# Quantum versus Classical Generative Modelling in Finance

**Brian Coyle**[*,§]**, Maxwell Henderson**[†]**, Justin Chan Jin Le**[†]**,
Niraj Kumar**[*]**, Marco Paini**[†]**, Elham Kashefi**[*,‡]

[*]School of Informatics, 10 Crichton Street, Edinburgh, United Kingdom, EH8 9AB.
[†]Rigetti Computing.
[‡]CNRS, LIP6, Sorbonne Université, 4 place Jussieu, 75005 Paris, France.

E-mail: [§]`brian.coyle@ed.ac.uk`

**Abstract.**　Finding a concrete use case for quantum computers in the near term is still an open question, with machine learning typically touted as one of the first fields which will be impacted by quantum technologies. In this work, we investigate and compare the capabilities of quantum versus classical models for the task of generative modelling in machine learning. We use a real world financial dataset consisting of correlated currency pairs and compare two models in their ability to learn the resulting distribution - a restricted Boltzmann machine, and a quantum circuit Born machine. We provide extensive numerical results indicating that the simulated Born machine always at least matches the performance of the Boltzmann machine in this task, and demonstrates superior performance as the model scales. We perform experiments on both simulated and physical quantum chips using the Rigetti forest platform, and also are able to partially train the largest instance to date of a quantum circuit Born machine on quantum hardware. Finally, by studying the entanglement capacity of the training Born machines, we find that entanglement typically plays a role in the problem instances which demonstrate an advantage over the Boltzmann machine.

## 1. Introduction

The prediction power of machine learning algorithms is limited by the quality of the datasets used to train the models. In the age of big data, possessing high-quality data can offer significant competitive advantage to institutions who utilize machine learning in their core business operations such as Facebook, Google and Amazon. However, for many organizations high-quality data can be scarce. This is because training data for industrial problems are often plagued by erroneous information, limited by privacy and over-fitting. Hence, high-quality data can be expensive or even impossible to obtain especially for machine learning applications at industrial scales. Synthetic data generation (SDG) bridges the gap for training better machine learning models when such data is not readily available. Rather than collecting raw data, SDG uses statistical methods, simulation modeling, and neural networks to generate a synthetic equivalent of the real-world data set (i.e. sample generation). SDG allows users to overcome data scarcity, avoid privacy issues, and overcome over-fitting problems at lower costs. This is achieved by SDG removing erroneous or mislabeled data, as each sample is generated from predefined parameters to produce clean and machine learning-ready datasets. SDG can also produce realistic data for unobserved scenarios to train more generalized models. In machine learning terms, SDG is typically achieved by *generative modelling* or distribution learning. Using quantum models for SDG has garnered interested due to the ease of generating data samples (alternatively, performing the 'inference' step) from a quantum distribution, whereas even the very act of sample generation can be difficult classically, which we elaborate on through the text.

In terms of quantum capabilities, we are now firmly in the noisy intermediate scale quantum (NISQ) [1] era, where we have access to small, error-prone quantum computers, but which are sufficiently powerful to be able to address problems which are not classically simulatable [2]. However, finding a useful application for such devices is a non-trivial task, with quantum chemistry [3, 4] or quantum optimization [5, 6] being the usual suspects for areas in which to search. Problems in finance have also proved to be a lucrative area of study, [7–10]. With each discovered use case, an argument is frequently required as to why such a problem could not have been tackled by purely classical methods. The primary approaches to gain an advantage with quantum computers study the computational time complexity in solving these problems. The claims of exponential speedups [11, 12] in these cases usually rely on the non-existence of unlikely relationships between computational complexity classes. However, simply solving the problem faster is not the only way in which quantum computers can gain victories. Alternatively, one can examine other relevant problem dimensions, such as accuracy of solution, which is the goal we aim for in this work.

We explore two different machine learning approaches for generating synthetic financial market data. One model is completely classical (although trained using simulated quantum methods): the restricted Boltzmann machines (RBM) and the other is completely quantum in nature: a quantum circuit Born machines (QCBM). This is similar to other recent works [13], which addressed financial problems with these two models and found that the Born machine has the capacity to outperform the Boltzmann machine, when it comes to generating synthetic data. In this work, we draw a similar conclusion by enforcing similar constraints on both models in order to draw a fair comparison.

In Section 2 we discuss the main ideas involved in generative modelling, and elaborate on the two models we use for this task. We also discuss the financial dataset we use for training. In Section 3, we detail the specific architectures for the Boltzmann and Born machines, namely the underlying graph structures and the circuit Ansätze for the QCBM. In Section 4, we describe the training protocols we use for each model and finally in Section 5 we detail the numerical results we find, and showcase examples where the Born machine outperforms the Boltzmann machine in learning the financial dataset. We present simulated and experimental results implemented on the Rigetti QPU [14] using Quantum Cloud Services (QCS™). Finally, we conclude in Section 6 and discuss future work.

## 2. Generative Modelling

Generative models are powerful machine learning models, which essentially aim to learn a probability distribution, denoted $\pi$, over some data (say vectors, $\boldsymbol{x}$), which is sampled from $\pi$, $\boldsymbol{x} \sim \pi(\boldsymbol{x})$. A typical use case is in classification tasks, where a generative model seeks to learn the joint distribution over data and labels, $y$, $\pi(\boldsymbol{x}, y)$. We assume the distribution in question is defined over the space of binary strings of length $n$, $\{0, 1\}^n$. A generative model can be typically parameterised by some parameters, $\boldsymbol{\theta}$, and are represented by an output 'model' distribution over the data, which is a function of those parameters, $p_{\boldsymbol{\theta}}(\boldsymbol{x})$. The goal of training a generative model is to force the model distribution as close as possible to the data distribution, relative to some measure. This is done by finding a suitable setting of the parameters, typically using some optimization routine. In practice however,
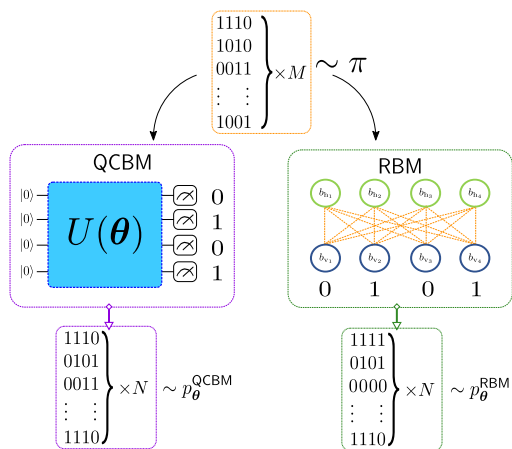
Figure 1: Illustration of synthetic data generation by a quantum circuit Born machine (QCBM) versus a restricted Boltzmann machine (RBM) for binary strings of length 4. Generative modelling involves learning a representation of an underlying distribution ($\pi$) from $M$ samples. The trained model can generate $N$ samples where $N$ can be larger than $M$. The samples associated to the QCBM are the binary results from measuring qubits in (typically) the computational basis and generate $p_{\boldsymbol{\theta}}^{\mathsf{QCBM}}$, whereas for the RBM, they are associated to configurations of the visible nodes, $p_{\boldsymbol{\theta}}^{\mathsf{RBM}}$. We judge the quality of the samples by the similarity of the generated distributions to $\pi$.

we typically do not have access to the true distributions (meaning their explicit probability density functions or otherwise) This is inevitably true when using implicit models ‡ like generative adversarial networks (GANs) [16] or quantum circuit distributions, which by their very nature are distributions which are not directly accessible [2] due to the classical intractability of them. In this work, we assume we have $N$ samples from the model distribution, $\{\boldsymbol{x}_i\}_{i=1}^N, \boldsymbol{x}_i \sim p_{\boldsymbol{\theta}}(\boldsymbol{x})$ and $M$ samples from the data distribution, $\{\boldsymbol{y}_j\}_{j=1}^M, \boldsymbol{y}_j \sim \pi(\boldsymbol{y})$.

Common use cases for generative models are in image generation, but they have also received interest from the quantum computing community, as the acceleration in training of generative models using quantum techniques was one of the early areas of interest in the field of quantum machine learning [17]. The focus of the area has shifted somewhat in recent years, from accelerating training and inference of *classical* models using quantum techniques, to the development of *completely new* models in the quantum world. One of the earliest examples of which is the *quantum Boltzmann machine* (QBM), which is a generalization of the classical Boltzmann machine (see

Section 2.2). This was followed by the introduction of Born machines [18] and *quantum circuit Born machines* (QCBMs) [19, 20], which sample from the fundamentally quantum distribution underlying a pure state of a quantum system. One of the most recent additions to this family are Hamiltonian based models and the variational quantum thermalizer (VQT) [21], which generalizes all of the above since it contains the distribution provided by a *mixed* quantum state as the underlying model. The latter is also a generalization of 'energy-based' models, of which the Boltzmann machine is an example. Furthermore, quantum generative models are some of the most promising applications for near term quantum computers since their nature aligns them closely with demonstrations of '*quantum supremacy*' [2] and such connections have recently been made [22, 23] with extensions into different architectures [24].

In this work we focus on two of these models in order to make a direct comparison and study any potential indication of quantum advantage for these models over purely classical generative models. We investigate a Born machine and a restricted Boltzmann machine (RBM) and make a thorough comparison between the two for a generative modelling task. We do this using a realistic dataset in a financial application, which facilitates a simple way to compare the models at differing scales. Our motivation is the work of [25] which showed the outperformance of an RBM over parametric models, for this dataset, which are the common tool used in the finance industry. This was subsequently followed by the subsequent outperformance of the RBM by a Born machine [13] on the same dataset. However, the degree to which this advantage was observable was not obvious. This research and [13] supplements the work of [26] which demonstrated a similar outperformance of an RBM by a QCBM, but for a different problem domain. Our work expands on the latter by running larger problem instances on simulators and physical hardware, using alternative training methods, and also using alternative methods of comparison of the models. Finally, drawing a comparison between a Born and Boltzmann machine is part of the goal of [18], in which they consider the problem from a mutual information point of view. They further conjecture that properties such as mutual information of the dataset, and entanglement entropy in the target problem, an/or model would be useful in determining problems where the Born machine could have superior performance over an RBM.

## 2.1. Born Machine

A Born machine [18] is a fundamentally quantum model, which achieves synthetic data generation

by generating samples according to Born's rule of quantum mechanics. The fundamentally non-classical nature of the model has provided motivation for why it can outperform classical models in at least its expressive power [22]. This expressive power translates in an ability to represent certain distributions efficiently which cannot be done by any classical model, for example, those utilized in a recent demonstration of quantum computational supremacy [2].

In the most common scenario, a binary sample, $\boldsymbol{x} \in \{0,1\}^n$, is generated from a quantum state, $\rho$, according to:

$$\boldsymbol{x} \sim p(\boldsymbol{x}) = \mathrm{Tr}\left(|\boldsymbol{x}\rangle\langle\boldsymbol{x}|\,\rho\right) \tag{1}$$

where $|\boldsymbol{x}\rangle\langle\boldsymbol{x}|$ is the projector onto the computational basis state described by $\boldsymbol{x}$. In order to obtain a trainable machine learning model, we parameterize the state: $\rho \to \rho_{\boldsymbol{\theta}}$. We also further consider the scenario where the parameterised state is a pure state, i.e. $\rho_{\boldsymbol{\theta}} := |\psi_{\boldsymbol{\theta}}\rangle\langle\psi_{\boldsymbol{\theta}}|$. In this case, the correlations present in the model will be of a purely quantum nature. The parameterised distribution is then:

$$\boldsymbol{x} \sim p_{\boldsymbol{\theta}}^{\mathsf{QCBM}}(\boldsymbol{x}) = |\langle\boldsymbol{x}|\psi_{\boldsymbol{\theta}}\rangle|^2 \tag{2}$$

Finally, if the state, $|\psi_{\boldsymbol{\theta}}\rangle$ is generated by a quantum circuit (as opposed to, for example, by a continuous time Hamiltonian evolution), the model is referred to as a *quantum circuit* Born machine [19, 20] (QCBM). In this form, the ease of performing inference becomes apparent: once trained, the parameterized quantum state prepared by a quantum circuit and then simply measured. The measurement results then constitute an (approximate) sample from the data distribution. Furthermore, utilizing quantum randomness as a sample generation mechanism this way relaxes the need to input randomness into the model as is usually done to build GANs. However, we mention that inputting randomness has been considered in the quantum case [27] as well, although the advantage of doing so has yet to be explored.

A generalization of the above can be achieved by relaxing the purity assumption of the underlying state, and doing so results in quantum Hamiltonian based models [21], which instead can carry *both* classical and quantum correlations.

In order to find a good fit to the data distribution, $\pi$, such that the model, $p_{\boldsymbol{\theta}}(\boldsymbol{x})$, can effectively generate synthetic data, an optimization routine is invoked to search over the space of possible states $|\psi_{\boldsymbol{\theta}}\rangle$. Since Born machines are implicit models, careful consideration must be given to the choice of optimization routine, since any optimizer must be able to effectively, and efficiently, deal with samples alone. One may consider quantum training procedures [28],

but more commonly the optimization procedure will be a fully classical routine. This makes these models hybrid quantum-classical in nature and therefore friendly to NISQ devices, only using the quantum resource when necessary.

## 2.2. Boltzmann Machine

Generalized Boltzmann machines (GBMs) are graphical models with powerful synthetic data-generation capabilities. While GBMs can vary significantly in terms of how they are applied to various problems and their particular architectures (see some example architectures in Figure 2), they all share some defining characteristics. The model architecture is defined by a graph $\mathsf{G}$, which consists of a set of edges, $\mathcal{E}$, and nodes (vertices) which we denote $\mathcal{N}$. Each edge, $e \in \mathcal{E}$ has a corresponding edge *weight*, $W_e \in \mathcal{W}$. In generality, the edge weights can also be self-loops (biases in standard Boltzmann machine terminology), or hyper-edges (edges connecting more than two nodes) as illustrated in Figure 2(c). Crucially, the nodes are typically partitioned into *visible* and *hidden* nodes, $\mathcal{N} = \mathcal{N}_{\mathrm{v}} \bigcup \mathcal{N}_{\mathrm{h}}$. The visible nodes directly model some aspect of the data distribution, while the hidden nodes are used for capturing features of the data, and are not tied to any particular aspect of it. As such the hidden nodes typically correspond directly to the expressive power of the model. Finally, a sample generated by the GBM is distributed according to the Boltzmann distribution:

$$\boldsymbol{x} \sim p_{\boldsymbol{\theta}}^{\mathsf{GBM}}(\boldsymbol{x}) = \frac{e^{-\beta E(\boldsymbol{x})}}{\mathcal{Z}}, \tag{3}$$

$p_{\boldsymbol{\theta}}^{\mathsf{GBM}}(\boldsymbol{x})$ is the probability to observe the visible nodes in some state $\boldsymbol{n}_{\mathrm{v}} = \boldsymbol{x}$, and describes the model distribution for the Boltzmann machine. $\boldsymbol{n}_{\mathrm{v}}$ is a particular state (corresponding to a binary vector) of the visible nodes in $\mathcal{N}_{\mathrm{v}}$. In this case, the model parameters are the weights of the machine, $\boldsymbol{\theta} = \mathcal{W}$. $E$ and $\mathcal{Z}$ is the model *energy* (defining an energy based model [29]) and *partition function* respectively, and are defined by:

$$E(\boldsymbol{v}) := -\sum_{e \in \mathcal{E}} W_e \prod_{v_i \in e} v_i, \tag{4}$$

$$\mathcal{Z} := \sum_{\boldsymbol{v}} e^{-\beta E(\boldsymbol{v})} \tag{5}$$

The notation, $v_i \in e$, refers to the nodes connected to edge $e$, and $\beta$ is an effective inverse temperature term. The sum in Equation (5) is taken over all possible binary vectors $\boldsymbol{v} \in \{0,1\}^n$. In this work, we focus specifically on the *restricted* version of the Boltzmann machine (RBM) corresponding to Figure 2(b), and we discuss this specification further in Section 3.2. In this case, we denote $p_{\boldsymbol{\theta}}^{\mathsf{GBM}}(\boldsymbol{x}) \to p_{\boldsymbol{\theta}}^{\mathsf{RBM}}(\boldsymbol{x})$ where the latter
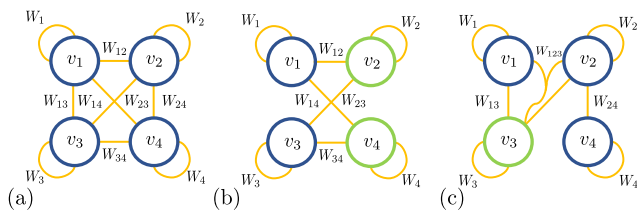
Figure 2: Visualizing various possible GBM architectures for a network with 4 nodes, such as a (a) fully-connected visible network, (b) restricted Boltzmann machine (RBM), and (c) partially-connected higher order Boltzmann machine. All visible nodes are shown in blue and hidden nodes are shown in green.
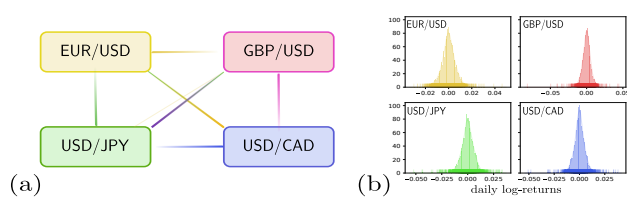


Figure 3: We use data generated from FX spot prices of the above currency pairs. The generative model aims to learn correlations between each pair based on a 16 bit binary representation. (a) The selection of currency pairs we use, and (b) the marginal distributions of the log-returns of each pair over a 20 year period. We aim to learn the joint distribution of subsets of the pair in this work.

distribution is generated by marginalizing over hidden units.

Finally, while all of the above is purely classical (in contrast to the Born machine, a GBM carries only classical correlations), the extension of the *model itself* into the quantum world has also been proposed in the quantum Boltzmann machine [30–33] as we mentioned above. In this framework, the energy function, Equation (4) is replaced by a quantum Hamiltonian and the model distribution in question is generated by sampling from the thermal state of this Hamiltonian, mimicking a Boltzmann distribution. This thermal state can be prepared either by quantum annealing [30] or by a gate based approach [34]. By introducing off-diagonal terms in this Hamiltonian, non-trivial quantum behavior can be exploited, and the model inherits some characteristics of a Born machine (i.e. some of the randomness originates from Born's rule).

In this work, however, we focus on the GBM as a completely classical object, which we detail in Section 3, however, we do leverage quantum inspired training methods which are discussed in Section 4. Furthermore, as mentioned we only study the RBM here, but we discuss the extension of the methods in this work to the more general Boltzmann machine structures in Section 6.

### 2.3. A Financial Dataset

In order to perform SDG, we require some dataset to learn. In this work, we focus on one of a financial origin, in particular one considered by [25]. This dataset comtains 5070 samples of daily log-returns of 4 currency pairs between $1999 - 2019$ (see Figure 3). In order to fit on the binary architecture of the Born and Boltzmann machines, the spot prices of each currency pair are converted to 16 bit binary values, resulting in samples of 64 bits long. This discretisation provides a convenient method for fitting various problem sizes onto models with different numbers of qubits or visible nodes for the Born machine or RBM respectively. In particular, we can tune both the number of currency pairs $(i)$, and the precision of each pair $(j)$ so the problem size is described by a tuple $(i, j)$. For example, as we revisit in Section 3, a 12 qubit Born machine can be tasked to learn the distribution of 4 currency pairs at 3 bits of precision, 3 pairs with 4 bits or 2 pairs at 6 bits of precision.

## 3. Model Structures

Here we provide specific details about the model architectures we choose to use, in order to derive as fair a comparison as possible. In the first instance, we choose to only train the bias terms in the RBM (the self-loops in Figure 2) for simplicity. We also fix the number of parameters in the Born machine by the number of layers, and then match the number of parameters in the RBM to this, since it is simpler to grow the number of RBM parameters by simply adding extra nodes.

### 3.1. Born Machine Ansatz

The Ansatz which we use for the QCBM is hardware efficient as we endeavor to run the model on real quantum hardware. We also restrict the number of parameters in the circuit to match the number used in the RBM, following [26]. We choose this hardware native approach to closely fit the structure of Rigetti's chip design (the structure of the `Aspen-7` and `Aspen-8` can be seen in Figure 4). Furthermore, we solely parameterize the single qubit unitaries to avoid compilation overheads arising out of two qubit unitary parameterization. If we were to do so, we could employ a similar strategy to [35], which uses 'blocks' of parameterized unitaries in such a way to enforce a linear scaling of the number of parameters with the number of qubits, when building a quantum classifier.
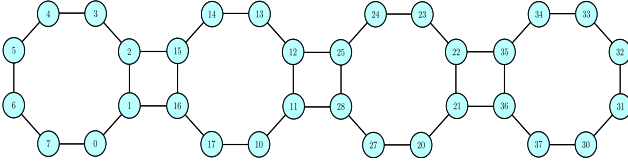
Figure 4: `Aspen 7`/`Aspen 8` 32 qubit chip designs. Note not all connections shown above are directly accessible on the chip itself.
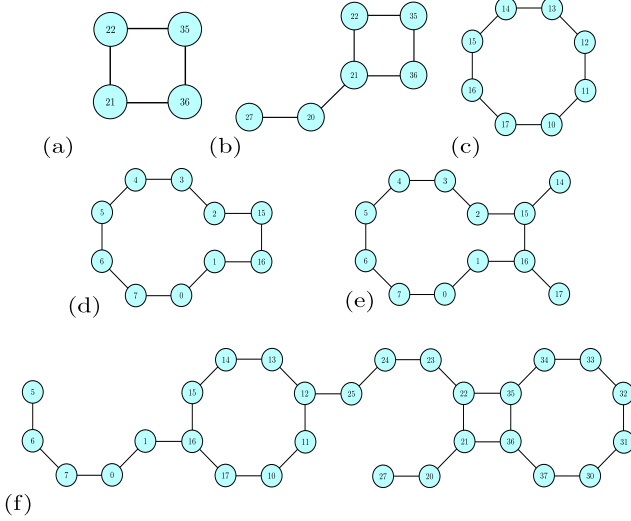


Figure 5: Select sublattices from the `Aspen-7` and `Aspen-8` chips, corresponding to different problem sizes. Figure shows the (a) `Aspen-7-4Q-C`, (b) `Aspen-7-6Q-C`, (c) `Aspen-7-8Q-C`, (d) 10 qubit `Aspen-8`, (d) 12 qubit `Aspen-8` and (f) `Aspen-7-28Q-A` sublattices. Using these topologies we can fit problems of size $(2, 2), (2, 3), (2, 4), (2, 5), (2, 6)$ and $(4, 7)$ respectively, where the notation $(i, j)$ indicates $i$ currency pairs, each described by $j$ bits of precision.

We run all experiments using the Rigetti `Aspen 7` and `Aspen 8` chips, which are designed to contain 32 qubits, however some qubits are not available. Each QPU can be divided into sublattices containing fewer qubits, some examples can be seen in Figure 5. The largest sublattice on the `Aspen-7` chip is the `Aspen-7-28Q-A` which contains 28 usable qubits (seen in Figure 5(f)).

For each of the lattices in Figure 5, we fit the native entanglement structure using CZ gates, and layers of single qubits rotations. For convenience, we use $R_y$ rotation gates as the single qubit gates, which have the decomposition $R_y(\theta) = R_x(\pi/2)R_z(\theta)R_x(-\pi/2)$, using the Rigetti native single qubit rotations. The first 'layer' contains only $R_y$ gates, and each layer thereafter consists of the hardware native CZ gates, plus a layer of $R_y$ gates. For the 4 qubit chip, `Aspen-7-4Q-C`, we illustrate this
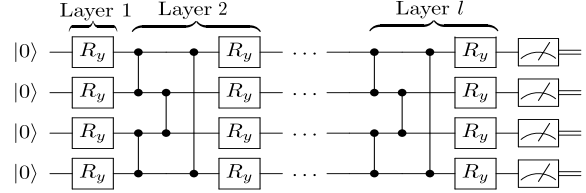


Figure 6: Hardware efficient circuit for the `Aspen-7-4Q-D`, with $l$ layers using the native entanglement structure native to the chip.
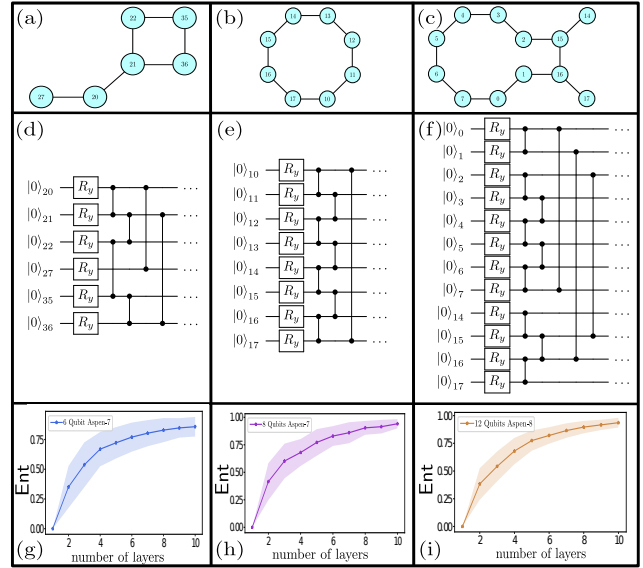


Figure 7: Hardware efficient circuits for 6, 8, 12 qubit Born machine ansatz. (a)-(c) show `Aspen-7-6Q-C`, `Aspen-7-8Q-C` from the `Aspen-7` chip and a 12 qubit sublattice from the `Aspen-8` chip which we consider. (d) - (f) illustrate the the entanglement structure in a single layer, which tightly matches the chip topology. (g) - (i) show the average entangling capability, Ent in (10) as a function of the number of layers in the circuit for each of the entangling structures shown in (d)-(f). Error bars show mean and standard deviation over 100 random parameter instances, $\{\boldsymbol{\theta}_i\}_{i=1}^{100}, \boldsymbol{\theta}_i^j \sim U(0, 2\pi)$, in the single qubit rotations. $U$ is the uniform distribution over the interval $[0, 2\pi]$.

in Figure 6. For the other sublattices in Figure 5, we illustrate the entanglement structure in Figure 7 for the first layer of the circuits. In this way, an $n$ qubit QCBM with $l$ layers will have $n \times l$ trainable parameters.

For the above circuits, we compute the average Meyer-Wallach [36] entanglement capacity, a measure of entanglement in quantum states proposed as a method of comparing different circuit Ansätze by [37]. This measure has been used in a similar context by [38] in order to draw connections between Ansatz structure and classification accuracy. The entanglement measure

$Q$ is defined, for a given input state $|\psi\rangle$ as:

$$Q(|\psi\rangle) := \frac{4}{n} \sum_{j=1}^{n} D(\iota_j(0)|\psi\rangle, \iota_j(1)|\psi\rangle) \qquad (6)$$

$$D(|\boldsymbol{u}\rangle, |\boldsymbol{v}\rangle) = \frac{1}{2} \sum_{i,j} |u_i v_j - u_j v_i|^2 \qquad (7)$$

where $D$ is a particular distance between two quantum states, $|\boldsymbol{u}\rangle := \sum_i u_i |i\rangle, |\boldsymbol{v}\rangle := \sum_j v_j |j\rangle$. This distance can be understood as the square of the area of the parallelogram created by vectors $|\boldsymbol{u}\rangle$ and $|\boldsymbol{v}\rangle$. The notation $\iota_j(b)$ is a linear map which acts on computational basis states as follows:

$$\iota_j(b)|b_1 \ldots b_n\rangle := \delta_{bb_j} \left| b_1 \ldots \hat{b}_j \ldots b_n \right\rangle \qquad (8)$$

where $\hat{\phantom{x}}$ indicates the absence of the $j^{th}$ qubit. For example, $\iota_2(0)|1001\rangle = |101\rangle$ However, to evaluate $Q$ for a quantum state, we instead use the equivalent formulation derived by [39], which involves computing the purities of each subsystem of the state $|\psi\rangle$:

$$Q(|\psi\rangle) = 2 \left(1 - \frac{1}{n} \sum_{k=1}^{n} \text{Tr}\left[\rho_k^2\right]\right) \qquad (9)$$

where $\rho_k := \text{Tr}_{\bar{k}}(|\psi\rangle\langle\psi|)$ is the partial trace over every one of the $n$ subsystem of $|\psi\rangle$ *except* $k$. This reformulation of $Q$ gives more efficient computation and operational meaning since the purity of a quantum state is efficiently computable. Given $Q$, we define [37] Ent as the average value of $Q$ over a set, $\mathcal{S}$ of $M$ randomly chosen parameter instances, $S := \{\boldsymbol{\theta}_i\}_{i=1}^{M}$:

$$\text{Ent} := \frac{1}{|S|} \sum_i Q(|\psi_{\boldsymbol{\theta}_i}\rangle) \qquad (10)$$

For the circuit Ansätze we choose, the value of Ent is plotted for a given number of layers in Figure 7.

### 3.2. Boltzmann Machine Structure

Given the above choice for a Born machine ansatz, we can build a corresponding restricted Boltzmann machine which has $n_{\text{v}} := |\mathcal{N}_{\text{v}}| = n$ visible nodes (where $n$ is the number of qubits) and $n_{\text{h}} := |\mathcal{N}_{\text{h}}| = nl - n = n \times (l-1)$ hidden nodes. To reiterate, we fix the RBM weights to have random values and only the local biases are trained. We revisit weight training in Appendix B.3.

### 4. Training Procedures

In order to fit the model distribution to the data, one need some means of comparing how close these two distributions are. Typically, this comes in the form of
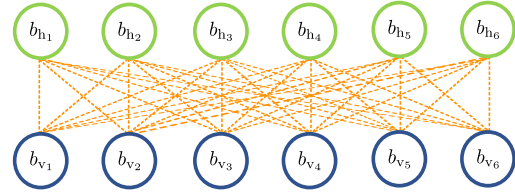


Figure 8: Restricted Boltzmann machine structure using 6 visible and hidden nodes with 12 parameters. Corresponds to the 6 qubit Born machine in Figure 7(d). Dotted lines indicate weights are not trainable but randomly chosen and fixed through training. Biases for visible and hidden nodes, $b_{\text{v}_i}, b_{\text{h}_i}$ correspond to the self-loop weights in Figure 2 which are trainable.

a *cost function*, $D(p_{\boldsymbol{\theta}}, \pi)$. In this work, we consider a variety of cost functions with which to compare both models we investigate.

This cost function is then minimized during the training procedure to find a setting of the parameters, $\boldsymbol{\theta}$ such that $D(p_{\boldsymbol{\theta}}, \pi)$ is as small as possible. Gradient descent (GD) is a common method to minimize such costs in machine learning as it finds the steepest direction of descent in the parameter landscape defined by, $\boldsymbol{\theta}$. GD proceeds with a number of 'epochs', where in each epoch $(t)$ the parameters are updated as follows:

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \Delta D(p_{\boldsymbol{\theta}^{(t)}}, \pi) \qquad (11)$$

$\Delta D(p_{\boldsymbol{\theta}}, \pi)$ is the update rule defining how each parameter should be updated, depending on the current value of $D$ and is negative since we wish to go downhill in the parameter landscape. The 'vanilla' form of gradient descent simply directly uses an update of the form $\Delta D(p_{\boldsymbol{\theta}}, \pi) = \eta \partial_{\boldsymbol{\theta}^{(t)}} D(p_{\boldsymbol{\theta}}, \pi)$, where $\eta$ is a *learning rate* and $\partial_{\boldsymbol{\theta}^{(t)}} D$ is the partial derivative of $D$ with respect to the current parameters. Computing this gradient efficiently can be a non-trivial procedure, and it is estimate given the data. More complicated update rules such as Adam [40] are also possible, which include terms like 'momentum' to the update rule, to improve convergence speed.

### 4.1. Born Machine Training

The primary cost function we choose to train the Born machine is the Sinkhorn divergence (SHD), a recently defined [41–43] method of distribution comparison, and related to *optimal transport* (OT) [44], which is known to be a relatively powerful metric between probability distributions.

$$D(p_{\boldsymbol{\theta}}, \pi) := \mathcal{L}_{\text{SHD}}^{\epsilon}(p_{\boldsymbol{\theta}}, \pi) :=$$
$$\text{OT}_{\epsilon}^c(p_{\boldsymbol{\theta}}, \pi) - \frac{1}{2}\text{OT}_{\epsilon}^c(p_{\boldsymbol{\theta}}, p_{\boldsymbol{\theta}}) - \frac{1}{2}\text{OT}_{\epsilon}^c(\pi, \pi) \quad (12)$$

$$\text{OT}^c_\epsilon(p_{\boldsymbol{\theta}}, \pi) := \min_{U \in \mathcal{U}(p_{\boldsymbol{\theta}}, \pi)}$$

$$\left( \sum_{\substack{(\boldsymbol{x}, \boldsymbol{y}) \\ \in \mathcal{X}^d \times \mathcal{Y}^d}} c(\boldsymbol{x}, \boldsymbol{y}) U(\boldsymbol{x}, \boldsymbol{y}) + \epsilon \text{KL}(U | p_{\boldsymbol{\theta}} \otimes \pi) \right) \quad (13)$$

where $\epsilon \geq 0$ is a regularisation parameter, and $\mathcal{U}(p_{\boldsymbol{\theta}}, \pi)$ is the set of all *couplings* between $p_{\boldsymbol{\theta}}$ and $\pi$, i.e. the set of all joint distributions, whose marginals with respect to $\boldsymbol{x}, \boldsymbol{y}$ are $p_{\boldsymbol{\theta}}(\boldsymbol{x}), \pi(\boldsymbol{y})$ respectively. $\text{KL}(U | p_{\boldsymbol{\theta}} \otimes \pi)$ is the Kullback-Leibler [45] divergence (also relative entropy) between the coupling, $U$, and a product distribution composed of the model and the data, $p_{\boldsymbol{\theta}} \otimes \pi$. The introduction of the entropy term smooths the problem, so that it becomes more easily solvable, as a function of $\epsilon$.

We use this cost function since we numerically found it to be the best choice, in terms of speed and accuracy of training. However, we provide a comparison to the *maximum mean discrepancy* (MMD) cost function, training with respect to an adversarial discriminator and a gradient free genetic algorithm in Appendix A.

As shown in [22] we can derive gradients of the Sinkhorn divergence, with respect to the given parameter, $\boldsymbol{\theta}_k$, since each parameterised gate we employ has the form $U(\theta) = \exp(i\theta/2\Sigma)$, where $\Sigma^2 = \mathbb{1}$. Using the parameter shift rule [46, 47], the gradient can be written as follows:

$$\frac{\partial \mathcal{L}^\epsilon_{\text{SHD}}(p_{\boldsymbol{\theta}}, \pi)}{\partial \boldsymbol{\theta}_k} = \sum_{\boldsymbol{x}} \frac{\partial \mathcal{L}^\epsilon_{\text{SHD}}(p_{\boldsymbol{\theta}}, \pi)}{\partial p_{\boldsymbol{\theta}}(\boldsymbol{x})} \frac{\partial p_{\boldsymbol{\theta}}(\boldsymbol{x})}{\partial \boldsymbol{\theta}_k} \quad (14)$$

$$= \frac{1}{2} \sum_{\boldsymbol{x}} \varphi(\boldsymbol{x}) \left( p_{\boldsymbol{\theta}_k^+}(\boldsymbol{x}) - p_{\boldsymbol{\theta}_k^-}(\boldsymbol{x}) \right) \quad (15)$$

$$= \frac{1}{2} \left( \mathop{\mathbb{E}}_{\boldsymbol{x} \sim p_{\boldsymbol{\theta}_k^+}} [\varphi(\boldsymbol{x})] - \mathop{\mathbb{E}}_{\boldsymbol{x} \sim p_{\boldsymbol{\theta}^-}} [\varphi(\boldsymbol{x})] \right) \quad (16)$$

The function $\varphi$ is defined[43] in order to ensure the gradient extends to the entire sample space, and is defined as follows for each sample, $\boldsymbol{x}$:

$$\varphi(\boldsymbol{x}) =$$
$$- \epsilon \text{LSE}^M_{k=1} \left( \log \left( \pi(\boldsymbol{y}^k) + \frac{1}{\epsilon} g(\boldsymbol{y}^k) - \frac{1}{\epsilon} C(\boldsymbol{x}, \boldsymbol{y}^k) \right) \right)$$
$$+ \epsilon \text{LSE}^N_{k=1} \left( \log \left( p_{\boldsymbol{\theta}}(\boldsymbol{x}^k) + \frac{1}{\epsilon} s(\boldsymbol{x}^k) - \frac{1}{\epsilon} C(\boldsymbol{x}, \boldsymbol{x}^k) \right) \right)$$
$$(17)$$

Therefore, one can compute the gradient by drawing samples from the distributions, $\hat{\boldsymbol{x}} \sim p_{\boldsymbol{\theta}^\pm}$, and computing the vector $\varphi(\boldsymbol{x})$, for each sample, $\boldsymbol{x}$. The functions $g$ and $s$ in Equation (17) are

optimal Sinkhorn potentials, arising from a primal-dual formulation of optimal transport. These are computed using the Sinkhorn algorithm, which gives the divergence its name [48]. $C(\boldsymbol{x}, \boldsymbol{y})$ is the optimal transport *cost matrix* derived from the cost function applied to all samples, $C_{ij}(\boldsymbol{x}^i, \boldsymbol{y}^j) = c(\boldsymbol{x}^i, \boldsymbol{y}^j)$ and $\text{LSE}^N_{k=1}(\boldsymbol{V}_k) = \log \sum_{k=1}^N \exp(\boldsymbol{V}_k)$ is a log-sum-exp reduction for a vector $\boldsymbol{V}$. For further details on how the functions, $g$ and $s$ are computed, see along with the Sinkhorn divergence and its gradient, see [22, 43].

### 4.2. Boltzmann Machine Training

For the RBM, we use the standard Boltzmann protocol of maximizing the log-likelihood function $\mathcal{L}$ §, i.e. the probability of generating vectors belonging to a training set $Y = \{\boldsymbol{y}\}$:

$$\mathcal{L}(Y, \boldsymbol{\theta}) = \log(p_{\boldsymbol{\theta}}(Y)) \quad (18)$$

where $\boldsymbol{\theta}$ are the Boltzmann machine model parameters and $p(\boldsymbol{y})$ is the probability of generating data vectors $\boldsymbol{y} \sim \pi(\boldsymbol{y})$. For a particular data vector $\boldsymbol{y} \in Y$, we can take the likelihood function $\mathcal{L}$ as our cost function as a function of the model parameters $\boldsymbol{\theta} = \mathcal{W} = \{W_e\}$ [49], which results in the gradient:

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{L}}{\partial W_e} = \langle v_e \rangle_\pi - \langle v_e \rangle_{p_{\boldsymbol{\theta}}} \quad (19)$$

wherein $W_e$ of (4) are the model parameters coupling their respective set of nodes $v_e := \prod_{v_i \in e} v_i$, $\langle v_e \rangle_\pi$ and $\langle v_e \rangle_{p_{\boldsymbol{\theta}}}$ are the expectation values of $v_e$ calculated from the data and model distributions respectively where $v_e$ is taken to be a random variable taking values $\{0, 1\}$.

As an example, consider an update to a (visible node) bias term (i.e. a self-loop in Figure 2), we have $W_e = b_{\text{v}_i}$ and also $v_e$ is simply $b_{\text{v}_i}$ since the edge connects only one node. Then the gradient is computed using the expectation value of the bias:

$$\frac{\partial \mathcal{L}}{\partial b_{\text{v}_i}} = \langle b_{\text{v}_i} \rangle_\pi - \langle b_{\text{v}_i} \rangle_{p_{\boldsymbol{\theta}}} \quad (20)$$

In this work, we use vanilla gradient descent as the update rule, but we note we also considered more complex update rules or optimizers such as Adam [40], and we found that this only improved the convergence speed, and not the final accuracy of the model.

The above is discussion has no quantum component, as the update rule and the model is completely classical. However, in order to actually compute the first and second order moment terms (using $M$ data

---

§ Equivalent to minimizing an empirical cost $\widetilde{D}(p_{\boldsymbol{\theta}}, \pi) = 1 - \mathcal{L}(Y, \boldsymbol{\theta})$. The maximization procedure adds an extra negative sign to the update rule.

and $N$ model-generated bitstrings) in (19), we require a method of generating samples from the RBM. Unlike the Born machine, sample generation from a Boltzmann machine is not a trivial matter. Typical approaches are based on Gibbs sampling, for example $k$-step contrastive divergence [50]. Here, we use a method called QxSQA, a GPGPU-Accelerated Simulated quantum annealer based on Path-Integral Monte Carlo (PIMC) [51]. This simulated QA has been shown to be useful for sampling Boltzmann-like distributions, and we have shown the ability to use this sampling to train large *quantum* generalized Boltzmann machines (QGBMs) for the purposes of generating synthetic data based on images [52] and financial data [13, 25] in previous research [53].

## 5. Results

Here we present the numerical results obtained above using the models and training methods detailed above. In particular, we focus of training using the Sinkhorn divergence with the Adam optimiser [40] and its analytic gradients for the Born Machine, and log-likelihood maximization using QxSQA for training the Boltzmann machine. In Appendix A, we revisit alternative training methods. As mentioned above, in the first instance, we also fix both models to only have trainable local parameters for simplicity. For the Born machine, this corresponds to only training single qubit unitaries, with the two qubit gates being unparameterised, and for the Boltzmann machine, this corresponds to training the biases of each node. We force each model to have the same number of trainable parameters in this way. The entanglement structure in the Born machine is fixed by the problem size, via the lattice topology, and the weights of the RBM are chosen to be random (but fixed) values on each instance. It is difficult to directly compare the connectivity of the models, howeveer we also experimented with randomly pruning the RBM weights to enforce the same number of connections as in the QCBM, but we found this did not affect performance significantly.

As a method of benchmarking the expressive power of each model in a fair way, we use an adversarial discriminator, and judge the performance relative to it. Specifically, we use a random forest discriminator from `scikit-learn` [54] with 1000 estimators. A higher discriminator error implies better performance, with an error of 50% indicates the discriminator can at best guess randomly when presented with a sample as to its origin - whether it came from the real data, or the model. Where error bars are shown, they correspond to the mean and standard deviations of the training over 5 independent runs. As the QCBM scales, the

classical simulation becomes a bottleneck and limits the number of runs which can be done.

In summary, we find the Born machine has the capacity to outperform the RBM as the precision of the currency pairs increases. In Figure 9, we use data from 2 currency pairs, at $2, 3, 4$ and $6$ bits of precision. We notice the Born machine outperforms the RBM around 4 bits (measured by a higher discriminator error), and still performs relatively well when run on the QPU. Similar bahaviour is observed for 3 currency pairs in Figure 11, which uses a precision of 2 and 4 bits, and with 4 pairs in Figure 12 for a precision of 2 and 3 bits. In Figure 13 we plot the entangling capability (defined by (9)) of the states generated by initial and final circuits learned via training. Curiously, we notice that in the problem instances in which the Born machine outperforms the Boltzmann machine (those with a higher level of precision), the trained circuits have a higher level of entanglement than those that do not, despite the data being completely classical in nature. This is especially prominent for 2 currency pairs in Figure 13(a), in which the training drives the entanglement capability at 2 and 3 bits of precision close to zero (even for increased numbers of layers), but it is significantly higher for 4 and 6 bits of precision, when the Born machine outperforms the RBM, as seen in Figure 9. Similar behaviour is seen for 3 currency pairs, but not as evident for 4 pairs. The latter effect is possibly correlated to the similar performance of both models for 4 currency pairs up to 3 bits of precision.

We are also able to somewhat successfully train the largest instance of a Born Machine to date in the literature, namely one consisting of 28 qubits on the Rigetti `Aspen-7` chip using the Sinkhorn divergence (whose topology is shown in Figure 5(e), and we find it performs surprisingly well. We show the performance of the 28 qubit model versus the a Boltzmann machine with 28 visible nodes, and a suitable number of hidden nodes to match the number of parameters in the Born machine in Figure 14. While the performance of the Born machine is significantly less than that of its counterpart, it is clear that the model is learning (despite hardware errors), up to a discriminator error of 20%. While this result seems to contradict the previous findings in this work, we emphasize that it does not, since we are not able to simulate the QCBM at this scale in a reasonable amount of time. We would not necessarily expect the Born machine to match the performance of the RBM *on hardware* at this scale for a number of reasons, the most likely cause for diminishing performance is quantum errors in the hardware. However we cannot rule out other factors, such as the Ansatz choice. We leave thorough investigation of improving hardware performance to future work, perhaps by including error mitigation [55]
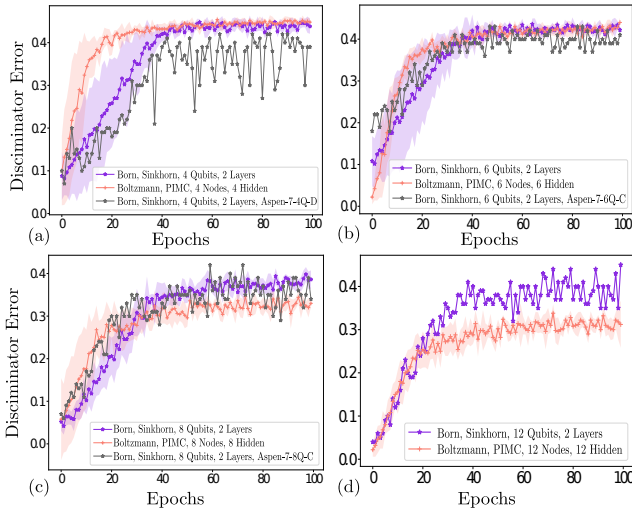
Figure 9: 2 currency pairs (specifically EUR/USD and GBP/USD) at (a) 2 bits, (b) 3 bits, (c) 4 bits and (d) 6 bits of precision. Correspondingly, we use a QCBM of 4, 6, 8 and 12 qubits using the Ansätze described above, and an RBM with the same numbers of visible nodes. The hidden units are scaled in each case to match 2 layers of the QCBM. Results when the QCBM is run on sublattices of the Aspen-7 QPU are shown in grey.
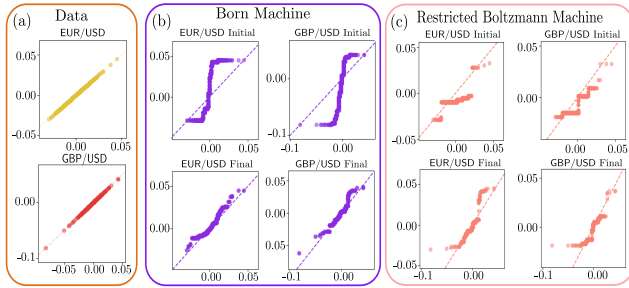


Figure 10: QQ Plots corresponding to Figure 9(d) of the marginal distributions of 2 currency pairs (EUR/USD and GBP/USD) at 6 bits of precision. The Born machine distributions (purple) and those generated by the Boltzmann machine in (pink). (a) shows the QQ plot for the marginal distribution of each currency pair with respect to itself as a benchmark. (b) Born machine initial (top panels) and final (bottom panels) marginal distributions for both pairs and similarly in (c) for the RBM. While not able to completely mimic the data due to the low number of parameters, the Born machine clearly produces a better fit.

to reduce errors, parametric compilation and active qubit reset [14, 56] to improve running time and other techniques.
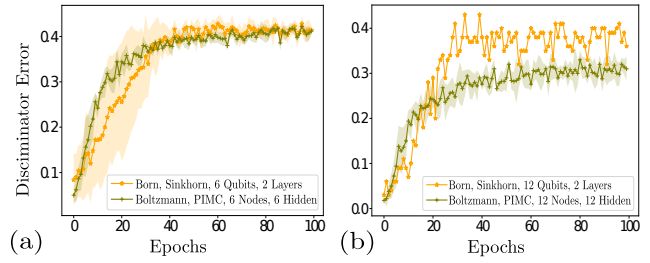


Figure 11: 3 currency pairs at (a) 2 bits and (b) 4 bits of precision, using a QCBM of 6 and 12 qubits and an RBM with the same numbers of visible nodes.
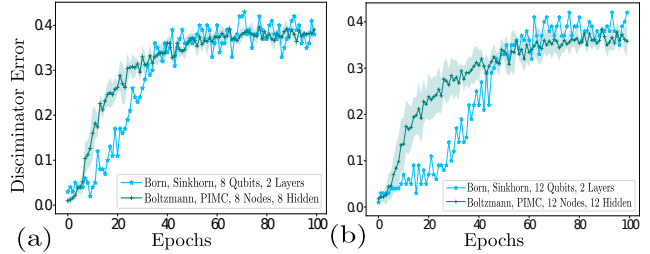


Figure 12: All 4 currency pairs at (a) 2 bits, (b) 3 bits of precision, using a QCBM of 8 and 12 qubits and RBM with the same numbers of visible nodes. We notice the RBM again performs similarly to the Born machine for 2 bits of precision, but begins to be outstripped by it at 3 bits.
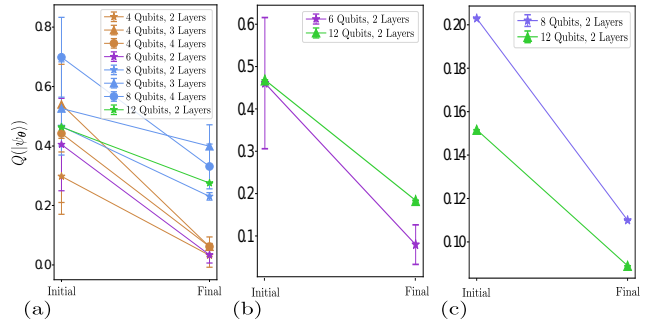


Figure 13: Meyer-Wallach entangling capability (9) for a random choice of parameters (Initial) and the trained parameters (Final) in the same circuit. Error bars represent mean and standard deviation over 5 independent training runs, where they are shown. The circuit ansätze used are those above in Figure 7 closely matching the corresponding chip topology. In each panel we see the circuits trained on (a) 2 currency pairs at 2, 3, 4, 6 bits of precision, (b) 3 currency pairs at 2 and 4 bits of precision and (c) 4 currency pairs at 2 and 3 bits of precision.

## 6. Discussion

In conclusion, we investigate and compare two different models when trained on a real-world financial dataset
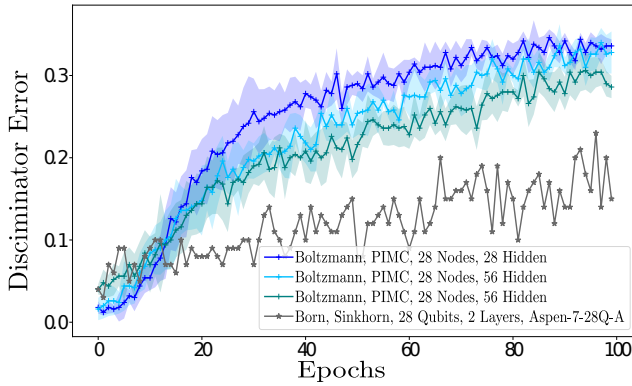
Figure 14: Random forest discriminator during training for a problem size of 4 currency pairs at 7 bits of precision, using 28 visible nodes in the Boltzmann machine and 28 qubits in the Born machine. The 28 qubit Born machine is run exclusively on the `Aspen-7-28Q-A` chip using 2 layers of the hardware efficient ansatz similar to those shown in Figure 7.

consisting of currency pairs at varying levels of precision. We chose a completely classical model in the restricted Boltzmann machine, and put it up against a completely quantum model in the form of a quantum circuit Born machine, in order to compare their relevant expressive powers, and supplement recent related work in this direction [13, 26]. As a benchmark of fairness, we fixed the models to have the same numbers of trainable parameters and found that the simulated Born machine always performed at least as well as the RBM, and in several cases outperformed it, measured relative to the accuracy of an adversarial discriminator. To complement this finding, we investigated the entangling capability of the circuits learned by the QCBM, and found a rough correlation between training towards higher levels of entanglement, and outperforming the classical model.

From this work, there are many possible avenues for exploration. The first, is improving the Born machine training speed by, for example, leveraging GPU accelerated computation of the cost functions, and also incorporating techniques to improve running time and execution on the QPU. Furthermore, to improve performance, one could consider variable structure Ansätze [57, 58] or quantum-specific optimizers [59–61] for the model and training. An alternative direction, is to enlarge the suite of classical model comparison to compare the Born machine to, in order to solidify any perceived advantage and extending the model into mixed states to potentially increase the expressive power [21]. Alternatively, one could investigate methods to divide the classical-quantum resources in the learning procedure [62].

## References

[1] John Preskill. Quantum Computing in the NISQ era and beyond. *Quantum*, 2:79, August 2018. Publisher: Verein zur Förderung des Open Access Publizierens in den Quantenwissenschaften.

[2] Frank Arute et. al. Quantum supremacy using a programmable superconducting processor. *Nature*, 574(7779):505–510, October 2019.

[3] Sam McArdle, Suguru Endo, Alán Aspuru-Guzik, Simon C. Benjamin, and Xiao Yuan. Quantum computational chemistry. *Rev. Mod. Phys.*, 92(1):015003, March 2020. Publisher: American Physical Society _eprint: 1808.10402.

[4] Frank Arute et. al. Hartree-Fock on a superconducting qubit quantum computer. *arXiv:2004.04174 [physics, physics:quant-ph]*, April 2020. arXiv: 2004.04174.

[5] E Farhi, J Goldstone, and S Gutmann. A Quantum Approximate Optimization Algorithm. *arXiv Prepr. arXiv1411.4028*, 2014. _eprint: 1411.4028.

[6] Frank Arute et. al. Quantum approximate optimization of non-planar graph problems on a planar superconducting processor. *arXiv Prepr. arXiv2004.04197*, 2020.

[7] Christa Zoufal, Aurélien Lucchi, and Stefan Woerner. Quantum Generative Adversarial Networks for learning and loading random distributions. *npj Quantum Inf.*, 5(1):103, December 2019. Publisher: Nature Research _eprint: 1904.00043.

[8] Sergi Ramos-Calderer, Adrián Pérez-Salinas, Diego García-Martín, Carlos Bravo-Prieto, Jorge Cortada, Jordi Planagumà, and José I. Latorre. Quantum unary approach to option pricing. *arxiv Prepr. arXiv1912.01618*, December 2019. _eprint: 1912.01618.

[9] Patrick Rebentrost, Brajesh Gupt, and Thomas R Bromley. Quantum computational finance: Monte Carlo pricing of financial derivatives. *Phys. Rev. A*, 98(2):22321, August 2018. Publisher: American Physical Society.

[10] Samuel Mugel, Carlos Kuchkovsky, Escolastico Sanchez, Samuel Fernandez-Lorenzo, Jorge Luis-Hita, Enrique Lizaso, and Roman Orus. Dynamic Portfolio Optimization with Real Datasets Using Quantum Processors and Quantum-Inspired Tensor Networks. *arXiv Prepr. arXiv2007.00017*, June 2020. _eprint: 2007.00017.

[11] Peter W Shor. Polynomial-Time Algorithms for Prime Factorization and Discrete Logarithms on a Quantum Computer. *SIAM J. Comput.*, 26(5):1484–1509, 1997. Place: Philadelphia, PA, USA Publisher: Society for Industrial and Applied Mathematics.

[12] Aram W. Harrow, Avinatan Hassidim, and Seth Lloyd. Quantum Algorithm for Linear Systems of Equations. *Phys. Rev. Lett.*, 103(15):150502, October 2009. Publisher: American Physical Society.

[13] Alexei Kondratyev. Non-Differentiable Learning of Quantum Circuit Born Machine with Genetic Algorithm. *Available SSRN 3569226*, April 2020. Publisher: Elsevier BV.

[14] Peter J Karalekas, Nikolas A Tezak, Eric C Peterson, Colm A Ryan, Marcus P da Silva, and Robert S Smith. A quantum-classical cloud platform optimized for variational hybrid algorithms. *Quantum Sci. Technol.*, 5(2):24003, April 2020. Publisher: {IOP} Publishing.

[15] Shakir Mohamed and Balaji Lakshminarayanan. Learning in Implicit Generative Models. *arXiv1610.03483 [cs, stat]*, October 2016.

[16] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Networks. *arXiv1406.2661 [cs, stat]*, June 2014.

[17] Nathan Wiebe, Ashish Kapoor, Christopher Granade, and Krysta M Svore. Quantum Inspired Training for Boltzmann Machines. *arXiv:1507.02642 [quant-ph]*, July 2015.

[18] Song Cheng, Jing Chen, and Lei Wang. Information Perspective to Probabilistic Modeling: Boltzmann Machines versus Born Machines. *Entropy*, 20(8):583, August 2018. Publisher: MDPI AG.

[19] Marcello Benedetti, Delfina Garcia-Pintos, Oscar Perdomo, Vicente Leyton-Ortega, Yunseong Nam, and Alejandro Perdomo-Ortiz. A generative modeling approach for benchmarking and training shallow quantum circuits. *npj Quantum Inf.*, 5(1):45, May 2019.

[20] Jin-Guo Liu and Lei Wang. Differentiable learning of quantum circuit Born machines. *Phys. Rev. A*, 98(6):62324, December 2018. _eprint: 1804.04168.

[21] Guillaume Verdon, Jacob Marks, Sasha Nanda, Stefan Leichenauer, and Jack Hidary. Quantum Hamiltonian-Based Models and the Variational Quantum Thermalizer Algorithm. *arXiv Prepr. arXiv1910.02071*, October 2019. _eprint: 1910.02071.

[22] Brian 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, July 2020.

[23] Ryan Sweke, Jean-Pierre Seifert, Dominik Hangleiter, and Jens Eisert. On the Quantum versus Classical Learnability of Discrete Distributions. *arXiv:2007.14451 [quant-ph]*, July 2020. arXiv: 2007.14451.

[24] Jirawat Tangpanitanon, Supanut Thanasilp, Ninnat Dangniam, Marc-Antoine Lemonde, and Dimitris G. Angelakis. Expressibility and trainability of parameterized analog quantum systems for machine learning applications. *arXiv Prepr. arXiv2005.11222*, May 2020. _eprint: 2005.11222.

[25] Alexei Kondratyev and Christian Schwarz. The Market Generator. *Available SSRN 3384948*, 2019.

[26] Javier Alcazar, Vicente Leyton-Ortega, and Alejandro Perdomo-Ortiz. Classical versus quantum models in machine learning: insights from a finance application. *Machine Learning: Science and Technology*, 2020.

[27] Jonathan Romero and Alan Aspuru-Guzik. Variational quantum generators: Generative adversarial quantum machine learning for continuous distributions. *arXiv:1901.00848 [quant-ph]*, January 2019. arXiv: 1901.00848.

[28] Guillaume Verdon, Jason Pye, and Michael Broughton. A Universal Training Algorithm for Quantum Deep Learning. *arXiv Prepr. arXiv1806.09729*, June 2018. _eprint: 1806.09729.

[29] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.

[30] Mohammad H. Amin, Evgeny Andriyash, Jason Rolfe, Bohdan Kulchytskyy, and Roger Melko. Quantum Boltzmann Machine. *Phys. Rev. X*, 8(2):21050, May 2018. Publisher: American Physical Society _eprint: 1601.02036.

[31] Maria Kieferova and Nathan Wiebe. Tomography and Generative Data Modeling via Quantum Boltzmann Training. *Phys. Rev. A*, 96(6), December 2017.

[32] Hai Jing Song, Tieling Song, Qi Kai He, Yang Liu, and D. L. Zhou. Geometry and symmetry

in the quantum Boltzmann machine. *Phys. Rev. A*, 99(4):042307, April 2019. Publisher: American Physical Society _eprint: 1808.04567.

[33] Nathan Wiebe and Leonard Wossnig. Generative training of quantum Boltzmann machines with hidden units. *arXiv Prepr. arXiv1905.09902*, May 2019. _eprint: 1905.09902.

[34] Guillaume Verdon, Michael Broughton, and Jacob Biamonte. A quantum algorithm to train neural networks using low-depth circuits. *arXiv Prepr. arXiv1712.05304*, December 2017. _eprint: 1712.05304.

[35] Maria Schuld, Alex Bocharov, Krysta M. Svore, and Nathan Wiebe. Circuit-centric quantum classifiers. *Phys. Rev. A*, 101(3):032308, March 2020. Publisher: American Physical Society.

[36] David A. Meyer and Nolan R. Wallach. Global entanglement in multiparticle systems. *J. Math. Phys.*, 43(9):4273–4278, September 2002. Publisher: American Institute of PhysicsAIP _eprint: 0108104.

[37] Sukin Sim, Peter D. Johnson, and Alán Aspuru-Guzik. Expressibility and Entangling Capability of Parameterized Quantum Circuits for Hybrid Quantum-Classical Algorithms. *Adv. Quantum Technol.*, 2(12):1900070, December 2019. Publisher: Wiley _eprint: 1905.10876.

[38] Thomas Hubregtsen, Josef Pichlmeier, and Koen Bertels. Evaluation of Parameterized Quantum Circuits: on the design, and the relation between classification accuracy, expressibility and entangling capability. *arXiv:2003.09887 [quant-ph]*, March 2020. arXiv: 2003.09887.

[39] Gavin K. Brennen. An observable measure of entanglement for pure states of multi-qubit systems. *arXiv:quant-ph/0305094*, November 2003. arXiv: quant-ph/0305094.

[40] Diederik P Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd Int. Conf. Learn. Represent. {ICLR} 2015, San Diego, CA, USA, May 7-9, 2015, Conf. Track Proc.*, 2015.

[41] Aaditya Ramdas, Nicolas Garcia, and Marco Cuturi. On Wasserstein Two Sample Testing and Related Families of Nonparametric Tests. *arXiv:1509.02237 [math, stat]*, September 2015.

[42] Aude Genevay, Gabriel Peyre, and Marco Cuturi. Learning Generative Models with Sinkhorn Divergences. In Amos Storkey and Fernando Perez-Cruz, editors, *Proc. Twenty-First Int. Conf. Artif. Intell. Stat.*, volume 84 of *Proceedings of {Machine} {Learning} {Research}*, pages 1608–1617, Playa Blanca, Lanzarote, Canary Islands, April 2018. PMLR.

[43] Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trouve, and Gabriel Peyré. Interpolating between Optimal Transport and MMD using Sinkhorn Divergences. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proc. Mach. Learn. Res.*, volume 89 of *Proceedings of {Machine} {Learning} {Research}*, pages 2681–2690. PMLR, April 2019.

[44] Cédric Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen {Wissenschaften}. Springer-Verlag, Berlin Heidelberg, 2009.

[45] S Kullback and R A Leibler. On Information and Sufficiency. *Ann. Math. Stat.*, 22(1):79–86, 1951.

[46] Kosuke Mitarai, Makoto Negoro, Masahiro Kitagawa, and Keisuke Fujii. Quantum Circuit Learning. *Phys. Rev. A*, 98(3):32309, March 2018. _eprint: 1803.00745.

[47] Maria Schuld, Ville Bergholm, Christian Gogolin, Josh Izaac, and Nathan Killoran. Evaluating analytic gradients on quantum hardware. *Phys. Rev. A*, 99(3):32331, March 2019.

[48] Richard Sinkhorn. A Relationship Between Arbitrary Positive Matrices and Doubly Stochastic Matrices. *Ann. Math. Stat.*, 35(2):876–879, June 1964.

[49] Asja Fischer and Christian Igel. Training restricted Boltzmann machines: An introduction. *Pattern Recognit.*, 47(1):25–39, January 2014. Publisher: Pergamon.

[50] Geoffrey E Hinton. A Practical Guide to Training Restricted Boltzmann Machines. In Grégoire Montavon, Geneviève B Orr, and Klaus-Robert Müller, editors, *Neural Networks: Tricks of the Trade*, volume 7700, pages 599–619. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.

[51] D. Padilha, S. Weinstock, and M. Hodson. QxSQA: GPGPU-Accelerated Simulated Quantum Annealer within a Non-Linear Optimization and Boltzmann Sampling Framework. In *2019 IEEE High Performance Extreme Computing Conference (HPEC)*, pages 1–8, 2019.

[52] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2323, 1998. ISBN: 0018-9219 _eprint: 1102.0183.

[53] Maxwell P. Henderson and Justin Chan Jin Le. Generation of industry-relevant synthetic data using simulated quantum annealing-trained Boltzmann machines. In *QTML - Quantum Tech. Mach. Learn.*, Daejeon, South Korea, 2019.

[54] F Pedregosa, G Varoquaux, A Gramfort, V Michel, B Thirion, O Grisel, M Blondel, P Prettenhofer, R Weiss, V Dubourg, J Vanderplas, A Passos, D Cournapeau, M Brucher, M Perrot, and E Duchesnay. Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.*, 12:2825–2830, 2011.

[55] Kathleen E. Hamilton and Raphael C. Pooser. Error-mitigated data-driven circuit learning on noisy quantum hardware. *arXiv:1911.13289 [quant-ph]*, November 2019. arXiv: 1911.13289.

[56] Robert S Smith, Michael J Curtis, and William J Zeng. A Practical Quantum Instruction Set Architecture. *arXiv:1608.03355 [quant-ph]*, August 2016.

[57] Lukasz Cincio, Yiğit Subaşı, Andrew T Sornborger, and Patrick J Coles. Learning the quantum algorithm for state overlap. *New J. Phys.*, 20(11):113022, November 2018.

[58] M. Cerezo, Kunal Sharma, Andrew Arrasmith, and Patrick J. Coles. Variational Quantum State Eigensolver. *arXiv:2004.01372 [quant-ph]*, April 2020. arXiv: 2004.01372.

[59] Jonas M. Kübler, Andrew Arrasmith, Lukasz Cincio, and Patrick J. Coles. An Adaptive Optimizer for Measurement-Frugal Variational Algorithms. *Quantum*, 4:263, May 2020. Publisher: Verein zur Förderung des Open Access Publizierens in den Quantenwissenschaften.

[60] Andrew Arrasmith, Lukasz Cincio, Rolando D. Somma, and Patrick J. Coles. Operator Sampling for Shot-frugal Optimization in Variational Algorithms. *arXiv:2004.06252 [quant-ph]*, April 2020. arXiv: 2004.06252.

[61] Wim Lavrijsen, Ana Tudor, Juliane Müller, Costin Iancu, and Wibe de Jong. Classical Optimizers for Noisy Intermediate-Scale Quantum Devices. *arXiv:2004.03004 [quant-ph]*, April 2020. arXiv: 2004.03004.

[62] Marco Paini and Amir Kalev. An approximate description of quantum states. *arXiv:1910.10543 [quant-ph]*, November 2019. arXiv: 1910.10543.

[63] Karsten M Borgwardt, Arthur Gretton, Malte J Rasch, Hans-Peter Kriegel, Bernhard Schölkopf, and Alex J Smola. Integrating structured biological data by Kernel Maximum Mean Discrepancy. *Bioinformatics*, 22(14):e49–e57, 2006.

[64] Arthur Gretton, Karsten M Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex J Smola. A Kernel Method for the Two-Sample-Problem. In B Schölkopf, J C Platt, and T Hoffman, editors, *Adv. Neural Information Processing Systems 19*, pages 513–520. MIT Press, 2007.

[65] Jonas M Kübler, Krikamol Muandet, and Bernhard Schölkopf. Quantum mean embedding of probability distributions. *Phys. Rev. Res.*, 1(3):33159, December 2019.

[66] Maria Schuld and Nathan Killoran. Quantum Machine Learning in Feature Hilbert Spaces. *Phys. Rev. Lett.*, 122(4), March 2019. _eprint: 1803.07128.

[67] Vojtěch Havlíček, Antonio D Córcoles, Kristan Temme, Aram W Harrow, Abhinav Kandala, Jerry M Chow, and Jay M Gambetta. Supervised learning with quantum-enhanced feature spaces. *Nature*, 567(7747):209–212, 2019.

[68] Seth Lloyd and Christian Weedbrook. Quantum Generative Adversarial Learning. *Phys. Rev. Lett.*, 121(4):040502, July 2018. Publisher: American Physical Society.

[69] Pierre-Luc Dallaire-Demers and Nathan Killoran. Quantum generative adversarial networks. *Phys. Rev. A*, 98(1):12324, July 2018. _eprint: 1804.08641.

[70] Abhinav Anand, Jonathan Romero, Matthias Degroote, and Alán Aspuru-Guzik. Experimental demonstration of a quantum generative adversarial network for continuous distributions. *arXiv:2006.01976 [quant-ph]*, June 2020. arXiv: 2006.01976.

## Appendix A. Alternative Training Methods

Here we provide numerical results illustrating the training of the Born machine using some alternative methods and cost functions, for small numbers of qubits.

### *Appendix A.1. Maximum Mean Discrepancy*

The first alternative method is derived by using a different cost function, the so-called maximum mean discrepancy (MMD). Like optimal transport, this defines a metric on the space of probability distributions, and from which, an efficient-to-compute method of comparison can be defined [63, 64]:

$$
D(p_{\boldsymbol{\theta}}, \pi) \coloneqq \mathcal{L}_{\mathrm{MMD}} \coloneqq
$$
$$
\underset{\substack{\mathbf{x} \sim p_{\boldsymbol{\theta}} \\ \mathbf{y} \sim p_{\boldsymbol{\theta}}}}{\mathbb{E}} (\kappa(\mathbf{x}, \mathbf{y})) + \underset{\substack{\mathbf{x} \sim \pi \\ \mathbf{y} \sim \pi}}{\mathbb{E}} (\kappa(\mathbf{x}, \mathbf{y})) - \underset{\substack{\mathbf{x} \sim p_{\boldsymbol{\theta}} \\ \mathbf{y} \sim \pi}}{2\mathbb{E}} (\kappa(\mathbf{x}, \mathbf{y})) \quad \text{(A.1)}
$$

This cost function was originally utilized for hypothesis testing [64], but has since found use in training generative models. In particular, it enabled the first approach to train a QCBM [20, 22] in a differentiable way.

The function $\kappa$ is a *kernel*, which enables a means of comparison on the support spaces, $\mathcal{X}, \mathcal{Y}$. For

this work, we choose the common Gaussian mixture kernel [20] for the MMD, which is universal, and hence enables the MMD to act as a faithful method of distribution comparison:

$$\kappa_G(\mathbf{x}, \mathbf{y}) := \frac{1}{c} \sum_{i=1}^{c} \exp\left(-\frac{||\mathbf{x} - \mathbf{y}||_2^2}{2\sigma_i}\right) \qquad \text{(A.2)}$$

The parameters, $\sigma_i$, are *bandwidths* which determine the scale at which the samples are compared, and $||\cdot||_2$ is the $\ell_2$ norm. Here we choose $\sigma = [0.25, 10, 1000]$, as in [20]. Typically, the kernel is a classical function, but quantum kernels can also be considered here [22, 65–67]

### *Appendix A.2. Adversarial Discriminator*

The second method we can choose to use is to not only use a discriminator as a benchmark, but also to train the model relative to it. As in the above cases, this is a gradient based approach, with the analytic gradient taken by differentiating the discriminator loss.

Adversarial training has become a popular and powerful way to train neural networks, originating with generative adversarial networks (GANs) [16]. GANs are composed of two machine learning components, a discriminator, $\mathcal{D}$, which attempts to predict if a sample $\boldsymbol{x}$ is from a data distribution or rather has been generated by a generator network $\mathcal{G}$ (in our notation, the generator network samples from a probability distribution, $p_{\boldsymbol{\theta}}$ and is either a Born machine or a Boltzmann machine). Generalizations of the GAN into the quantum domain have also been considered [7, 27, 68–70]. The generator attempts to minimize the following loss:

$$\mathcal{L}_{\mathcal{G}} := \mathbb{E}\left[\log\left(1 - \mathcal{D}\left(\boldsymbol{x}\right)\right)\right] \approx \frac{1}{N} \sum_{i=1}^{N} \left[\log\left(1 - \mathcal{D}\left(\boldsymbol{x}_i\right)\right)\right] \tag{A.3}$$

where $\mathcal{D}\left(\boldsymbol{x}\right)$ is the probability that a discriminator, $\mathcal{D}$, guesses that $\boldsymbol{x}$ is from the real data set. The approximation to the expectation value is taken over $N$ generated samples in practice. In order to train the generator with respect to this cost function (taken with respect to a specific discriminator, $\mathcal{D}$), gradient descent can be used to minimize (A.3), with the gradient given by:

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}\left[\log\left(1 - \mathcal{D}\left(\boldsymbol{x}_i\right)\right)\right] =$$
$$\nabla_{\boldsymbol{\theta}} \sum_{\boldsymbol{x}} p_{\boldsymbol{\theta}}(\boldsymbol{x})\left[\log\left(1 - \mathcal{D}\left(\boldsymbol{x}_i\right)\right)\right] =$$
$$\sum_{\boldsymbol{x}} \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\boldsymbol{x})\left[\log\left(1 - \mathcal{D}\left(\boldsymbol{x}_i\right)\right)\right] \tag{A.4}$$

If we again assume the generator network is a Born machine, composed of quantum gates of the form

$U(\theta) = \exp(i\theta/2\Sigma)$, then using the parameter shift rule as for the Sinkhorn divergence above (14), we get:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{G}} = \frac{1}{2} \sum_{\boldsymbol{x}} p_{\boldsymbol{\theta}^+} \left[\log\left(1 - \mathcal{D}\left(\boldsymbol{x}_i\right)\right)\right] -$$
$$p_{\boldsymbol{\theta}^-} \left[\log\left(1 - \mathcal{D}\left(\boldsymbol{x}_i\right)\right)\right] \quad \text{(A.5)}$$

These expectation values can be evaluated by sampling from the parameter shifted circuit distributions, $p_{\boldsymbol{\theta}^\pm}$ as usual. Correspondingly, while a generator is trying to minimize the above cost, (A.3), the discriminator can also be trained for a number of sub-steps to become better at identifying false samples. This can be done by using gradient *ascent* to maximize the following cost:

$$\mathcal{L}_{\mathcal{D}} \approx \frac{1}{M} \sum_{j=1}^{M} \log \mathcal{D}\left(\boldsymbol{y}_j\right) + \frac{1}{N} \sum_{i=1}^{N} \log\left(1 - \mathcal{D}\left(\boldsymbol{x}_i\right)\right) \tag{A.6}$$

where the latter term is the same as in (A.4), and the former represents the probability that $\mathcal{D}$ is able to correctly identify true data samples, $\boldsymbol{y} \sim \pi$. The gradient of (A.6) can be computed similarly.

In this work, we implement the training laid out in this section with two slight variations to note:

(i) We used a slightly revised version of Eq. A.5 which dropped the 0.5 and log components (simply used $1 - \mathcal{D}$ in both terms). As we were still using Adam as the update optimizer, we believe that asserting this overall should not pose any major impact to performance. Moreover, the adversarial approach was still slower compared to the Sinkhorn divergence, and therefore did not garner increased focus in this work.

(ii) As the modeling problem in this work was extremely small, we choose to simply re-train a new discriminator model from scratch at every generative model training iteration, with corresponding test set error of the discriminator being recorded and used as a primary metric in this work. Similarly, a different discriminator model was used for calculating model parameter updates every training iteration while using adversarial training.

### *Appendix A.3. Genetic Algorithm*

Finally, we use a gradient free approach in a genetic algorithm, since this was also used to train a QCBM on this same dataset [13]. One could also choose one of the many gradient free optimisers from `scikit-learn`, as has been also done for Born machines [19]. A simplified version of a genetic algorithm was implemented in [13] due to the low number of parameters in a 12 qubit Born machine. We found that this method was significantly slower than the gradient based methods we discuss above.
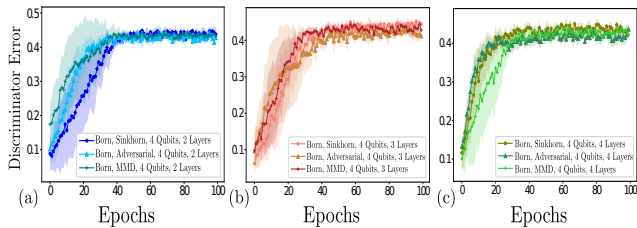
Figure A1: Training with the Sinkhorn divergence, the MMD, and the adversarial discriminator for (a) 2, (b) 3, (c) 4 layers of the hardware efficient ansatz for 4 qubits.
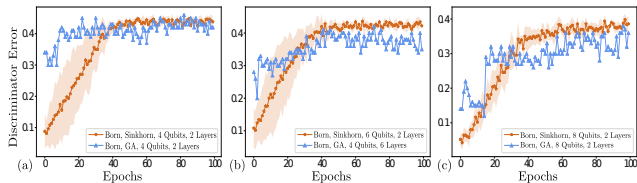


Figure A2: Training a QCBM with $4, 6$ and $8$ qubits using the Sinkhorn divergence (gradient based) versus a genetic algorithm (gradient free).

## Appendix B. Alternative Model Structures

Here we showcase the effect of using alternative model structures for the QCBM and the RBM.

### Appendix B.1. Differing numbers of Born machine layers

For completeness, in Figure B1, we show the effect of alternating the number of layers of the hardware efficient Ansätze, shown in Figure 7 for 4 and 8 qubits. In particular, we notice that increasing the number of layers does not have a significant impact, at least at these scales, except perhaps in convergence speed of the training. It is likely however, that at larger scales, increased parameter numbers would be required to improve performance.

### Appendix B.2. Differing numbers of Boltzmann hidden nodes

We also demonstrate the effect of changing the number of hidden nodes in the Boltzmann machine in Figure B2, where we have $4, 8$ and $28$ visible nodes. Again, we observe that an increasing number of hidden nodes (and by extension, number of parameters) does not substantially affect the performance of the model, in fact it can hinder it, at least when training only biases of the Boltzmann machine. In particular, it does not substantially alter the final accuracy achieved by the model. We also noticed similar behavior when also training the weights of the Boltzmann machine.
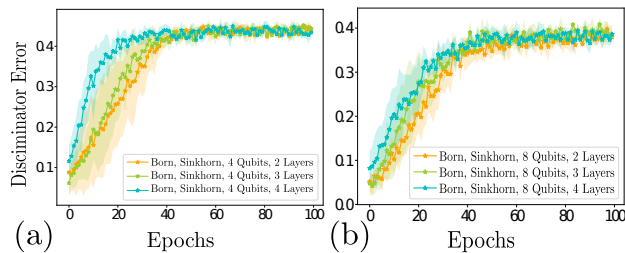


Figure B1: $2, 3$ and 4 layers of the hardware efficient ansatz for (a) 4 and (b) 8 qubits. Models are trained on 2 currency pairs at 2 and 4 bits of precision respectively. No major advantage observed for using an increasing number of layers, except perhaps in convergence speed, suggesting that 2 layers is sufficient for these problem instances.
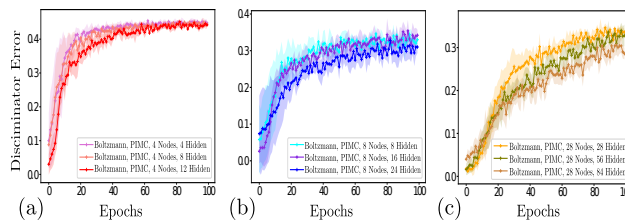


Figure B2: Increasing number of hidden nodes for RBMs with (a) 4 (b) 8 (c) 28 visible nodes. Enlarging the hidden space for the RBM again did not impact significantly for these problem sizes and in particular would not give a performance boost to outperform the QCBM.

### Appendix B.3. Weight training of Boltzmann machine

Finally, we compare the effect of weight training of the Boltzmann machine to training the bias terms alone. For the problem instances where the Boltzmann machine was able to converge to the best discriminator accuracy (i.e. in the small problem instances), we find training the weights has the effect of increasing convergence speed, and also increased accuracy where training the biases only was insufficient to achieve high discriminator error. Interestingly, we note that the Born machine still outperforms the 8 and 12 visible node RBMs, even when the weights are also trained, and this does not seem to majorly affect the performance. However, training the weights does make a large difference for 28 nodes, as seen in Figure B3(c), so again further investigation is needed in future work of this phenomenon.
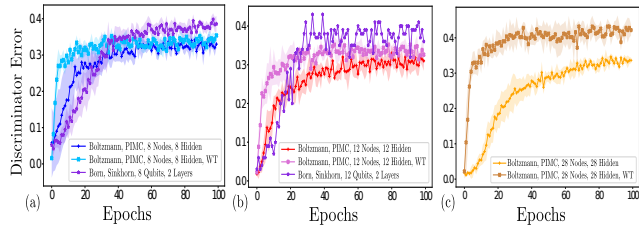
Figure B3: Weight training (WT) on the Boltzmann machine along with the node biases. We compare (a) 8, (b) 12 and (c) 28 visible node RBMs along with the corresponding Born machine. The latter uses 4 currency pairs, while the others use 2, as in the text above.