

1           **BAYESIAN INFERENCE IN QUANTUM GENERATIVE**  
2           **MODELLING\***

3           R. RAJPAL , J. ZORYK , AND U. KARADAG

4       **Abstract.** Quantum computing has the potential to revolutionize various fields, but current  
5 quantum hardware is prone to noise and errors that can limit the performance of quantum al-  
6 gorithms. This report explores using Bayesian inference to estimate the distribution of quantum  
7 circuit parameters in the Quantum Circuit Born Machine (QCBM) used for generative modeling.  
8 Bayesian inference allows uncertainty quantification and incorporation of prior knowledge. However,  
9 computing gradients in quantum circuits can be costly for large datasets due to the high-dimensional  
10 parameter space, so stochastic gradient approximations are often used. While this reduces compu-  
11 tational cost, it introduces gradient noise that can bias the parameter samples. To address this,  
12 we investigate using the Stochastic Gradient Adaptive Langevin Thermostat (SGAdLT) to reduce  
13 the bias from stochastic gradient noise. Experiments are conducted on a Gaussian mixture model  
14 and the Bars-and-Stripes dataset using simulated quantum circuits. The effect of the temperature  
15 parameter  $\beta$ , gradient batch size, and artificial gradient noise are analyzed. Results show SGAdLT  
16 helps explore local cost minima at lower  $\beta$  values compared to standard SGLD. This work demon-  
17 strates the potential of combining Bayesian inference with quantum machine learning to enhance the  
18 robustness of variational quantum algorithms in the presence of hardware noise and computational  
19 limitations. Further research directions are also discussed.

20      **1. Motivation.** Quantum computing has the potential to revolutionize various  
21 fields, including cryptography, drug development, and optimization. However, the  
22 current state of quantum hardware is still prone to noise and errors, which can limit the  
23 performance and reliability of quantum algorithms. Variational quantum algorithms,  
24 such as the Quantum Circuit Born Machine (QCBM) [1] used for generative modeling,  
25 offer a promising approach to mitigate these issues by optimizing the parameters of  
26 the quantum circuit to minimise the impact of noise and improve robustness.

27      Bayesian inference provides a principled framework for estimating the distribu-  
28 tion of quantum circuit parameters, allowing for the quantification of uncertainty and  
29 the incorporation of prior knowledge [2]. By sampling from the posterior distribu-  
30 tion of circuit parameters using sampling methods, such as Langevin Monte Carlo  
31 or Hamiltonian Monte Carlo (HMC), we can explore the complex probability land-  
32 scape and identify optimal parameter configurations that enhance the performance  
33 and reliability of the quantum circuit.

---

\*RSCAM Group Project, Supervisor: Prof. Benedict Leimkuhler

34        However, computing gradients in variational quantum algorithms can be costly,  
35 especially for large datasets. To address this challenge, stochastic gradient approx-  
36 imations are often employed [2], where a subset of the data is used to estimate the  
37 gradient. While this approach reduces computational overhead, it introduces gradi-  
38 ent noise, which can ultimately bias the samples and lead to suboptimal parameter  
39 estimates.

40        To overcome this limitation, various debiasing strategies have been proposed in  
41 the literature such as: stepsize decay and Multilevel Monte-Carlo [3]. These tech-  
42 niques aim to correct the bias introduced by stochastic gradient approximations and  
43 improve the quality of the samples. In this project, we use the Stochastic Gradient  
44 Adaptive Langevin Thermostat (SGAdLT) [4], also known as the stochastic gradient  
45 Nosé-Hoover thermostat, as our debiasing strategy.

46        SGAdLT is an adaptive sampling method that adjusts the noise level of the sam-  
47 pler based on the observed gradient noise. By incorporating a thermostat variable  
48 that controls the temperature of the system, SGAdLT can effectively counteract the  
49 bias introduced by stochastic gradients [4]. The method dynamically adapts the noise  
50 level to ensure that the samples are drawn from the correct target distribution, even  
51 in the presence of gradient noise.

52        The motivation for this project is to explore the effectiveness of SGAdLT in  
53 the context of generative quantum models, specifically focusing on its application  
54 in Quantum Circuit Born Machine (QCBM). By investigating the performance of  
55 SGAdLT in reducing the bias caused by stochastic gradient approximations, we aim  
56 to improve the quality of samples, leading to more accurate and robust parameter  
57 estimates.

58        This research has the potential to enhance the reliability and performance of  
59 quantum models in the presence of hardware noise and computational limitations.  
60 By developing more robust Bayesian inference methods for quantum circuits, we can  
61 contribute to the advancement of quantum computing and its practical applications  
62 in various domains.

63        **2. Background.**

64        **2.1. Generative modelling.** Generative AI has seen an explosion in applica-  
65 tions over the past decades with examples such as: LLMs like ChatGPT, Diffusion  
66 Models in Images, Audio Generation models etc [5]. The premise of generative model-

67     ling is that we have samples  $X$  from an unknown distribution  $P$  i.e.,  $X \sim P$  with the  
 68     aim to generate samples  $Y$  that look like  $X$  by sampling  $Y \sim Q$  from some tuneable  
 69     distribution  $Q$ .

70     The data distribution  $P$  represents the true, underlying distribution from which  
 71     the observed data samples  $X$  are drawn. In real-world scenarios, the data distribution  
 72     is often unknown and can only be approximated using the available data samples. The  
 73     goal of generative modelling is to learn a model that can generate new samples that  
 74     closely resemble the data distribution.

75     The generated distribution  $Q$  is a parametric distribution that the generative  
 76     model learns to approximate the data distribution  $P$ . The parameters of the gener-  
 77     ated distribution are adjusted during the training process to minimize the difference  
 78     between the generated samples and the data samples. Common choices for the gen-  
 79     erated distribution include Gaussian distributions, Bernoulli distributions, or more  
 80     complex distributions like mixture models and implicit distributions defined by neu-  
 81     ral networks. To quantify the similarity in distributions, practitioners often use a cost  
 82     function or a scoring rule such as Maximum Mean Discrepancy (MMD), where it is 0  
 83     when distributions  $P$  and  $Q$  are identical. See Appendix C for more detail on MMD.

84       **2.2. Bayesian Learning.** Consider parameters  $\boldsymbol{\theta}$  as random variables with a  
 85     posterior distribution  $\pi(\boldsymbol{\theta}|\mathcal{D})$ , prior  $p(\boldsymbol{\theta})$ , likelihood  $f(\mathcal{D}|\boldsymbol{\theta})$  and data  $\mathcal{D}$ . Using Bayes  
 86     formula:

$$87 \quad (2.1) \quad \pi(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta})f(\mathcal{D}|\boldsymbol{\theta}) = p(\boldsymbol{\theta}) \exp(-\beta C(\boldsymbol{\theta}))$$

88     where  $\beta$  is the reciprocal temperature and  $C(\boldsymbol{\theta})$  is a cost function such as MMD.

89     The goal becomes sampling  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N \sim \pi(\boldsymbol{\theta}|\mathcal{D})$  to approximate the predictive  
 90     distribution with a Bayesian Model Average (BMA). See Appendix B for a visual  
 91     representation of BMA.

$$92 \quad (2.2) \quad p_{\text{BMA}}(y|\mathcal{D}) = \int_{\boldsymbol{\theta} \in \mathbb{R}^n} p(y|\mathcal{D}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} \approx \frac{1}{N} \sum_{i=1}^N p(y|\boldsymbol{\theta}_i)$$

93     Working in log-space, we use gradient-based stepping methods (e.g., Langevin  
 94     Monte Carlo) to sample from the posterior  $\pi(\boldsymbol{\theta}|\mathcal{D})$ :

$$95 \quad \nabla_{\boldsymbol{\theta}} \log \pi(\boldsymbol{\theta}|\mathcal{D}) = \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}) - \beta \nabla_{\boldsymbol{\theta}} C(\boldsymbol{\theta}).$$

96        The reciprocal temperature  $\beta$  plays a crucial role in shaping the posterior distribution  
97        and controlling the exploration-exploitation trade-off during sampling. When  
98         $\beta$  is high (i.e., the temperature is low), the posterior distribution becomes sharper  
99        and more peaked around the modes, as shown in Figure D.1(a). This can be useful  
100      for focusing inference on specific regions of interest or for identifying modes of the  
101      distribution with high confidence, a process known as annealing [6]. Annealing allows  
102      the sampling algorithm to gradually converge to the most probable regions of the  
103      parameter space as  $\beta$  increases.

104        On the other hand, when  $\beta$  is low (i.e., the temperature is high), the posterior  
105      distribution becomes smoother and more diffuse, as illustrated in Figure D.1(b). This  
106      can make the distribution more tractable for sampling algorithms by reducing the  
107      number of local maxima and making the distribution less peaked, a process called  
108      tempering [7]. Tempering allows for better exploration of the parameter space and  
109      can help prevent the sampling algorithm from getting stuck in local optima.

110        In the context of Stochastic Gradient Langevin Dynamics (SGLD) and Stochastic  
111      Gradient Adaptive Langevin Thermostats (SGAdLT), which will be introduced  
112      later, the choice of  $\beta$  can be interpreted as controlling the balance between the stochastic  
113      gradient noise and the Langevin dynamics. A higher  $\beta$  value will result in  
114      a more deterministic update rule, emphasizing exploitation of the current best solutions,  
115      while a lower  $\beta$  value will introduce more stochasticity, promoting exploration  
116      of the parameter space.

117        **2.3. Quantum Circuit Born Machine (QCBM).** Quantum Circuit Born  
118      Machine (QCBM), a special subclass of Parameterized Quantum Circuits (PQCs),  
119      can be used as an Ansatz <sup>1</sup> for generative modeling. A PQC consists of a sequence of  
120      parameterized quantum gates applied to a set of qubits. By adjusting the parameters  
121      of the gates, the PQC can generate different quantum states.

122        In a QCBM, the PQC is used to generate a probability distribution over classical  
123      bit strings. According to the Born rule, the probability of measuring a particular bit  
124      string  $x$  is given by  $p_{\theta}(x) = |\langle x | \psi_{\theta} \rangle|^2$ , where  $|\psi_{\theta}\rangle$  is the quantum state generated  
125      by the PQC with parameters  $\theta$ . By sampling from this distribution, the QCBM can

---

<sup>1</sup>Ansatz: Assuming a particular structure for the solution and optimizing based on that. E.g. Assume a polynomial ansatz  $y(x) = ax^2 + bx + c$  for the solution of the ODE  $\frac{d^2y}{dx^2} - 4y = 0$  and optimize w.r.t  $a, b, c$ .

126 generate classical data. Figure 2.1 shows an example QCBM with 8 qubits and a  
 127 circuit depth of 3.

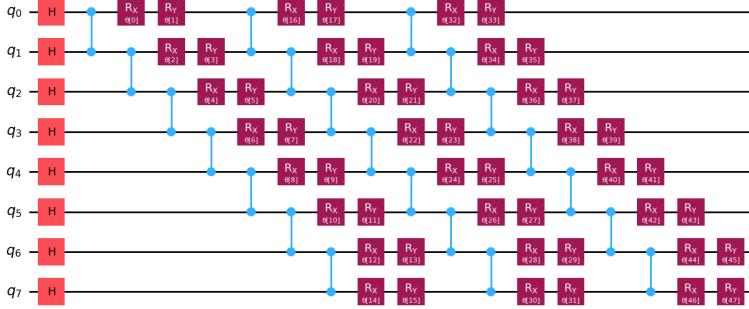


FIG. 2.1. *Quantum Circuit Born Machine with 8 Qubits and a Circuit Depth of 3.*

128 QCBMs have several potential advantages over classical generative models. First,  
 129 they can efficiently represent certain probability distributions that are difficult to  
 130 model classically, such as those with strong correlations or entanglement between  
 131 variables. Second, they may be more expressive than classical models with a similar  
 132 number of parameters, due to the exponential size of the Hilbert space.

133 However, training QCBMs can be challenging due to the following issues:

134 Reasons why  $\beta$  from 2.2 may be helpful:

- 135 • Exponentially many local minima [8]
- 136 • Barren Plateau Phenomenon [9]

137 To address these challenges, we propose using Bayesian learning to train QCBMs.  
 138 By introducing a prior distribution over the parameters and updating it based on the  
 139 observed data, we can potentially avoid getting stuck in local minima and escape  
 140 barren plateaus. The posterior distribution can be approximated using Langevin  
 141 dynamics, as discussed in Section 2.2.

142 Fig 2.2 shows an overview of our quantum Bayesian learning approach. The  
 143 QCBM, denoted by  $U(\theta)$ , generates samples from a probability distribution parame-  
 144 terized by  $\theta$ . The prior distribution  $p(\theta)$  is updated based on the observed data  $\mathcal{D}$  to  
 145 obtain the posterior distribution  $\pi(\theta|\mathcal{D})$ . Langevin dynamics is used to sample from  
 146 the posterior and update the parameters of the QCBM.

147 The gradients of the cost function  $C(\theta)$  required for Langevin dynamics can  
 148 be computed using the parameter-shift rule, which is a finite-difference method for

149 estimating gradients on quantum circuits. This may add a deterministic bias in the  
 150 estimated gradients. However, for the purposes of this project, we assume that this  
 151 is exact as it is proven to be unbiased in [1]. See Appendix E for details.

152 In summary, QCBMs are a promising approach for generative modeling on quan-  
 153 tum computers, but training can be challenging due to the complex cost landscape.  
 154 By combining QCBMs with Bayesian learning, we aim to develop a more robust and  
 155 effective training method.

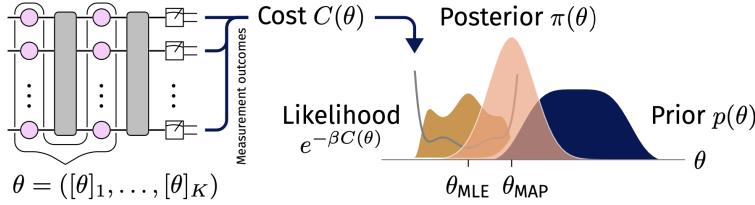


FIG. 2.2. Overview of quantum Bayesian learning for QCBMs [2].

156 **3. Stochastic Gradient Langevin Dynamics (SGLD).** Langevin Dynamics  
 157 is used to sample the posterior  $\pi(\theta|D)$ . The parameter space of QCBM is quite  
 158 high-dimensional, especially as the number of qubits increase. The QCBM in Fig 2.1  
 159 has parameter  $\theta$  that is  $2^8 = 64$ -dimensional. In generative modelling, we often have  
 160 many samples of data  $D$ . It can become quite computationally costly to compute the  
 161 gradient  $\nabla_{\theta} f(\theta)$ . Thus, it is beneficial to compute stochastic gradients.

162 SGLD [10] is based on the following stochastic differential equation:

$$163 \quad d\theta = -\frac{1}{\gamma} \nabla \tilde{C}(\theta) dt + \sqrt{\frac{2}{\gamma\beta}} dW$$

$$164 \quad = -\frac{1}{\gamma} \left( \frac{1}{|S|} \sum_{i \in S} \nabla C_i(\theta) \right) dt + \sqrt{\frac{2}{\gamma\beta}} dW$$

165 where  $S \subseteq 1, \dots, n$  is a random subset of indices.

- 166 •  $\theta$ : circuit parameters  
 167 •  $\gamma$ : friction coefficient, inverse step-size  
 168 •  $\tilde{C}(\theta)$ : stochastic gradient of the cost function  
 169 •  $\beta$ : reciprocal temperature  
 170 •  $dW$ : standard Wiener process

171       **4. Gradient Noise Model (GNM).** One of the biggest problems about the  
 172 modelling approach here is that since it is computationally expensive the generative  
 173 model utilizes the Stochastic Gradient at each step. This introduces bias to the  
 174 approximations for the gradients. Below Fig 4.1 illustrates the component-wise (mar-  
 175 ginal) distributions of bias vector. The methodology we used to conduct this analysis  
 176 was to take 2000 stochastic gradients with a batch size of 200 for each of the 64 com-  
 177 ponents then we obtain the bias value for that stochastic gradient by subtracting the  
 178 best approximation for the real gradient (i.e. using all samples to approximate the  
 179 gradient):

180       
$$h_i = \{(\nabla \tilde{C}_{|S_j|}(\boldsymbol{\theta}) - \nabla C(\boldsymbol{\theta}))_i\}_{j=1}^{2000} \text{ for } i = 1, \dots, 64.$$

181       where  $h_i$  is the bias vector for the respective component. The histograms just  
 182 illustrate the distributions of  $h_1, \dots, h_{64}$ .

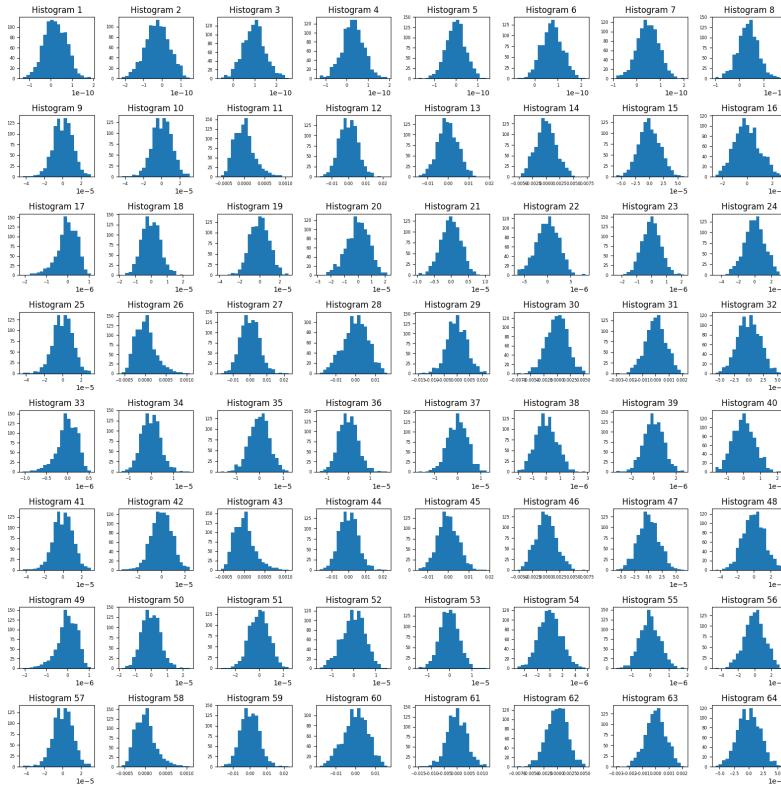


FIG. 4.1. *Gradient Noise Model for a Bimodal Gaussian with 5000 samples and batch size = 200. The component-wise biases appear to resemble the Gaussian distribution, even though some histograms are positively/negatively skewed or seem to have considerably different variations.*

183       **5. Stochastic Gradient Adaptive Langevin Thermostat (SGAdLT).**

184       **5.1. Introduction.** The Ad-Langevin Thermostat [4] debiases under the **as-**  
 185 **sumption of Gaussian batch noise with constant covariance.**

186       It is based on the following coupled system of SDEs:

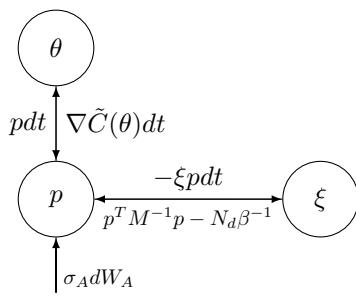
187       (5.1)                    $d\theta = \mathbf{M}^{-1} \mathbf{p} dt,$

188       (5.2)                    $d\mathbf{p} = \nabla \tilde{C}(\theta) dt - \xi \mathbf{p} dt + \sigma_A \sqrt{\mathbf{M}} d\mathbf{W}_A,$

189       (5.3)                    $d\xi = \mu^{-1} [\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} - N_d \beta^{-1}] dt$

- 190       •  $\theta$ : circuit parameters
- 191       •  $\mathbf{p}$ : momentum variable
- 192       •  $\mathbf{M}$ : mass matrix
- 193       •  $\nabla \tilde{C}(\theta)$ : stochastic gradient of the cost function
- 194       •  $\xi$ : auxillary variable
- 195       •  $\sigma_A$ : additive noise
- 196       •  $d\mathbf{W}_A$ : standard Wiener process
- 197       •  $\mu$ : thermal mass
- 198       •  $N_d$ : number of degrees of freedom
- 199       •  $\beta = (k_B T)^{-1}$ : reciprocal temperature

200       **5.2. Physical Interpretation.** Below is a diagram representing the interac-  
 201       tions between the variables.



203       Key interactions:

- 205       •  $\xi$  acts as the friction on  $\mathbf{p}$ ,
- 206       • Additive noise,  $\sigma_A$ , injects heat to the system,
- 207       • The dynamics of  $\xi$  are driven by the thermostat (kinetic energy difference)  
 208        $(\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} - N_d \beta^{-1})$ .

209      Negative feedback loop enables adaptive noise dissipation.

210      **5.3. Splitting the Ad-Langevin Thermostat.** Since stochastic systems in  
 211      most of the cases cannot be solved “exactly,” splitting methods are often adopted in  
 212      practice. For instance here, the vector field of the Ad-Langevin can be split into four  
 213      pieces which are denoted as “A,” “B,” “O,” and “D” [4]:

$$\begin{aligned} d \begin{bmatrix} \theta \\ p \\ \xi \end{bmatrix} &= \underbrace{\begin{bmatrix} M^{-1}p \\ 0 \\ 0 \end{bmatrix}}_A dt + \underbrace{\begin{bmatrix} 0 \\ -\nabla \tilde{C}(\theta) + \sigma M^{1/2} R \\ 0 \end{bmatrix}}_B dt \\ &+ \underbrace{\begin{bmatrix} 0 \\ -\xi pdt + \sigma_A M^{1/2} dW_A \\ 0 \end{bmatrix}}_O + \underbrace{\begin{bmatrix} 0 \\ 0 \\ \mu^{-1}(p^T M^{-1} p - N_d \beta^{-1}) \end{bmatrix}}_D dt \end{aligned}$$

214      In this scheme we note that  $R$  appears as opposed to the SDE representation  
 215      from 5.1 to clear any confusion, here  $R$  is a vector of i.i.d. standard normal random  
 216      variables [4]. For our experiments, we chose to use the **BADODAB** scheme as used  
 217      in the Bayesian Logistic Regression example from [4].

219     **6. Experiments.** (Code can be found [here](#))

220     **6.1. Datasets.** We use the BARS-and-Stripes [11] Dataset and a Gaussian Mix-  
221     ture model. We used the noiseless simulator Qujax [12] to simulate our experiments.

222     **6.1.1. Gaussian Mixture Model.** Our first empirical data distribution is of  
223     sample size 5000 and is a non-symmetric bimodal Gaussian with dimensionality = 1.

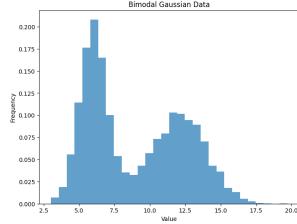


FIG. 6.1. *Bimodal Gaussian*

224     For our experiments on this dataset, unless otherwise mentioned, the default  
225     hyperparameters are listed in Appendix F.

226     **6.1.2. BARS-and-Stripes (BAS).** The Bars-and-Stripes (BAS) dataset is a  
227     synthetic dataset commonly used to benchmark generative models, particularly in the  
228     quantum computing domain [11]. The dataset consists of images with either vertical  
229     bars or horizontal stripes on a grid of pixels, as seen in Figure 6.2. For an  $n \times m$  grid,  
230     the total number of valid BAS patterns is  $N_{BAS} = 2^n + 2^m - 2$ . To generate samples  
231     from the BAS dataset, we define a target distribution  $p(x)$  where  $p(x_i) = 1/N_{BAS}$   
232     if  $x_i$  is a valid BAS image, and 0 otherwise. The prior distribution is assumed to  
233     be uniform, meaning that all valid BAS patterns are expected to be generated with  
234     equal probability. The BAS patterns are encoded into the qubits, with each data  
235     qubit representing a pixel of the BAS image. By applying a series of quantum gates  
236     and sampling the circuit parameters, the QCBM can generate samples that resemble  
237     the target BAS distribution.

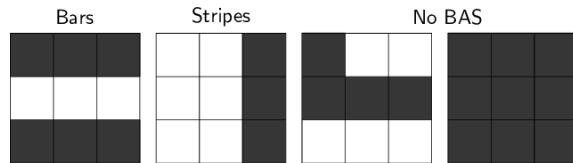


FIG. 6.2. *Bars and Stripes Dataset*

238        **6.1.3. Experiment 1: Expressivity of Circuit.** The expressive power of a  
 239 generative learning machine directly determines how well the generated distribution  
 240 can match the target distribution. The stronger the expressive power, the smaller the  
 241 dissimilarity between the two distributions will be. In this experiment, we investigated  
 242 the relationship between the expressivity of the quantum circuit by varying the num-  
 243 ber of qubits  $Q \in \{2, 4, 6, 8, 10\}$  and the circuit depth  $D \in \{1, 2, 3, 4, 5\}$ . The quantum  
 244 circuits are composed of  $D$  blocks, where each block implements a unitary operator  
 245  $U(\theta_i)$ . A unitary operator  $U(\theta) = \prod_{i=1}^D U(\theta_i)$  is applied to  $Q$  input qubits. The  
 246 experiment was conducted with the following hyper-parameters: number of training  
 247 steps  $N = 5000$ , step size  $h = 0.1$ , and reciprocal temperature  $\beta = 100$ . The initial  
 248 parameters of the quantum circuit were randomly initialized within a small range  
 249  $(-0.001/\pi, 0.001/\pi)$  to avoid vanishing or exploding gradients. Figure A.2 displays  
 250 the Predictive Distributions with a burn-in of 1000. We observe that for  $Q = 2$ , there  
 251 is no visible Predictive Distribution  $p(X)$ , while for  $Q = 4$ , the  $p(x)$  ends at  $x = 15$ .  
 252 This is because increasing  $Q$  increases the  $p(x)$  space sampling grid points by a factor  
 253 of  $2^Q$ , namely  $p(x) \rightarrow [0, 2^Q - 1]$ , given that the grid points are evenly spaced. How-  
 254 ever, due to limitations of the Qujax package used for the simulation, increasing the  
 255 number of qubits beyond a certain point may not further improve expressivity. To  
 256 overcome this, we propose mapping the qubits to a targeted sub-domain of the range,  
 257 such as  $[4, 20]$ , to allow for a more focused generative model of the target distribution.  
 258 Regarding the circuit depth  $D$ , we observe that increasing it improves expressivity but  
 259 also increases the computational cost per iteration, since the number of parameters  
 260 grows by  $|\theta| = 2Q(D + 1)$ .

261       **6.1.4. Experiment 2: Adding artificial noise  $\sigma$  to gradients.** Assume  
 262        $\tilde{C}(\boldsymbol{\theta}) = C(\boldsymbol{\theta}) + \alpha\boldsymbol{\varepsilon}$  where  $\alpha$  is multiplicative factor and  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I})$ . We assume in  
 263       this case that the noise induced by taking stochastic gradients is equivalent to adding  
 264       standard Gaussian noise. This gradient noise model is Gaussian with identity covari-  
 265       ance which satisfies the assumption of Adaptive Langevin Thermostat 5. We would  
 266       hence like to use this experiment as a proof of concept for the Adaptive Langevin  
 267       Thermostat. First of all, take a look at Fig A.3. The marginal distribution of noises  
 268       are all standard normal Gaussians i.e., equal scale equal to one and zero mean. As an  
 269       extreme case, we run experiments with  $\alpha = 10$  to demonstrate the benefit the ther-  
 270       mostat has on sampling that overdamped Langevin dynamics does not offer. This is  
 271       an extremely high gradient noise regime. Referring to Figures A.4, A.6, it is apparent  
 272       that SGLD's cost curve (Maximum Mean Discrepancy) is unable to sample from a  
 273       region of low cost (or high probability). Thus, the predictive distribution computed is  
 274       quite flat. It is nearly a uniform predictive distribution which is quite poor. Interest-  
 275       ingly though, it is able to generate two small modes similar to our data distribution.  
 276       On the other hand, quite pleasantly, the cost curve appears to be sampling from a  
 277       region of low cost (high probability) and attains a fairly accurate predictive distribu-  
 278       tion. It does however overfit to some extent as it contains more than two peaks as  
 279       indicated in Figures A.5, A.7. It is likely overfitting to the kinks in the data distribu-  
 280       tion as indicated in the histogram as opposed to learning the underlying signal. This  
 281       is perhaps due to the low temperature i.e., high value of  $\beta = 100$  which has the effect  
 282       of annealing the distribution resulting in overconfident predictions. We will look into  
 283       this effect of  $\beta$  in the following experiment.

284       **6.1.5. Experiment 3: The effect of  $\beta$ .** As indicated in Fig D.1, increasing  $\beta$   
285 has the effect of annealing the distribution while decreasing  $\beta$  has the effect of tem-  
286 pering the distribution. This forms a tradeoff between exploration and exploitation as  
287 discussed in Sec 2.2. Referring to Fig A.8 and A.9, we observe the loss curve of SGLD  
288 (row 1) and SGAdLT (row 2) where we take stochastic gradients of batch size 200.  
289 Note that we are no longer adding artificial noise to the gradients. The gradient noise  
290 model now looks like Fig 4.1 which looks approximately Gaussian with non-constant  
291 covariance, thus violating the assumptions of Adaptive Langevin Thermostat. How-  
292 ever, according to [13], it is possible for the method to work for even non-Gaussian and  
293 non-constant covariance with a minor modification. We decide to proceed with the  
294 experiments despite the violation in assumption by varying  $\beta \in \{1, 100, 1000, 10000\}$ .

295       As we increase  $\beta$ , SGLD begins to converge very quickly to nearly zero cost for  
296  $\beta = 10000$ , with very little variation. This implies that it acts as an optimizer for very  
297 large  $\beta$ . For  $\beta = 1$ , which is the standard temperature regime, it is unable to sample  
298 from a region of high probability. Note the corresponding predictive distributions in  
299 Fig A.10 and A.11. For  $\beta = 1$  and 100, the predictive distribution appears to be  
300 nearly uniform. The predictive distribution for  $\beta = 1000$  begins to show some, albeit  
301 poor, fitting to the data distribution. At  $\beta = 10000$ , it appears to give an excellent  
302 predictive distribution. However, this is most likely because it is sampling just around  
303 the maximum a posteriori estimate  $\theta_{MAP}$ . This is an insufficient exploration of the  
304 parameter space  $\theta$ . We cannot rely on it to provide us with accurate uncertainty  
305 estimates. When we baseline against Stochastic Gradient Adaptive Langevin Ther-  
306 mostat, we observe that for  $\beta = 1$ , it is unable to sufficiently explore a local basin  
307 of attraction, thus yielding a uniform predictive distribution. However, for  $\beta = 100$   
308 and higher, we observe that it begins to explore to some extent around the local min-  
309 ima, thus yielding reasonable predictive distributions. However, these appear to be  
310 overfitting to the data distribution. This is a standard effect of annealing. Moreover,  
311 we have chosen an arbitrary set of default parameters for  $\mu^{-1}, \sigma_A, etc.$  It would be  
312 beneficial to first find the accurate values for these hyper-parameters. On a positive  
313 note however, we observe that SGAdLT is able to begin exploring local minima for  
314 lower  $\beta$  than for SGLD. This is a promising result!

315        **6.1.6. Experiment 4: Varying batch size to see how the GNM varies.**

316        As indicated in Sec 4, the smaller our batch size, the noisier our stochastic gradients  
317        are. The gradient noise model in Fig 4.1 is of batch size 200 for a specific configuration  
318        of parameters  $\theta_{init}$ . We wish to study what kind of distribution is formed by the errors  
319        of the stochastic gradients and how they vary with batch size. The Adaptive Langevin  
320        Thermostat method assumes zero-mean Gaussian noise with constant covariance. We  
321        hypothesize that the scale of the error distributions should get smaller as the batch  
322        size increases since there is less uncertainty and hence less variance. We varied our  
323        batch size  $S \in \{100, 200, 500, 1000, 2000\}$ .

324        To analyze the gradient noise model, we plot the marginal histograms with 2000  
325        samples each. Due to our quantum circuit having 64 parameters to optimize over, we  
326        only present batch size 100 and 2000 to represent the disparity in Fig ???. The rest can  
327        be found in Experiment 4 notebook of the Github repository. After assessing the grids  
328        thoroughly, we were unable to find a discernible pattern in the change of the scales  
329        as the batch size varied. It is unclear to us why this is the case. Some of the scales  
330        of the parameters remain the same, some increase or decrease arbitrarily irrespective  
331        of whether the batch size is decreasing or increasing. However, we must recall that  
332        these marginals only represent the joint distribution of the gradient noise if they are  
333        all independent. In the case of Gaussian distributions, being uncorrelated implies  
334        independence of the marginals. First of all, it is unclear whether this distribution  
335        is Gaussian. If we make the assumption that it is Gaussian based on a qualitative  
336        visual assessment of the histograms, the arbitrary change of scale with batch size  
337        indicates that the components of the gradient noise are correlated. Thus, it would  
338        make more sense to validate whether the joint distribution of the gradient noise model  
339        is Gaussian with tests like the Kolmogorov-Smirnov test. However, this is expensive  
340        in high-dimensions (i.e., 64) and due to the time constraint, we leave this for future  
341        work. The stabilizing thermostat term  $\xi$  can also be studied more extensively to  
342        derive key insights.

343        We must also note that many of the scales have order of magnitudes in the range  
344        of  $[-6, -3]$ , which is very small. This means that the noise introduced by the stochastic  
345        gradients is quite small. This is most likely because our dataset is 1-dimensional. The  
346        bias introduced by the stochastic gradients may thus be very small as well.

347       **6.1.7. Experiment 5: The effect of stepsize  $h$ .** In this experiment, we in-  
348 vestigated the effect of the stepsize  $h$  on the performance of SGLD and SGAdLT.  
349 We varied  $h \in \{0.01, 0.05, 0.1, 0.5, 1.0\}$  while keeping other hyper-parameters fixed  
350 at: number of qubits  $Q = 8$ , circuit depth  $D = 3$ , reciprocal temperature  $\beta = 100$ ,  
351 training steps  $N = 5000$ , and batch size = 200. Figure ?? shows the MMD loss  
352 curves for SGLD and SGAdLT at different step sizes. For smaller  $h$  (0.01 and 0.05),  
353 both methods exhibit slower convergence. As  $h$  increases to 0.1, convergence speed  
354 improves. However, large step sizes (0.5 and 1.0) lead to instability and fluctuations,  
355 more pronounced in SGLD. The choice of  $h$  in SGLD and SGAdLT controls the bal-  
356 ance between stochastic gradient noise and Langevin dynamics, similar to the role of  
357  $\beta$ . A smaller  $h$  promotes exploitation of the current best solutions, while a larger  $h$   
358 encourages exploration of the parameter space. In the context of sampling,  $h$  controls  
359 the "energy landscape" of the posterior distribution. Smaller  $h$  values result in a more  
360 rugged landscape with well-defined minima, emphasizing the exploitation of local op-  
361 tima. Larger  $h$  values lead to a smoother landscape, allowing for better exploration  
362 and preventing the algorithm from getting stuck in suboptimal regions. SGAdLT gen-  
363 erally achieves lower MMD values and faster convergence than SGLD across different  
364 step sizes due to its adaptive noise adjustment based on the observed gradient noise.  
365 The optimal  $h$  depends on the problem and may require tuning. Further work may  
366 include using an adaptive step size control.

367       **6.1.8. Experiment 6: Thermal Mass  $\mu$  vs Additive noise  $\sigma_A$ .** Previously,  
368       we arbitrarily set  $\mu^{-1} = \sigma_A = 1$ . We perform a grid search to determine the optimal  
369       hyperparameter configuration  $(\mu^{-1}, \sigma_A)$  at  $\beta = 100$ . Thermal mass  $\mu$  regulates how  
370       particles interact with the heat bath, affecting equilibration rate. Additive noise  
371        $\sigma_A$  refers to random fluctuations added to the particle dynamics to mimic thermal  
372       fluctuations from the environment. Dynamically adjusting these two hyperparameters  
373       optimizes sampling efficiency and convergence in molecular dynamics simulations. We  
374       vary  $(\mu^{-1}, \sigma_A) \in \{0.01, 0.1, 1, 10\}^2$ .

375       Figure A.18 contains the cost curves and Bayesian model averages in a grid for-  
376       mat. Based on the cost curves, it appears that as  $\sigma_A$  increase while  $\mu^{-1}$  is kept  
377       constant (Row 1), there is more noise in the system, thus resulting in an inability to  
378       explore a region of high probability. As we increase  $\mu^{-1}$  (i.e., decrease the thermal  
379       mass), we coerce exploration of a region of low cost. There is thus a tradeoff between  
380       exploration and exploitation similar to that of  $\beta$ . We argue that for our purposes,  
381       the configurations  $(\mu^{-1}, \sigma_A) = (0.1, 1.0)$  or  $(1.0, 1.0)$  achieves the desired tradeoff for  
382       exploration while remaining in a region of low cost. We then compare their Bayesian  
383       Model Averages and conclude that the configuration  $(\mu^{-1}, \sigma_A) = (0.1, 1.0)$  is more  
384       favorable than  $(\mu^{-1}, \sigma_A) = (1.0, 1.0)$  because it overfits to kinks in the histogram. We  
385       proceed to use this configuration in future experiments. This decision, we admit, is  
386       somewhat subjective and may benefit from further objective and empirical analysis.

387       **6.1.9. Experiment 7: Preconditioning using mass matrix  $M$ .** From Ex-  
 388       periment 4, we know that the Gradient noise model has different variance along each  
 389       component and are possibly correlated. To remedy this, we can precondition the  
 390       algorithm using mass matrix  $M$  [14]. The choice of  $M$  can have drastic effects on  
 391       convergence and the quality of the samples, particularly for noise distributions which  
 392       are very anisotropic. A typical choice of  $M$  is the covariance of the noise distribution  
 393       in the convex setting. The effect of  $M^{-1}$  is to make the dynamics affine invariant. In  
 394       summary, we should expect Langevin Dynamics to be most efficient when applied to  
 395       roughly isotropic distributions with roughly Gaussian tail behavior. We decide to use  
 396       a batch size of 200 and approximate its stochastic gradient noise distribution with a  
 397       diagonal matrix  $M$  where

$$M = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_{64} \end{bmatrix}$$

399       where  $\sigma_i$  is the variance of histogram  $h_i$  for batch size 200 from Fig 4.1.

400       We set  $\mu^{-1} = 0.1$  and  $\sigma_A = 1.0$  based on Experiment 6 and re-run our experiment  
 401       for the bimodal Gaussian with  $\beta = 100$ . Regrettably, this did not yield favorable  
 402       results as indicated in Fig A.19. As indicated by the figure, it never explores a region  
 403       (basin) of high probability (low cost). It is unclear why this might be the case.  
 404       Potential reasons may include:

- 405       • The diagonal approximation to the covariance matrix of the noise distribution  
 406       is a crude one, due to the correlation between components.
- 407       • The scales of the histograms are of order of magnitude between  $[-6, -3]$  which  
 408       may result in numerical errors when performing operations with the inverted  
 409       matrix such as  $M^{-1}p$ .
- 410       • We assumed that the gradient generated by the parameter shift rule in Ap-  
 411       pendix E is exact, which may not be the case, resulting in erroneous variances  
 412        $\sigma_i$ .

413       **6.2. BARS-and-STRIPES.** To evaluate the performance of the QCBM on the  
414       BAS dataset, we conducted experiments on 2x2 and 3x3 grid sizes. The QCBM was  
415       trained using the SGAdLT algorithm with a step size of  $h = 0.5$  and 20,000 training  
416       steps. The reciprocal temperature  $\beta$  was set to 1000, the inverse thermal mass  $\mu^{-1}$   
417       was set to 10, and the additive noise parameter  $\sigma_a$  was set to 0.1.

418       For the 2x2 case, the total number of valid BAS patterns is  $N_{BAS} = 6$ , while  
419       for the 3x3 case,  $N_{BAS} = 14$ . The goal of the QCBM is to learn the probability  
420       distribution of the valid BAS patterns and generate samples that closely resemble the  
421       target distribution.

422       After training, the QCBM parameters corresponding to the minimum MMD cost  
423       were selected as the final parameters. Using these parameters, we generated samples  
424       from the QCBM and visualized them in a 2x2 grid (Figure A.22) and a 3x3 grid  
425       (Figure A.25). The corresponding probabilities of each generated sample were also  
426       calculated and displayed.

427       The results showed that the QCBM, trained using SGAdL, was able to generate  
428       valid BAS patterns in some cases, while in other cases, the generated patterns were  
429       close to the valid ones. However, further hyperparameter analysis and tuning could  
430       potentially improve the results for this dataset.

431       Upon observing the MMD loss plot, we noticed a rapid descent at the beginning,  
432       followed by a relative plateauing for the rest of the plot. The MMD did not reach  
433       close to zero as desired, indicating that the QCBM may have encountered challenges  
434       in fully capturing the target distribution. One hypothesis for this behavior is that the  
435       sampling particle may have gotten trapped in a local minimum during the training  
436       process.

437       To address this issue, a potential future work could involve implementing a par-  
438       allel tempering algorithm [15]. In this approach, multiple sampling particles are run  
439       simultaneously with different hyperparameters. After a certain number of iterations,  
440       the hyperparameters are swapped between the particles. The desired behavior of this  
441       algorithm is that some particles can jump between different basins in the loss land-  
442       scape, while other particles can explore the local minima within each basin. This  
443       approach could help the QCBM escape local minima and improve its ability to learn  
444       the target distribution.

445       **7. Conclusions.** In this report, we have reviewed the application of Bayesian  
446 inference to generative modelling using QCBMs. Our primary focus was to investigate  
447 the effectiveness of the SGAdLT in addressing the existing challenges posed by the  
448 stochastic gradient noise in these models. By conducting extensive experiments on  
449 different types of datasets such as the simpler Bimodal Gaussian data and the more  
450 complex BARS-and-STRIPES dataset, we aimed to provide valuable insights into the  
451 performance and utility of SGAdLT as an alternative to sampling methods used prior  
452 in the Quantum Machine Learning space.

453       The results from our experiments demonstrated that SGAdLT exhibits promising  
454 capabilities in mitigating the bias introduced by stochastic gradients. We observed  
455 that SGAdLT can effectively explore local cost minima at lower  $\beta$  values compared to  
456 the SGLD, the recently proposed sampling algorithm for this type of model. These  
457 findings highlight the potential of SGAdLT to improve the sampling process and  
458 enhance the overall performance of Quantum Generative Models in general.

459       Throughout the experimentation phase, we also analysed the impact of various  
460 hyper-parameters on the performance and generalisation of SGAdLT. We mainly in-  
461 vestigated the effects of the reciprocal temperature ( $\beta$ ) and gradient batch size  $S$ . Our  
462 study of the hyper-parameters followed a key theme: "exploitation versus exploration  
463 while sampling", by tuning the hyper-parameters in a specific configuration one can  
464 drive SGAdLT to explore high probability posterior landscapes (exploitation), how-  
465 ever if the task requires the parameter space to be extensively explored one can also  
466 choose a configuration of hyper-parameters to achieve the desired effect.

467       It is important to note that in our experiments, we employed uniform priors for  
468 the Bayesian inference/learning process. However as suggested in [2] exploring the  
469 use of Laplace could potentially come with further enhancements to the performance  
470 of SGAdLT as it was a great choice to use in conjunction with SGLD. We assume  
471 the main benefit of incorporating the Laplace Prior to our model would be better  
472 generalisation and regularisation by imposing sparsity.

473       Moreover, while our experiments were conducted on simulated quantum circuits  
474 via the use of Qujax [12], it is critical for us to extend this research to actual quan-  
475 tum hardware. Running our model on real quantum devices would allow us to assess  
476 the performance and scalability in the presence of hardware noise and see how well  
477 SGAdLT works for mitigating noise that does not fit into the prior assumptions (non-

478 Gaussian, non-constant covariance). This understanding could lead to the develop-  
479 ment of more bespoke variations to use for quantum machine learning but before that  
480 we might also benefit from exploring the use of modified variations of Ad-Langevin  
481 Thermostat, such as the one proposed by Sekkat and Stoltz [13].

482 Another note about noise is that it is an important aspect of the work that war-  
483 rants further research is the gradient noise model. Our experiments highlighted the  
484 significance of understanding the characteristics and behaviour of gradient noise in  
485 QCBMs. By conducting a more-in depth analysis of the GNM, we can gain valuable  
486 insights into its impact on the convergence and stability of SGAdLT. These algo-  
487 rithmic advancements may help in addressing specific challenges encountered in quantum  
488 generative modeling, such as the barren plateau phenomenon and the presence of local  
489 minima in the cost landscape.

490 In conclusion, our research showcases the immense potential of using a Bayesian  
491 framework within quantum machine learning, particularly in the context of quan-  
492 tum generative modeling using QCBMs. The application of SGAdLT demonstrates  
493 promising results in enhancing the robustness and performance of these models in the  
494 presence of stochastic gradient noise. However, this work also highlights the need for  
495 further research and exploration in several key areas, including the incorporation of  
496 Laplace priors, simulation on actual quantum hardware, in-depth analysis of the gra-  
497 dient noise model, and the investigation of modified algorithms. By addressing these  
498 aspects, we can continue to push the boundaries of quantum generative modeling and  
499 unlock its full potential for various practical applications. As quantum computing  
500 technologies continue to advance, the integration of Bayesian methods with quantum  
501 machine learning will hopefully play a pivotal role in shaping the future of this exciting  
502 field.

- [1] Jin-Guo Liu and Lei Wang. Differentiable learning of quantum circuit born machines. *Physical Review A*, 98(6), December 2018.
- [2] Samuel Duffield, Marcello Benedetti, and Matthias Rosenkranz. Bayesian learning of parameterised quantum circuits. *Machine Learning: Science and Technology*, 4(2):025007, Apr 2023.
- [3] Neil K. Chada, Benedict Leimkuhler, Daniel Paulin, and Peter A. Whalley. Unbiased kinetic langevin monte carlo with inexact gradients, 2023.
- [4] Benedict Leimkuhler and Xiaocheng Shang. Adaptive thermostats for noisy gradient systems. *SIAM Journal on Scientific Computing*, 38(2):A712–A736, Jan 2016.
- [5] Maria Teresa Baldassarre, Danilo Caivano, Berenice Fernandez Nieto, Domenico Gigante, and Azzurra Ragone. The social impact of generative ai: An analysis on chatgpt. In *Proceedings of the 2023 ACM Conference on Information Technology for Social Good*, pages 363–373, 2023.
- [6] Aram W. Harrow and Annie Y. Wei. *Adaptive Quantum Simulated Annealing for Bayesian Inference and Estimating Partition Functions*, page 193–212. Society for Industrial and Applied Mathematics, January 2020.
- [7] Nanyang Ye, Zhanxing Zhu, and Rafal K. Mantiuk. Langevin dynamics with continuous tempering for training deep neural networks, 2017.
- [8] Zhan Yu, Qiuhan Chen, Yuling Jiao, Yinan Li, Xiliang Lu, Xin Wang, and Jerry Zhijian Yang. Provable Advantage of Parameterized Quantum Circuit in Function Approximation. 10 2023.
- [9] Jarrod R. McClean, Sergio Boixo, Vadim N. Smelyanskiy, Ryan Babbush, and Hartmut Neven. Barren plateaus in quantum neural network training landscapes. *Nature Communications*, 9(1), November 2018.
- [10] Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ICML’11, page 681–688, Madison, WI, USA, 2011. Omnipress.
- [11] Yuxuan Du, Min-Hsiu Hsieh, Tongliang Liu, and Dacheng Tao. Expressive power of parametrized quantum circuits. *Physical Review Research*, 2(3), July 2020.
- [12] Samuel Duffield, Gabriel Matos, and Melf Johannsen. qujax: Simulating quantum circuits with JAX. *Journal of Open Source Software*, 8(89):5504, September 2023.
- [13] Inass Sekkat and Gabriel Stoltz. Quantifying the mini-batching error in bayesian inference for adaptive langevin dynamics, 2023.
- [14] James M Flegal, Kshitij Khare, and Yan Zhou. Preconditioned langevin monte carlo for bayesian inference. *arXiv preprint arXiv:1905.11503*, 2019.
- [15] Rohitash Chandra, Konark Jain, Ratneel V. Deo, and Sally Cripps. Langevin-gradient parallel tempering for bayesian neural learning. *Neurocomputing*, 359:315–326, 2019.

## Appendix A. Figures for Experiments.

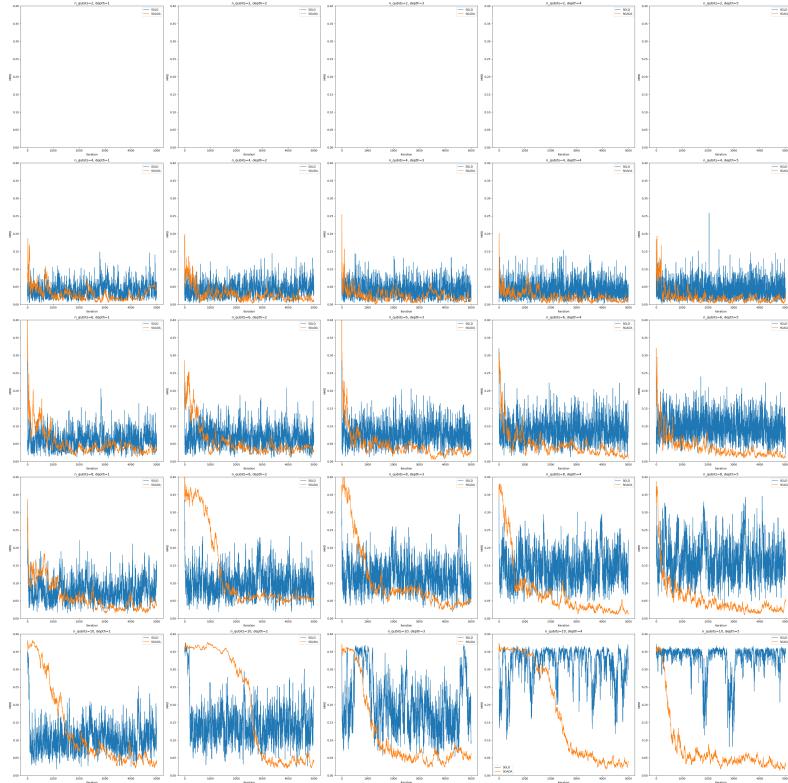


FIG. A.1. Maximum Mean Discrepancy (MMD) convergence plots for SGGLD (Blue) and SGADA (Orange) algorithms across different numbers of qubits (Rows, 2-10) and circuit depths (columns 1-5).

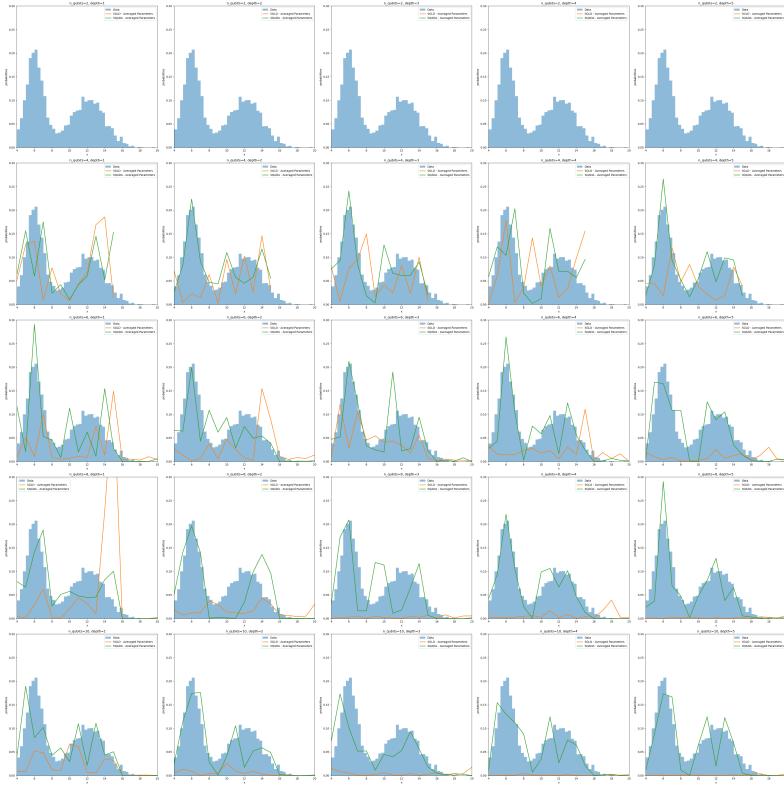


FIG. A.2. Probability distributions of the generated samples for SGGLD (Orange) and SGADA (Green) algorithms across different numbers of qubits( Rows, 2-10) and circuit depths (columns 1-5).

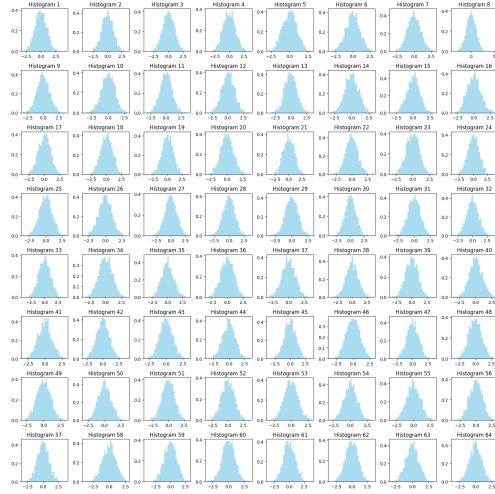


FIG. A.3. Gradient noise model for artificial noise

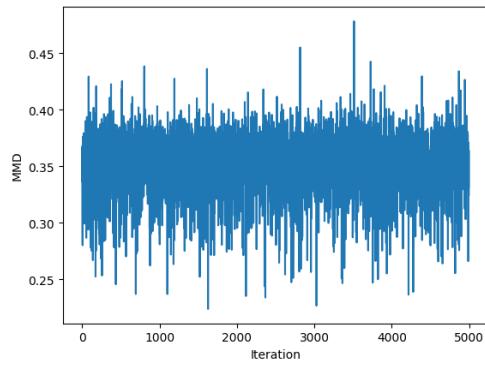


FIG. A.4. *Cost Curve for SGLD*

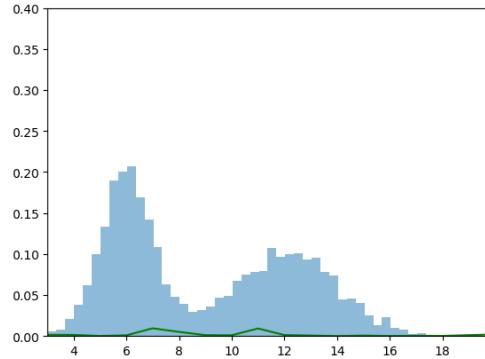


FIG. A.5.  $p(y|\mathcal{D}) \approx \frac{1}{N} \sum_{i=1}^N p(y|\boldsymbol{\theta}_i^{SGLD})$

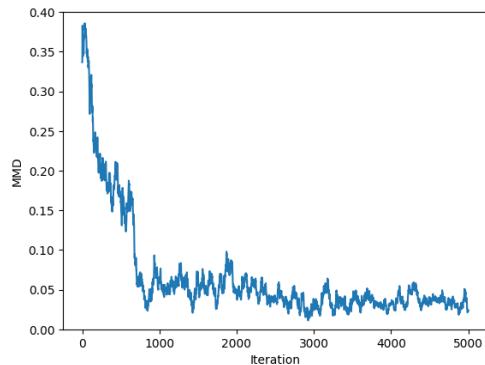


FIG. A.6. *Cost Curve for SGAdLT*

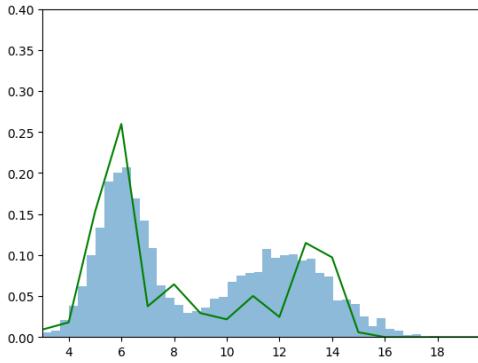


FIG. A.7.  $p(y|\mathbf{D}) \approx \frac{1}{N} \sum_{i=1}^N p(y|\boldsymbol{\theta}_i^{SGAdLT})$

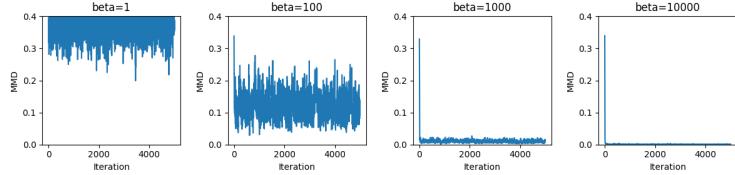


FIG. A.8. (a) Cost Curves for SGLD

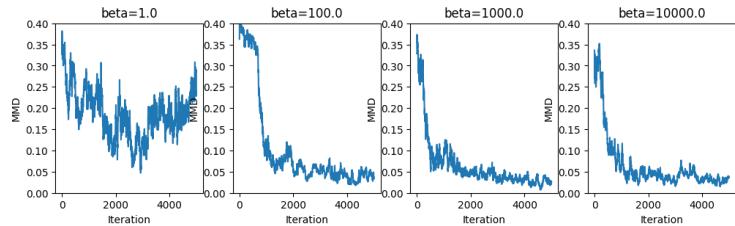


FIG. A.9. (b) Cost Curves for SGAdLT

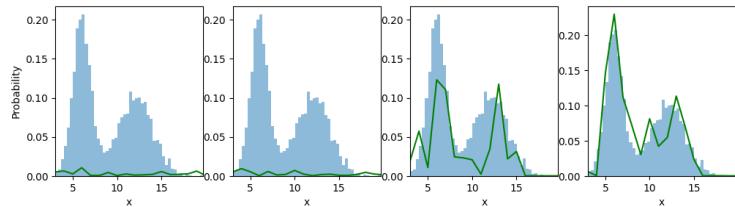


FIG. A.10. Predictive Distribution  $p(y|\mathbf{D})$  for SGLD

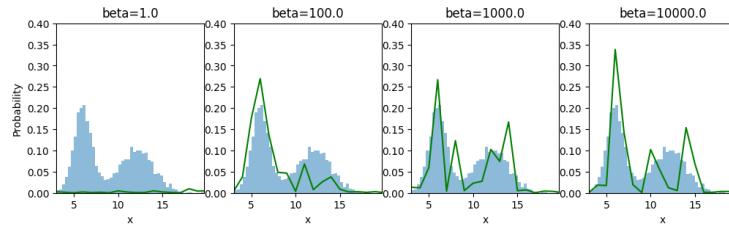


FIG. A.11. Predictive Distribution  $p(y|\mathcal{D})$  for SGAdLT

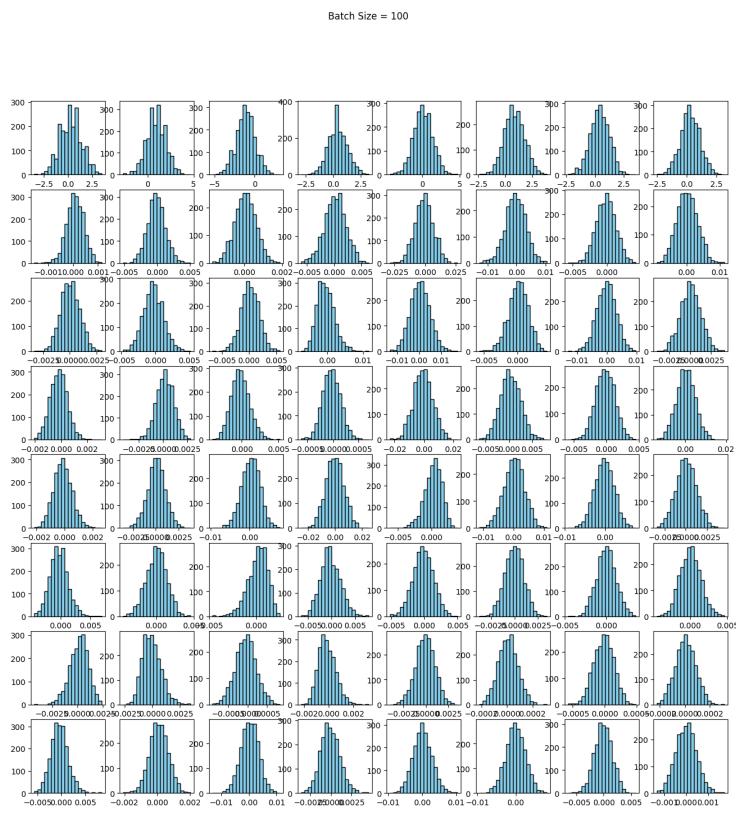


FIG. A.12. Batch Size = 100

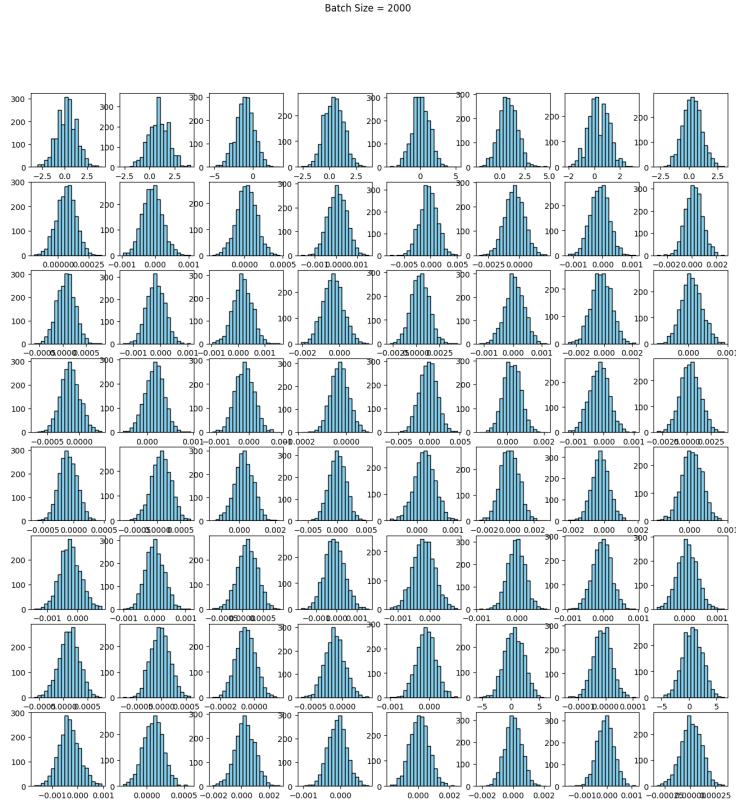


FIG. A.13. Batch Size = 2000

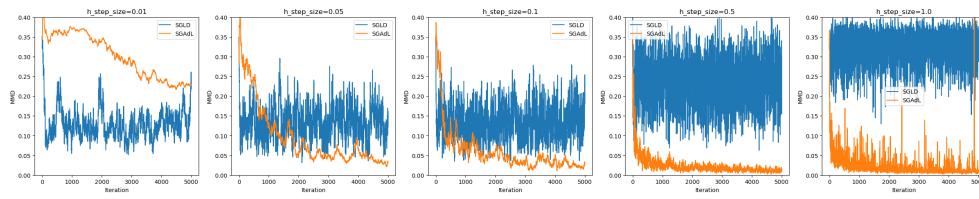


FIG. A.14. Maximum Mean Discrepancy (MMD) convergence plots for SGLD (Blue) and SGADA (Orange) algorithms across different range of step size  $h$ , (0.01, 0.05, 0.1, 0.5, 1)

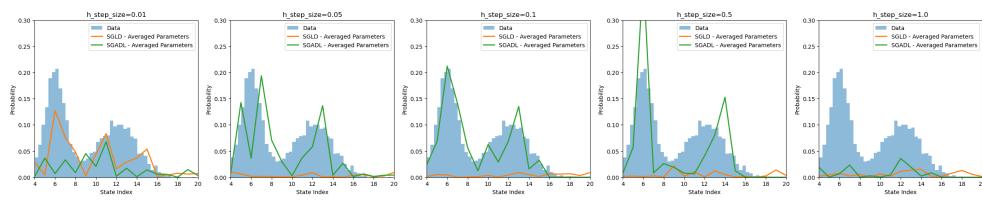


FIG. A.15. Probability distributions of the generated samples for SGLD (Orange) and SGADA (Green) algorithms across different range of step size  $h$ , (0.01, 0.05, 0.1, 0.5, 1)

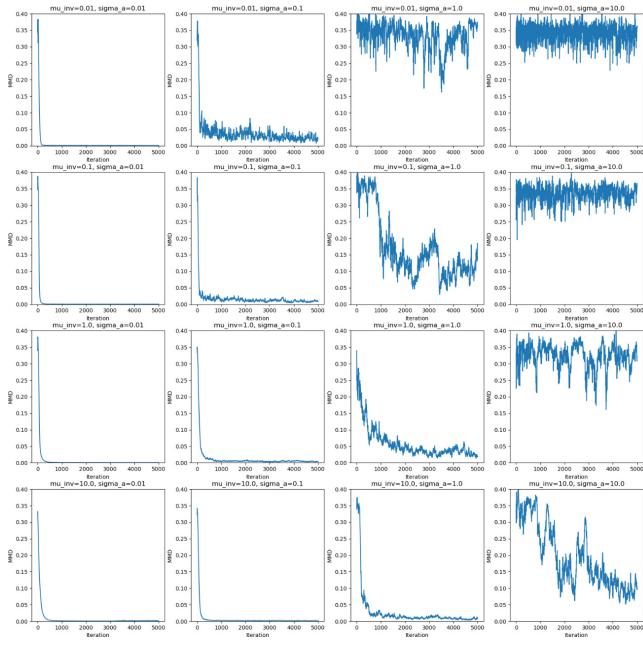


FIG. A.16. Maximum Mean Discrepancy Curves

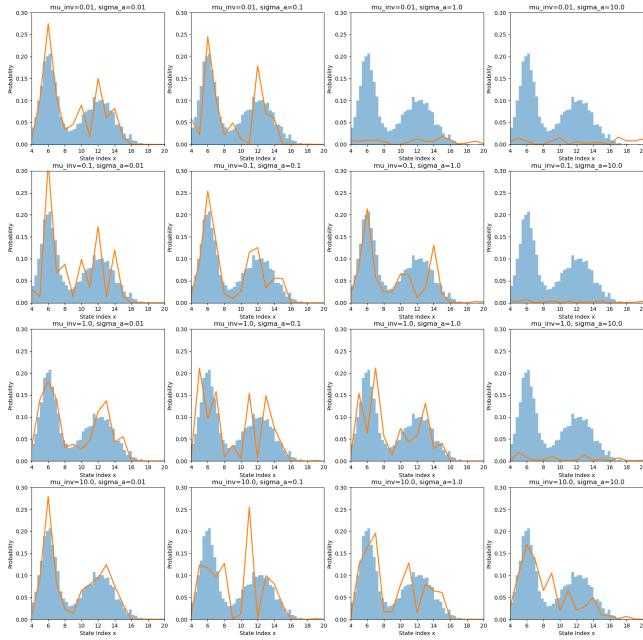


FIG. A.17. Bayesian Model Averages

FIG. A.18. Thermal Mass  $\mu$  vs Additive Noise  $\sigma_A$  for {0.01, 0.1, 1, 10}

*This manuscript is for review purposes only.*

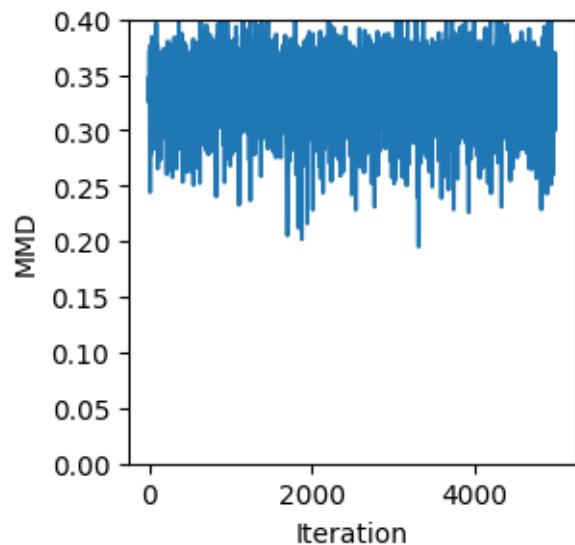


FIG. A.19. *Loss curve for preconditioned SGAdLT*

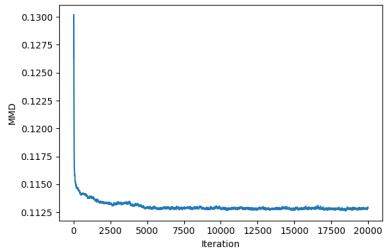


FIG. A.20. *MMD plot*

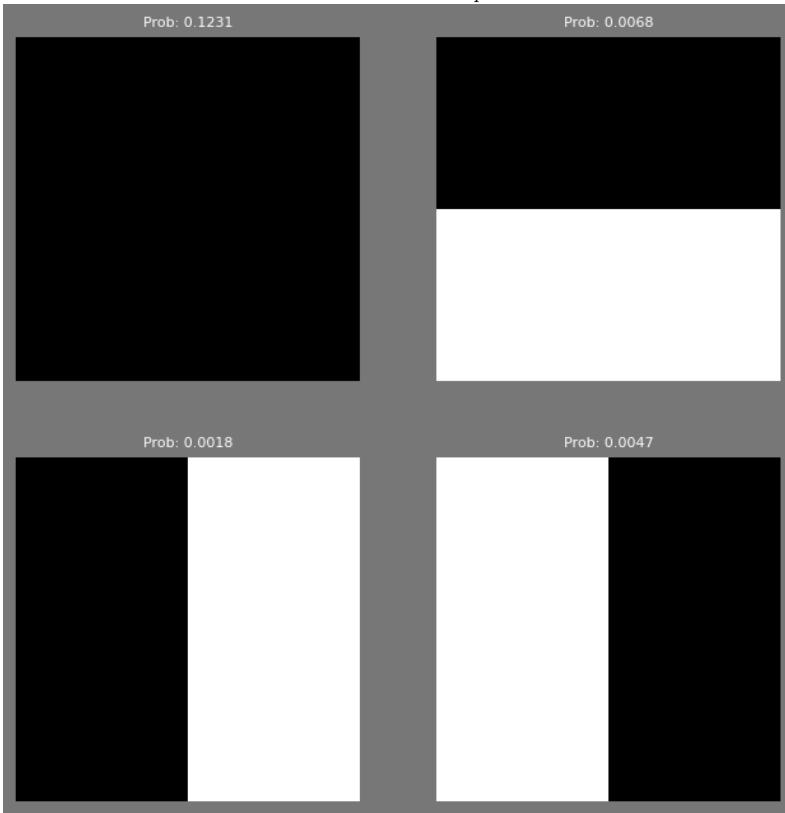


FIG. A.21. *Generated samples with probabilities*

FIG. A.22. *Comparison of MMD plot and generated samples with probabilities for the 2x2 Bars and Stripes dataset using SGAdL.* (a) The MMD plot shows the convergence of the SGAdLT algorithm over iterations. (b) The generated samples demonstrate the learned probability distribution of the Bars and Stripes patterns.

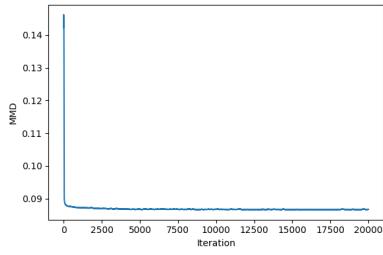


FIG. A.23. *MMD plot*

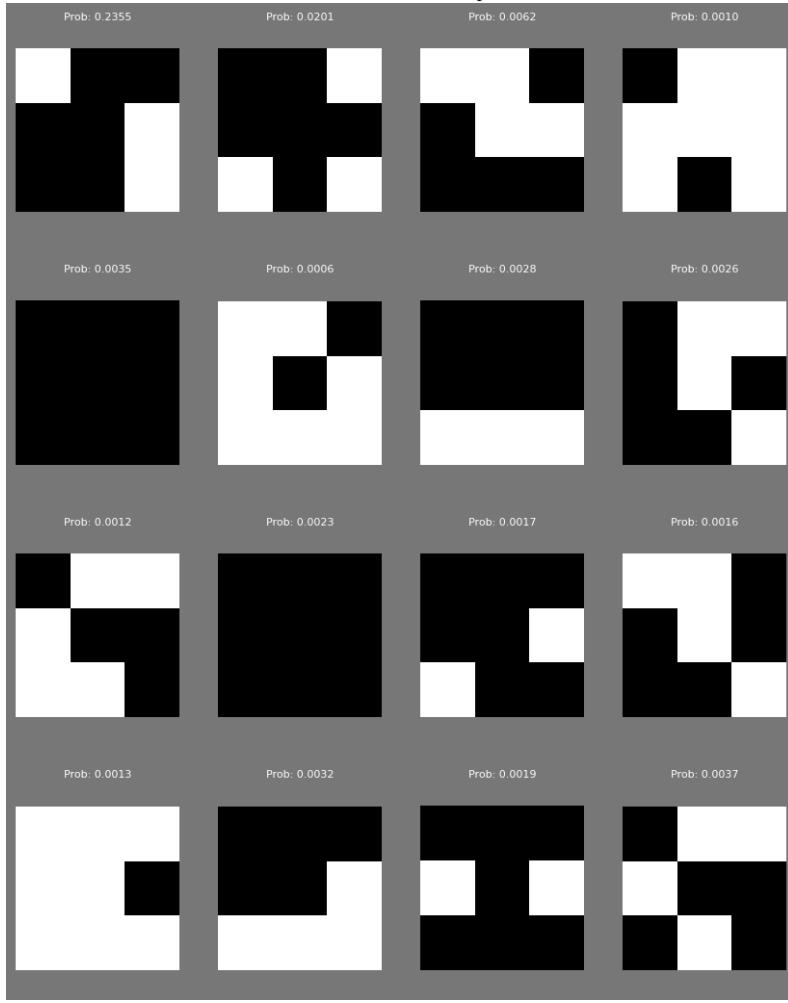


FIG. A.24. *Generated samples with probabilities*

FIG. A.25. *Comparison of MMD plot and generated samples with probabilities for the 3x3 Bars and Stripes dataset using SGAdL. (a) The MMD plot shows the convergence of the SGAdL algorithm over iterations. (b) The generated samples demonstrate the learned probability distribution of the Bars and Stripes patterns.*

542

## Appendix B. Bayesian Model Averaging.

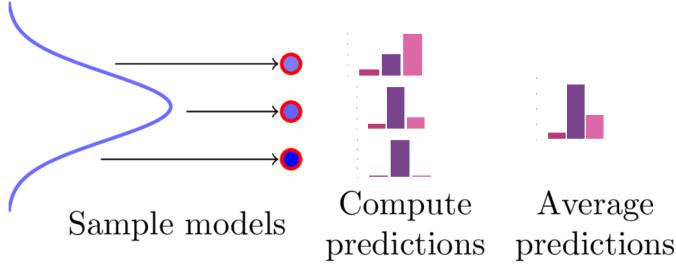


FIG. B.1. *Bayesian Model Averaging (BMA) procedure. The models  $w$  are sampled  $N$  times leading to predictions  $p(y|w_i, x)$  which are uniformly averaged to give  $p_{BMA}(y|x, D)$*

543

## Appendix C. Maximum Mean Discrepancy.

544

The Maximum Mean Discrepancy (MMD) is used as a cost function by computing the difference in means of the distributions  $P$  and  $Q$  lifted to feature space  $H_k$ .

$$\begin{aligned} MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{H_k}^2 \\ &= \langle \mu_P, \mu_P \rangle_{H_k} + \langle \mu_Q, \mu_Q \rangle_{H_k} - 2\langle \mu_P, \mu_Q \rangle_{H_k} \\ &= E_P[k(X, X')] + E_Q[k(Y, Y')] + E_{P,Q}[k(X, Y)] \end{aligned}$$

547

where  $k(\mathbf{z}, \mathbf{z}') = e^{-\frac{(\mathbf{z}-\mathbf{z}')^2}{2\sigma^2}}$  is Gaussian kernel with bandwidth  $\sigma$ . We chose this to be the median of the dataset.

549

## Appendix D. Reciprocal Temperature $\beta$ .

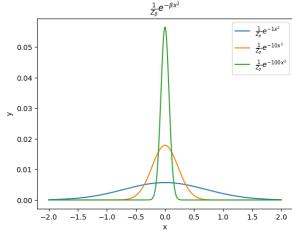


FIG. D.1. *Effect of  $\beta$*

550

## Appendix E. Parameter-Shift Rule for Gradient Estimation.

551

To perform gradient-based optimization of QCBMs, we need to compute the gradients of the loss function with respect to the circuit parameters. However, unlike

553 classical neural networks, quantum circuits do not allow for direct back-propagation  
 554 of gradients through the gates. Instead, we can use the parameter-shift rule [2], which  
 555 is a finite-difference method for estimating gradients on quantum circuits.

556 Consider a parameterized quantum circuit  $U(\theta)$  that consists of a sequence of  
 557 parameterized gates  $U_i(\theta_i)$ , where  $\theta = (\theta_1, \dots, \theta_n)$  is the vector of parameters. The  
 558 parameter-shift rule states that the gradient of the expectation value of an observable  
 559  $O$  with respect to a parameter  $\theta_i$  can be estimated as:

560 (E.1) 
$$\frac{\partial \langle O \rangle}{\partial \theta_i} = \frac{1}{2} (\langle O \rangle_{\theta_i + \frac{\pi}{2}} - \langle O \rangle_{\theta_i - \frac{\pi}{2}}),$$

561 where  $\langle O \rangle_{\theta_i \pm \frac{\pi}{2}}$  denotes the expectation value of  $O$  when the parameter  $\theta_i$  is shifted  
 562 by  $\pm \frac{\pi}{2}$ .

563 To estimate the gradients using the parameter-shift rule, we need to evaluate the  
 564 expectation values  $\langle O \rangle_{\theta_i \pm \frac{\pi}{2}}$  for each parameter  $\theta_i$ . This can be done by running the  
 565 quantum circuit twice for each parameter, with the corresponding parameter shifted  
 566 by  $\pm \frac{\pi}{2}$ , and measuring the observable  $O$ . The gradients can then be computed using  
 567 the finite-difference formula. The parameter-shift rule provides a way to estimate  
 568 gradients on quantum circuits without the need for ancilla qubits or complex quantum  
 569 operations.

570 **Appendix F. Hyperparameters.**

571 Unless otherwise specified, the default hyperparameters are:

- 572 • Number of Qubits = 8
- 573 • Circuit Depth = 3
- 574 • Cost Function: MMD
- 575 • Prior: Uniform
- 576 • Mass Matrix  $M = I_{64}$
- 577 • Batch Size  $S = 200$
- 578 • Degrees of Freedom  $N_d = 1$
- 579 • Additive Noise  $\sigma_A = 1.0$
- 580 • Thermal Mass  $\mu = 1.0$
- 581 • Number of samples = 5000
- 582 • Burn-in Period = 1000
- 583 • Stepsize  $h = 0.1$