





Quantum Generative Adversarial Autoencoders: Learning latent representations for quantum data generation

Naipunnya Raj ^{*}, Rajiv Sangle ^{*}, Avinash Singh , and Krishna Kumar Sabapathy 
Quantum Lab, Fujitsu Research of India

In this work, we introduce the Quantum Generative Adversarial Autoencoder (QGAA), a quantum model for generation of quantum data. The QGAA consists of two components: (a) Quantum Autoencoder (QAE) to compress quantum states, and (b) Quantum Generative Adversarial Network (QGAN) to learn the latent space of the trained QAE. This approach imparts the QAE with generative capabilities. The utility of QGAA is demonstrated in two representative scenarios: (a) generation of pure entangled states, and (b) generation of parameterized molecular ground states for H_2 and LiH . The average errors in the energies estimated by the trained QGAA are 0.02 Ha for H_2 and 0.06 Ha for LiH in simulations upto 6 qubits. These results illustrate the potential of QGAA for quantum state generation, quantum chemistry, and near-term quantum machine learning applications.

I. INTRODUCTION

Over the past decade, machine learning has undergone transformative advancements, primarily fueled by the development of sophisticated deep learning architectures and training methodologies. In parallel, Quantum Machine Learning (QML) has emerged as a field dedicated to exploring how quantum algorithms and quantum computing platforms can be utilized to process, model, and extract meaningful insights from data [9, 14, 65], and also generate new data [26, 59]. While efforts in QML primarily focused on leveraging quantum computing to accelerate classical machine learning tasks [19, 34], a significant and increasingly important direction involves the development of quantum models that operate directly on quantum data [7, 9, 41]. These models, tailored specifically to quantum data, are essential for realizing the full potential of quantum technologies, enabling applications in quantum information processing that are intractable with classical methods [25].

A notable model within QML for handling quantum data is the Quantum Autoencoder (QAE), which draws inspiration from its classical counterpart, the Autoencoder (AE) [5, 58]. QAE has been applied to demonstrate how quantum circuits can be trained to compress quantum states, with applications to quantum simulation and quantum information [13, 29, 42, 44, 57]. Further developments extend these architectures to the denoising of entangled quantum states under realistic noise models [1, 10, 62, 63], along with proposals for error mitigation strategies tailored to Noisy Intermediate-Scale Quantum (NISQ) devices [46, 66]. Practical realizations of QAE in quantum hardware, such as nitrogen-vacancy centers, demonstrated robust compression and the preservation of entanglement, while significantly lengthening the coherence times of Bell states [67].

Despite these promising applications, it is important to note that the QAE, on its own, does not possess generative capabilities. It is fundamentally designed to learn a low-dimensional latent space that captures the essential features of a given dataset [57]. However, this compression process is not inherently generative since it does not allow direct access to the quantum latent space. In comparison, the classical Variational Autoencoder (VAE) enables direct access to the latent space of an Autoencoder by regularizing the latent space with a known distribution. This allows sampling from the latent space to generate new data samples [4]. There is no straightforward natural analogue to access the quantum latent space of a QAE for quantum generative tasks.

Alongside these developments, the intersection of classical generative models and QML has catalyzed significant interest in the formulation and development of quantum generative models, designed to learn and generate quantum data using quantum processes [61, 64]. Notable examples include the Quantum Generative Adversarial Network (QGAN) [22, 35, 41, 47], Quantum Boltzmann Machine [2], Quantum Circuit Born Machine [20], Quantum Transformers [33] and Quantum Generative Diffusion Model [18]. These quantum generative models are formulated to learn and synthesize quantum states. In particular among these, works on the QGAN [22, 35, 41, 47] establish theoretical foundations for quantum adversarial learning for quantum generative tasks, inspired by the formalism of classical Generative Adversarial Networks (GAN) [26].

To address the lack of generative capability of QAE, we propose a novel QML architecture, the Quantum Generative Adversarial Autoencoders (QGAA). It combines the compression power of the QAE with the generative capabilities of the QGAN. We demonstrate how the adversarial training framework of QGAN can be used to learn the latent quantum representation encoded by a trained QAE, imparting the model with generative capabilities. The key contributions of this work are as follows:

1. We present the theoretical framework of quantum adversarial learning to learn the representation of

^{*} These two authors contributed equally.

Emails: naipunnya.raj@fujitsu.com ; rajiv.sangle@fujitsu.com

the quantum latent space of a trained QAE. The adversarial learning approach gives the ability to directly generate new quantum latent states.

2. We demonstrate two applications:

- (a) Learning the latent representation of a set of 2-qubit entangled states.
- (b) Learning the ground state energy profiles of the parameterized molecular Hamiltonians of H_2 (4 qubits) and LiH (6 qubits).

These tasks allow us to evaluate the performance of QGAA while also highlighting practical challenges encountered during implementation.

A key advantage of the QGAA approach lies in its potential to reduce quantum resource requirements by compressing quantum states into a lower-dimensional latent space, while simultaneously enabling enhanced generative capabilities through adversarial learning.

The paper is organized as follows: Section II introduces the QAE, which finds the compressed representation of quantum states, and QGAN, a quantum generative model. Section III details the quantum adversarial framework for learning the latent space representation of a trained QAE. Sections IV and V demonstrate the implementation of this architecture for two different quantum-native generative tasks. Section VI concludes the article with insights gained from this work and suggests future research directions.

II. BACKGROUND

In this section, we provide an overview of the two QML models that form the foundation of this work: (i) the Quantum Autoencoder for the compression stage and (ii) the Quantum Generative Adversarial Network for the generative training stage.

A. Quantum Autoencoder

Autoencoder (AE) is a classical machine learning model designed to compress data into a low-dimensional latent space and reconstruct the compressed data to its original structure [37]. Building on the principles of the AE, the Quantum Autoencoder (QAE) was introduced as a quantum framework for efficient compression and reconstruction of quantum states [57]. An overview of classical AE is provided in Appendix A. This section details the architecture and the training process, designed to minimize information loss during compression and reconstruction of quantum states by QAE, as illustrated in Figure 1.

The core concept of a QAE is learning an optimal compressed or latent form of a given set of input quantum

states $\Gamma := \{\sigma_K\}$. Each state, $\sigma_K \in \Gamma$, is uniquely characterized by some label, $K \in \Lambda$. Here, $\Lambda := \{K\}$ denotes the set of all valid labels that define Γ .

Architecture: The QAE consists of three key components. The encoder is a parametrized unitary transformation

$$U_E(\vec{\theta}_E) : \mathcal{H}_A \longrightarrow \mathcal{H}_L \otimes \mathcal{H}_T, \quad (1)$$

acting on input states $\sigma_K \in \mathcal{D}(\mathcal{H}_A)$, where $\mathcal{D}(\mathcal{H})$ denotes the set of density operators on the Hilbert space \mathcal{H} . Here, $\vec{\theta}_E$ represents the trainable parameters, $\mathcal{H}_A \cong (\mathbb{C}^2)^{\otimes n}$ is the n -qubit Hilbert space of the input, $\mathcal{H}_L \cong (\mathbb{C}^2)^{\otimes \ell}$ is the latent subspace of dimension 2^ℓ , and $\mathcal{H}_T \cong (\mathbb{C}^2)^{\otimes (n-\ell)}$ is the trash subspace of dimension $2^{n-\ell}$. Tracing out the trash subsystem yields the compressed latent state

$$\eta_K = \text{Tr}_T[U_E(\vec{\theta}_E) \sigma_K U_E^\dagger(\vec{\theta}_E)] \in \mathcal{D}(\mathcal{H}_L). \quad (2)$$

The decoder is a second parametrized unitary

$$U_D(\vec{\phi}_D) : \mathcal{H}_L \otimes \mathcal{H}_T \longrightarrow \mathcal{H}_A, \quad (3)$$

governed by trainable parameters $\vec{\phi}_D$. It acts on the latent state η_K , together with an initialized trash register, to produce the reconstructed output state $\rho_K \in \mathcal{D}(\mathcal{H}_A)$. The unitaries $U_E(\vec{\theta}_E)$ and $U_D(\vec{\phi}_D)$ are typically realized as parameterized quantum circuits, composed of single-qubit rotation gates (with trainable angles) and fixed entangling gates.

Training: The objective is to minimize the loss function that quantifies the distance between the original and reconstructed quantum states in a given data ensemble.

In practice, the SWAP test [6, 12, 17] is employed to estimate the overlap between two quantum states. If $\sigma_K = |\psi_K\rangle\langle\psi_K|$ and $\rho_K = |\phi_K\rangle\langle\phi_K|$ are pure states, then the fidelity of the two states is equal to the overlap evaluated using the SWAP test, i.e., $\mathcal{F}(\sigma_K, \rho_K) = \text{SWAP}(\sigma_K, \rho_K) = |\langle\psi_K|\phi_K\rangle|^2$. In this case, the results of the SWAP test between $\{\sigma_K\}$ and $\{\rho_K\}$ are used to evaluate the QAE loss function \mathcal{L}_{QAE} .

If the standard SWAP test is applied to two mixed states, the outcome corresponds to the overlap $\text{Tr}(\rho\sigma)$ rather than the fidelity. In other words, the SWAP test provides a direct estimate of fidelity only for pure states. For mixed states, the fidelity must be evaluated analytically. A detailed definition of the SWAP operator, the associated protocol, and the analytical expressions for fidelity in both pure and mixed state settings are provided in Appendix B.

For the purpose of the demonstrative examples discussed in this work, we have used the analytical expression of fidelity to define the cost function of the QAE,

$$\mathcal{L}_{\text{QAE}} = \mathbb{E}_{\sigma_K \in \{\sigma_K\}} [1 - \mathcal{F}(\sigma_K, \rho_K)]. \quad (4)$$

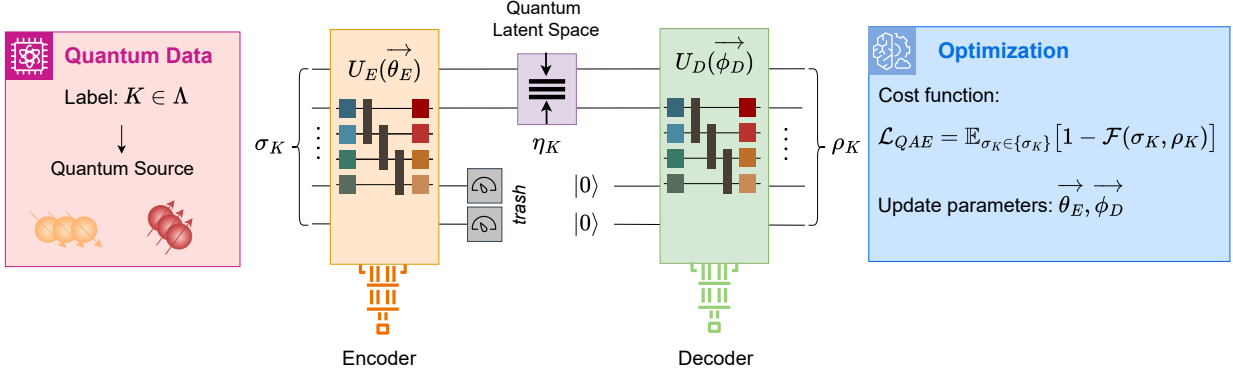


FIG. 1: **Compression stage using QAE:** The input quantum states σ_K is indexed by some label $K \in \Lambda$ which uniquely characterizes the state σ_K . Here, $\Lambda := \{K\}$ denotes the set of all valid labels that define $\{\sigma_K\}$. The encoder applies a parametrized unitary $U_E(\vec{\theta}_E)$ to each input state σ_K , generating an entangled intermediate state. A designated subset of qubits referred to as the *trash*, is then traced out to yield the compressed latent state η_K . The decoder subsequently applies a second parametrized unitary $U_D(\vec{\phi}_D)$ to reconstruct the output state ρ_K . The parameters $\vec{\theta}_E$ and $\vec{\phi}_D$ are trained to minimize the reconstruction loss, \mathcal{L}_{QAE} .

The loss function, \mathcal{L}_{QAE} , is minimized using classical optimization techniques to find parameter values of the encoder, $U_E(\vec{\theta}_E)$, and decoder, $U_D(\vec{\phi}_D)$, such that the reconstructed state, ρ_K , closely matches its corresponding input state, σ_K . The optimal parameters, $\vec{\theta}_E^*$ and $\vec{\phi}_D^*$ obtained after training, are used to evaluate the performance of the QAE on unseen test data by computing metrics such as fidelity. The classical Variational Autoencoder (VAE) extends the standard AE by enabling sampling from a probabilistic latent space, thereby combining compression with generative modeling, as detailed in Appendix C. In contrast, the QAE lacks generative capabilities and does not allow sampling of new quantum states from its learned latent space. This absence of a fully realized quantum analogue motivates the integration of a QGAN to address the generative limitation.

B. Quantum Generative Adversarial Network

The notion of Quantum Generative Adversarial Learning [41] was formalized taking inspiration from the classical GAN [26]. An overview of the classical GAN is provided in Appendix D.

Quantum mechanics is inherently probabilistic and therefore the notion of generation of quantum data is different from that of classical data. Given access to a quantum source, R , described by some density matrix, σ , the quantum generative task is defined as training a quantum agent - the generator, G , to generate ρ , whose statistics are equal to that of σ , i.e., $\rho = \sigma$ [41]. Here, statistics refers to the measurement outcomes of σ and

ρ corresponding to any tomographically complete set of observables.

Architecture: In a QGAN, the training of G is enabled by another quantum agent - the discriminator, D , whose objective is to evaluate whether the input quantum state supplied to it is σ (*real data*) from R or ρ (*fake data*) generated by G . Both G and D can be implemented as quantum circuits parameterized by $\vec{\theta}_g$ and $\vec{\theta}_d$ respectively.

D consists of an input register, where either σ (*real data*) or ρ (*fake data*) can be loaded. Apart from the input data register, D also consists of a probe qubit in some reference state such as $|0\rangle$. The evolution dynamics of $D(\vec{\theta}_d)$ cause the input state, σ or ρ , to interact with the reference probe. After this joint evolution of the input state and the reference probe, measuring the probe qubit in the Pauli-Z basis, $\{|0\rangle\langle 0|, |1\rangle\langle 1|\}$, is equivalent to the action of a binary Positive Operator-Valued Measure (POVM) $\{\hat{T} \equiv \hat{T}(\vec{\theta}_d), \hat{F} \equiv \hat{F}(\vec{\theta}_d)\}$ on the input state corresponding to the decisions of the input being *real* (+1) or *fake* (-1) [22, 27, 41, 56]. Naturally, for any $\vec{\theta}_d$, \hat{T} and \hat{F} are positive semi-definite operators satisfying the completeness condition $\hat{T} + \hat{F} = \mathbb{I}$.

D implements its objective of distinguishing between the states σ (*real data*) and ρ (*fake data*) by maximizing the probability cost function \mathcal{L}_{QGAN} [22, 41],

$$\rho \equiv G(\vec{\theta}_g) \text{ and } \hat{T} \equiv \hat{T}(\vec{\theta}_d), \quad (5a)$$

1. Conditional QGAN

$$\min_{\vec{\theta}_g} \max_{\vec{\theta}_d} \mathcal{L}_{\text{QGAN}} := \frac{1}{2} \left[1 + \left\{ \text{Tr}(\hat{T}\sigma) - \text{Tr}(\hat{T}\rho) \right\} \right]. \quad (5b)$$

$\mathcal{L}_{\text{QGAN}}$ is designed in accordance with the quantum state discrimination protocol [27, 56] designed to maximize the probability of distinguishing between the quantum states σ and ρ . The upper bound on $\mathcal{L}_{\text{QGAN}}$ is given by the theory of Helstrom measurement or the minimum-error distinguishing measurement [27, 56].

Training: The quantum generative adversarial learning game is set up as follows:

1. The parameters, $\vec{\theta}_g$ and $\vec{\theta}_d$, are initialized, and G generates an initial state, ρ_0 . The corresponding $\mathcal{L}_{\text{QGAN}}$ is evaluated with the initial value of the operator $\hat{T} \equiv \hat{T}(\vec{\theta}_d)$.
2. The strategy of D is to optimize $\vec{\theta}_d$ to produce a corresponding \hat{T}' that maximizes $\mathcal{L}_{\text{QGAN}}$ keeping the parameters of G fixed.
3. The strategy of G is to optimize $\vec{\theta}_g$ to generate ρ' such that $\text{Tr}\{\hat{T}'\rho'\}$ gets closer to $\text{Tr}\{\hat{T}'\sigma\}$ and minimizes $\mathcal{L}_{\text{QGAN}}$ keeping the parameters of D fixed.
4. Subsequent iterations of Step 2 and Step 3 correspond to the max-min optimization of $\mathcal{L}_{\text{QGAN}}$. Both the space of density matrices, $\{\rho(\vec{\theta}_g)\}$, and the space of positive semi-definite operators, $\{\hat{T}(\vec{\theta}_d)\}$, are convex and compact. Therefore, optimization over these spaces is possible and a unique Nash equilibrium exists exactly at $\rho(\vec{\theta}_g^*) = \sigma$ for some optimal parameters $\vec{\theta}_g^*$ [41].
5. Upon convergence to the unique fixed-point of the game, D either distinguishes between σ and $\rho(\vec{\theta}_g^*)$ as good as a fair-coin toss [41] or identically concludes that the state is *real* no matter whether the input state is *real* or *generated/fake* [22].

However, the following assumptions are made in this framework:

1. First, it is assumed that both the quantum information processors G and D have enough *capacity* [47] or *expressibility* [60] such that they can approximate any arbitrary function or transformation using their parameters $\vec{\theta}_g$ and $\vec{\theta}_d$ respectively.
2. Second, it is assumed that an efficient optimization scheme exists that can drive the algorithm to *converge* to the unique Nash equilibrium.

The highlighted assumptions are detailed in Appendix E.

As an augmentation over QGAN, the quantum source, R , can have an additional description such that a label, $K \in \Lambda$, can be supplied to it as an input, conditioning $R \equiv R(K)$ with the density matrix representation σ_K . Similar to the notation defined in Section II A, $\Lambda := \{K\}$ denotes the set of all valid labels that can be supplied to R to produce the set of quantum states $\Gamma := \{\sigma_K\}$. The QGAN formalism described so far can also be extended to train the generator, G , to generate a *fake state*, ρ_K , conditioned on K , such that, $\rho_K = \sigma_K \forall K \in \Lambda$ at the Nash equilibrium of the quantum adversarial game [22].

Similar to the QGAN training, the quantum adversarial game for the conditional QGAN is as follows:

1. The parameters $\vec{\theta}_g$ and $\vec{\theta}_d$ of the generator, G , and the discriminator, D , respectively are initialized. Conditioned on the label, K , the real source, $R \equiv R(K)$, produces, σ_K , and the initial strategies of G and D can be described by ρ_K and \hat{T}_K respectively where,

$$\rho_K \equiv G(K, \vec{\theta}_g) \text{ and } \hat{T}_K \equiv \hat{T}(K, \vec{\theta}_d). \quad (6a)$$

Based on Equation 5b, the probability of distinguishing between σ_K and ρ_K using the discriminating strategy \hat{T}_K is evaluated as,

$$\mathcal{L}_K := \frac{1}{2} \left[1 + \left\{ \text{Tr}(\hat{T}_K \sigma_K) - \text{Tr}(\hat{T}_K \rho_K) \right\} \right]. \quad (6b)$$

The cost function, $\mathcal{L}_{\text{QGAN}}$, for the adversarial training is defined as,

$$\mathcal{L}_{\text{QGAN}} = \frac{1}{N} \sum_{K \in \Lambda} \mathcal{L}_K, \quad (6c)$$

the average of the probabilities of discrimination corresponding to different labels, $K \in \Lambda$ [22].

2. The strategy of D is to optimize $\vec{\theta}_d$ to produce \hat{T}'_K that maximizes $\mathcal{L}_{\text{QGAN}}$ keeping the parameters of G fixed.
3. The strategy of G is to optimize $\vec{\theta}_g$ to generate ρ'_K such that $\text{Tr}\{\hat{T}'_K \rho'_K\}$ gets closer to $\text{Tr}\{\hat{T}'_K \sigma_K\}$ and minimizes $\mathcal{L}_{\text{QGAN}}$ while keeping the parameters of D fixed.
4. Subsequent iterations of Step 2 and Step 3 correspond to the max-min optimization of $\mathcal{L}_{\text{QGAN}}$ in Equation 6c. A unique Nash equilibrium exists for some optimal parameters $\vec{\theta}_g^*$ such that $\rho_K(\vec{\theta}_g^*) = \sigma_K \forall K \in \Lambda$ [22, 41].

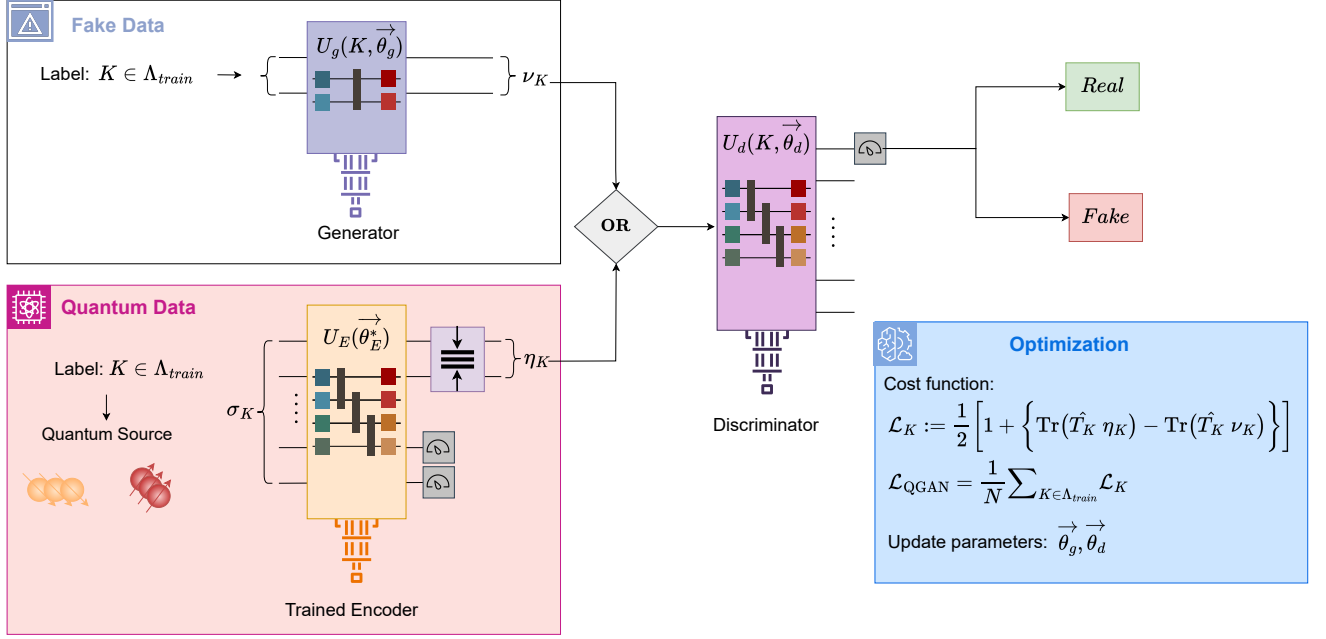


FIG. 2: **QGAA adversarial learning stage** – The states $\{\eta_K\}$ obtained from the *trained encoder* $U_E(\vec{\theta}_E^*)$ of the QAE are the *real data* for training the QGAN. The training label set, Λ_{train} , is the set of all labels whose corresponding $\{\sigma_K\}$ are used to train the QGAA. The *fake data* $\{\nu_K\}$, conditioned on $K \in \Lambda_{train}$, is generated by the generator $U_g(K, \vec{\theta}_g)$. In a training iteration, the discriminator $U_d(K, \vec{\theta}_d)$ estimates the likelihood of the input state (η_K or ν_K) supplied to it being *real*. The objective of the discriminator is to optimize $\vec{\theta}_d$ to correctly estimate η_K as *real* and ν_K as *fake* by maximizing the cost function \mathcal{L}_{QGAN} . Whereas, the objective of the generator is to optimize $\vec{\theta}_g$ to generate ν_K whose likelihood of being classified as *fake* is minimized. The adversarial training implements these objectives via the min-max optimization of the cost function \mathcal{L}_{QGAN} by the generator and the discriminator respectively. The Nash equilibrium of the adversarial training is reached when the discriminator can no longer distinguish between η_K and ν_K . At this point, the latent space representation of the *trained encoder* has been learned since $\nu_K(\vec{\theta}_g^*) = \eta_K \forall K \in \Lambda_{train}$ for some optimal parameters $\vec{\theta}_g^*$. If Λ_{train} is chosen to be a representative sample of the larger set Λ , then it is expected that $\nu_K(\vec{\theta}_g^*) = \eta_K \forall K \in \Lambda$ as well.

5. Upon convergence to the unique fixed-point of the game, D either distinguishes between σ_K and $\rho_K(\vec{\theta}_g^*) \forall K \in \Lambda$ as good as a fair-coin toss [41] or identically concludes that the state is *real* no matter whether the input state is *real* or *fake* (i.e. generated) [22].

This framework automatically requires R to be controllable with respect to K . Such a structure also necessitates that G is equipped with a suitable feature map that can encode any label, $K \in \Lambda$, and that the space of density matrices $\{\rho_K(\vec{\theta}_g)\}$ contains the *unique fixed-point* $\rho_K = \sigma_K \forall K \in \Lambda$ for the optimal parameters $\vec{\theta}_g^*$. For quantum generative tasks such as those described in Sections IV and V, the encoded label is a continuous variable. In such cases, a representative sample $\Lambda_{train} \in \Lambda$ can be chosen whose corresponding states $\{\sigma_K\}_{train}$ can be used for the adversarial training.

Section III describes how this formalism of conditional

QGAN can be leveraged to confer generative capability to the traditional QAE that is designed to only perform quantum state compression and reconstruction.

III. ADVERSARIAL FORMALISM FOR LEARNING QUANTUM LATENT SPACE

The work detailed in this section draws inspiration from the classical Adversarial Autoencoder [45], where, adversarial training is used to regularize the latent space of a classical Autoencoder. Appendix F describes this classical machine learning architecture.

In this work, we propose using quantum adversarial learning to model the latent space representation, $\{\eta_K\}$, of a trained QAE. $\{\eta_K\}$ essentially captures the core features of the original quantum data $\Gamma := \{\sigma_K\}$, and the adversarial learning approach enables direct access to $\{\eta_K\}$.

With the *trained encoder* of a QAE as the real source

R , we demonstrate that, under the assumptions described in Section II B, the quantum generative adversarial learning approach can be used to directly generate quantum latent states, $\{\nu_K = \eta_K\}$, corresponding to the labels, $K \in \Lambda$. Then the generated latent states, $\{\nu_K\}$, can be supplied to the corresponding *trained decoder* to generate the reconstructed quantum states, $\{\xi_K = \rho_K\}$. If the trained QAE performs perfect reconstruction, i.e., $\{\rho_K = \sigma_K\}$, then $\{\xi_K = \sigma_K\}$.

The following quantum adversarial game to achieve this is detailed in Algorithm 1 and depicted in Figure 2:

1. The parameters $\vec{\theta}_g$ and $\vec{\theta}_d$ of the generator, G , and the discriminator, D , respectively are initialized, and a representative sample of training labels, $\Lambda_{train} := \{K\}_{train} \in \Lambda$ is defined, where, $N = |\Lambda_{train}|$, is the size of the training label set. The real source, $R \equiv R(K)$, comprises the *trained encoder*, $U_E(\vec{\theta}_E^*)$, which can produce the compressed representation η_K , corresponding to the *real state*, σ_K , conditioned on the label K . Similarly, G can generate a *fake state*, ν_K , conditioned on the label, K , and the initial strategy of D can be described by the corresponding positive semi-definite operator \hat{T}_K ,

$$\nu_K \equiv G(K, \vec{\theta}_g) \text{ and } \hat{T}_K \equiv \hat{T}(K, \vec{\theta}_d). \quad (7a)$$

Similar to Equation 6b, the probability of distinguishing between η_K and ν_K using the discriminating strategy \hat{T}_K is evaluated as,

$$\mathcal{L}_K := \frac{1}{2} \left[1 + \left\{ \text{Tr}(\hat{T}_K \eta_K) - \text{Tr}(\hat{T}_K \nu_K) \right\} \right]. \quad (7b)$$

Similar to Equation 6c, the cost function, \mathcal{L}_{QGAN} , for the adversarial training is defined as,

$$\mathcal{L}_{QGAN} = \frac{1}{N} \sum_{K \in \Lambda_{train}} \mathcal{L}_K, \quad (7c)$$

the average of the probabilities of discrimination corresponding to different labels, $K \in \Lambda_{train}$ [22].

2. The strategy of D is to optimize $\vec{\theta}_d$ to produce \hat{T}'_K that maximizes \mathcal{L}_{QGAN} keeping the parameters of G fixed.
3. The strategy of G is to optimize $\vec{\theta}_g$ to generate ν'_K such that $\text{Tr}\{\hat{T}'_K \nu'_K\}$ gets closer to $\text{Tr}\{\hat{T}'_K \eta_K\}$ and minimizes \mathcal{L}_{QGAN} while keeping the parameters of D fixed.
4. Subsequent iterations of Step 2 and Step 3 correspond to the max-min optimization of \mathcal{L}_{QGAN} in Equation 7c. A unique Nash equilibrium exists for some optimal parameters $\vec{\theta}_g^*$ such that $\nu_K(\vec{\theta}_g^*) = \eta_K \forall K \in \Lambda_{train}$ [22, 41].

Algorithm 1: Adversarial training for learning the latent space of a Quantum Autoencoder.

- 1 **Input Quantum Data:** $\{\eta_K\}$
Result: $\nu_K = \eta_K \forall K \in \Lambda$
 - 2 $U_E^* \equiv U_E(\vec{\theta}_E^*)$;
 - 3 The Generator $U_g(K, \vec{\theta}_g)$ generates $\nu_K(\vec{\theta}_g)$;
 - 4 The Discriminator $U_d \equiv U_d(K, \vec{\theta}_d)$ evaluates the probabilities $\text{Prob}(+1|\eta_K)$ and $\text{Prob}(+1|\nu_K)$ corresponding to the probe qubit determining the input state (η_K or ν_K) to be *real* (+1) ;
 - 5 Initialize the parameters $\vec{\theta}_g$ and $\vec{\theta}_d$, *iter* = 1;
 - 6 Define maximum number of iterations: *max_iter* ;
 - 7 **while** *iter* ≤ *max_iter* **do**
 - 8 **while** $K \in \Lambda_{train}$ **do**
 - 9 $\eta_K := \text{Tr}_{trash} \left\{ U_E^* \sigma_K U_E^{*\dagger} \right\}$
 - 10 $\langle Z_{real} \rangle :=$
 $\text{Tr} \left[\left\{ U_d(\eta_K \otimes |0\rangle \langle 0|_{probe}) U_d^\dagger \right\} Z_{probe} \right]$
 - 11 $\therefore \text{Prob}(\text{probe} = +1 | \eta_K) = \frac{1 + \langle Z_{real} \rangle}{2}$
 - 12 $\langle Z_{fake} \rangle :=$
 $\text{Tr} \left[\left\{ U_d(\nu_K \otimes |0\rangle \langle 0|_{probe}) U_d^\dagger \right\} Z_{probe} \right]$
 - 13 $\therefore \text{Prob}(\text{probe} = +1 | \nu_K) = \frac{1 + \langle Z_{fake} \rangle}{2}$
 - 14 **Success Probability as cost function:**
 $\mathcal{L}_K := \frac{1}{2} \left[1 + \{ \text{P}(+1|\eta_K) - \text{P}(+1|\nu_K) \} \right]$
 - 15 $\mathcal{L}_{QGAN} = \frac{1}{N} \sum_{K \in \Lambda_{train}} \mathcal{L}_K$
 - 16 **→ Discriminator's Strategy:**
 Optimize and Update $\vec{\theta}_d$ to maximize \mathcal{L}_{QGAN} ;
 - 17 **Compute \mathcal{L}_{QGAN} using updated $\vec{\theta}_d$;**
 - 18 **→ Generator's Strategy:**
 Optimize and Update $\vec{\theta}_g$ to minimize \mathcal{L}_{QGAN} ;
 - 19 *iter* += 1
-

5. Upon convergence to the unique fixed-point of the game, D *either* distinguishes between η_K and $\nu_K(\vec{\theta}_g^*) \forall K \in \Lambda_{train}$ as good as a fair-coin toss [41] *or* identically concludes that the state is *real* no matter whether the input state is *real* or *fake* (i.e. generated) [22].

After the model has been trained to converge to its Nash equilibrium, G is capable of directly generating latent states corresponding to any label, $K \in \Lambda_{train}$. Since Λ_{train} is chosen to be a representative sample of the larger set Λ , it is expected that $\nu_K(\vec{\theta}_g^*) = \eta_K \forall K \in \Lambda$ as well. The generated latent state, ν_K , can be fed into the *trained decoder*, $U_D(\vec{\theta}_D^*)$, of the QAE to reconstruct

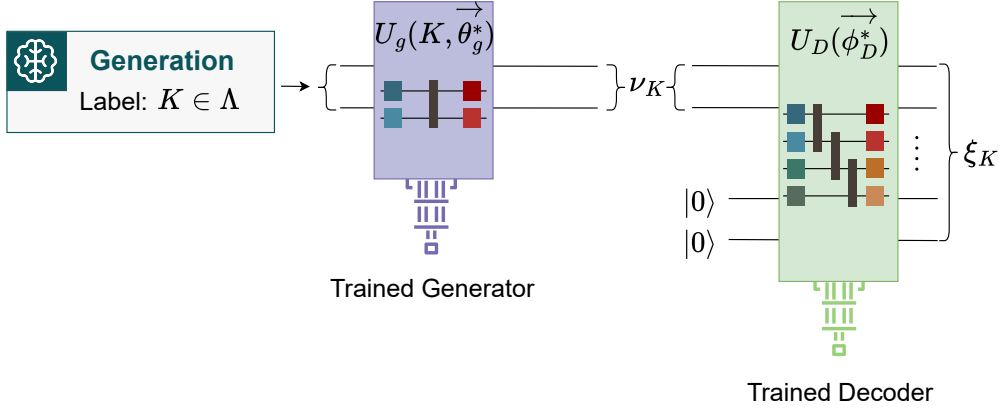


FIG. 3: **QGAA generation stage**— Upon learning the latent space representation, the *trained generator* $U_g(K, \vec{\theta}_g^*)$ is capable of directly generating the latent state $\nu_K = \eta_K$ corresponding to any label $K \in \Lambda$.

The generated latent state ν_K is fed into the *trained decoder* $U_D(\vec{\theta}_D^*)$ of the QAE to reconstruct $\xi_K = \rho_K$. If the trained QAE performs perfect reconstruction, i.e., $\{\rho_K = \sigma_K\}$, then $\{\xi_K = \sigma_K\}$. Thus, the quantum adversarial protocol is able to give generative capabilities to a QAE by learning the representation of its latent space.

$\xi_K = \rho_K$ as depicted in Figure 3. Further, in case of perfect reconstruction capabilities of the trained QAE, $\xi_K = \sigma_K$. Thus, the quantum adversarial protocol is able to give generative capabilities to a QAE by learning the representation of its latent space.

Through informative applications in the following Sections IV and V, we investigate this formalism to learn the latent space representations of trained QAEs and generate quantum data directly from the latent spaces. We show that the proposed adversarial learning scheme is capable of generating new quantum states (corresponding to labels outside the training dataset) exhibiting a particular property, such as entanglement and ground states of a parameterized molecular Hamiltonian.

IV. APPLICATION EXAMPLE 1: LEARNING LATENT REPRESENTATION OF ENTANGLED STATES

This section demonstrates the application of the adversarial formalism to learn the latent space representation of entangled states. The following 2-qubit entangled state is prepared conditioned on the label information $K = (k_0, k_1)$:

$$|\psi_K\rangle = \cos\left(\frac{k_0}{2}\right) |00\rangle + e^{ik_1} \sin\left(\frac{k_0}{2}\right) |11\rangle, \quad (8a)$$

$$\sigma_K = |\psi_K\rangle \langle \psi_K|. \quad (8b)$$

Here, the label, K , is further comprised of sub-labels, k_0 and k_1 , that characterize σ_K .

$|\psi_K\rangle$ is created by first preparing the single qubit state $|\chi_K\rangle$ conditioned on the sub-labels k_0 and k_1 :

$$\begin{aligned} |\chi_K\rangle &= RZ(k_1) RY(k_0) |0\rangle \\ &= \cos\left(\frac{k_0}{2}\right) |0\rangle + e^{ik_1} \sin\left(\frac{k_0}{2}\right) |1\rangle. \end{aligned} \quad (9a)$$

Applying the CX gate on another qubit q_1 (in the state $|0\rangle$) with qubit q_0 (in the state $|\chi_K\rangle$) as the control, generates the 2-qubit entangled state $|\psi_K\rangle$ as depicted in Figure 4:

$$|\psi_K\rangle = \text{CX} |\chi_K\rangle |0\rangle. \quad (9b)$$

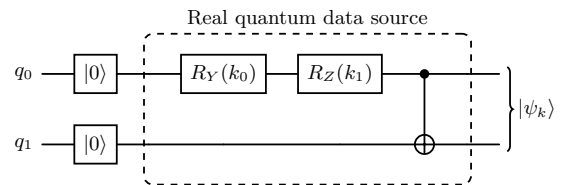


FIG. 4: Quantum circuit for generating 2-qubit parameterized entangled states $|\psi_K\rangle$ conditioned on the label information $K = (k_0, k_1) \in \Lambda$ supplied as input.

Therefore, the set $\{|\chi_K\rangle\}$ is a 1-qubit latent space representation of $\{\sigma_K\}$ since the reversible operation of a CX gate on $|\psi_K\rangle$ can “compress” it to a corresponding 1-qubit representation, $|\chi_K\rangle$. It is worthwhile to note that this 1-qubit latent space representation of $\{\sigma_K\}$ is not unique. A QAE upon iterative training can converge to a 1-qubit latent representation different than the one in Equation 9a.

For ease of visualization in the following discussion, we consider the set of training labels, Λ_{train} , to be

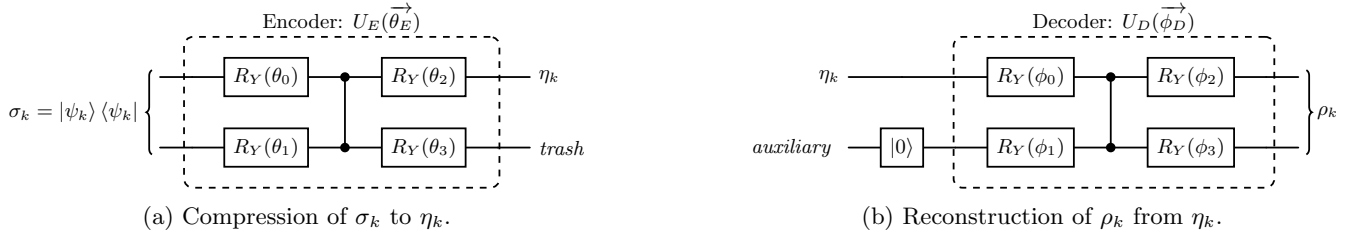


FIG. 5: Parameterized quantum circuits of (a) the encoder to compress 2-qubit entangled states $\{\sigma_k\}$ to a 1-qubit representation $\{\eta_k\}$ by discarding the processed *trash qubit*, and (b) the decoder to reconstruct of 2-qubit entangled states $\{\rho_k\}$ from $\{\eta_k\}$. The parameters $\vec{\theta}_E$ and $\vec{\phi}_D$ of the encoder and the decoder respectively are optimized during the QAE training to achieve $\mathcal{F}(\sigma_k, \rho_k) = 1$.

the set of labels $\{K\}_{train}$, whose corresponding $\{\sigma_K = |\psi_K\rangle\langle\psi_K|\}$ comprises all the states with entanglement entropy greater than 97%. The corresponding $|\chi_K\rangle$ are depicted on the Bloch sphere in Figure 6. The entanglement entropy of $\{|\psi_K\rangle\}$ corresponding to $k_0 = 0.5\pi \pm 0.06\pi$ and $k_0 = 0.5\pi$ is 0.97 and 1.00 respectively.

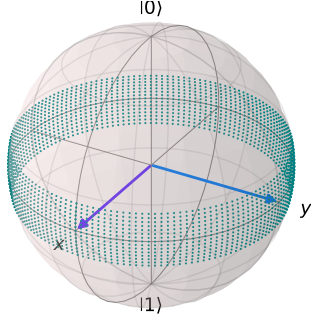


FIG. 6: 1-qubit pure states $\{|\chi_K\rangle\}$ on the Bloch sphere corresponding to $k_0 \in \left[\frac{\pi}{2} - 0.06\pi, \frac{\pi}{2} + 0.06\pi\right]$ and $k_1 \in [0, 2\pi)$. The two states marked by Bloch vectors correspond to $K = (0.5\pi, 0)$ and $K = (0.5\pi, 0.5\pi)$.

A. Training the QAE for state compression

The QAE is optimized according to Section II A to compress the 2-qubit entangled states $\sigma_K = |\psi_K\rangle\langle\psi_K|$ to a 1-qubit latent representation η_K . Figure 5 depicts the ansatzes used for the encoder and the decoder of the QAE. Figure 7 depicts the decreasing loss function of the QAE upon training using the COBYLA optimizer.

The QAE achieves perfect reconstruction of the input states $\{\sigma_K\}$ with fidelity of state reconstruction converging to $\mathcal{F}(\sigma_K, \rho_K) = 1$. The compressed states $\{\eta_K\}$ have a purity equal to 1. Therefore, the states $\{\eta_K\}$ have a pure state representation $\{\eta_K = |\gamma_K\rangle\langle\gamma_K|\}$ and lie on the surface of the Bloch sphere as depicted in Figure 8. Following are the optimal parameters of the *trained encoder* and the *trained decoder* that achieve these results:

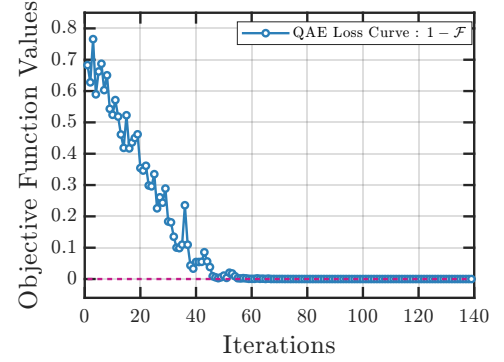


FIG. 7: Plot depicting the minimization of the QAE cost function (Equation 4) for optimal compression (from 2 qubits to 1 qubit) and reconstruction of entangled states (Equation 8a).

$$\vec{\theta}_E^* = (0.5\pi, \pi, 1.18\pi, 0.6\pi), \quad (10a)$$

$$\vec{\phi}_D^* = (0.4\pi, 0.5\pi, 0, 1.5\pi). \quad (10b)$$

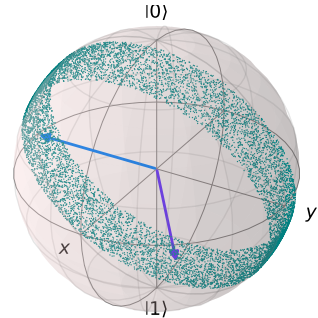


FIG. 8: 1-qubit latent space $\{|\gamma_K\rangle\}$ generated by the *trained encoder* upon compression of 2-qubit entangled states $\{\sigma_K\}$. The two states marked by Bloch vectors correspond to $K = \{0.5\pi, 0\}$ and $K = \{0.5\pi, 0.5\pi\}$.

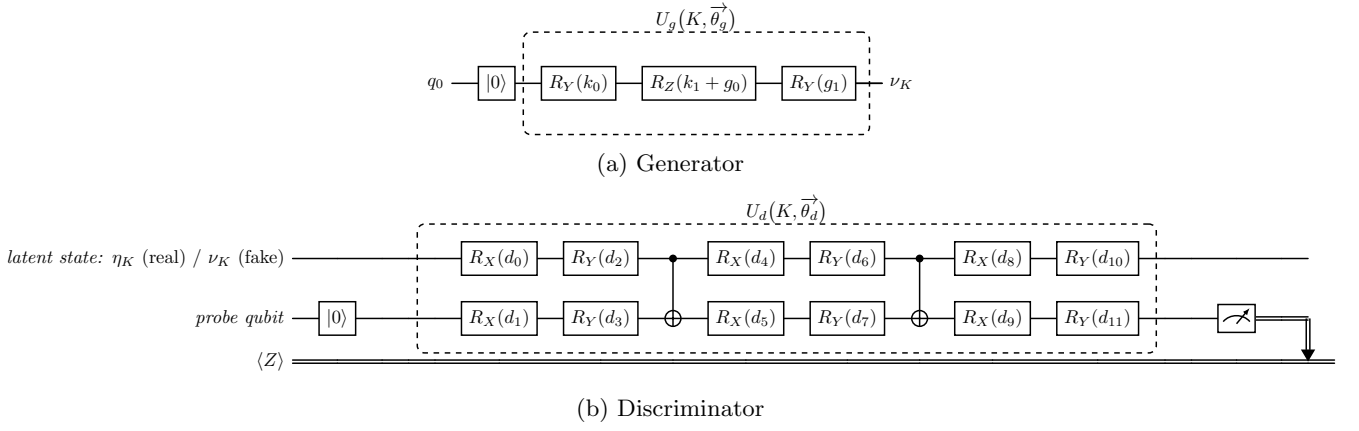


FIG. 9: Parameterized quantum circuits of (a) the generator and (b) the discriminator for the adversarial learning of 1-qubit latent representations of entangled states obtained using a trained QAE. The generator generates a 1-qubit state ν_K conditioned on the label information $K = (k_0, k_1)$, and (b) the discriminator evaluates the probability of the input state (either η_K or ν_K) being *real*. The parameters $\vec{\theta}_g$ and $\vec{\theta}_d$ of the generator and the discriminator respectively undergo min-max optimization during the adversarial training.

Upon closer numerical inspection of the latent states in Figure 8, it can be shown that the *trained encoder* $U_E(\vec{\theta}_E^*)$ effectively implements the following transformation on the state $|\chi_K\rangle$ in Equation 9a:

$$|\gamma_K\rangle = RY(-(\pi - 0.3)) RZ(\pi) |\chi_K\rangle. \quad (11)$$

Equation 11 demonstrates the non-uniqueness of the latent representation since $\eta_K = |\gamma_K\rangle\langle\gamma_K|$ obtained upon iteratively training the QAE is different than the 1-qubit representation $|\chi_K\rangle$ described in Equation 9a.

The two states $|\gamma_K\rangle$ corresponding to $k_0 = \pi/2$ and $k_1 = \pm\pi/2$ are eigenstates of the effective transformation in Equation 11. For the purpose of visualizing the adversarial learning protocol in the following discussion, it is helpful to keep track of the evolution of one of these eigenstates.

B. Adversarial Learning of QAE Latent Space

We now implement a QGAN as explained in Section III to learn the representation of the latent space $\{\eta_K\}$ depicted in Figure 8.

Taking insights from the effective transformation of the *trained encoder* (Equation 11), we design a suitable single-qubit Euler-decomposition ansatz,

$$U_g(K, \vec{\theta}_g) = RY(g_1) RZ(g_0 + k_1) RY(k_0), \quad (12)$$

as the generator with parameters $\vec{\theta}_g = (g_0, g_1)$.

Figure 9a depicts the circuit for the generator that encodes the label K and generates a *fake state* $\nu_K = |\nu_K\rangle\langle\nu_K|$. Figure 9b depicts the structure of the discriminator ansatz $U_d(K, \vec{\theta}_d)$ with parameters $\vec{\theta}_d = (d_0, \dots, d_{11})$.

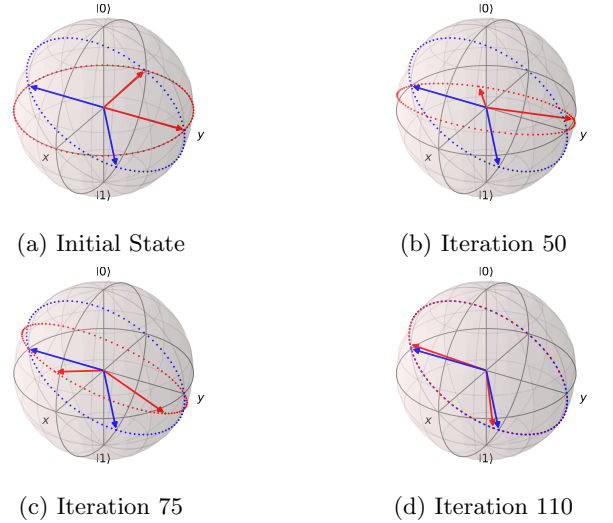


FIG. 10: Evolution of the learned latent space for $k[0] = \pi/2$ (red points) towards their target latent space (blue points) during adversarial training. Red points show states produced at a particular iteration of training $\vec{\theta}_g$, while blue points indicate the target states using optimal parameters $\vec{\theta}_g^*$. States corresponding to $k[1] = 0$ and $k[1] = \pi/2$ marked by Bloch vectors for ease of visualization.

Note: Ideally, the discriminator should be equipped with additional qubit(s) that encode the label information K conditioned on which the input state has been created [22]. This is to ensure that the generator does not *cheat* in the adversarial game by generating an input state of an incorrect label to confuse the discriminator. For our examples, we ensure that such a scenario of cheating in the adversarial game does not take place and that

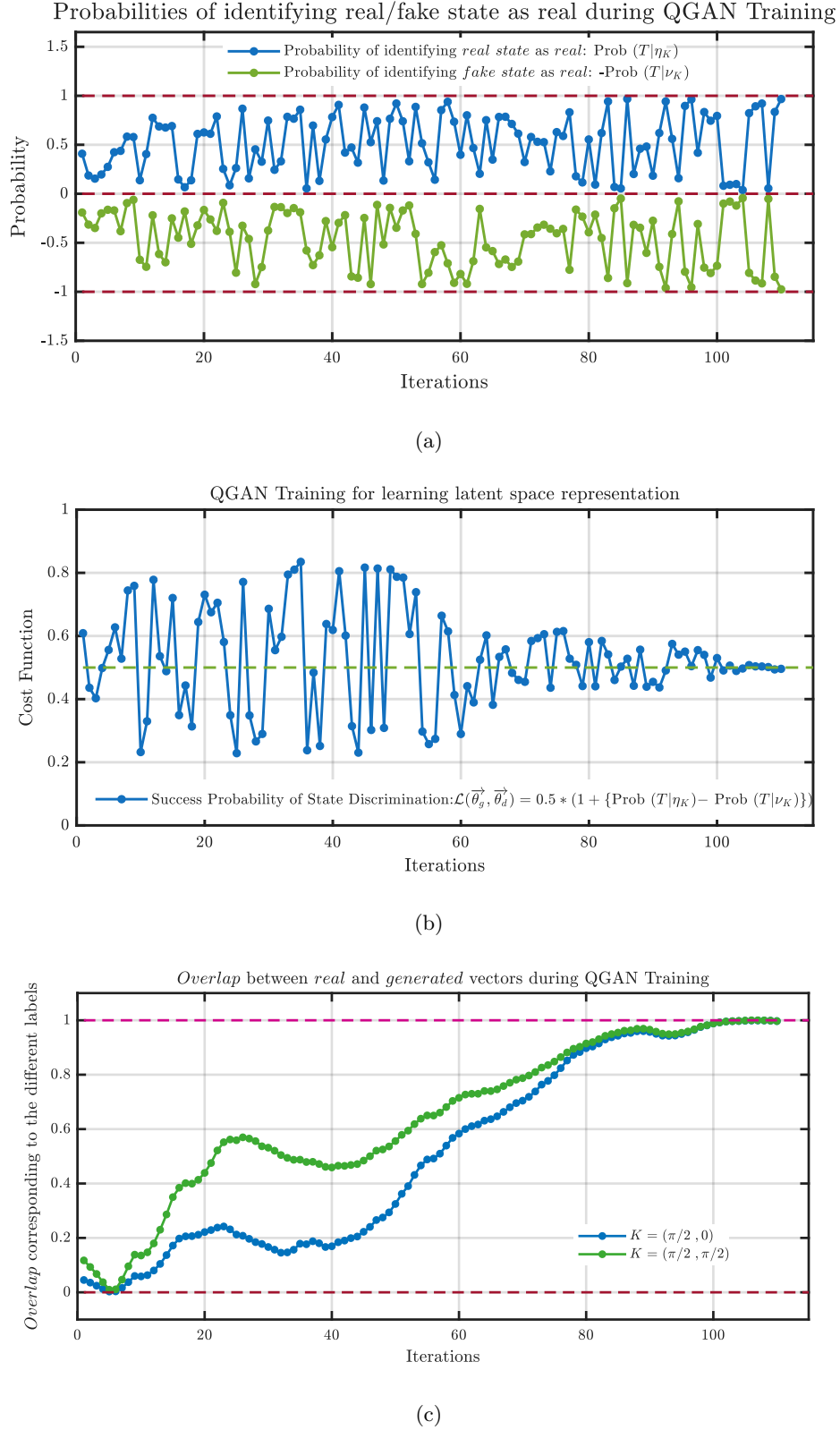


FIG. 11: (a) Probabilities of identifying the real/fake state as *real* during QGAN training. The curves clearly demonstrate adversarial behaviour. At the final 110th iteration of the training, the discriminator identifies both the real and the fake states as *real* with probability very close to 1. Therefore, at this stage, the generator is able to generate fake states which the discriminator cannot distinguish from real states; (b) This results in the cost function (Equation 7b) to converge close to 0.5, implying that the probability of the discriminator to distinguish between real and fake states is as good as a fair-coin toss; (c) Further, the *overlap* between the Bloch vectors of the real (η_k) and the generated (ν_k) state, corresponding to any label K , converges to 1 at the end of the QGAN training. This implies that the generator has indeed learnt the latent representation and can generate latent states $\nu_K = \eta_K$.

the information of the *actual label* is inherently present in the input state and hence implicitly supplied to the discriminator. This assumption reduces the qubit count of the discriminator and makes our instructive examples convenient for demonstrating adversarial learning.

The *trained encoder* (Figure 5a) with optimal parameters $\vec{\theta}_E^*$ (Equation 10a) serves as the source of *real data* $\{\eta_K\}$ for the discriminator, whereas the generator (Figure 9a) serves as the source of *fake data* $\{\nu_K\}$. A well-trained adversarial strategy should guide the parameters $\vec{\theta}_g$ of the generator toward their optimal value $\vec{\theta}_g^*$ where $\{\nu_K = \eta_K\}$. This would mean that the model has been able to learn and directly generate the latent space representation $\{\eta_K\}$ of the *trained* QAE.

Training QGANs is a challenging task and requires careful considerations. The initial choice of parameters also greatly determines the trajectory of QGAN training. Good starting points help the optimizer reach the required minimum much easily. After observing the QGAN training behaviour across multiple initializations, we choose the starting parameters as $\vec{\theta}_g = (\pi/2, 0)$ and $\vec{\theta}_d = \vec{0}$ with additional noise sampled from a standard normal distribution added to both. This injected noise ensures variability in the initialization, promoting more robust training. We use two Adam optimizers - one each for the generator and the discriminator respectively. The learning rate for the generator is chosen as 10^{-1} and that of the discriminator as 10^{-2} . The evolution of the learning of the latent space can be observed in Figure 10.

Figure 11 shows the evolution of various metrics over the course of the QGAN training to learn the latent space of the trained QAE. The adversarial behaviour is observed in the probability curves in Figure 11a. The training is stopped at the 110th iteration where the following two observations are simultaneous met:

1. the discriminator identifies both the *real* and the *fake* states as real with probability very close to 1 (Figure 11a).
2. the cost function $\mathcal{L}_{\text{QGAN}}$ converges close to 0.5 (Figure 11b) meaning that the discriminator can distinguish between the *real* and the *fake* states no better than a fair-coin toss.

Both these observable quantities together are useful to decide a stopping criteria for the QGAN training. Further, as a conclusive observation, Figures 10d, and 11c show that the overlap between the *real* and the *fake* states is indeed close to 1.0 which is in line with the observations in Figures 11a and 11b.

The trained generator now has the ability to reproduce the latent representation in Figure 8 by directly generating latent states $\{\nu_K = \eta_K\}$ conditioned on the label K using the following optimal parameters:

$$\vec{\theta}_g^* = (1.126\pi, -0.457\pi). \quad (13)$$

As depicted in Figure 3, a generated latent state ν_K can be processed using the *trained* decoder (with the optimal parameters $\vec{\phi}_D^*$ in Equation 10b) to generate 2-qubit entangled states $\xi_K = \sigma_K$ in Equation 8a.

This example was primarily motivated to showcase the adversarial learning aspect of this framework. The generative aspect of this framework can be appreciated better in the following example.

V. APPLICATION EXAMPLE 2: GENERATING GROUND STATE ENERGY PROFILES OF MOLECULAR HAMILTONIANS

This section investigates the utility of the proposed QGAA model to generate the ground state energy profiles of the parameterized molecular Hamiltonians of the Hydrogen molecule (H_2) and Lithium Hydride (LiH). The details of these Hamiltonians are provided in Appendix G. A molecular Hamiltonian is dependent upon various parameters such as the coordinates of the constituent nuclei and electrons. Each unique set of parameters corresponds to a specific molecular configuration, which in turn has an associated ground state. The fermionic Hamiltonians of both H_2 and LiH are parameterized by the interatomic distance r between their constituent atoms. When the fermionic Hamiltonians are mapped to the qubit or the spin versions, the parameter r is implicitly encoded in the real coefficients $c_i \equiv c_i(r)$ associated with each Pauli term in the spin Hamiltonian which in general has the following form:

$$\mathcal{H}(r) = \sum_i c_i P_i \quad (14)$$

where each $P_i \in \{I, X, Y, Z\}^{\otimes n}$ is a tensor product of n -Pauli operators and $c_i \in \mathbb{R}$.

The objective is to generate the ground state energy landscape $E(r)$ of the ground state of the target molecules as a function of the interatomic distance r , using a set of known ground states which serve as the training data for the QGAA model as depicted in Figure 12. We outline the experimental setup and describe the evaluation metrics employed for assessing the different components of the proposed model.

A. Training the QAE for Compression of Molecular Ground States

According to the QGAA formalism, a QAE is first trained to compress the training data into a lower-dimensional latent space. The sparsity in the density matrices, as observed in Figures 13a and 14a for the H_2 and LiH molecule respectively, captures the possibility of a reduced latent representation of the sparse molecular ground states using QAE [57].

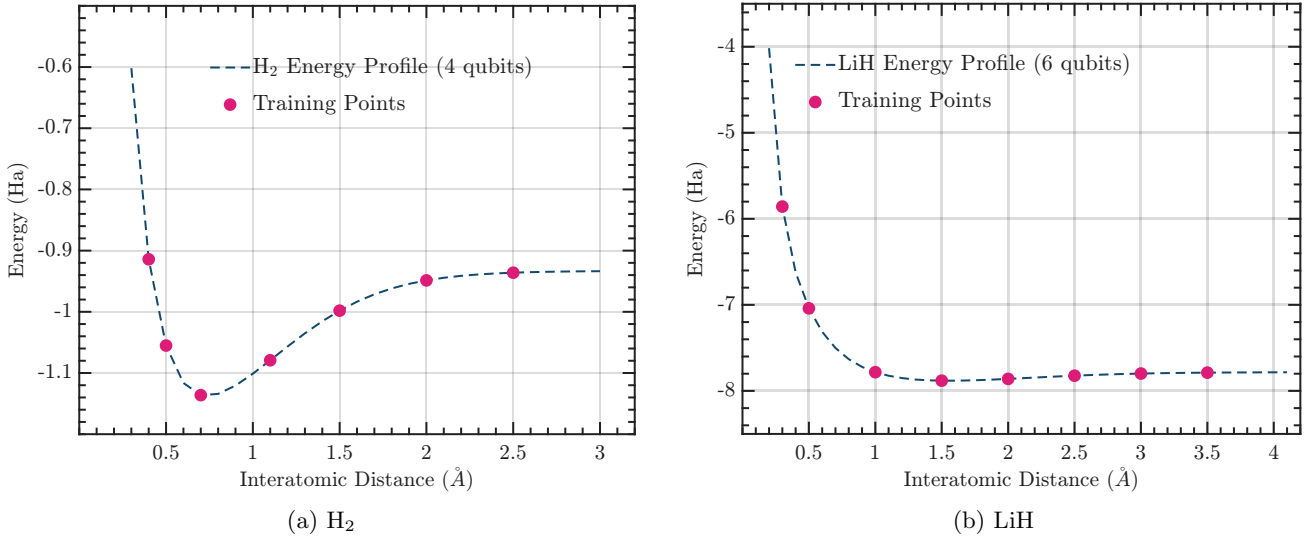


FIG. 12: Ground state energy profiles across different interatomic distances for the molecules (a) H₂ and (b) LiH. The red points indicate the ground states used as training data for the QAE. The QAE is trained to compress molecular ground states of H₂ from 4 qubits to 1 qubit, and of LiH from 6 qubits to 4 qubit. The performance metrics of the trained QAE are depicted in Table I.

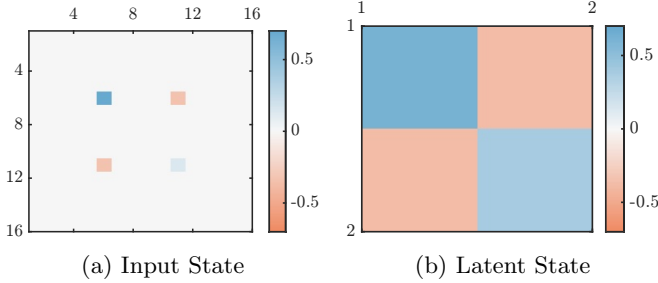


FIG. 13: Heatmaps of the real component of (a) the 4-qubit density matrix of the H₂ ground state at $r = 1.5$ Å and (b) its 1-qubit latent representation from a trained QAE.

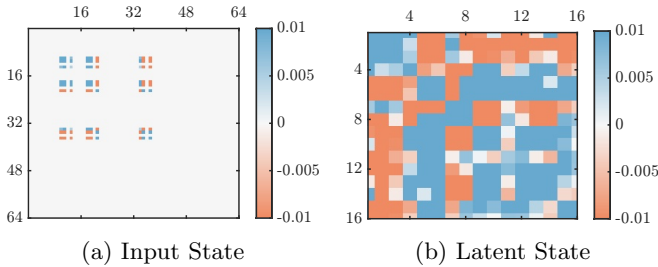


FIG. 14: Heatmaps of the real component of (a) the 6-qubit density matrix of the LiH ground state at $r = 1.5$ Å and (b) its 4-qubit latent representation from a trained QAE.

The training dataset $\{\sigma_r\}$ of the QAE comprises 7 ground states for the case of H₂ and 8 ground states for the case of LiH. Figure 12 depicts the true ground

state energy profiles of H₂ and LiH. The datapoints corresponding to the training ground states used to train the QAE are also marked in Figure 12. The label r denotes the interatomic distance for each data point. The ground states $\{\sigma_r\}$ used for training the QAE can be prepared on quantum hardware using algorithms such as the Variational Quantum Eigensolver (VQE) [51]. For our implementation, we selected the compression rates of the molecular ground states as 4 to 1 qubit for the case of H₂ and 6 to 4 qubits for LiH, as detailed in Appendix G.

The general structure of the ansatz used to implement the encoder and decoder comprises layers of R_X and R_Y rotation gates followed by entangling $CNOT$ gates. The exact form of these ansatzes is depicted in Figures 24 and 25, corresponding to H₂ and LiH respectively, in Appendices H and I.

The parameters of both the encoder and decoder circuits in the QAE are initialized randomly prior to training. To minimize the QAE cost function, the gradient-free COBYLA optimizer [53] is employed for the case of H₂, whereas for the gradient-based ADAM optimizer [36] is employed for the case of LiH.

Table I reports the performance metrics of the QAE averaged over the label r , namely the fidelity of reconstruction, $\mathcal{F}(\sigma_r, \rho_r)$, and the error in the energy, $|\Delta E(r)| = |\text{Tr}\{\mathcal{H}(r)\sigma_r\} - \text{Tr}\{\mathcal{H}(r)\rho_r\}|$. The results depict that the QAE ansatz structure used for the case of H₂ enables the reconstruction of the molecular ground states with high average fidelity, achieving $\mathcal{F}(\sigma_r, \rho_r) = 0.99$ and an energy reconstruction error of approximately $|\Delta E(r)| \sim 0.1$ mHa. However, in comparison, the performance of the QAE is suboptimal for the compression and reconstruction of the LiH molecular ground states

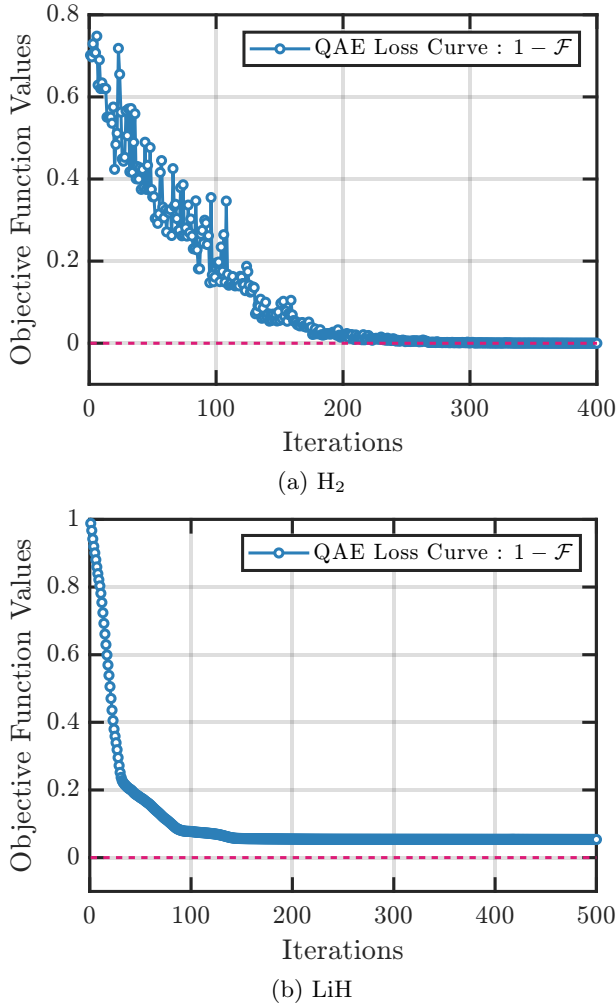


FIG. 15: Convergence of the QAE cost function (Equation 4) during training for optimal compression and reconstruction of the molecular ground states of H_2 and LiH . The H_2 QAE is optimized using the gradient-free COBYLA method, while the LiH QAE is trained using the ADAM optimizer.

even with an increased number of layers, with an average fidelity of $\mathcal{F}(\sigma_r, \rho_r) = 0.94 \pm 0.03$ and reconstruction error of $\langle |\Delta E(r)| \rangle = 0.02 \pm 0.01$. The heatmaps of the compressed latent state representations learned by the trained QAE at a bond length of $r = 1.5 \text{ \AA}$ and $r = 2.0 \text{ \AA}$ for H_2 and LiH are shown in Figures 13b and 14b, respectively. Figure 15 shows the convergence of the QAE cost function in Equation 4 during training, illustrating the optimization process for achieving optimal compression and reconstruction of the molecular ground states of H_2 and LiH .

The encoder of the trained QAE maps each input ground state to a corresponding latent representation. We show that the QAE learns pure-state latent representations of the form $\{\eta_r = |\gamma_r\rangle \langle \gamma_r|\}$ for the training states $\{\sigma_r\}$, and enables the reconstruction of the correspond-

Molecule	Compression Rate	$\langle \mathcal{F}(\sigma_r, \rho_r) \rangle$	$\langle \Delta E(r) \rangle$ (Hartree)
H_2	$4 \rightarrow 1$ qubits	0.99	~ 0.0001
LiH	$6 \rightarrow 4$ qubits	0.94 ± 0.03	0.02 ± 0.01

TABLE I: Average fidelity error over the ground states (\mathcal{F}) upon compression and reconstruction using the QAE trained on ground states of H_2 and LiH . The average absolute error in the energy of the reconstructed states is also reported.

ing output states $\{\rho_r\}$ from these latent representations. We note certain caveats in the Quantum Autoencoder formalism, which are discussed in detail in Appendix J. Additional experiments on the LiH Quantum Autoencoder are presented in Appendix K.

As described in Section III, the adversarial training formalism is leveraged in the following section to learn the latent space representation of the trained QAE to generate ground states $\{\xi_r\}$ beyond the training set.

B. Adversarial Learning of QAE Latent Space for Generation of Molecular Ground States

The latent states $\{\eta_r\}$ of the *trained* QAE are the *real data* for the QGAN whereas the states $\{\nu_r\}$ generated by the trainable generator conditioned on the label r are the *fake data*. Assuming that the chosen parametrized ansatz for the generator can approximate the required solution, the aim of the adversarial training is to find optimal parameters $\vec{\theta}_g^*$ such that the *trained* decoder transforms the generated state ν_r to a state ξ_r that closely approximates the ground state of the molecular Hamiltonian $\mathcal{H}(r)$.

The architectures employed for the generator in the cases of H_2 and LiH are shown in Figures 17a and 17b, respectively. The architecture used for the discriminator is depicted in Figure 18. The parameters of the generator and the discriminator are randomly initialized and the QGAN is trained in accordance with Algorithm 1. The adversarial training employs two separate ADAM optimizers [36] with learning rates of 0.1 and 0.01 for the generator and discriminator, respectively, for the case of H_2 . Whereas, for the case of LiH , learning rate schedulers are incorporated that reduce the generator's learning rate by a factor of 0.75 every 100 iterations and the discriminator's learning rate by the same factor every 250 iterations. This adaptive learning strategy is intended to improve convergence.

Table II reports the average fidelity of the generated states, $\mathcal{F}(\sigma_r, \xi_r)$, and the error in the energy, $|\Delta E(r)| = |\text{Tr}\{\mathcal{H}(r)\sigma_r\} - \text{Tr}\{\mathcal{H}(r)\xi_r\}|$, evaluated over the interatomic distance r . The learned ground state energy profile is shown for the H_2 molecule in Fig. 16a and for the LiH molecule in Fig. 16b. The results demonstrate

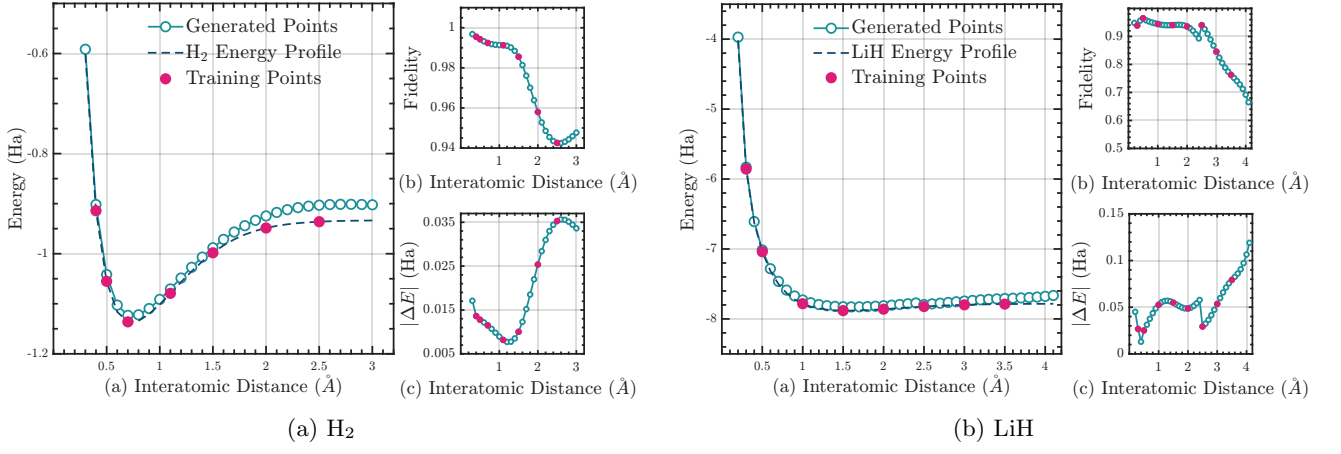


FIG. 16: Learned ground state energy profiles across different interatomic distances for the molecules (a) H_2 and (b) LiH . The red points indicate the ground states used as training data for the QGAN. The average fidelity of the generated states is 0.97 for H_2 and 0.88 for LiH . Panels (b) and (c) in each sub-figure show the fidelity and absolute error for the QGAN-generated states, respectively.

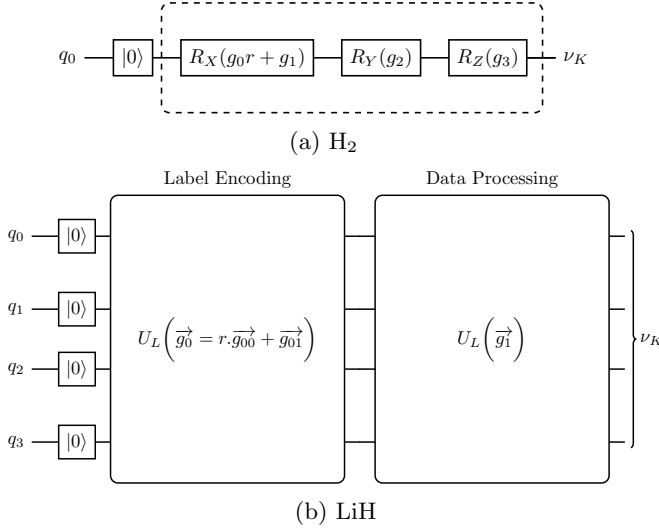


FIG. 17: Structure of the generator $U_g(K, \vec{\theta}_g)$ for (a) H_2 and (b) LiH - comprising a label encoding block followed by a data processing block [15] - with the internal structure of U_L detailed in Figure 26 in Appendix I.

that the trained QGAN closely approximates the true energy landscape, achieving an average state fidelity of $\mathcal{F}(\sigma_r, \xi_r) = 0.97 \pm 0.02$ and an energy reconstruction error of $|\Delta E(r)| = 0.02 \pm 0.01$ Ha for the case of H_2 . In contrast, the results for LiH are comparatively suboptimal, primarily due to limitations in the trained encoder, as reported in Table I. These discrepancies highlight the influence of both model architecture and training configuration in the QAE and QGAN frameworks on performance outcomes, especially in the context of scaling to larger molecular systems. We compare the QGAA with

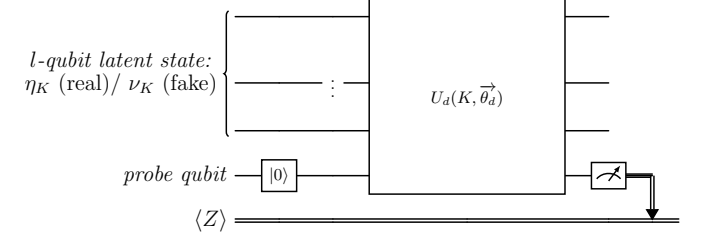


FIG. 18: Structure of the discriminator $U_d(K, \vec{\theta}_d)$, which distinguishes between real latent states (η_k) and fake latent states (ν_k). The ansatz structure to implement $U_d(K, \vec{\theta}_d)$ for the case of H_2 is the same as depicted in Figure 9b, and for the case of LiH is detailed in Figure 25 in Appendix I.

a baseline QGAN trained directly on the 4-qubit ground states of H_2 , with results presented in Appendix L, highlighting the utility of using the QAE to achieve more stable and efficient training of the QGAN.

Molecule	$\langle \mathcal{F}(\sigma_r, \xi_r) \rangle$	$\langle \Delta E(r) \rangle$ (Hartree)
H_2	0.97 ± 0.02	0.02 ± 0.01
LiH	0.88 ± 0.09	0.06 ± 0.02

TABLE II: Average fidelity error (\mathcal{F}) of the QGAA model, i.e., the QGAN trained on the compressed latent space of the ground states of H_2 and LiH . The average absolute error in the energy of the generated states is also reported.

VI. CONCLUSION AND FUTURE OUTLOOK

In this work, we proposed the formalism of the *Quantum Generative Adversarial Autoencoder* (QGAA) for generating quantum states with desired features.

By leveraging the quantum adversarial learning approach, we offer a way to directly access the latent space of a trained QAE, and thus impart the QAE with generative capabilities which it inherently lacks. Through the following illustrative examples, we demonstrated that the QGAA can successfully capture and reconstruct the quantum latent space of the QAE:

1. Learning the latent space of entangled states compressed by a QAE to analyze the adversarial training behaviour consistent with the theory of QGAN.
2. Energy profiling in quantum chemistry: learning the ground state energy landscape of the parameterized Hamiltonians, H_2 and LiH .

The examples implemented in Sections IV and V are demonstrative of the potential for applying QGAA to generative tasks where quantum state compression can be employed for efficient utilization of resources in QML algorithms.

In our examples, we demonstrate the use of a QGAN to learn the latent representation of the QAE, yielding two key benefits. First, it endows the QAE with generative capabilities by providing direct access to its latent space through the QGAN’s generator. Second, the inherent ability of the QAE to compress quantum data into a lower dimensional latent space enables a more resource efficient implementation of the QGAN. Thus, the QGAA has a bidirectional advantage of enhancing the QAE with generative functionality while simultaneously reducing the resource costs for a QGAN.

Notably, the training procedure only requires access to quantum states, and the resource requirements for training a QGAA on a compressed subspace are significantly lesser than that for training a QGAN on the original larger space of the quantum states.

Although the theory of quantum adversarial learning guarantees the existence of a unique Nash equilibrium, we encountered practical challenges in steering the model toward this optimal solution during training. Nevertheless, our results show that even with standard optimization techniques, it is possible to reasonably learn the latent space representation and generate quantum states that closely approximate the target states. In this context, various QGAN variants have been proposed in the literature, which could be integrated into our framework to stabilize training dynamics and improve optimization efficiency [16, 23, 24].

Even upon learning the latent representation approximately and not exactly, QGAA has utility in quantum machine learning applications such as warm-starting another quantum algorithm like the VQE with a better guess of the initial state [15]. The generated ground

states may serve as useful initial points that could potentially lead to more efficient convergence to the optimal solution [15].

We also highlighted the nuances and the assumptions in the QGAA framework and are explained in Appendices J and E. The framework of QGAA proposed in this work thus paves a new direction for exploring the efficient execution of quantum generative tasks.

VII. OPEN QUESTIONS AND CHALLENGES

As a QML model, the QGAA has several open questions and avenues for future investigation. The following are key aspects to be explored:

1. Scalability and Training: The implementations explored in this study were limited to small-scale quantum systems with the largest being 6 qubits for the LiH molecule. Scaling the QGAA to larger-qubit systems is challenging both for the QAE and the QGAN. Overcoming the issues posed by vanishing-gradients [8, 55] is critical for resource-efficient and stable training of QML models of increasing system sizes. Particularly for the QAE component, the SWAP test (employed to estimate the overlap of the reconstructed state) is computationally expensive to evaluate with scaling system size. The utility of methods such as approximate fidelity estimation using shadow tomography techniques can be explored to address the issues associated with fidelity-based loss functions [30].
2. Hardware implementation: While our experiments were conducted on simulators, deploying the proposed model on real quantum hardware introduces additional considerations such as gate noise, decoherence, and limited qubit connectivity. These hardware-induced imperfections can significantly affect the performance of both the QAE and QGAN components. Fine-tuning the model for useful performance on quantum devices requires the incorporation of noise-resilient circuit designs, error mitigation techniques, and hardware-aware optimization strategies to ensure reliable implementation under realistic conditions.

VIII. CODE AVAILABILITY

The code developed to obtain the results in this paper will be made available at a later date.

ACKNOWLEDGMENTS

The authors acknowledge the role of Fujitsu Research in supporting this project. We are also grateful to our

colleagues - Ruchira Bhat, Rahul Bhowmick, Harsh Wadhwa, Aritra Sarkar, and Quoc Hoan Tran for their insights and valuable feedback that helped develop this work. The authors extend their immense gratitude to Yasuhiro Endo, Hirotaka Oshima and Shintaro Sato, as

well as the entire Robust Quantum Computing Department at Fujitsu Limited for their strategic and technical support. A preliminary version of this paper was presented at the Fujitsu-IISc Quantum Workshop in Bangalore, India, and a poster presentation is scheduled for TQC 2025, to be held from September 15–19.

-
- [1] Tom Achache, Lior Hoeshe, and John Smolin. Denoising quantum states with quantum autoencoders – theory and applications. *ArXiv*, 2020. URL <https://arxiv.org/abs/2012.14714>.
 - [2] Mohammad H. Amin, Evgeny Andriyash, Jason Rolfe, Bohdan Kulchytskyy, and Roger Melko. Quantum boltzmann machine. *Phys. Rev. X*, 8:021050, May 2018. doi: 10.1103/PhysRevX.8.021050. URL <https://link.aps.org/doi/10.1103/PhysRevX.8.021050>.
 - [3] Richard Jozsa and. Fidelity for mixed quantum states. *Journal of Modern Optics*, 41(12):2315–2323, 1994. doi: 10.1080/09500349414552171. URL <https://doi.org/10.1080/09500349414552171>.
 - [4] Omar Bacarreza, Thorin Farnsworth, Alexander Makarovskiy, Hugo Wallner, Tessa Hicks, Santiago Sempere-Llagostera, John Price, Robert J. A. Francis-Jones, and William R. Clements. Quantum latent distributions in deep generative models. *ArXiv*, 2025. URL <https://arxiv.org/abs/2508.19857>.
 - [5] Dor Bank, Noam Koenigstein, and Raja Giryes. Autoencoders. *ArXiv*, 2020. URL <https://api.semanticscholar.org/CorpusID:260535551>.
 - [6] Adriano Barenco, André Berthiaume, David Deutsch, Artur Ekert, Richard Jozsa, and Chiara Macchiavello. Stabilization of quantum computations by symmetrization. *SIAM Journal on Computing*, 26(5):1541–1557, 1997. doi:10.1137/S0097539796302452. URL <https://doi.org/10.1137/S0097539796302452>.
 - [7] Ruchira V Bhat, Rahul Bhowmick, Avinash Singh, and Krishna Kumar Sabapathy. Meta-learning of gibbs states for many-body hamiltonians with applications to quantum boltzmann machines. *ArXiv*, 2025. URL <https://arxiv.org/abs/2507.16373>.
 - [8] Rahul Bhowmick, Harsh Wadhwa, Avinash Singh, Tania Sidana, Quoc Hoan Tran, and Krishna Kumar Sabapathy. Enhancing variational quantum algorithms by balancing training on classical and quantum hardware. *ArXiv*, 2025. URL <https://arxiv.org/abs/2503.16361>.
 - [9] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd. Quantum machine learning. *Nature*, 549(7671):195–202, 2017. doi: 10.1038/nature23474. URL <https://doi.org/10.1038/nature23474>.
 - [10] Dmytro Bondarenko and Polina Feldmann. Quantum autoencoders to denoise quantum data. *Phys. Rev. Lett.*, 124:130502, Mar 2020. doi: 10.1103/PhysRevLett.124.130502. URL <https://link.aps.org/doi/10.1103/PhysRevLett.124.130502>.
 - [11] Paolo Braccia, Filippo Caruso, and Leonardo Banchi. How to enhance quantum generative adversarial learning of noisy information. *New Journal of Physics*, 23(5): 053024, may 2021. doi:10.1088/1367-2630/abf798. URL <https://dx.doi.org/10.1088/1367-2630/abf798>.
 - [12] Harry Buhrman, Richard Cleve, John Watrous, and Ronald de Wolf. Quantum fingerprinting. *Phys. Rev. Lett.*, 87:167902, Sep 2001. doi: 10.1103/PhysRevLett.87.167902. URL <https://link.aps.org/doi/10.1103/PhysRevLett.87.167902>.
 - [13] Chenfeng Cao and Xin Wang. Noise-assisted quantum autoencoder. *Phys. Rev. Appl.*, 15:054012, May 2021. doi:10.1103/PhysRevApplied.15.054012. URL <https://link.aps.org/doi/10.1103/PhysRevApplied.15.054012>.
 - [14] M. Cerezo, A. Arrasmith, R. Babbush, S. C. Benjamin, S. Endo, K. Fujii, J. R. McClean, K. Mitarai, X. Yuan, L. Cincio, and P. J. Coles. Variational quantum algorithms. *Nature Reviews Physics*, 3(9):625–644, 2021. doi: 10.1038/s42254-021-00348-9.
 - [15] Alba Cervera-Lierta, Jakob S. Kottmann, and Alán Aspuru-Guzik. Meta-variational quantum eigensolver: Learning energy profiles of parameterized hamiltonians for quantum simulation. *PRX Quantum*, 2:020329, May 2021. doi:10.1103/PRXQuantum.2.020329.
 - [16] Shouvanik Chakrabarti, Yiming Huang, Tongyang Li, Soheil Feizi, and Xiaodi Wu. *Quantum wasserstein GANs*. Curran Associates Inc., Red Hook, NY, USA, 2019. URL <https://dl.acm.org/doi/abs/10.5555/3454287.3454896>.
 - [17] P Chamorro-Posada and J C Garcia-Escartin. The switch test for discriminating quantum evolutions. *Journal of Physics A: Mathematical and Theoretical*, 56(35):355301, aug 2023. doi:10.1088/1751-8121/acecc5. URL <https://dx.doi.org/10.1088/1751-8121/acecc5>.
 - [18] Chuangtao Chen, Qinglin Zhao, MengChu Zhou, Zhimin He, Zhili Sun, and Haozhen Situ. Quantum generative diffusion model: A fully quantum-mechanical model for generating quantum state ensemble. *ArXiv*, 2024. URL <https://arxiv.org/abs/2401.07039>.
 - [19] El Cherrat, Iordanis Kerenidis, Natansh Mathur, Jonas Landman, Martin Strahm, and Yun Li. Quantum vision transformers. *Quantum*, 8:1265, 02 2024. doi:10.22331/q-2024-02-22-1265.
 - [20] Brendan Coyle, Daniel Mills, Vincent Danos, and Elham Kashefi. The born supremacy: quantum advantage and training of an ising born machine. *npj Quantum Information*, 6(1):60, 2020. doi:10.1038/s41534-020-00288-9. URL <https://doi.org/10.1038/s41534-020-00288-9>.
 - [21] M. Cramer, M. B. Plenio, S. T. Flammia, R. Somma, D. Gross, S. D. Bartlett, O. Landon-Cardinal, D. Poulin, and Y.-K. Liu. Efficient quantum state tomography. *Nature Communications*, 1:149, 2010. doi: 10.1038/ncomms1147. URL <https://doi.org/10.1038/ncomms1147>.
 - [22] Pierre-Luc Dallaire-Demers and Nathan Killoran. Quantum generative adversarial networks. *Phys. Rev. A*,

- 98:012324, Jul 2018. doi:10.1103/PhysRevA.98.012324. URL <https://link.aps.org/doi/10.1103/PhysRevA.98.012324>.
- [23] Giacomo De Palma, Milad Marvian, Dario Trevisan, and Seth Lloyd. The quantum wasserstein distance of order 1. *IEEE Transactions on Information Theory*, 67(10):6627–6643, 2021. doi:10.1109/TIT.2021.3076442.
- [24] Giacomo De Palma, Tristan Klein, and Davide Pastorello. Classical shadows meet quantum optimal mass transport. *J. Math. Phys.*, 65(9):092201, 2024. doi:10.1063/5.0178897.
- [25] Raghavendra M Devadas and Sowmya T. Quantum machine learning: A comprehensive review of integrating ai with quantum computing for computational advancements. *MethodsX*, 14:103318, 2025. ISSN 2215-0161. doi:<https://doi.org/10.1016/j.mex.2025.103318>. URL <https://www.sciencedirect.com/science/article/pii/S2215016125001645>.
- [26] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Commun. ACM*, 63(11):139–144, October 2020. ISSN 0001-0782. doi:10.1145/3422622. URL <https://doi.org/10.1145/3422622>.
- [27] C. W. Helstrom. Quantum detection and estimation theory. *Journal of Statistical Physics*, 1(2):231–252, 1969. doi:10.1007/BF01007479. URL <https://doi.org/10.1007/BF01007479>.
- [28] Ling Hu, Shu-Hao Wu, Weizhou Cai, Yuwei Ma, Xianghao Mu, Yuan Xu, Haiyan Wang, Yipu Song, Dong-Ling Deng, Chang-Ling Zou, and Luyan Sun. Quantum generative adversarial learning in a superconducting quantum circuit. *Science Advances*, 5(1):eaav2761, 2019. doi:10.1126/sciadv.aav2761. URL <https://www.science.org/doi/abs/10.1126/sciadv.aav2761>.
- [29] Chang-Jiang Huang, Hailan Ma, Qi Yin, Jun-Feng Tang, Daoyi Dong, Chunlin Chen, Guo-Yong Xiang, Chuan-Feng Li, and Guang-Can Guo. Realization of a quantum autoencoder for lossless compression of quantum data. *Phys. Rev. A*, 102:032412, Sep 2020. doi:10.1103/PhysRevA.102.032412. URL <https://link.aps.org/doi/10.1103/PhysRevA.102.032412>.
- [30] Hsin-Yuan Huang, Richard Kueng, and John Preskill. Predicting many properties of a quantum system from very few measurements. *Nature Physics*, 16:1050–1057, 2020. doi:10.1038/s41567-020-0932-7. URL <https://doi.org/10.1038/s41567-020-0932-7>.
- [31] K. Huang, Z. A. Wang, C. Song, W. Zhang, H. Li, Y. Liu, H. Liu, L. Sun, G. Guo, and G. Guo. Quantum generative adversarial networks with multiple superconducting qubits. *npj Quantum Information*, 7:165, 2021. doi:10.1038/s41534-021-00503-1. URL <https://doi.org/10.1038/s41534-021-00503-1>.
- [32] Igor O. Sokolov. IBM Quantum Challenge 2021 - Exercise 5 solution, 2021. URL <https://github.com/qiskit-community/ibm-quantum-challenge-2021/blob/main/solutions%20by%20authors/ex5/ex5-solution.ipynb>.
- [33] Yuichi Kamata, Quoc Hoan Tran, Yasuhiro Endo, and Hirotaka Oshima. Molecular quantum transformer. *ArXiv*, 2025. URL <https://arxiv.org/abs/2503.21686>.
- [34] Amir Khoshaman, Walter Vinci, Brandon Denis, Evgeny Andriyash, Hossein Sadeghi, and Mohammad H Amin. Quantum variational autoencoder. *Quantum Science and Technology*, 4(1):014001, sep 2018. doi:10.1088/2058-9565/aada1f. URL <https://dx.doi.org/10.1088/2058-9565/aada1f>.
- [35] Leeseok Kim, Seth Lloyd, and Milad Marvian. Hamiltonian quantum generative adversarial networks. *Phys. Rev. Res.*, 6:033019, Jul 2024. doi:10.1103/PhysRevResearch.6.033019. URL <https://link.aps.org/doi/10.1103/PhysRevResearch.6.033019>.
- [36] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. URL <https://api.semanticscholar.org/CorpusID:6628106>.
- [37] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013. URL <https://api.semanticscholar.org/CorpusID:216078090>.
- [38] Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Found. Trends Mach. Learn.*, 12(4):307–392, November 2019. ISSN 1935-8237. doi:10.1561/22000000056. URL <https://doi.org/10.1561/22000000056>.
- [39] Alistair Letcher, Stefan Woerner, and Christa Zoufal. Tight and Efficient Gradient Bounds for Parameterized Quantum Circuits. *Quantum*, 8:1484, September 2024. ISSN 2521-327X. doi:10.22331/q-2024-09-25-1484. URL <https://doi.org/10.22331/q-2024-09-25-1484>.
- [40] Yeong-Cherng Liang, Yu-Hao Yeh, Paulo E M F Mendonça, Run Yan Teh, Margaret D Reid, and Peter D Drummond. Quantum fidelity measures for mixed states. *Reports on Progress in Physics*, 82(7):076001, jun 2019. doi:10.1088/1361-6633/ab1ca4. URL <https://dx.doi.org/10.1088/1361-6633/ab1ca4>.
- [41] Seth Lloyd and Christian Weedbrook. Quantum generative adversarial learning. *Phys. Rev. Lett.*, 121:040502, Jul 2018. doi:10.1103/PhysRevLett.121.040502. URL <https://link.aps.org/doi/10.1103/PhysRevLett.121.040502>.
- [42] H. Ma, G. J. Mooney, I. R. Petersen, C. Ferrie, and R. Harper. Quantum autoencoders using mixed reference states. *npj Quantum Information*, 10:86, 2024. doi:10.1038/s41534-024-00872-3. URL <https://doi.org/10.1038/s41534-024-00872-3>.
- [43] H. Ma, G. J. Mooney, I. R. Petersen, T. Tighe, A. van der Reest, S. Gorman, and S. C. Benjamin. Quantum autoencoders using mixed reference states. *npj Quantum Information*, 10:86, 2024. doi:10.1038/s41534-024-00872-3. URL <https://doi.org/10.1038/s41534-024-00872-3>.
- [44] Hailan Ma, Chang-Jiang Huang, Chunlin Chen, Daoyi Dong, Yuanlong Wang, Re-Bing Wu, and Guo-Yong Xiang. On compression rate of quantum autoencoders: Control design, numerical and experimental realization. *Automatica*, 147:110659, 2023. ISSN 0005-1098. doi:<https://doi.org/10.1016/j.automatica.2022.110659>. URL <https://www.sciencedirect.com/science/article/pii/S0005109822005234>.
- [45] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *ArXiv*, 2015. URL <https://arxiv.org/abs/1511.05644>.
- [46] Wai-Keong Mok, Hui Zhang, Tobias Haug, Xianshu Luo, Guo-Qiang Lo, Zhenyu Li, Hong Cai, M. S. Kim, Ai Qun Liu, and Leong-Chuan Kwek. Rigorous noise

- reduction with quantum autoencoders. *AVS Quantum Science*, 6(2):023803, 05 2024. ISSN 2639-0213. doi: 10.1116/5.0192456. URL <https://doi.org/10.1116/5.0192456>.
- [47] Murphy Yuezhen Niu, Alexander Zlokapa, Michael Broughton, Sergio Boixo, Masoud Mohseni, Vadim Smelyanskiy, and Hartmut Neven. Entangling quantum generative adversarial networks. *Phys. Rev. Lett.*, 128:220505, Jun 2022. doi: 10.1103/PhysRevLett.128.220505. URL <https://link.aps.org/doi/10.1103/PhysRevLett.128.220505>.
- [48] A. D. Patel. The quantum density matrix and its many uses. *Journal of the Indian Institute of Science*, 103: 401–417, 2023. doi:10.1007/s41745-023-00406-4. URL <https://doi.org/10.1007/s41745-023-00406-4>.
- [49] Saahil Patel, Benjamin Collis, William Duong, Daniel Koch, Massimiliano Cutugno, Laura Wessing, and Paul Alsing. Information loss and run time from practical application of quantum data compression. *Physica Scripta*, 98(4):045111, mar 2023. doi:10.1088/1402-4896/acc492. URL <https://dx.doi.org/10.1088/1402-4896/acc492>.
- [50] Alex Pepper, Nora Tischler, and Geoff J. Pryde. Experimental realization of a quantum autoencoder: The compression of qutrits via machine learning. *Phys. Rev. Lett.*, 122:060501, Feb 2019. doi: 10.1103/PhysRevLett.122.060501. URL <https://link.aps.org/doi/10.1103/PhysRevLett.122.060501>.
- [51] A. Peruzzo, J. McClean, P. Shadbolt, M.-H. Yung, X.-Q. Zhou, P. J. Love, A. Aspuru-Guzik, and J. L. O’Brien. A variational eigenvalue solver on a photonic quantum processor. *Nature Communications*, 5:4213, 2014. doi: 10.1038/ncomms5213. URL <https://doi.org/10.1038/ncomms5213>.
- [52] M. Pivoluska and M. Plesch. Implementation of quantum compression on ibm quantum computers. *Scientific Reports*, 12:5841, 2022. doi:10.1038/s41598-022-09881-8. URL <https://doi.org/10.1038/s41598-022-09881-8>.
- [53] M. J. D. Powell. *A Direct Search Optimization Method That Models the Objective and Constraint Functions by Linear Interpolation*, pages 51–67. Springer Netherlands, Dordrecht, 1994. ISBN 978-94-015-8330-5. doi: 10.1007/978-94-015-8330-5_4. URL https://doi.org/10.1007/978-94-015-8330-5_4.
- [54] Qiskit Nature Development Team. Qiskit Ecosystem: Mapping to the qubit space, 2024. URL https://qiskit-community.github.io/qiskit-nature/tutorials/06_qubit_mappers.html.
- [55] M. Ragone, B. N. Bakalov, F. Sauvage, M. Rosenkranz, F. Ticozzi, G. Carleo, and M. Dalmonte. A lie algebraic theory of barren plateaus for deep parameterized quantum circuits. *Nature Communications*, 15:7172, 2024. doi:10.1038/s41467-024-49909-3. URL <https://doi.org/10.1038/s41467-024-49909-3>.
- [56] J. M. Renes. Mark m. wilde: Quantum information theory. *Quantum Information Processing*, 13(3):587–590, 2014. doi:10.1007/s11128-013-0690-x. URL <https://doi.org/10.1007/s11128-013-0690-x>.
- [57] Jonathan Romero, Jonathan P Olson, and Alan Aspuru-Guzik. Quantum autoencoders for efficient compression of quantum data. *Quantum Science and Technology*, 2(4): 045001, aug 2017. doi:10.1088/2058-9565/aa8072. URL <https://dx.doi.org/10.1088/2058-9565/aa8072>.
- [58] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In *Parallel Distributed Processing, Volume 1: Explorations in the Microstructure of Cognition: Foundations*. The MIT Press, 07 1986. ISBN 9780262291408. doi: 10.7551/mitpress/5236.003.0012. URL <https://doi.org/10.7551/mitpress/5236.003.0012>.
- [59] Sandeep Singh Sengar, Affan Bin Hasan, Sanjay Kumar, and Fiona Carroll. Generative artificial intelligence: a systematic review and applications. *Multimedia Tools and Applications*, 84(21):23661–23700, 2025. doi: 10.1007/s11042-024-20016-1. URL <https://doi.org/10.1007/s11042-024-20016-1>.
- [60] Sukin Sim, Peter D. Johnson, and Alán Aspuru-Guzik. Expressibility and entangling capability of parameterized quantum circuits for hybrid quantum-classical algorithms. *Advanced Quantum Technologies*, 2(12):1900070, 2019. doi:https://doi.org/10.1002/qute.201900070. URL <https://advanced.onlinelibrary.wiley.com/doi/abs/10.1002/qute.201900070>.
- [61] Jinkai Tian, Xiaoyu Sun, Yuxuan Du, Shanshan Zhao, Qing Liu, Kaining Zhang, Wei Yi, Wanrong Huang, Chaoyue Wang, Xingyao Wu, Min-Hsiu Hsieh, Tongliang Liu, Wenjing Yang, and Dacheng Tao. Recent advances for quantum neural networks in generative learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):12321–12340, 2023. doi: 10.1109/TPAMI.2023.3272029.
- [62] Quoc Hoan Tran, Shinji Kikuchi, and Hirotaka Oshima. Variational denoising for variational quantum eigensolver. *Phys. Rev. Res.*, 6:023181, May 2024. doi:10.1103/PhysRevResearch.6.023181. URL <https://link.aps.org/doi/10.1103/PhysRevResearch.6.023181>.
- [63] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning*, ICML ’08, page 1096–1103, New York, NY, USA, 2008. Association for Computing Machinery. ISBN 9781605582054. doi:10.1145/1390156.1390294. URL <https://doi.org/10.1145/1390156.1390294>.
- [64] Adit Vishnu, Abhay Shastri, Dhruva Kashyap, and Chiranjib Bhattacharyya. Do-em: Density operator expectation maximization. *ArXiv*, 2025. URL <https://arxiv.org/abs/2507.22786>.
- [65] Yunfei Wang and Junyu Liu. A comprehensive review of quantum machine learning: from nisq to fault tolerance. *Reports on Progress in Physics*, 87(11):116402, oct 2024. doi:10.1088/1361-6633/ad7f69. URL <https://dx.doi.org/10.1088/1361-6633/ad7f69>.
- [66] Xiao-Ming Zhang, Weicheng Kong, Muhammad Usman Farooq, Man-Hong Yung, Guoping Guo, and Xin Wang. Generic detection-based error mitigation using quantum autoencoders. *Phys. Rev. A*, 103:L040403, Apr 2021. doi: 10.1103/PhysRevA.103.L040403. URL <https://link.aps.org/doi/10.1103/PhysRevA.103.L040403>.
- [67] Feifei Zhou, Yu Tian, Yumeng Song, Chudan Qiu, Xiangyu Wang, Mingti Zhou, Bing Chen, Nanyang Xu, and Dawei Lu. Preserving entanglement in a solid-spin system using quantum autoencoders. *Applied Physics Letters*, 121(13):134001, 09 2022. ISSN 0003-6951. doi: 10.1063/5.0120060. URL <https://doi.org/10.1063/5.0120060>.

Appendix A: Autoencoders

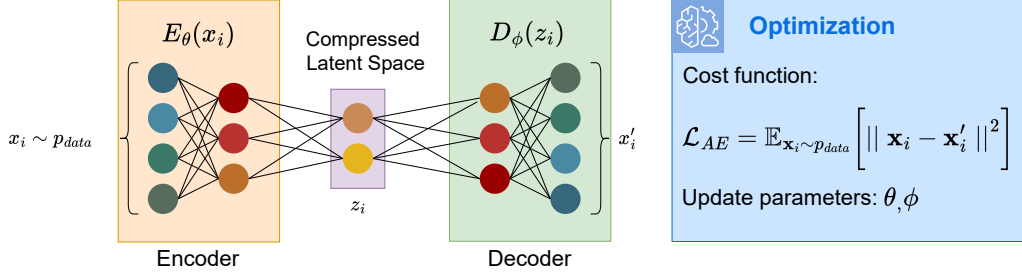


FIG. 19: Schematic of the Autoencoder (AE). The training set is denoted as $\mathbf{X}_{\text{train}} = \{\mathbf{x}_i\}_{i=1}^n$, where each input $\mathbf{x}_i \in \mathbb{R}^d$ is a d -dimensional vector sampled from the data distribution p_{data} . The encoder $E_{\theta} : \mathbb{R}^d \mapsto \mathbb{R}^m$, with parameters θ , compresses each input into a latent vector $\mathbf{z}_i \in \mathbb{R}^m$ with $m < d$. The decoder $D_{\phi} : \mathbb{R}^m \mapsto \mathbb{R}^d$, with parameters ϕ , maps \mathbf{z}_i back to the reconstructed output $\mathbf{x}'_i \in \mathbb{R}^d$. Training minimizes the reconstruction loss $\mathcal{L}_{\text{AE}} = \mathbb{E}_{\mathbf{x}_i \sim p_{\text{data}}} [\|\mathbf{x}_i - \mathbf{x}'_i\|^2]$, ensuring that \mathbf{x}'_i remains close to \mathbf{x}_i . While effective for compression, the AE does not possess generative capability.

An autoencoder (AE) [58] is a machine learning architecture designed to compress input data into a latent representation, as depicted schematically in Figure 19. AE has been applied in various domains such as dimensionality reduction, classification, anomaly detection, and denoising [63]. Consider a training set $\mathbf{X}_{\text{train}} = \{\mathbf{x}_i\}_{i=1}^n$ sampled from some data distribution p_{data} , where n is the number of training points and each $\mathbf{x}_i \in \mathbb{R}^d$ is a d -dimensional input vector. The architecture consists of two classical neural networks. The encoder $E_{\theta} : \mathbb{R}^d \mapsto \mathbb{R}^m$ with trainable parameters θ maps an input vector \mathbf{x}_i to a compressed latent representation $\mathbf{z}_i = E_{\theta}(\mathbf{x}_i) \in \mathbb{R}^m$, where $m < d$ ensures compression. The decoder $D_{\phi} : \mathbb{R}^m \mapsto \mathbb{R}^d$ with parameters ϕ maps this latent vector back to the input space, reconstructing the original data point as $\mathbf{x}'_i = D_{\phi}(\mathbf{z}_i)$. Training is carried out by minimizing the reconstruction loss, defined as the average squared error between the input and its reconstruction,

$$\min_{\theta, \phi} \mathcal{L}_{\text{AE}} := \mathbb{E}_{\mathbf{x}_i \sim p_{\text{data}}} [\|\mathbf{x}_i - \mathbf{x}'_i\|^2], \quad (\text{A1})$$

which ensures that the reconstructed outputs $\mathbf{x}'_i = D_{\phi}(E_{\theta}(\mathbf{x}_i))$ remain close to the original inputs \mathbf{x}_i . The AE is a deterministic model that effectively compresses data into a lower-dimensional latent space but lacks generative capability. The Variational Autoencoder (VAE) [37, 38], explained in the Appendix C, extends this framework by enabling both compression and the generation of new samples not present in the training set $\mathbf{X}_{\text{train}}$.

Appendix B: Metric Used to evaluate the QAE Cost function

The SWAP test is a standard quantum subroutine used to estimate the overlap between two quantum states σ_K, ρ_K . Given two states, the quantity of interest is the Hilbert–Schmidt inner product

$$\text{SWAP}(\sigma_K, \rho_K) := \text{Tr}(\sigma_K \rho_K). \quad (\text{B1})$$

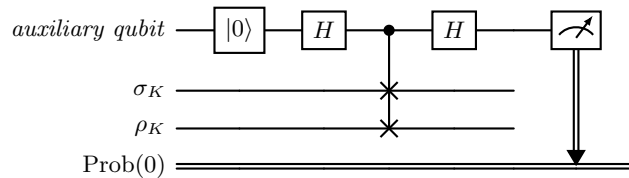


FIG. 20: SWAP Test for evaluating overlap between the quantum states σ_K and ρ_K . Prob(0) is the probability of measuring the auxiliary qubit in the state $|0\rangle$.

The protocol, illustrated in Figure 20, introduces an auxiliary qubit, initialized in the state $|0\rangle$, which controls a SWAP operation between the two registers containing σ_K and ρ_K . After applying Hadamard gates before and

after the controlled-SWAP, the auxiliary qubit is measured in the computational basis. The probability of observing outcome $|0\rangle$ is given by,

$$\text{Prob}(0) = \frac{1}{2}(1 + \text{Tr}(\sigma_K \rho_K)). \quad (\text{B2})$$

Rearranging, one obtains

$$\text{SWAP}(\sigma_K, \rho_K) = \text{Tr}(\sigma_K \rho_K) = 2\text{Prob}(0) - 1. \quad (\text{B3})$$

Thus, the SWAP test provides an operational way to estimate the overlap between two density operators by repeated sampling of the ancilla measurement outcome.

Fidelity: For two general (possibly mixed) states ρ_k and σ_k , the *Uhlmann–Jozsa fidelity* [3, 40] is

$$\mathcal{F}(\sigma_k, \rho_k) = \left(\text{Tr} \sqrt{\sqrt{\sigma_k} \rho_k \sqrt{\sigma_k}} \right)^2, \quad (\text{B4})$$

which reduces to $|\langle \psi_k | \phi_k \rangle|^2$ for two pure states and to $\langle \psi_k | \rho_k | \psi_k \rangle$ when $\sigma_k = |\psi_k\rangle \langle \psi_k|$ is pure.

Key properties:

- $0 \leq \mathcal{F} \leq 1$
- $\mathcal{F} = 1$ if and only if the states are identical
- Fidelity is symmetric: $\mathcal{F}(\rho_k, \sigma_k) = \mathcal{F}(\sigma_k, \rho_k)$
- Fidelity is invariant under unitary transformation U such that $\mathcal{F}(\rho_k, \sigma_k) = \mathcal{F}(U\rho_k U^\dagger, U\sigma_k U^\dagger)$.

Appendix C: Variational Autoencoders

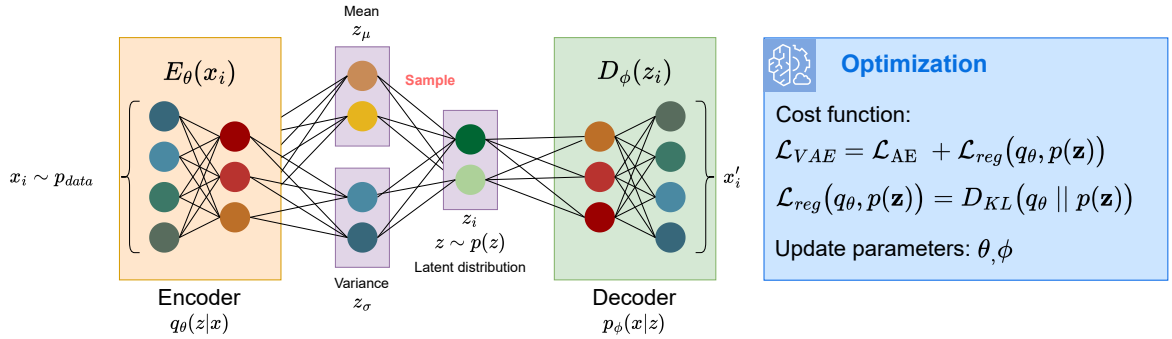


FIG. 21: Schematic of the Variational Autoencoder (VAE). An input datapoint $\mathbf{x}_i \in \mathbb{R}^d$ from the training set is mapped by the encoder E_θ into the parameters of a latent probability distribution $q_\theta(\mathbf{z}|\mathbf{x}_i)$ over latent variables $\mathbf{z}_i \in \mathbb{R}^\ell$, where $\ell < d$. Unlike a deterministic Autoencoder, which encodes each input into a single latent vector, the VAE samples \mathbf{z}_i from this distribution. The latent distribution is regularized to match a chosen prior $p(\mathbf{z})$, typically the standard normal $\mathcal{N}(\mathbf{0}, \mathbb{I})$, ensuring that the latent space is continuous and well-structured. The decoder D_ϕ then maps a sampled latent vector \mathbf{z}_i back to the data space, producing a reconstruction $\mathbf{x}'_i = D_\phi(\mathbf{z}_i)$. This probabilistic formulation enables the VAE to both reconstruct training data and generate novel data points \mathbf{x}'_i by sampling $\mathbf{z}_i' \sim p(\mathbf{z})$.

Unlike an Autoencoder (AE), which is a deterministic model that maps input data into a discrete latent representation, a Variational Autoencoder (VAE) [37, 38] is a probabilistic model, as shown schematically in Figure 21. Instead of encoding the latent variables \mathbf{z} of training data as fixed points, the VAE encodes them as a probability distribution. Specifically, the encoder outputs an approximate posterior distribution $q_\theta(\mathbf{z}|\mathbf{x})$, rather than a single latent vector, while the prior distribution over latent variables is denoted as $p(\mathbf{z})$. This ensures that the latent space has a continuous probabilistic structure. Once trained, the model can sample \mathbf{z} from $p(\mathbf{z})$ and generate new datapoints through the decoder. The architecture of a VAE resembles that of an AE, consisting of an encoder E_θ that maps an input $\mathbf{x}_i \in \mathbf{X}_{\text{train}}$ to a latent distribution $q_\theta(\mathbf{z}|\mathbf{x}_i)$, and a decoder D_ϕ that reconstructs or generates data as $\mathbf{x}'_i = D_\phi(\mathbf{z}_i)$, where \mathbf{z}_i is sampled from $q_\theta(\mathbf{z}|\mathbf{x}_i)$. Training involves not only minimizing the reconstruction loss \mathcal{L}_{AE} (Equation A1),

but also regularizing the latent distribution $q_{\theta}(\mathbf{z}|\mathbf{x})$ so that it approximates a known prior $p(\mathbf{z})$, typically chosen as the standard normal distribution $\mathcal{N}(\mathbf{0}, \mathbb{I})$. The resulting loss function is

$$\min_{\theta, \phi} \mathcal{L}_{\text{VAE}} := \mathcal{L}_{\text{AE}} + \mathcal{L}_{\text{reg}}(q_{\theta}(\mathbf{z}|\mathbf{x}), p(\mathbf{z})), \quad (\text{C1})$$

where the regularization term \mathcal{L}_{reg} measures the divergence between the approximate posterior $q_{\theta}(\mathbf{z}|\mathbf{x})$ and the prior $p(\mathbf{z})$, defined using the Kullback–Leibler (KL) divergence as

$$\mathcal{L}_{\text{reg}}(q_{\theta}(\mathbf{z}|\mathbf{x}), p(\mathbf{z})) = D_{\text{KL}}(q_{\theta}(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})). \quad (\text{C2})$$

At convergence, $q_{\theta}(\mathbf{z}|\mathbf{x})$ aligns with $p(\mathbf{z})$, ensuring that the latent space is continuous and structured. As a result, new latent samples $\mathbf{z}_{i'} \sim p(\mathbf{z})$, which do not correspond to any encoded training datapoint, can be passed through the trained decoder to generate novel datapoints $\mathbf{x}_{i'} = D_{\phi}(\mathbf{z}_{i'})$. Thus, unlike deterministic AE which is limited to reconstructing training data, VAE possesses true generative capability by directly sampling from the latent space to create new data beyond the training set.

Appendix D: Generative Adversarial Networks

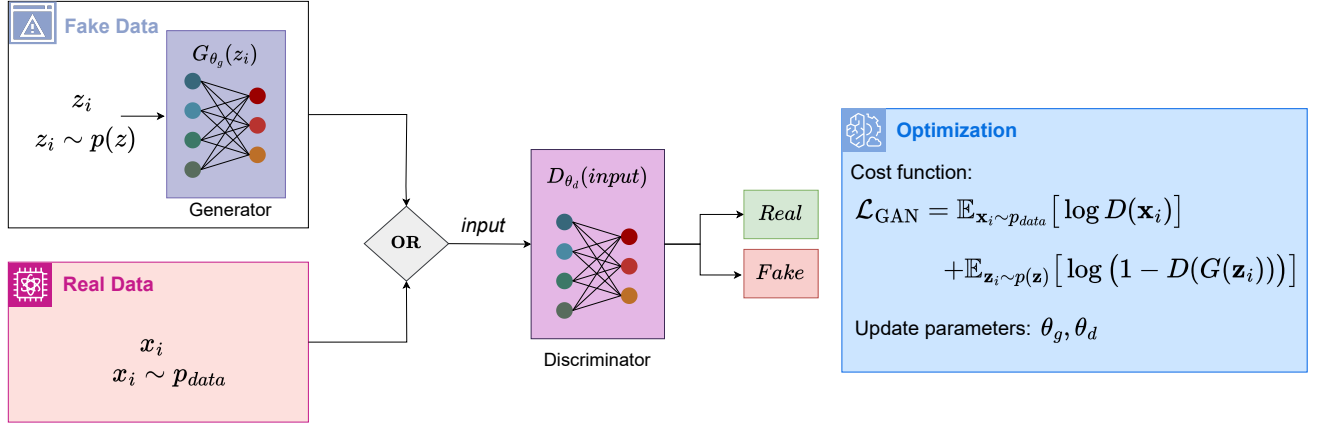


FIG. 22: Schematic of the Generative Adversarial Network (GAN). The min-max optimization of \mathcal{L}_{GAN} by the generator and the discriminator characterizes the adversarial training. At the Nash equilibrium of the training, the generator generates *fake* data samples $G(\mathbf{z}_i)$ that the discriminator cannot distinguish from the *real* data samples $\mathbf{x}_i \sim p_{\text{data}}(\mathbf{x})$.

A Generative Adversarial Network (GAN) [26] is a generative framework where two neural networks — a generator G and a discriminator D with parameters θ_g and θ_d respectively — are trained in an adversarial manner, as shown schematically in Figure 22. Conditioned on some prior entropy source $p(\mathbf{z})$, the objective of G is to produce samples of synthetic data $G(\mathbf{z}_i \sim p(\mathbf{z}))$ that is indistinguishable from samples of the training data $\mathbf{x}_i \sim p_{\text{data}}(\mathbf{x})$. On the other hand, the objective of D is to correctly identify with high probability whether the input data supplied to it is from the real distribution p_{data} or from the prior $p(\mathbf{z})$ to the generated/fake distribution. In other words, G generates synthetic data, and D outputs has a probabilistic binary outcome.

The way the training is setup is such that G and D compete against each other in an adversarial game until they reach the Nash equilibrium or the fixed-point. In general, $D(\cdot)$ outputs a score corresponding to the input data and the maximizing objective of $D(\cdot)$ is to assign a higher score to real data samples $\mathbf{x}_i \sim p_{\text{data}}$ and a lower score to generated/fake data samples $G(\mathbf{z}_i)$ where $\mathbf{z}_i \sim p(\mathbf{z})$. The minimizing objective of $G(\cdot)$ is to generate samples $G(\mathbf{z}_i)$ such that $D(\cdot)$ assigns them a higher score. This min-max optimization of the two competing objectives is captured by the cost function $\mathcal{L}_{\text{GAN}}(G, D)$ in Equation D1.

$$\min_{\theta_g} \max_{\theta_d} \mathcal{L}_{\text{GAN}} := \mathbb{E}_{\mathbf{x}_i \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x}_i)] + \mathbb{E}_{\mathbf{z}_i \sim p(\mathbf{z})} [\log (1 - D(G(\mathbf{z}_i)))] \quad (\text{D1})$$

The Nash equilibrium of this adversarial game is reached when the features of the generated data are close to that of the real data, at which point the discriminator can no longer discriminate between the real and generated/fake data

[26].

Appendix E: Caveats in the Quantum Generative Adversarial Networks formalism

The assumptions mentioned in Section II B with regards to the generator, G , and the discriminator, D , are detailed further as follows :

1. *First, it is assumed that both the quantum information processors G and D have enough capacity [47] or expressibility [60] such that they can approximate any arbitrary function or transformation using their parameters $\vec{\theta}_g$ and $\vec{\theta}_d$ respectively.*

In practice, this requirement is assumed to be fulfilled if the structure of G and D is a deep quantum feature map that is *expressive enough* for G and D to achieve their respective optimal strategies. These feature maps are often chosen heuristically and typically comprise blocks of rotation gates followed by entangling gates. It can be reasoned that ansatzes with such a structure are *highly versatile* [48] and hence are able to create a large variety of feature maps based on different parameterizations, some of which can approximate the desired transformation. This ansatz structure can be easily generalized to higher qubits and can be made *deep* or *shallow* based on the number of repeating blocks composed. However, the effectiveness of the chosen ansatz remains dependent on the problem and the physical realization of the gates.

2. *Second, it is assumed that an efficient optimization scheme exists that can drive the algorithm to converge to the unique Nash equilibrium.*

Although the theory of Helstrom measurement helps in setting up the quantum adversarial game, it does not give any guarantees on the convergence to the Nash equilibrium. Analytically, at a given $(i)th$ iteration of the game the *optimal strategy* of D is to optimize $\vec{\theta}_d$ such that operator $\hat{T}^{(i)}$ projects onto the positive eigenspace of $(\sigma - \rho^{(i-1)})$ thus maximizing $\mathcal{L}_{\text{QGAN}}$ [56]. In practice, this is difficult to achieve since the state σ is not known apriori. That is the exact reason why the adversarial protocol is being executed so that the statistics of σ could be learned. However, let's assume that D can somehow achieve this particular analytical solution in practice. Then in the next $(i+1)th$ iteration, the analytical *optimal strategy* of G is to optimize $\vec{\theta}_g$ such that the form of the density matrix of the state $\rho^{(i+1)}$ is a pure state that projects onto the largest eigenvalue of $\hat{T}^{(i)}$ thus minimizing $\mathcal{L}_{\text{QGAN}}$. This, however, is not useful for learning mixed-states in general [11].

Such *exact optimization* at every iteration can potentially lead to *mode collapse* or oscillations between some quantum states [47] and *limit cycles* [11]. One approach towards addressing *mode collapse* is by using the square of $\mathcal{L}_{\text{QGAN}}$ as the cost function instead [35]. Other efforts include the development of quantum extensions of the classical Wasserstein GAN [16, 23, 24].

Further, deep feature maps are argued to be powerful enough to approximate arbitrary transformations, and besides the QGAN-specific challenges discussed so far, vanishing gradients or barren plateaus is another significant challenge in training such deep feature maps [39]. Although a unique fixed-point exists in the quantum adversarial game [41], efficient optimization techniques to converge to the fixed-point is a challenging and open problem. Studies training fully-quantum GANs to generate quantum states are limited to 1 – 3 qubits both on simulators [11, 22, 35, 47] and real devices [28, 31].

Appendix F: Adversarial Autoencoder

The original framework of the VAE requires that the form of $p(\mathbf{z})$ be known, as mentioned previously in Appendix C. $p(\mathbf{z})$ can be quantified and described by measurable quantities of the distribution like mean and variance in the case of a normal distribution. In practice, these quantities are evaluated during training (Figure 21) to regularize the latent distribution $q_\theta(\mathbf{z}|\mathbf{x})$ to the known distribution $p(\mathbf{z})$. However, the Adversarial Autoencoder (AAE) [45] was developed to show that the adversarial formalism of a GAN can also be employed to regularize the latent distribution $q_\theta(\mathbf{z}|\mathbf{x})$ with the target distribution $p(\mathbf{z})$, as shown schematically in Figure 23. Here the encoder E_θ acts as the generator that generates a sample $E_\theta(\mathbf{x}_i) \sim q_\theta(\mathbf{z}|\mathbf{x})$ of the latent space given input \mathbf{x}_i . Another neural network D is used as a discriminator to distinguish between samples from $p(\mathbf{z})$ (real) and $q_\theta(\mathbf{z}|\mathbf{x})$ (fake). The advantage of using a GAN

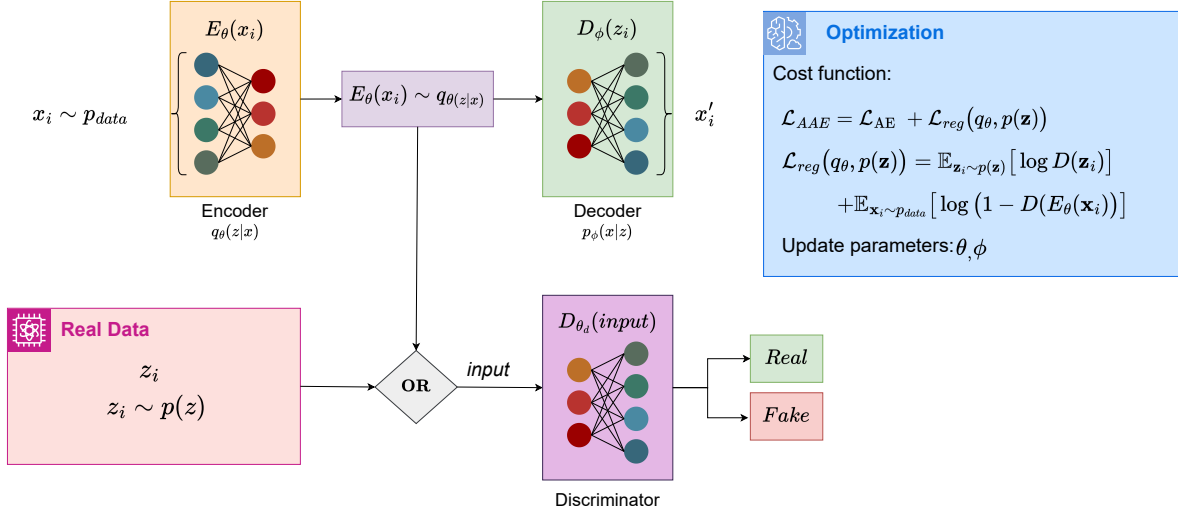


FIG. 23: Schematic of the Adversarial Autoencoder (AAE). Unlike the VAE in Figure 21, the regularization of the latent space in this case is performed using adversarial training instead of evaluating the KL-divergence.

for VAE latent space regularization is that the complete form of the target distribution $p(\mathbf{z})$ is not required to be known. Access to $p(\mathbf{z})$ with the ability to sample from it is enough to regularize $q_\theta(\mathbf{z}|\mathbf{x})$ using the adversarial protocol. Therefore $\mathcal{L}_{reg}(q_\theta(\mathbf{z}|\mathbf{x}), p(\mathbf{z}))$ (Equation C2) takes the form represented in Equation D1 as depicted in Equation F1:

$$\min_{E_\theta} \max_D \mathcal{L}_{reg}(q_\theta(\mathbf{z}|\mathbf{x}), p(\mathbf{z})) := \mathbb{E}_{\mathbf{z}_i \sim p(\mathbf{z})} [\log D(\mathbf{z}_i)] + \mathbb{E}_{\mathbf{x}_i \sim p_{data}(\mathbf{x})} [\log (1 - D(E_\theta(\mathbf{x}_i)))] \quad (\text{F1})$$

Appendix G: Molecular Hamiltonians

1. Hydrogen molecule: H_2

The two hydrogen atoms in the H_2 molecule are separated by interatomic distance r and each atom is associated with a $1s$ atomic orbital. Therefore, the H_2 molecule has two molecular orbitals. The two molecular orbitals correspond to four spin orbitals and hence can be described by 4 qubits. We use the Jordan-Wigner fermion-to-qubit mapping to arrive at the general form of the qubit hamiltonian (Equation G1) corresponding to the H_2 molecule [54].

$$\begin{aligned} H(r) = & c_0 I + c_1(Z_0 + Z_2) + c_2(Z_1 + Z_3) + c_3(Z_1 Z_0 + Z_3 Z_2) \\ & + c_4(Z_2 Z_0) + c_5(Z_3 Z_0 + Z_1 Z_2) + c_6 Z_3 Z_1 \\ & + c_7(Y_3 Y_2 Y_1 Y_0 + X_3 X_2 Y_1 Y_0 + Y_3 Y_2 X_1 X_0 + X_3 X_2 X_1 X_0) \end{aligned} \quad (\text{G1})$$

The coefficients $\{c_i\}$ in Equation G1 are a function of the interatomic distance r . The constant nuclear repulsion energy term is contained in the constant factor c_0 . However, using alternative mapping techniques that exploit the symmetries in the structure of $H(r)$, it is possible for $H(r)$ to have a simple 1-qubit representation as well: $H(r) = d_0 I_0 + d_1 Z_0 + d_2 X_0$ [54]. Therefore, H_2 system can be simulated by using just 1-qubit for finding the ground state energy profile (blue curve in Figure 12a).

Alternatively, one can also observe that the density matrices of the 4-qubit ground states of H_2 have a spare 1-qubit representation (an example depicted in Figure 13a) and can therefore be compressed using a QAE. The sparsity in the ground state density matrix is due to the symmetries in the H_2 Hamiltonian (Equation G1), allowing for a compression to a single qubit.

2. Lithium Hydride: LiH

In the Lithium Hydride (LiH) molecule, the electron from the H atom occupies the $1s$ atomic orbital, while the LiH atom contributes three electrons: one in the $n = 1$ shell (1 orbital) and two in the $n = 2$ shell (4 orbitals). This results in a total of $1 + (1 + 4) = 6$ molecular orbitals, which correspond to 12 spin orbitals. Consequently, when applying a fermion-to-qubit mapping such as the Jordan-Wigner transformation, simulating the LiH molecule requires 12 qubits.

However, the size of this large system can be reduced to 6 qubits by the employing following strategies:

1. The core electrons in the $1s$ shell of Li do not contribute significantly to the chemistry of the LiH molecule. Therefore only the contribution from the valence shell of Li can be considered. This eliminates 2 qubits corresponding to the core $1s$ orbital.
2. The LiH molecule has rotational symmetry about one of the axes (say the z -axis). Therefore, along with freezing the inner $1s$ orbital, the contributions from its p_x and p_y orbitals can also be eliminated. This results in the further reduction of $(1 + 1) = 2$ molecular orbitals or 4 spin orbitals (or qubits).

Further, the Parity fermion-to-qubit mapping can be applied with two-qubit reduction to bring down the qubit count to 4 qubits with certain loss in accuracy in the energy estimated. Different sequence of approximations and fermion-to-qubit mapping based on symmetries can also be applied to arrive at the 4 qubit reduction [32]. Further elimination of qubits is not possible and this can be verified computationally [32].

Appendix H: Ansatzes used for the example of H_2

In this appendix, we present the structure of parametrized quantum circuits employed for the compression and reconstruction of the molecular ground states of H_2 . The choice of ansatz plays a critical role in determining the accuracy and efficiency of variational quantum algorithms. The specific ansatz circuits used in our simulations are illustrated in Figure 24.

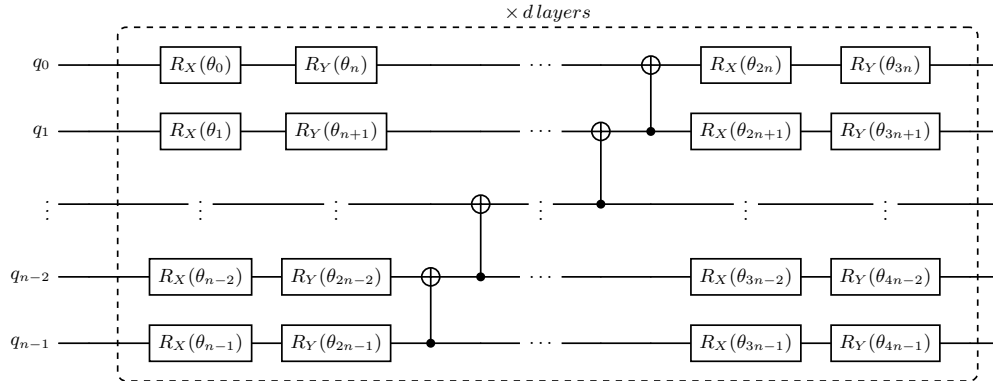


FIG. 24: Structure of the ansatz used to implement the encoder/decoder to compress/reconstruct molecular ground states of H_2 . Qubits 0 to $l - 1$ of such an ansatz for both the encoder and the decoder refer to the l -qubits of the latent state. The remaining l to $n - 1$ qubits refer to the *trash/auxiliary* qubits corresponding to the encoder/decoder. Each layer of this n -qubit ansatz of this structure has $4n$ tunable parameters for the rotation gates and $n - 1$ CX gates. For the case of H_2 , $l = 1$, $n = 4$ and $d = 1$.

Appendix I: Ansatzes used for the example of LiH

In this appendix, we present the quantum circuit ansatzes employed for the simulation of the Lithium Hydride (LiH) molecule. Figure 25 shows the structure of the ansatz used to implement the encoder and decoder for compressing and reconstructing the molecular ground states of LiH. In addition, Figure 26 illustrates the ansatz structure used to implement $U_L(\theta_L)$ in Figure 17b.

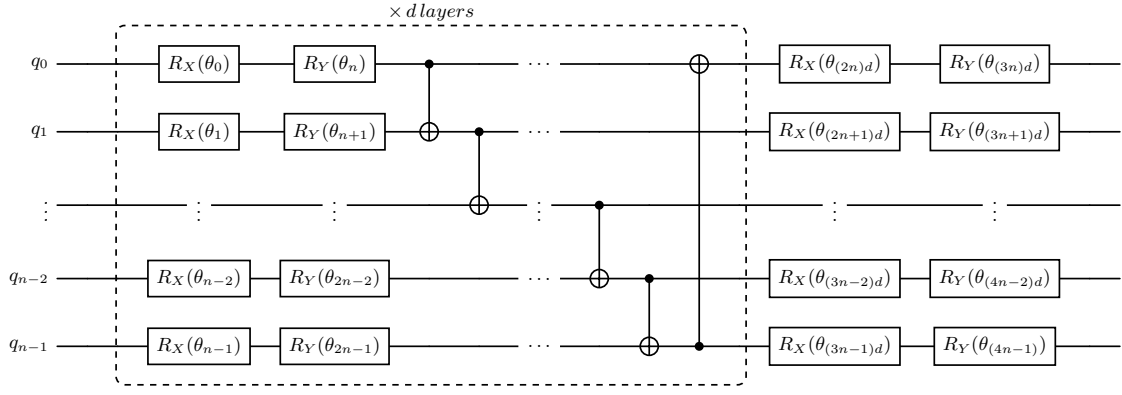


FIG. 25: Structure of the ansatz used to implement the encoder/decoder to compress/reconstruct molecular ground states of LiH. Qubits 0 to $l - 1$ of such an ansatz for both the encoder and the decoder refer to the l -qubits of the latent state. The remaining l to $n - 1$ qubits refer to the *trash/auxiliary* qubits corresponding to the encoder/decoder. Each layer of this n -qubit ansatz of this structure has $2n$ tunable parameters for the rotation gates and n circular entangling CX gates. For the encoder/decoder for the case of LiH, $l = 4$, $n = 6$ and $d = 3$. The same ansatz structure is also used to implement the discriminator $U_D(K, \vec{\theta}_D)$ in Figure 18 with $n = l + 1$ and $d = 3$.

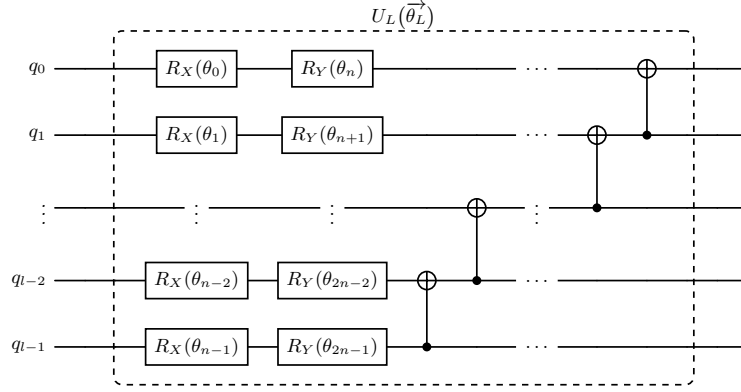


FIG. 26: The ansatz structure used to implement $U_L(\vec{\theta}_L)$ in Figure 17b. Each block of $U_L(\vec{\theta}_L)$ has $2l$ tunable parameters for the rotation gates and $l - 1$ CX gates. For the case of LiH, $l = 4$.

Appendix J: Caveats in the Quantum Autoencoder formalism

The compression capabilities of QAEs are subject to fundamental constraints such as structural properties, entanglement, and entropy of the input states. Theoretical investigations have demonstrated that QAEs can achieve lossless compression of high-dimensional quantum states into a lower-dimensional latent space [29, 44]. A necessary condition for such lossless compression is that the number of linearly independent input states does not exceed the dimensionality of the target latent space [44].

Despite their potential, QAEs face several challenges in scalability and implementation. Current experimental demonstrations are restricted to small systems due to hardware limitations such as noise, shallow circuit depths, and limited qubit connectivity, all of which reduce reconstruction fidelity [43, 49, 50, 52]. QAE performance is sensitive to design choices including circuit depth, compression ratio, and the choice of parameterized unitaries [49]. Additionally, large-scale deployment of QAEs may require deep quantum circuits and exponential classical resources for tasks like fidelity or entropy estimation (e.g., via SWAP tests [12] or full state tomography [21]), resulting in significant overhead that exceeds current NISQ-era capabilities. Classical shadow methods offer a potential means to partially alleviate these challenges by enabling more efficient estimations; however, their effectiveness in this context has yet to be thoroughly tested [30]. Table III summarizes existing implementations on quantum hardware and simulators, along with their corresponding compression rates.

	References	Compression Rate
HARDWARE	Realization of a quantum autoencoder for lossless compression of quantum data. [29]	$2 \rightarrow 1$
	Experimental Realization of a Quantum Autoencoder: The Compression of Qutrits via Machine Learning. [50]	Qutrits \rightarrow Qubits
	Implementation of quantum compression on IBM quantum computers. [52]	$3 \rightarrow 2$
	Information loss and run time from practical application of quantum data compression. [49]	$4 \rightarrow 1$
	Quantum autoencoders using mixed reference states. [42]	$4 \rightarrow 2$
SIMULATION	Variational Denoising for Variational Quantum Eigensolver. [62]	Up to 14 qubits
	Quantum autoencoders for efficient compression of quantum data. [57]	Up to 8 qubits

TABLE III: Summary of Quantum Autoencoder Studies and Their Compression Capabilities

Appendix K: Additional QAE experiments on LiH

We additionally consider a necessary condition for perfect quantum state compression: the number of linearly independent input states must not exceed the dimensionality of the target latent space [29, 44]. In the case of pure quantum states, each individual state is rank-one. However, when dealing with a dataset comprising multiple pure states, the effective rank is determined by the span of these states, which corresponds to the rank of the density matrix formed by the ensemble of pure states. In contrast, if the dataset contains mixed states with some prior distribution, the relevant rank is that of the averaged density matrix describing the ensemble. This rank captures both the intrinsic mixedness of each state and the statistical mixture induced by the prior, and it sets the minimum latent dimension required for lossless compression.

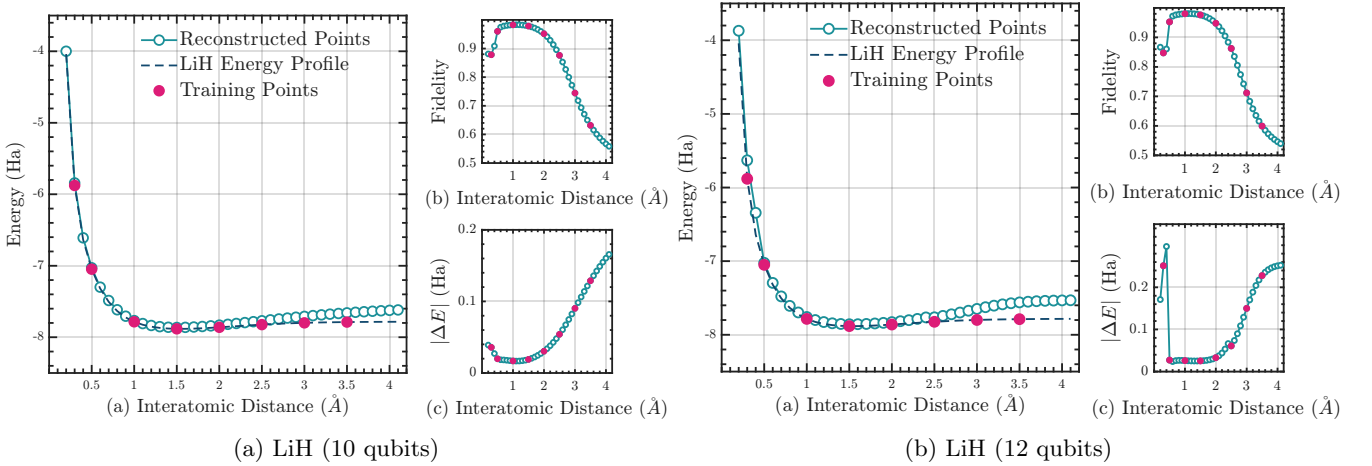


FIG. 27: Ground state energy profile across different interatomic distances of the molecules. The red points denote the ground states that were selected as training data for the QAE. The key takeaway is that, by learning from these representative configurations, the QAE can ideally reduce the ground states of LiH molecules from 12 or 10 qubits to a significantly smaller latent space while retaining essential physical information.

Considering the LiH molecule, where the dataset is generated by varying the interatomic distance from 0.2 to 4.2 Hartree, in steps of 0.1 Hartree. Each configuration yields a ground state, represented as a pure state. Assuming an equal probability distribution over all configurations, the ensemble can be used to construct a density matrix. The rank of this density matrix reflects the number of linearly independent pure states in the dataset. We compute the rank of the input state ensemble for Hamiltonians of 12, 10, and 6 qubits, obtaining ranks of 18, 7, and 6, respectively.

These results suggest that the latent space can be represented using 5 qubits for the 12-qubit Hamiltonian, and 3 qubits for both the 10- and 6-qubit Hamiltonians for lossless compression. While this approach provides a guideline for selecting the number of qubits in the latent space, the exact computation of the rank scales exponentially with the size of the input Hilbert space. This renders the method computationally intensive and challenging to scale for larger systems.

Compressing LiH ground State: We employ QAE to compress the ground states of the LiH molecule for systems initialized with 12 and 10 qubits. Table IV and Figure 27 present the reconstructed energy profiles along with the corresponding fidelity and reconstruction error for both cases.

Molecule	Compression Rate	$\langle \mathcal{F}(\sigma_r, \rho_r) \rangle$	$\langle \Delta E(r) \rangle$
LiH	10 \rightarrow 4 qubits	0.842 ± 0.148	0.062 ± 0.049 Ha
LiH	12 \rightarrow 4 qubits	0.826 ± 0.159	0.111 ± 0.093 Ha

TABLE IV: Performance Metrics of QAE for compression and reconstruction of Molecular Ground States.

Appendix L: QGAA vs QGAN: Training the QGAN directly on the 4 qubit ground states of H_2

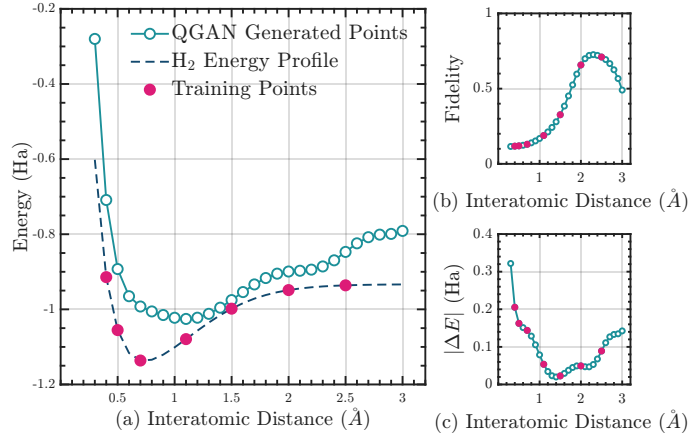


FIG. 28: Learned energy profiles after adversarial training on 4 qubits input space (without compression using QAE). Each plot also includes the fidelity and reconstruction error of the generated states as a function of interatomic distance. The QGAN learns the energy profile whereas the QGAA achieves better precision.

In this section, we compare the performance of the Quantum Generative Adversarial Network (QGAN) when trained directly on the 4-qubit ground states of the hydrogen molecule (H_2) with the Quantum Generative Autoencoder Architecture (QGAA), which is trained on a compressed latent representation of these ground states. While the QGAA leverages an autoencoding step to reduce the dimensionality of quantum states before adversarial training, the QGAN attempts to learn the full ground state distribution without compression. This comparison highlights the trade-offs between direct generative modeling of high-dimensional quantum states and the efficiency gains achieved by incorporating an autoencoder into the training pipeline. As shown in Fig. 28, we present the learned energy profile obtained by training the QGAN on 4-qubit ground states, i.e., without employing the autoencoding step used in the QGAA.