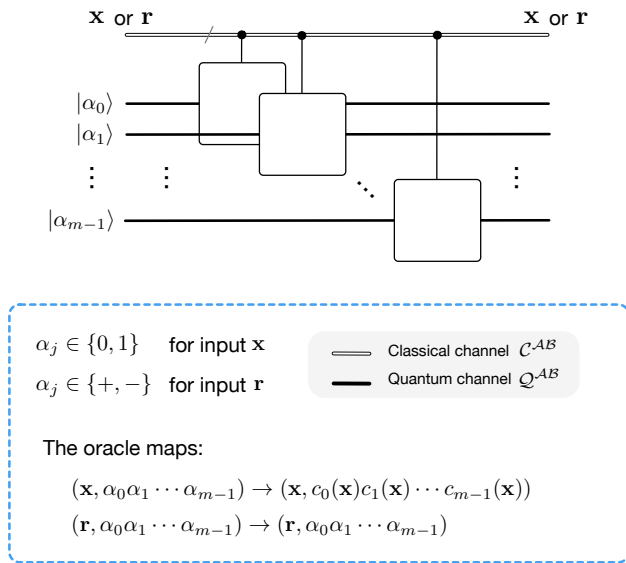


FIG. 5. Schematic of the oracle for a N-class output.

reduction, we can apply the condition for secure PAC learning (Eq. (7) of our main paper), as it is, to each trial performed for identifying $h_i(\mathbf{x})$. However, in this case, the learning time is increased as we should prepare

the dataset to train 2^m classifiers and the secure PAC learning condition should be satisfied for every 2^m trials.

2. A strategy of single-machine approach

Another useful approach is to solve a single optimization problem that trains many binary classifiers *simultaneously*; this approach is akin to the so-called “single machine approach” [19]. To apply this approach, we should consider an oracle that, given an input $\mathbf{x} \in \{0, 1\}^n$, outputs the corresponding label $\mathbf{y} \in \{0, 1\}^m$ for all 2^m classes, for example, by employing an arbitrary function $h : \{0, 1\}^n \rightarrow \{0, 1\}^m$. This is possible by allowing m qubits conditioned by the same \mathbf{x} -input channels (see Fig. 5). More specifically, in this generalization, the oracle performs the following mapping

$$(\mathbf{x}, \alpha_0 \alpha_1 \cdots \alpha_{m-1}) \rightarrow (\mathbf{x}, c_0(\mathbf{x}) c_1(\mathbf{x}) \cdots c_{m-1}(\mathbf{x})) \quad (\text{C2})$$

for the learning (i.e., for $\alpha_0 \alpha_1 \cdots \alpha_{m-1} \in \{0, 1\}^m$) and the mapping

$$(\mathbf{r}, \alpha_0 \alpha_1 \cdots \alpha_{m-1}) \rightarrow (\mathbf{r}, \alpha_0 \alpha_1 \cdots \alpha_{m-1}) \quad (\text{C3})$$

for the security check (i.e., for $\alpha_0 \alpha_1 \cdots \alpha_{m-1} \in \{+, -\}^m$). The learner (Alice, here) can identify the oracle’s output by measuring each returning qubit and construct the training samples for the learning. In this strategy, our theorems and the secure PAC learning condition can be applied to the states of an arbitrary number of qubits. Note that in our analysis, the states $\hat{\rho}_s$ and $\hat{\rho}_\mathcal{E}^{(k)}$ comprise an arbitrary number of qubits. The rules [R.1] and [R.2] derived for practical use of our protocol are applicable to each qubit measurement outcome. However, in this case, $M_b(\epsilon, \delta)$ is expected to increase as a higher model complexity, $|\mathcal{H}|$, would be imposed for large m . Furthermore, $M_c(\epsilon, \delta)$ decreases since η_c increases for large m ; specifically, we have [34]

$$\eta_c = 1 - \max F(\hat{\rho}_0(\mathbf{x})^{\otimes m}, \hat{\rho}_\mathcal{E}^{\otimes m}) = \frac{1}{(2m+4)}. \quad (\text{C4})$$

Consequently, the region $|M_c(\epsilon, \delta) - M_b(\epsilon, \delta)|$ that satisfies the secure PAC learning narrows as m increases; in other words, the security condition becomes more stringent. Therefore, there exists a *trade-off* between the two aforementioned approaches. Note that $|M_c(\epsilon, \delta) - M_b(\epsilon, \delta)| \geq 0$ is always satisfied along with the no-broadcasting theorem with the condition $(\eta_c \geq \eta^{\mathcal{A}}) \wedge (\eta_c \geq \eta^{\mathcal{E}})$ in Eq. (4) of our main paper.

[1] A. W. Harrow, A. Hassidim, and S. Lloyd, Physical review letters **103**, 150502 (2009).

[2] P. Reberstrost, M. Mohseni, and S. Lloyd, Physical review letters **113**, 130503 (2014).

[3] S. Lloyd, M. Mohseni, and P. Reberstrost, Nature Physics **10**, 631 (2014).

[4] M. Schuld, I. Sinayskiy, and F. Petruccione, Physical Review A **94**, 022342 (2016).

- [5] I. Kerenidis and A. Prakash, arXiv preprint arXiv:1603.08675 (2016).
- [6] J. Biamonte, P. Wittek, N. Pancotti, P. Rebentrost, N. Wiebe, and S. Lloyd, *Nature* **549**, 195 (2017).
- [7] M. Schuld, I. Sinayskiy, and F. Petruccione, *Contemporary Physics* **56**, 172 (2015).
- [8] M. Barreno, B. Nelson, A. Joseph, and J. Tygar, *Machine Learning* **81**, 121 (2010).
- [9] B. Nelson, B. I. P. Rubinstein, L. Huang, A. D. Joseph, S. J. Lee, S. Rao, and J. D. Tygar, *J. Mach. Learn. Res.* **13**, 1293 (2012).
- [10] J. Bang, S.-W. Lee, and H. Jeong, *Quantum Information Processing* **14**, 3933 (2015).
- [11] Y.-B. Sheng and L. Zhou, *Science Bulletin* **62**, 1025 (2017).
- [12] N. Liu and P. Rebentrost, *Physical Review A* **97**, 042315 (2018).
- [13] A. W. Harrow, arXiv preprint arXiv:2004.00026 (2020).
- [14] H. Barnum, C. M. Caves, C. A. Fuchs, R. Jozsa, and B. Schumacher, *Physical Review Letters* **76**, 2818 (1996).
- [15] G.-F. Dang and H. Fan, *Physical Review A* **76**, 022323 (2007).
- [16] J.-S. Lee, J. Bang, S. Hong, C. Lee, K. H. Seol, J. Lee, and K.-G. Lee, *Physical Review A* **99**, 012313 (2019).
- [17] W. Song, M. Wieśniak, N. Liu, M. Pawłowski, J. Lee, J. Kim, and J. Bang, arXiv preprint arXiv:1905.05751 (2019).
- [18] J. Preskill, *Quantum* **2**, 79 (2018).
- [19] R. Rifkin and A. Klautau, *Journal of machine learning research* **5**, 101 (2004).
- [20] M. J. Kearns, U. V. Vazirani, and U. Vazirani, *An introduction to computational learning theory* (MIT Press, Cambridge, MA, 1994).
- [21] P. Langley, *Elements of machine learning* (Morgan Kaufmann, San Francisco, CA, 1995).
- [22] V. Scarani, H. Bechmann-Pasquinucci, N. J. Cerf, M. Dušek, N. Lütkenhaus, and M. Peev, *Reviews of modern physics* **81**, 1301 (2009).
- [23] J.-Y. Hu, B. Yu, M.-Y. Jing, L.-T. Xiao, S.-T. Jia, G.-Q. Qin, and G.-L. Long, *Light: Science & Applications* **5**, e16144 (2016).
- [24] A. Chefles and S. M. Barnett, *Journal of Physics A: Mathematical and General* **31**, 10097 (1998).
- [25] H. Barnum, J. Barrett, M. Leifer, and A. Wilce, *Physical review letters* **99**, 240501 (2007).
- [26] R. Jozsa, *Journal of modern optics* **41**, 2315 (1994).
- [27] L. G. Valiant, *Communications of the ACM* **27**, 1134 (1984).
- [28] D. Angluin and D. K. Slonim, *Machine Learning* **14**, 7 (1994).
- [29] C. Ciliberto, M. Herbster, A. D. Ialongo, M. Pontil, A. Rocchetto, S. Severini, and L. Wossnig, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* **474**, 20170551 (2018).
- [30] V. Giovannetti, S. Lloyd, and L. Maccone, *Phys. Rev. Lett.* **100**, 160501 (2008).
- [31] V. Giovannetti, S. Lloyd, and L. Maccone, *Phys. Rev. A* **78**, 052310 (2008).
- [32] P. Gupta, A. Agrawal, and N. K. Jha, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* **25**, 2317 (2006).
- [33] T. Toffoli, in *International Colloquium on Automata, Languages, and Programming* (Springer, New York, 1980), p. 632.
- [34] L. Chen and Y.-X. Chen, *Physical Review A* **75**, 062322 (2007).
- [35] S. Hanneke, *The Journal of Machine Learning Research* **17**, 1319 (2016).