# Towards efficient and generic entanglement detection

Jue Xu[*] and Qi Zhao[†]

(Dated: October 17, 2022)

Detection of entanglement is an indispensable step to practical quantum computation and communication. In this work, we propose an end-to-end, machine learning assisted entanglement detection protocol. In this protocol, an entanglment witness for a generic entangled state is obtained by classical machine learning with a synthetic dataset which consists of classical features of states and their labels. In actual experiments, classical features of a state, that is expectation values of a set of Pauli observables, are estimated by sample-efficient methods such as classical shadow.

## I. INTRODUCTION

Entanglement [1] is the key ingredient of quantum computation [2], quantum communication, and quantum cryptography [3]. However, decoherence is inevitable in real-world, which means the interaction between a quantum system and classical environment would significantly affect entanglement quality and diminish quantum advantage. So, for practical purpose, it is essential to benchmark (characterize) entanglement structures of certain target states in actual (real) experiments. The goal of this paper is to find an efficient and generic way to achieve it. Machine learning (ML) is a powerful tool for such purpose. Many ML techniques including quantum machine learning models [4] have been proposed for classification tasks in physics, such as classification of phases and prediction of ground states [5] [6].

Background entanglement, many-body, structure, Existing method, conventional suffer from some problems, ML. address their problems our work address, performance: one or two sentence , summary consider different types of physical noises, ML, witness, + entanglement structure (certain partition, depth, separability (intactness))

Assume we would like to distinguish an entangled state incluing its 'vicinity' (proximity) from undesired states (e.g., all separable states), our method derive such a classifier by fitting a synthetic dataset randomly sampled states with their labels (entangled or not). Specifically, our pipeline starts from evaluation of expectations of $n$-qubit Pauli observables of a target state. The set of expectation values that serves as classical features of the target state, together with its label, consist of a data point of a dataset. Then, a classical ML classifier is obtained by training with this dataset. With the trained classifier at hand, it is expected that brand new samples from real experiments can be classified with high accuracy, where classical features of quantum states are estimated by classical shadow method [7] with affordable samples complexity.

This paper is organized as follows: in Section II, we briefly present necessary definitions about entanglement structures and mainstream entanglement detection methods; Section III demonstrates our end-to-end protocol including two parts: learning an entanglement witness from synthetic data and estimating classical features of states from experiments; at last, numerical simulation results are discussed in Section IV.

## II. PRELIMINARIES

**Notation 1.** If no ambiguity, we omit the tensor products between subsystems and the hats on operators for readability, e.g., $|\psi_A\rangle|\psi_B\rangle \equiv |\psi_A\rangle \otimes |\psi_B\rangle$ and $XIZ \equiv \hat{X} \otimes \mathbb{1} \otimes \hat{Z}$.

**Notation 2.** Denote $O_\sigma \in \{I, X, Y, Z\}^{\otimes n}$ for a Pauli observable. Denote $\mathbf{x}_{\rho,\boldsymbol{\sigma}} := (\mathrm{Tr}(\rho O_{\sigma_1}), \ldots, \mathrm{Tr}(\rho O_{\sigma_M}))$ for expectations of $M$ Pauli observables with respect to the state $\rho$ where $\boldsymbol{\sigma} \subseteq \{I, X, Y, Z\}^n$. Denote vectors $\mathbf{x}$, $\boldsymbol{\sigma}$, $\mathbf{w}$ by boldface font.

### A. Entanglement structures

Large scale entanglement involving multiple particles maybe the main resource for quantum advantages in quantum computation and communication. Roughly, we say a quantum state is *entangled* if it is not fully separable, i.e., the state cannot be written as the tensor product of all subsystems. However, the simple statement 'the state is entangled' would allow that only two of the particles are entangled while the rest is in a product state. So, the more interesting entanglement property is bipartite separability. Consider a system partitioned into two subsystems $\mathcal{H}_A \otimes \mathcal{H}_B$, where each has dimension $d_A$ and $d_B$ respectively.

**Definition 1** (bi-separable)**.** A pure state $|\psi\rangle$ is bipartite (bi-)separable if it can be written as a tensor product form $|\psi_{\mathrm{bi}}\rangle = |\phi_A\rangle \otimes |\phi_B\rangle$. A mixed state $\rho$ is separable if and only if it can be written as a convex combination of pure bi-separable states, i.e., $\rho_{\mathrm{bi}} = \sum_i p_i |\psi_i\rangle\langle\psi_i|_{\mathrm{bi}}$ with probability distribution $\{p_i\}$. The set of all bi-separable states is denoted as $\mathcal{S}_{\mathrm{bi}}$.

On the contrary, if a state is not a convex combination of any (partition) biseparable states, it means that all

subsystems are indeed entangled with each other. This is the strongest form of entanglement, formally

**Definition 2** (GME)**.** If a state is not in $\mathcal{S}_{\mathrm{bi}}$, it possesses genuine multipartite entanglement.

There is another restricted way for generalizing bi-separability to mixed states: if it is a mixing of pure bi-separable states with the same partition $\mathcal{P}_2$, and we denote the state set as $\mathcal{S}_{\mathrm{bi}}^{\mathcal{P}_2}$. It is practically interesting to study entanglement structure under certain partition, because it naturally indicates the quantum information processing capabilities among a real geometric configuration. We have a definition concerning partitions

**Definition 3** (full entanglement)**.** A state $\rho$ possesses full entanglement if it is outside of the separable state set $\mathcal{S}_{\mathrm{bi}}^{\mathcal{P}_2}$ for any partition, that is, $\forall \mathcal{P}_2 = \left\{ A, \bar{A} \right\}, \rho \notin \mathcal{S}_{\mathrm{bi}}^{\mathcal{P}_2}$.

For a state with full entanglement, it is possible to prepare it by mixing bi-separable states with different bipartitions, so full entanglement is weaker than GME but still useful.

### B. Entanglement detection

After introducing the definitions about entanglement, the next basic question is how to determine entanglement and its computational complexity. Despite its clear definitions, entanglement detection for a general state is a highly non-trivial problem. For a general review on this subject, we refer readers to [8]. The most widely studied problem in this area maybe bi-separability.

**Problem 1** (separability)**.** Given a state $\rho$ in its density matrix repsentation, to determine if it is bi-separable.

#### 1. Hardness of separability

It is not hard to prove that if a state is bi-separable, then it must have positive partial transpose (PPT), that is, the partially transposed (PT) density matrix $\rho_{AB}^{\mathsf{T}_A}$ is PSD [9]. By contrapositive, we have a criterion for entanglement, that is

**Theorem 1** (PPT criterion)**.** *If the smallest eigenvalue of partial transpose $\rho_{AB}^{\mathsf{T}_A}$ is negative (NPT), then the state is entangled (cannot be bi-separable) with respect to the partition $\mathcal{P} = \{A, B\}$.*

We should mention that PPT criterion is a necessary and sufficient condition only for separability of low-dimensional systems (when $d_A d_B \leq 6$) [10]. Therefore, no general solution for the separability problem is known. Then, a natural question is whether it is possible to solve separability approximately. By relaxing the defintion (promise a gap), a reformulation of separability in the theoretic computer science language is

**Problem 2** (Weak membership problem for separability)**.** Given a density matrix $\rho$ with the promise that either (i) $\rho \in \mathcal{S}_{\mathrm{bi}}$ or (ii) $\|\rho - \rho_{\mathrm{bi}}\| \geq \epsilon$ with certain norm, decide which is the case.

Unfortunately, even we are given the complete information about a state and promised a gap, it is still hard to determine separability approximately by classical computation. Weak membership problem for separability is NP-Hard for $\epsilon = 1/\operatorname{poly}(D)$ with respect to Euclidean norm and trace norm [11]. [12] [13] while there exists a quasipolynomial-time algorithm with respect to $\|\cdot\|_{\mathrm{LOCC}}$ (and $\|\cdot\|_2$?) [14]. quantum hardness ... [15] A notable example is the widely-used and powerful criteria called $k$-symmetric extension hierarchy based on SDP [16], which is computationally intractable with growing $k$. Whereas, these hardness results does not rule out the possibility to solve it efficiently with stronger promise (approximation) or by machine learning (heuristic) techniques powered by data.

A related but different problem setting is how to determine bi-separable given copies of an unknown state (from experiments) rather than its density matrix. Since the input to this problem is quantum data (states), direct estimation of reduced density matrix's spectrum by quantum circuits is a good option (without fully recovering density matrix). For example, multivariate trace $\operatorname{Tr}(\rho_A^m)$ encodes the entanglement information (e.g., purity, negativity, and entanglement entropy) of $\rho_{AB}$ where $\rho_A$ is the reduced density matrix [17] [18] [19]. The multivariate trace can be estimated by constant depth quantum circuits [20] [21], but this line of work is still based on PPT criterion (not iff).

#### 2. Entanglement witness based on fidelity

The problem we study in this paper is another variant:

**Problem 3** (entanglement detection)**.** Given an unknown state $\rho$ (from experiments) is promised either (i) $\rho \in \mathcal{S}_{\mathrm{bi}}$ or (ii) in proximity of a target $|\psi_{\mathrm{tar}}\rangle$ (i.e., possesses 'useful' entanglement such as GME, full entanglement, depth ...) [22], determine which is the case.

The typical scenario for this problem is one aims to prepare a pure entangled state $|\psi_{\mathrm{tar}}\rangle$ in experiments and would like to detect (verify) it as true multipartite entangled. While the preparation is not perfect, it is reasonable to assume that the prepared mixed state $\rho_{\mathrm{pre}}$ is in the proximity of the target state, that is, $|\psi_{\mathrm{tar}}\rangle$ undergoes noise channels restricted to white noise, bit/phase-flip error, or random local unitary.

This problem can be expected to be solved more efficiently, because we have a much stronger promise than the separability problem. The usual method is constructing an observable $W$ called entanglement witness such that

$$\operatorname{Tr}(W \rho_{\mathrm{bi}}) \geq 0 \text{ and } \operatorname{Tr}(W |\psi_{\mathrm{tar}}\rangle\langle\psi_{\mathrm{tar}}|) < 0 \qquad (1)$$

Eq. (1) means that the witness $W$ has a positive expectation value on all separable states, hence a negative expectation value implys the presence of entanglement (GME). For every entangled state, a witness can always be constructed, but no entanglement witness works for all entangled states [23]. Bell (CHSH) inequalities that were originally proposed to rule out local hidden variable models, can be regarded as the oldest entanglement witness [24]. A Bell inequality is a linear combination of Pauli observables $W_{\mathrm{Bell}} := \mathbf{w}_{\mathrm{Bell}} \cdot \mathbf{O}_{\mathrm{Bell}}$ such that only entangled states $\rho$ have $|\mathrm{Tr}(\rho W_{\mathrm{Bell}})|$ greater than a threshold [25].

While various methods for constructing an entanglement witness exist, the most common one is based on the fidelity of a state to the target (pure entangled) state

$$W_\psi = \alpha \mathbb{1} - |\psi_{\mathrm{tar}}\rangle\langle\psi_{\mathrm{tar}}| \qquad (2)$$

where $\alpha$ is the smallest constant such that for every product state $\mathrm{Tr}(\rho W) \geq 0$. For instance, assume the target state is $|\mathrm{GHZ}\rangle$, the maximal overlap between GHZ and bi-separable states is $1/2$, such that the witness Eq. (2) with $\alpha = 1/2$ certifies tripartite entanglement [26]. This kind of fidelity witness is projector-based witness [27]. In order to effectly measure a witness in an experiment, it is preferable to decompose the projector term into a sum of locally measurable observables. For graph states (stabilizer states), a witness can be constructed by very few local measurement settings (tradeoff between robustness and meaurement efficiency) [28] [29] [30], while for non-stabilizer cases (e.g., W state), more careful analysis is required [31] [32]. Related experiments: photonic implementation with a few qubits (generation, verification) [33]; fully entangled graph state (ring of 16 qubits) IBM by measuring negativity [34]; optical lattice (homogeneous, restricted measurement, detect GME, nonstabilizer) [35];

## III. END-TO-END ENTANGLEMENT DETECTION PROTOCOL

### A. Motivation: Beyond fidelity witness

The most common robustness measure of fidelity witness is the tolerence of white noise

$$\rho' = (1 - p_{\mathrm{noise}}) |\psi_{\mathrm{tar}}\rangle\langle\psi_{\mathrm{tar}}| + p_{\mathrm{noise}}/2^n \mathbb{1} \qquad (3)$$

where the limit of (maximal) $p_{\mathrm{noise}}$ indicates the robustness of the witness. For example, the maximally-entangled Bell state can maximally violate the CHSH inequality, but Bell states mixed with white noise don't violate the CHSH inequality when $1 - 1/\sqrt{2} < p_{\mathrm{noise}} < 2/3$, despite they are still entangled in this regime.

For GHZ and W states mixed with white noise, we can analytically compute the white noise threshold for NPT (implies bipartite entanglement): when $p_{\mathrm{noise}} < 0.8$, GHZ states cannot be bi-separable with respect to

any partition (that is full entanglement). However, the conventional fidelity witness only detects GME when $p_{\mathrm{noise}} < 1/2 \cdot (1 - 1/2^n)^{-1} \approx 1/2$ (cf. Table II). So, it would be practically interesting to have a witness for this white noise regime.

Other than white noise, more realistic noise happened in (photonic) experiments is coherent noise, e.g., local rotations. Take GHZ state as an example, unconscious phase accumulation and rotation on the first control qubit can be modeled as [36]

$$|\mathrm{GHZ}(\phi,\theta)\rangle = \cos\theta\,|0\rangle^{\otimes n} + e^{\mathrm{i}\phi}\sin\theta\,|1\rangle^{\otimes n}. \qquad (4)$$

In certain noise regime (see Fig. 3 in [36]), $|\mathrm{GHZ}(\phi,\theta)\rangle$ cannot be detected by conventional fidelity witness because coherent noises diminish the fidelity but not change entanglement property.
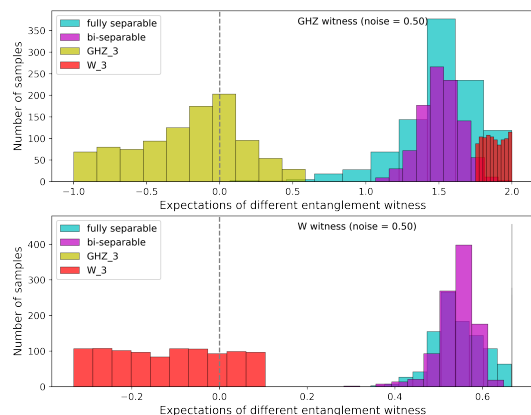


FIG. 1: Entanglement cannot be detected by fidelity witness (GHZ state with coherent noise, W state with large white noise)

To formally characterize the cases beyond fidelity witness, Weilenmann et. al [37] coined the term *unfaithful states* which systematically analyze 2-qudit entangled state mixed with white noise that cannot be detected by fidelity witness. They found that for $d \geq 3$ that almost all states in the Hilbert space are unfaithful. Subsequently, Güthe et. al [38] [39] give a formal definition:

**Definition 4** (unfaithful state)**.** A 2-qudit state $\rho_{AB}$ is faithful if and only if there are local unitary transformations $U_A$ and $U_B$ such that

$$\langle\phi^+|U_A \otimes U_B \rho_{AB} U_A^\dagger \otimes U_B^\dagger|\phi^+\rangle > \frac{1}{d}. \qquad (5)$$

Consequently, they found a necessary and sufficient condition for 2-qubit unfaithfulness, determined by the spectrum of

$$\mathcal{X}_2(\rho_{AB}) = \rho_{AB} - \frac{1}{2}(\rho_A \otimes \mathbb{1} + \mathbb{1} \otimes \rho_B) + \frac{1}{2}\mathbb{1} \otimes \mathbb{1}, \qquad (6)$$

i.e., a 2-qubit state $\rho_{AB}$ is faithful if and only if the maximal eigenvalue of $\mathcal{X}_2(\rho_{AB})$ is larger than $1/2$. We can see in Fig. 2, even for 2-qubit systems, nonnegligible portion of states are unfaithful but still entangled (NPT).
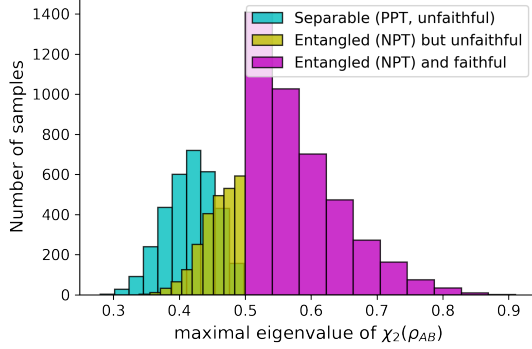


FIG. 2: Unfaithfulness of 2-qubit states determined by the maximal eigenvalue of Eq. (6).

Although there are variants of witness [36], such as nonlinear witness [40] and post-processing [41], designed to remedy the shortcomings of conventional fidelity witness respectively, it would be meaningful in practice to find a generic method to construct witnesses for entanglement detection. Machine learning techniques satisfy the needs well because supervised learning can be regarded as a powerful nonlinear post-processing tool.

### B. Training a generic witness via SVM

The introductory task of classical machine learning (ML) is the binary classification, such as cat/dog images classification. In (supervised) learning, the input to a ML algorithm is a (training) dataset $\left\{\mathbf{x}^{(i)}, y^{(i)}\right\}_{i=1}^{m}$ consists of $m$ data points, where each is a pair of feature vector $\mathbf{x}$ and its label $y$ (either 0 or 1). For example, the feature $\mathbf{x}$ of an image is a flatten vector of all pixel values and the label $y = 0(1)$ for cat (dog) respectively. It is clear to see separability or entanglement detection are exactly such binary classification problems where each quantum state has a binary label, such as entangled/bi-separable. The features $\mathbf{x}$ of a quantum state $\rho$ can be the entries of its density matrix, or more realistically, the expectation values of certain observables.

With the surge of ML research, ML algorithms have been proposed for classification tasks related to entanglement. Lu et. al [42] trained a (universal) separability classifier by classical neural network where features are the entries of density matrices. For the similar purpose, Ma and Yung [43] generalized Bell inequalities to a Bell-like ansatz $W_{\mathrm{ml}} := \mathbf{w}_{\mathrm{ml}} \cdot \mathbf{O}_{\mathrm{Bell}}$ where the optimal weights $\mathbf{w}_{\mathrm{ml}}$ are obtained via a neural network. And they found the tomographic ansatz

$$W_{\mathrm{ml}} := \mathbf{w}_{\mathrm{ml}} \cdot \mathbf{O}_{\sigma} \,, \ \forall \sigma \in \{I, X, Y, Z\}^n \qquad (7)$$

has better performance, also required [44] for a universal separability classifier. It is worth noting that training a universal classifier for high-dimensional systems is hard if the gap between two state sets is small (weak promise).
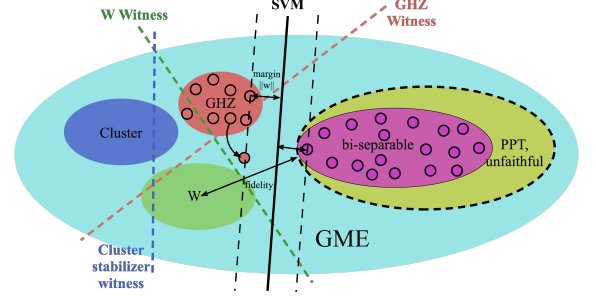


FIG. 3: Schematic diagram for entanglement detection methods: entanglement witnesses for different states are depicted by colored dash lines (hyperplane). SVM with the linear kernel (hyperplane). PPT criterion (non-linear, one-side) ...

In our paper, we focus on the problem entanglement detection with training data. In other words, we derive the entanglement witness for certain target states with desired entanglement structure by fitting a synthetic dataset.

**Problem 4** (Learn an entanglement witness). learn a witness for entanglement detection of $\rho$ from data

- **Input**: synthetic data consist of density matrices $\rho$ with corresponding labels $y$

- **Output**: a (minimal) classifier $\mathbf{w}_{\mathrm{ml}}$ with high accuracy

Learn an entanglement witness problem has also been studied by classical ML [32] [45], but by a different technique called Support Vector Machine (SVM) [46]. The features $\mathbf{x} := \mathrm{Tr}(\rho O_{\sigma})$ of a state is a vector of expectations of Pauli observables. An SVM constructs a hyperplane $\langle W_{\mathrm{svm}} \rangle = \mathbf{w}_{\mathrm{svm}} \cdot \mathbf{x}$ that clearly delineates two kinds of states, see Fig. 3. Both SVM witness and conventional fidelity witness is a weighted sum of observables (features) which represents a hyperplane in the state (feature) space. The SVM witness is more flexible because coefficients are optimized (automatically derived) from training. This method only requires local (Pauli) measurements even when the target state is a non-stabilizer state, such as W state (normally need nonlocal measurements).

---

**Algorithm III.1:** train a witness via kernel SVM

---

**input** : states with labels: $\left\{\left(\rho^{(i)}, y^{(i)}\right)\right\}_{i=1}^{m}$
**output:** a classifier $\mathbf{w}_{\mathrm{ml}}$, $\boldsymbol{\sigma}_{\mathrm{ml}}$

**1** Evaluate Pauli observables $\mathbf{x}_{\rho,\boldsymbol{\sigma}}^{(i)} := \mathrm{Tr}\left(\rho^{(i)} O_{\boldsymbol{\sigma}}\right)$, $\forall i$
**2** **for** $j = 1, 2, \ldots, 4^n$ **do**
**3**     **while** *accuracy not high enough* **do**
**4**        randomly select $j$ features $\tilde{\mathbf{x}}_i$ from $\mathbf{x}_i$, $\forall i$
**5**        accuracy, classifier = SVM($\left\{\left(\tilde{\mathbf{x}}^{(i)}, y^{(i)}\right)\right\}_{i}^{m}$)
**6**     **return** classifier $\mathbf{w}_{\mathrm{ml}}$

---

A key drawback (constrain) of conventional witnesses is its linearity. Despite the non-linear witness [40] proposed, its implementation in experiments is more challenging than linear one. The good news is, within the framework of SVM, non-linearity can be easily achieved by the kernel method. We focus on kernel methods rather than neural networks (also non-linear), not only because of its clear geometric interpretation, but also its equivalence to neural network in terms of neural tangent kernel [47]. The advantages of SVM: (1) the training of an SVM is convex; if a solution exists for the given target state and ansatz, the optimal SVM will be found. (2) this SVM formalism allows for the programmatic removal of features, i.e., reducing the number of experimental measurements and copies (samples).

| | # observables | weights | promise |
|---|---|---|---|
| fidelity witness | few local | fixed | strongest |
| Bell (CHSH) inequality | constant | fixed | weak |
| tomographic classifier | $4^n - 1$ | trained | weakest |
| SVM (kernel) witness | $\ll 4^n - 1$ | trained | strong |

TABLE I: Comparison of CHSH inequality, fidelity witness, and ML witness ansatz.

However, these previous (prior) ML witnesses only consider the robustness to white noise and cannot be directly applied to experiments. In numerical simulation, we can efficiently evaluate classical features by direct calculation, but, in actual experiments, entries of a density matrix are not explicitly known. Instead, we need to estimate the observables by many measurements, which we will discuss in next section.

### C. Sample-efficient expectation estimation methods

The brute force approach to fully characterize a state in an experiment is quantum state tomography [48] [49]. With a recovered densitry matrix, we can directly calculate classcial features or separability measures, but full tomography is experimentally and computationally demanding. Even if adaptive or collective measurements (and post-processing) allowed [50], rigorous analysis [51] [52] proved that $\Omega(D^2/\epsilon^2)$ measurements (copies?) are required for recovering a $D \times D$ density matrix with error tolerence $\epsilon$ measured in trace distance.

Now that full tomography is not practical for large systems, a natural question is whether it is possible to extract a bunch of information about a state without fully recovering it. The answer is yes. Many interesting properties of a quantum system are often linear functions of the underlying density matrix $\rho$, such as classical features $x_{\rho,\sigma} = \mathrm{Tr}(\rho O_\sigma)$ for entanglement witness [53]. This enables the possibility to *shadow tomography* [54].

**Problem 5** (shadow tomography). Aaronson's formulation ($\mathbb{P}[E_i \text{ accept } \rho] =? \mathrm{Tr}(E_i\rho)$)

- **Input**: copies of an unknown $D$-dimensional state $\rho$ and $M$ known 2-outcome measurements $\{E_1, \ldots, E_M\}$
- **Output**: estimate $\mathbb{P}[E_i \text{ accept } \rho]$ to within additive error $\epsilon$, $\forall i \in [M]$, with $\geq 2/3$ success probability.

Though it is proved that shadow tomography can be implemented in a samples-efficient (copies) manner, $\tilde{\mathcal{O}}(\log^4 M \cdot \log D \cdot \epsilon^{-4})$ copies [54] [55], Aaronson's shadow tomography procedure is very demanding in terms of quantum hardware. So, Huang et. al [7] introduce classical shadow (CS) method that is more friendly to experiments. In our pipeline, we focus on the classical shadow method and its variants.

A classical shadow is a succinct classical description of a quantum state, which can be extracted by performing reasonably simple single-copy measurements on a reasonably small number of copies of the state. The classical shadow attempts to approximate this expectation value by an empirical average over $R$ independent samples, much like Monte Carlo sampling approximates an integral.

$$o_i = \mathrm{Tr}(O_i \rho_{\mathrm{cs}}) \text{ obeys } \mathbb{E}[o] = \mathrm{Tr}(O_i \rho) \tag{8}$$

classical shadows are based on random Clifford measurements and do not depend on the structure of the concrete witness in question. In contrast, direct estimation crucially depends on the concrete witness in question and may be considerably more difficult to implement.

---

**Algorithm III.2:** estimate features by CS

---

**input** : samples of $\rho$ and $O_{\boldsymbol{\sigma}_{\mathrm{ml}}}$
**output:** estimation of $\mathbf{x}_{\rho,\boldsymbol{\sigma}_{\mathrm{ml}}} := \mathrm{Tr}(\rho O_{\boldsymbol{\sigma}_{\mathrm{ml}}})$

**1** **for** $i = 1, 2, \ldots, R$ **do**
**2**    $\rho \mapsto U\rho U^\dagger$      // apply a random unitary
**3**    $|b\rangle \in \{0,1\}^n$      // measurement outcome
**4**    $\rho_{\mathrm{cs}} = \mathcal{M}^{-1}\left(U^\dagger |b\rangle\langle b| U\right)$   // $\mathcal{M}$ quantum channel
**5** $\mathrm{CS}(\rho, R) = \{\rho_{\mathrm{cs}_1}, \ldots, \rho_{\mathrm{cs}_R}\}$    // classical shadow
   // estimate features for SVM from classical shadow
**6** **return** $\mathbf{x}_{\rho,\boldsymbol{\sigma}_{\mathrm{ml}}} = \textsc{MedianOfMeans}(\mathrm{CS}(\rho, R))_{\boldsymbol{\sigma}_{\mathrm{ml}}}$

---

Given a quantum state $\rho$, a classical shadow is created by repeatedly performing a simple procedure: Apply a unitary transformation $\rho \mapsto U\rho U^\dagger$, and then measure all the qubits in the computational basis $|\mathbf{b}\rangle \in \{|0\rangle, |1\rangle\}^{\otimes n}$. Its classical shadow (snapshots) $\rho_{\mathrm{cs}}$ (a density matrix) can be reconstructed

$$\rho_{\mathrm{cs}} := \mathcal{M}^{-1}\left(U^\dagger |\mathbf{b}\rangle\langle\mathbf{b}| U\right) \tag{9}$$

where $\mathcal{M}$ is a quantum channel that depends on the ensemble of random unitary transformation... . The algorithm is summarized in Algorithm. III.2. The number of times this procedure is repeated is called the size of the classical shadow. Classical shadows with size of order $\log(M)$ suffice to predict $M$ target functions $\{O_1, \ldots, O_M\}$.

The classical shadow size required to accurately approximate all reduced $k$-body density matrices scales exponentially in subsystem size $k$ $\Omega(\log(M)3^k/\epsilon^2)$ [7], but is independent of the total number of qubits $n$. [56] The derandomized variant of classical shadow [57] is the refinement of the original randomized protocol, but not necessarily guarantees better performance for global observables (involving all subsystems). noise-resilient variant [58] ... experiments of classical shadow and related comparison [59]; detect entanglement by estimating $p_3$-PPT with classical shadow [60]. Similar to the PPT condition, the $p_3$-PPT condition (without full tomography) applies to mixed states and is completely independent of the state in question [60].

The task of estimating expectation value can also be achieved efficiently by machin learning with training data [61] [62] [63] [6] [64]. Huang et. al rigorously show that, for any quantum process $\mathcal{E}$, observables $O$, and distribution $\mathcal{D}$, and for any quantum ML model, one can always design a classical ML model achieving a similar average prediction error such that $N_C$ (number of experiments?) is larger than $N_Q$ by at worst a small polynomial factor. In contrast, for achieving accurate prediction on all inputs $\mathrm{Tr}(\rho O_\sigma), \forall \sigma \in \{I, X, Y, Z\}^n$, exponential quantum advantage is possible. [65] [66] [67]

## IV. NUMERICAL SIMULATION

### A. Dataset preparation and states generation

We generate quantum state samples, construct quantum circuits, and manipulate quantum objects numerically by QuTiP library [68] [69]. We generate multipartite entangled states (synthetic data) including: Bell states, 3-qubit GHZ and W states, 4-qubit graph (1D cluster) state, see Fig. 4 for examples. In contrast to entangled states, we generate random separable states for different number of qubits by tensoring random density matrices of subsystems. For example, 2-qubit: bipartite $\rho_A \otimes \rho_B$ where $\rho_A$ and $\rho_B$ are random density matrices (sampled by Haar measure); 3-qubit pure states: $\rho_A \otimes \rho_{BC}$, $\rho_C \otimes \rho_{AB}$, and $\rho_B \otimes \rho_{AC}$. For different noise channels: white noise according to Eq. (3), coherent noise according to Eq. (4).

### B. Classification accuracy and comparison

For the machine learning part, we make use of scikit-learning package [70] to train SVM with the radial basis
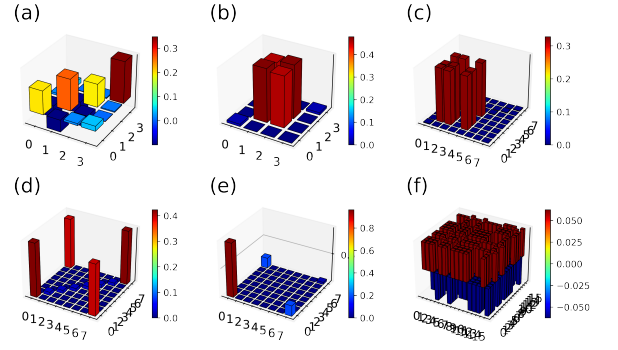


FIG. 4: Data praparation (real part): (a) random 2-qubit density matrix; (b) Bell state with white noise; (c) 3-qubit W state with white noise; (d) 3-qubit GHZ state with white noise; (e) GHZ state with coherent noise; (f) 4-qubit linear cluster.

function (RBF) kernel.

Fig. 5 shows the two-dimensional embedding of 2-qubit states (feature space). The colored shade indicates the decision boundary of our trained classifier (ML witness), which exhibits that two kinds of data points are clearly classified.
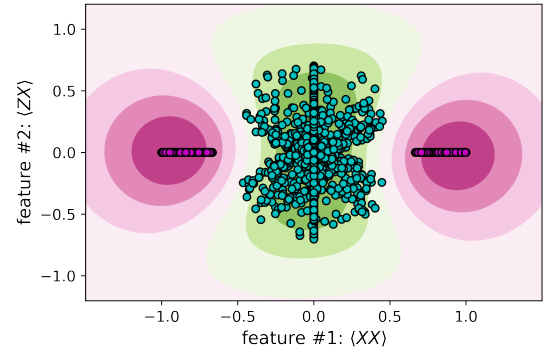


FIG. 5: two-dimensional embedding (feature space): green dots represent the separable states, while pink one represent entangled Bell states mixed with white noise.

Fig. 6 shows that the SVM witness can classify the states that cannot be detected by conventional fidelity witness, where the noise is randomly (uniform) sampled from $[0, p_{\mathrm{noise}}]$.

dataset size for training: $10^3$ for 3-qubit case; more qubits [TODO]...

## V. CONCLUSION AND DISCUSSION

Related experiments: .... Possible directions for future research: (1) rigorous proof for dataset size and number of features (required for high training accuracy) scaling with the system size; (2) better kernel options such as
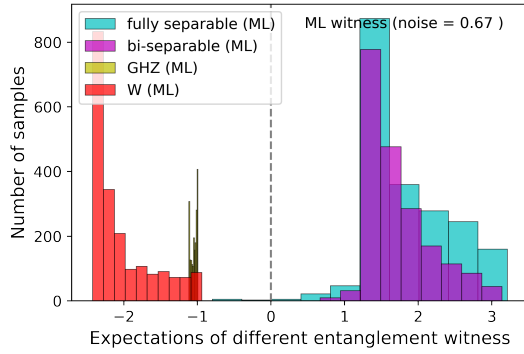
FIG. 6: ML witness for the states cannot be detect by fidelity witness (GHZ state with coherence noise, and W state with large white noise)

graph kernel, quantum kernel, shadow kernel and neural tangent kernel; (3) quantum machine learning for estimating all classical features (tomography) efficiently; (4) if we have all classical features, is it possible to train a universal classifier or with weaker promise; (5) can we estimate concurrence... by quantum circuit; (6) quantum complexity for separability

[1] R. Horodecki, P. Horodecki, M. Horodecki, and K. Horodecki, Rev. Mod. Phys. **81**, 865 (2009), arXiv:quant-ph/0702225.

[2] H. J. Briegel, D. E. Browne, W. Dür, R. Raussendorf, and M. V. den Nest, Nature Phys **5**, 19 (2009), arXiv:0910.1116.

[3] F. Xu, X. Ma, Q. Zhang, H.-K. Lo, and J.-W. Pan, Rev. Mod. Phys. **92**, 025002 (2020), arXiv:1903.09051.

[4] I. Cong, S. Choi, and M. D. Lukin, Nat. Phys. **15**, 1273 (2019), arXiv:1810.03787 [cond-mat, physics:quant-ph].

[5] J. Carrasquilla and R. G. Melko, Nature Phys **13**, 431 (2017), arXiv:1605.01735.

[6] H.-Y. Huang, R. Kueng, G. Torlai, V. V. Albert, and J. Preskill, Science **377**, eabk3333 (2022), arXiv:2106.12627.

[7] H.-Y. Huang, R. Kueng, and J. Preskill, Nat. Phys. **16**, 1050 (2020), arXiv:2002.08953 [quant-ph].

[8] O. Gühne and G. Toth, Physics Reports **474**, 1 (2009), arXiv:0811.2803 [cond-mat, physics:physics, physics:quant-ph].

[9] A. Peres, Phys. Rev. Lett. **77**, 1413 (1996), arXiv:quant-ph/9604005.

[10] M. Horodecki, P. Horodecki, and R. Horodecki, Physics Letters A **223**, 1 (1996), arXiv:quant-ph/9605038.

[11] L. Gurvits, Classical deterministic complexity of Edmonds' problem and Quantum Entanglement (2003), arXiv:quant-ph/0303055.

[12] L. M. Ioannou, Quantum Inf. Comput. **7**, 335 (2007), arXiv:quant-ph/0603199.

[13] A. C. Doherty, P. A. Parrilo, and F. M. Spedalieri, Phys. Rev. A **69**, 022308 (2004), arXiv:quant-ph/0308032.

[14] F. G. Brandão, M. Christandl, and J. Yard, in *Proc. 43rd Annu. ACM Symp. Theory Comput. - STOC 11* (ACM Press, San Jose, California, USA, 2011) p. 343, arXiv:1011.2751 [quant-ph].

[15] G. Gutoski, P. Hayden, K. Milner, and M. M. Wilde, Theory of Comput. **11**, 59 (2015), arXiv:1308.5788 [quant-ph].

[16] M. Navascues, M. Owari, and M. B. Plenio, Phys. Rev. A **80**, 052306 (2009), arXiv:0906.2731 [quant-ph].

[17] A. K. Ekert, C. M. Alves, D. K. L. Oi, M. Horodecki, P. Horodecki, and L. C. Kwek, Phys. Rev. Lett. **88**, 217901 (2002), arXiv:quant-ph/0203016.

[18] P. Horodecki and A. Ekert, Phys. Rev. Lett. **89**, 127902 (2002), arXiv:quant-ph/0111064.

[19] The well-known identity (related to the replica trick originating in spin glass theory)

$$\mathrm{Tr}(U^{\pi}(\rho_1 \otimes \cdots \otimes \rho_m)) = \mathrm{Tr}(\rho_1 \cdots \rho_m) \qquad (10)$$

where the RHS is the multivariate trace and $U^{\pi}$ is a unitary representation of the cyclic shift permutation.

[20] S. Johri, D. S. Steiger, and M. Troyer, Phys. Rev. B **96**, 195136 (2017), arXiv:1707.07658.

[21] Y. Quek, M. M. Wilde, and E. Kaur, Multivariate trace estimation in constant quantum depth (2022), arXiv:2206.15405 [hep-th, physics:quant-ph].

[22] For fidelity witness, promise that the state is either (1) fidelity $\|\rho_{\mathrm{bi}} - |\psi_{\mathrm{tar}}\rangle\langle\psi_{\mathrm{tar}}|\| \geq \alpha$; (2) fidelity $< \alpha$ implies $\rho \in \mathcal{S}_{\mathrm{bi}}$.

[23] T. Heinosaari and M. Ziman, *The Mathematical Language of Quantum Theory: From Uncertainty to Entanglement*, 1st ed. (Cambridge University Press, 2011).

[24] B. M. Terhal, Physics Letters A **271**, 319 (2000), arXiv:quant-ph/9911057.

[25] The CHSH inequality: $\mathbf{O}_{\mathrm{CHSH}} = (\mathbb{1}, ab, ab', a'b, a'b')$ with $a = Z, a' = X, b = (X - Z)/\sqrt{2}, b = (X + Z)/\sqrt{2}$ and $\mathbf{w}_{\mathrm{CHSH}} = (\pm 2, 1, -1, 1, 1)$.

[26] A. Acin, D. Bruss, M. Lewenstein, and A. Sanpera, Phys. Rev. Lett. **87**, 040401 (2001), arXiv:quant-ph/0103025.

[27] M. Bourennane, M. Eibl, C. Kurtsiefer, S. Gaertner, H. Weinfurter, O. Guehne, P. Hyllus, D. Bruss, M. Lewenstein, and A. Sanpera, Phys. Rev. Lett. **92**, 087902 (2004), arXiv:quant-ph/0309043.

[28] G. Toth and O. Guehne, Phys. Rev. Lett. **94**, 060501 (2005), arXiv:quant-ph/0405165.

[29] G. Tóth and O. Gühne, Phys. Rev. A **72**, 022340 (2005).

[30] Y. Zhou, Q. Zhao, X. Yuan, and X. Ma, npj Quantum Inf **5**, 83 (2019).

[31] Y. Zhang, Y. Tang, Y. Zhou, and X. Ma, Phys. Rev. A **103**, 052426 (2021), arXiv:2012.07606 [quant-ph].

[32] E. Y. Zhu, L. T. H. Wu, O. Levi, and L. Qian, Machine Learning-Derived Entanglement Witnesses (2021), arXiv:2107.02301 [quant-ph].

[33] H. Lu, Q. Zhao, Z.-D. Li, X.-F. Yin, X. Yuan, J.-C. Hung, L.-K. Chen, L. Li, N.-L. Liu, C.-Z. Peng, Y.-C. Liang, X. Ma, Y.-A. Chen, and J.-W. Pan, Phys. Rev. X 8, 021072 (2018).

[34] Y. Wang, Y. Li, Z.-q. Yin, and B. Zeng, npj Quantum Inf 4, 46 (2018), arXiv:1801.03782.

[35] Y. Zhou, B. Xiao, M.-D. Li, Q. Zhao, Z.-S. Yuan, X. Ma, and J.-W. Pan, npj Quantum Inf 8, 1 (2022).

[36] Y. Zhou, Phys. Rev. A 101, 012301 (2020), arXiv:1907.11495 [quant-ph].

[37] M. Weilenmann, B. Dive, D. Trillo, E. A. Aguilar, and M. Navascués, Phys. Rev. Lett. 124, 200502 (2020), arXiv:1912.10056 [quant-ph].

[38] O. Gühne, Y. Mao, and X.-D. Yu, Phys. Rev. Lett. 126, 140503 (2021), arXiv:2008.05961 [quant-ph].

[39] G. Riccardi, D. E. Jones, X.-D. Yu, O. Gühne, and B. T. Kirby, Exploring the relationship between the faithfulness and entanglement of two qubits (2021), arXiv:2102.10121 [quant-ph].

[40] O. Gühne and N. Lütkenhaus, Phys. Rev. Lett. 96, 170502 (2006).

[41] Y. Zhan and H.-K. Lo, Detecting Entanglement in Unfaithful States (2021), arXiv:2010.06054 [quant-ph].

[42] S. Lu, S. Huang, K. Li, J. Li, J. Chen, D. Lu, Z. Ji, Y. Shen, D. Zhou, and B. Zeng, Phys. Rev. A 98, 012315 (2018), arXiv:1705.01523 [quant-ph].

[43] Y.-C. Ma and M.-H. Yung, npj Quantum Inf 4, 34 (2018), arXiv:1705.00813 [quant-ph].

[44] D. Lu, T. Xin, N. Yu, Z. Ji, J. Chen, G. Long, J. Baugh, X. Peng, B. Zeng, and R. Laflamme, Phys. Rev. Lett. 116, 230501 (2016), arXiv:1511.00581 [quant-ph].

[45] S. V. Vintskevich, N. Bao, A. Nomerotski, P. Stankus, and D. A. Grigoriev, Classification of four-qubit entangled states via Machine Learning (2022), arXiv:2205.11512 [quant-ph].

[46] C. Cortes and V. Vapnik, Mach Learn 20, 273 (1995).

[47] A. Jacot, F. Gabriel, and C. Hongler, Neural Tangent Kernel: Convergence and Generalization in Neural Networks (2020), arXiv:1806.07572 [cs, math, stat].

[48] J. Altepeter, E. Jeffrey, and P. Kwiat, in Advances In Atomic, Molecular, and Optical Physics, Vol. 52 (Elsevier, 2005) pp. 105–159.

[49] Informally, quantum state tomography refers to the task of estimating complete description (density matrix) of an unknown $D$-dimensional state $\rho$ within error tolerence $\epsilon$, given the ability to prepare and measure copies of $\rho$.

[50] Intermediate between independent measurements and unrestricted (also called "collective" or "entangled") measurements are adaptive measurements in which the copies of $\rho$ are measured individually, but the choice of measurement basis can change in response to earlier measurements.

[51] J. Haah, A. W. Harrow, Z. Ji, X. Wu, and N. Yu, IEEE Trans. Inform. Theory , 1 (2017).

[52] R. O'Donnell and J. Wright, in Proc. Forty-Eighth Annu. ACM Symp. Theory Comput. (ACM, Cambridge MA USA, 2016) pp. 899–912.

[53] Nonlinear functions: entropy; multivariate functions: $\mathrm{Tr}(\rho_1 \cdots \rho_m)$, ?? $\mathrm{Tr}(\rho\rho')$, quadratic $\mathrm{Tr}(O\rho_i \otimes \rho_j)$, fidelity $F(\rho, \rho')$.

[54] S. Aaronson, in Proc. 50th Annu. ACM SIGACT Symp. Theory Comput., STOC 2018 (Association for Computing Machinery, New York, NY, USA, 2018) pp. 325–338, arXiv:1711.01053.

[55] Full tomography: additive error $\epsilon \ll 1/D$.

[56] Known fundamental lower bounds state that classical shadows of exponential size (at least) $T = \Omega(2^n/\epsilon^2)$ are required to $\epsilon$-approximate $\rho$ in trace distance.

[57] H.-Y. Huang, R. Kueng, and J. Preskill, Phys. Rev. Lett. 127, 030503 (2021), arXiv:2103.07510 [quant-ph].

[58] S. Chen, W. Yu, P. Zeng, and S. T. Flammia, PRX Quantum 2, 030348 (2021), arXiv:2011.09636 [quant-ph].

[59] T. Zhang, J. Sun, X.-X. Fang, X.-M. Zhang, X. Yuan, and H. Lu, Experimental quantum state measurement with classical shadows (2021), arXiv:2106.10190 [physics, physics:quant-ph].

[60] A. Elben, R. Kueng, H.-Y. Huang, R. van Bijnen, C. Kokail, M. Dalmonte, P. Calabrese, B. Kraus, J. Preskill, P. Zoller, and B. Vermersch, Phys. Rev. Lett. 125, 200501 (2020), arXiv:2007.06305 [cond-mat, physics:quant-ph].

[61] X. Gao and L.-M. Duan, Nat Commun 8, 662 (2017), arXiv:1701.05039 [cond-mat, physics:quant-ph].

[62] G. Torlai, G. Mazzola, J. Carrasquilla, M. Troyer, R. Melko, and G. Carleo, Nature Phys 14, 447 (2018), arXiv:1703.05334.

[63] H.-Y. Huang, M. Broughton, M. Mohseni, R. Babbush, S. Boixo, H. Neven, and J. R. McClean, Nat Commun 12, 2631 (2021), arXiv:2011.01938 [quant-ph].

[64] Y. Zhu, Y.-D. Wu, G. Bai, D.-S. Wang, Y. Wang, and G. Chiribella, Flexible learning of quantum states with generative query neural networks (2022), arXiv:2202.06804 [quant-ph].

[65] H.-Y. Huang, R. Kueng, and J. Preskill, Phys. Rev. Lett. 126, 190505 (2021), arXiv:2101.02464 [quant-ph].

[66] $\mathcal{O}(\log(M/\delta)\epsilon^{-4})$ copies of the unknown quantum state $\rho$. ($M = 4^n$ implies linear copy for full tomography???).

[67] The required amount of training data scales badly with $\epsilon$. This unfortunate scaling is not a shortcoming of the considered ML algorithm, but a necessary feature.

[68] J. R. Johansson, P. D. Nation, and F. Nori, Computer Physics Communications 184, 1234 (2013), arXiv:1110.0573.

[69] B. Li, S. Ahmed, S. Saraogi, N. Lambert, F. Nori, A. Pitchford, and N. Shammah, Quantum 6, 630 (2022), arXiv:2105.09902 [quant-ph].

[70] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and É. Duchesnay, J. Mach. Learn. Res. 12, 2825 (2011).

[71] M. Hein, W. Dür, J. Eisert, R. Raussendorf, M. V. den Nest, and H.-J. Briegel, Entanglement in Graph States and its Applications (2006), arXiv:quant-ph/0602096.

## Appendix A: Definitions

**Definition 5** (density matrix). A quantum (mixed) state $\rho$ can be represented by a density matrix which is a Hermitian, PSD operator (matrix) of trace one. If the rank of $\rho$ is 1, then the state is a pure state $\rho \equiv |\psi\rangle\langle\psi|$.

**Definition 6** (POVM). A positive-operator valued measurement (POVM) $M$ consists of a set of positive operators that sum to the identity operator $\mathbb{1}$. When a measurement $M = \{E_1, \ldots, E_k\}$ is applied to a quantum state $\rho$, the outcome is $i \in [k]$ with probability $p_i = \mathrm{tr}(\rho E_i)$. observables ... $\mathbb{E}[x] \equiv \langle O_x \rangle := \mathrm{tr}(O_x \rho)$

**Definition 7** (PSD). A matrix (operator) is positive, semidefinite (PSD) if all its eigenvalues are non-negative.

**Definition 8** (reduced density matrix). reduced density matrix $\rho_A = \mathrm{Tr}_B(\rho_{AB})$

**Definition 9** (partial transpose). [10] The partial transpose (PT) operation - acting on subsystem $A$ - is defined as

$$|k_A, k_B\rangle\langle l_A, l_B|^{\mathsf{T}_A} := |l_A, k_B\rangle\langle k_A, l_B| \tag{A1}$$

where $\{|k_A, k_B\rangle\}$ is a product basis of the joint system $\mathcal{H}_{AB}$.

**Definition 10** (Schmidt measure). Consider the following bipartite pure state, written in Schmidt form:

$$|\psi\rangle = \sum_i^r \sqrt{\lambda_i} \, |\phi_i^A\rangle \otimes |\phi_i^B\rangle \tag{A2}$$

where $\{|\phi_i^A\rangle\}$ is a basis for $\mathcal{H}_A$ and $\{|\phi_i^A\rangle\}$ for $\mathcal{H}_B$. The strictly positive values $\sqrt{\lambda_i}$ in the Schmidt decomposition are its *Schmidt coefficients*. The number of Schmidt coefficients, counted with multiplicity, is called its *Schmidt rank*, or Schmidt number. Schmidt measure is minimum of $\log_2 r$ where $r$ is number of terms in an expansion of the state in product basis. (Schmidt rank ?? $\mathrm{SR}^A(\psi) = \mathrm{rank}(\rho_\psi^A)$)

**Definition 11** (entropy). In quantum mechanics (information), the von Neumann *entropy* of a density matrix is $H_N(\rho) := -\mathrm{Tr}(\rho \log \rho) = -\sum_i \lambda_i \log(\lambda_i)$; In classical information (statistical) theory, the Shannon entropy of a probability distribution $P$ is $H_S(P) := -\sum_i P(x_i) \log P(x_i)$.

**Definition 12** (entanglement entropy). The bipartite *von Neumann entanglement entropy* $S$ is defined as the von Neumann entropy of either of its reduced density matrix $\rho_A$. For a pure state $\rho_{AB} = |\Psi\rangle\langle\Psi|_{AB}$, it is given by

$$E(\Psi_{AB}) = S(\rho_A) = -\mathrm{Tr}(\rho_A \log \rho_A) = -\mathrm{Tr}(\rho_B \log \rho_B) = S(\rho_B) \tag{A3}$$

where $\rho_A$ and $\rho_B$ are the reduced density matrix for each partition. With Schmidt decomposition (Eq. (A2)), the entropy of entanglement is simply $-\sum_i p_i^2 \log(p_i)$. the $n$th Renyi entropy, $S_n = \frac{1}{n-1} \log(R_n)$ where $R_n = \mathrm{Tr}(\rho_A^n)$

**Example 1.** The Schmidt measure for any multi-partite GHZ states is 1, because there are just two terms. Schmidt measure for 1D, 2D, 3D-cluster state is $\lfloor \frac{N}{2} \rfloor$. Schmidt measure of tree is the size of its minimal vertex cover[??]. other entanglement measures...

**Definition 13** (fidelity). Given a pair of states (target $\rho$ and prepared $\rho'$), Uhlmann fidelity $F(\rho, \rho') := \mathrm{Tr}\left(\sqrt{\sqrt{\rho}\rho'\sqrt{\rho}}\right) \equiv \left\|\sqrt{\rho}\sqrt{\rho'}\right\|_1$, where $\sqrt{\rho}$ dentoes the positive semidefinite square root of the operator $\rho$. (infidelity $1 - F(\rho, \rho')$) For any mixed state $\rho$ and pure state $|\psi\rangle$, $F(\rho, |\psi\rangle\langle\psi|) = \sqrt{\langle\psi|\rho|\psi\rangle} \equiv \sqrt{\mathrm{Tr}(\rho \, |\psi\rangle\langle\psi|)}$ which can be obtained by the Swap-test[?]. linear fidelity or overlap $F(\rho, \rho') := \mathrm{tr}(\rho\rho')$.

**Notation 3** (norm). Schatten p-norm $\|x\|_p := (\sum_i |x_i|^p)^{1/p}$. Euclidean norm $l_2$ norm; Spectral (operator) norm $\|\mathbf{x}\|_\infty$; Trace norm $\|A\|_{\mathrm{Tr}} \equiv \|A\|_1 := \mathrm{Tr}(|A|) \equiv \mathrm{Tr}\left(\sqrt{A^\dagger A}\right)$, $|A| := \sqrt{A^\dagger A}$, $p = 1$; Frobenius norm $\|A\|_F := \sqrt{\mathrm{Tr}(A^\dagger A)}$, $p = 2$; Hilbert-Schmidt norm $\|A\|_{HS} := \sqrt{\sum_{i,j} A_{ij}^2} =? \sum_{i \in I} \|Ae_i\|_H^2$; Hilbert-Schmidt inner product $\langle A, B \rangle_{\mathrm{HS}} := \mathrm{Tr}(A^\dagger B)$, Frobenius inner product $\langle A, B \rangle_{\mathrm{F}} := \mathrm{Tr}(A^\dagger B)$? (in finite-dimensionala Euclidean space, the HS norm is identical to the Frobenius norm) Although the Hilbert-Schmidt distance is arguably not too meaningful, operationally, one can use Cauchy-Schwarz to relate it to the very natural trace distance. shadow norm ...

**Definition 14** (distance). For mixed states, trace distance $d_{\mathrm{tr}}(\rho, \rho') := \frac{1}{2}\|\rho - \rho'\|_1$. For pure states, $d_{\mathrm{tr}}(|\psi\rangle, |\psi'\rangle) := \frac{1}{2}\||\psi\rangle\langle\psi| - |\psi'\rangle\langle\psi'|\|_1 = \sqrt{1 - |\langle\psi|\psi'\rangle|^2}$. fidelity and trace distance are related by the inequalities

$$1 - F \leq D_{\mathrm{tr}}(\rho, \rho') \leq \sqrt{1 - F^2} \tag{A4}$$

variation distance of two distribution $d_{var}(p, p') := \frac{1}{2}\sum_i |p_i - p_i'| = \frac{1}{2}\|p - p'\|_1$. $l_2$ distance ... Hellinger distance ... HS distance $D_{\mathrm{HS}}(\rho, \rho') := \|\rho - \rho'\|_{\mathrm{HS}} = \sqrt{\mathrm{Tr}((\rho - \rho')^2)}$

**Definition 15** (stabilizer)**.** An observable $S_k$ is a stabilizing operator of an $n$-qubit state $|\psi\rangle$ if the state $|\psi\rangle$ is an eigenstate of $S_k$ with eigenvalue 1, A stabilizer set $S = \{S_1, \ldots, S_n\}$ consisting of $n$ mutually commuting and independent stabilizer operators is called the set of stabilizer "generators".

Many highly entangled $n$-qubit states can be uniquely defined by $n$ stabilizing operators which are locally measurable, i.e., they are products of Pauli matrices. A stabilizer $S_i$ is an $n$-fold tensor product of $n$ operators chosen from the one qubit Pauli operators $\{\mathbb{1}, X, Y, Z\}$.

**Example 2** (GHZ)**.** For GHZ state: $|\text{GHZ}\rangle := \frac{1}{\sqrt{2}}(|0\rangle^{\otimes n} + |1\rangle^{\otimes n})$, the projector based witness $W_{\text{GHZ}_3} = \mathbb{1}/2 - |\text{GHZ}\rangle\langle\text{GHZ}|$ requires four measurement settings. For three-qubit GHZ state [28], the local measurement witness

$$W_{\text{GHZ}_3} := \frac{3}{2}\mathbb{1} - X^{(1)}X^{(2)}X^{(3)} - \frac{1}{2}\Big(Z^{(1)}Z^{(2)} + Z^{(2)}Z^{(3)} + Z^{(1)}Z^{(3)}\Big) \tag{A5}$$

This witness requires the measurement of the $\left\{\hat{\sigma}_x^{(1)}, \hat{\sigma}_x^{(2)}, \hat{\sigma}_x^{(3)}\right\}$ and $\left\{\hat{\sigma}_z^{(1)}, \hat{\sigma}_z^{(2)}, \hat{\sigma}_z^{(3)}\right\}$ settings. For $n$-qubit case, detect genuine $n$-qubit entanglement close to $\text{GHZ}_n$

$$W_{\text{GHZ}_n} = (n-1)\mathbb{1} - \sum_{k=1}^{n} S_k(\text{GHZ}_n) \tag{A6}$$

where $\hat{S}_k$ is the stabilizer ... [29]

**Definition 16** (cluster state)**.** 1D four qubits

$$\left|\psi_4^{1D}\right\rangle = \frac{1}{2}(|+00+\rangle + |+01-\rangle + |-10+\rangle - |-11-\rangle) \tag{A7}$$

The entanglement in a graph state is related to the topology of its underlying graph [71].

**Remark 1.** LU, LC equivalence, local operations and classical communication (LOCC),

**Definition 17** (graph state)**.** Given a simple graph (undirected, unweighted, no loop and multiple edge) $G = (V, E)$, a graph state is constructed as from the initial state $|+\rangle^{\otimes n}$ corresponding to $n$ vertices. Then, apply controlled-Z gate to every edge, that is $|G\rangle := \prod_{(i,j)\in E} \mathsf{cZ}_{(i,j)} |+\rangle^{\otimes n}$ with $|+\rangle := (|0\rangle + |1\rangle)/\sqrt{2}$.

| | $|\text{GHZ}_3\rangle$ | $|W_3\rangle$ | $|CL_3\rangle$ | $|\psi_2\rangle$ | $|\mathcal{D}_{2,4}\rangle$ | $|\text{GHZ}_n\rangle$ | $|W_n\rangle$ | $|G_n\rangle$ |
|---|---|---|---|---|---|---|---|---|
| maximal overlap $\alpha$ | 1/2 | 2/3 | 1/2 | 3/4 | 2/3 | 1/2 | $(n-1)/n$ | 1/2 |
| maximal $p_{\text{noise}}$ | 4/7 | 8/21 | 8/15 | 4/15 | 16/45 | $1/2 \cdot (1-1/2^n)^{-1}$ | $1/n \cdot (1-1/2^n)^{-1}$ | $1/2 \cdot (1-1/2^n)^{-1}$ |
| # local measurements | 4 | 5 | 9 | 15 | 21 | $n+1$ | $2n-1$ | depend on graphs |

TABLE II: Results on local decompositions of different entanglement witnesses for different states. [8]

**Appendix B: Machine learning background**

Notations: The (classical) training data (for supervised learning) is a set of $m$ data points $\left\{(\mathbf{x}^{(i)}, y^{(i)})\right\}_{i=1}^{m}$ where each data point is a pair $(\mathbf{x}, y)$. Normally, the input (e.g., an image) $\mathbf{x} := (x_1, x_2, \ldots, x_d) \in \mathbb{R}^d$ is a vector where $d$ is the number of *features* and its *label* $y \in \Sigma$ is a scalar with some discrete set $\Sigma$ of alphabet/categories. For simplicity and the purpose of this paper, we assume $\Sigma = \{-1, 1\}$ (binary classification).

**1. Support vector machine**

SVM is a typical supervised learning algorithm for classification. Taking the example of classifying cat/dog images, supervised learning means we are given a dataset in which every image is labeled either a cat or a dog such that we can find a function classifying new images with high accuracy. More precisely, the training dataset is a set of pairs of features X and their labels y. In the image classification case, features are obtained by transforming all pixels of an image into a vector. In SVM, we want to find a linear function, that is a hyperplane which separates cat data from dog data. So, the prediction label is given by the sign of the inner product (projection) of the hyperplane and the feature vector. We can observe that the problem setting of image classification by SVM is quite analogous to entanglement detection, where input data are quantum states now and the labels are either entangled or separable.

**Definition 18** (SVM). Given a set of (binary) labeled data, support vector machine (SVM) is designed to find a hyperplane (a linear function) such that maximize the margin between two partitions...

$$\max_{\mathbf{w}} \|\mathbf{w}\|^2 \text{ s.t. } \forall i, y^{(i)} \cdot (\mathbf{w} \cdot \mathbf{x} + b) \geq 1. \tag{B1}$$

Lagrange multipliers $\alpha$

$$L = \frac{1}{2}\|\mathbf{w}\|^2 - \sum_i^m \alpha^{(i)}\left(\mathbf{w} \cdot \mathbf{x}^{(i)} + b\right) + \sum_i^m \alpha^{(i)} \tag{B2}$$

[TODO]

### a. kernel method

However, note that SVM is only a linear classifier. while most real-world data, such as cat/dog images and entangled/separable quantum states are not linearly separable. For example, with this two dimension dataset, we are unable to find a hyperplane to separate red points from the purple points very well. Fortunately, there is a very useful tool called kernel method or kernel trick to remedy this drawback. The main idea is mapping the features to a higher dimensional space such that they can be linearly separated in the high dimensional feature space. Just like this example, two dimensional data are mapped to the three dimensional space. Now, we can easily find the separating plane. With SVM and kernel methods, we expect to find a generic and flexible way for entanglement detection. kernel

**Definition 19** (kernel). In general, the kernel function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ measures the similarity between two input data points by an inner product

$$k(\mathbf{x}, \mathbf{x}') := \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \tag{B3}$$

If the input $\mathbf{x} \in \mathbb{R}^d$ (conventional machine learning task, e.g., image classification), the feature map $\phi(\mathbf{x}) : \mathbb{R}^d \to \mathbb{R}^n$ ($d < n$) from a low dimensional space to a higher dimensional space. The corresponding kernel (Gram) matrix $\mathbf{K}$ is PSD.

**Example 3** (kernels). Some common kernels: the polynomial kernel $k_{\text{poly}}(\mathbf{x}, \mathbf{x}') := (1 + \mathbf{x} \cdot \mathbf{x}')^q$ with feature map $\phi(\mathbf{x})$ ... The Gaussian kernel $k_{\text{gaus}}(\mathbf{x}, \mathbf{x}') := \exp\left(-\gamma\|\mathbf{x} - \mathbf{x}'\|_2^2\right)$ with an infinite dimensional feature map $\phi(\mathbf{x})$. An important feature of kernel method is that kernels can be computed efficiently without evaluating feature map (might be infinite dimension) explicitly.