
Thermodynamic Bayesian Inference

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A fully Bayesian treatment of complicated predictive models (such as deep neural
2 networks) would enable rigorous uncertainty quantification and the automation
3 of higher-level tasks including model selection. However, the intractability of
4 sampling Bayesian posteriors over many parameters inhibits the use of Bayesian
5 methods where they are most needed. Thermodynamic computing has emerged as
6 a paradigm for accelerating operations used in machine learning, such as matrix
7 inversion, and is based on the mapping of Langevin equations to the dynamics
8 of noisy physical systems. Hence, it is natural to consider the implementation
9 of Langevin sampling algorithms on thermodynamic devices. In this work we
10 propose electronic analog devices that sample from Bayesian posteriors by real-
11 izing Langevin dynamics physically. Circuit designs are given for sampling the
12 posterior of a Gaussian-Gaussian model and for Bayesian logistic regression, and
13 are validated by simulations. It is shown, under reasonable assumptions, that the
14 time-complexity of sampling the Gaussian-Gaussian posterior is sublinear in di-
15 mension. These results highlight the potential to accelerate Bayesian inference
16 with thermodynamic computing.

17 1 Introduction

18 Bayesian statistics has proved an effective framework for making predictions under uncertainty [1,
19 2, 3, 4, 5, 6], and it is central to proposals for automating machine learning [7]. Bayesian methods
20 enable uncertainty quantification by incorporating prior knowledge and modeling a distribution over
21 the parameters of interest. Popular machine learning methods that employ this approach include
22 Bayesian linear and non-linear regression [8], Kalman filters [9], Thompson sampling [2], continual
23 learning [10, 11], and Bayesian neural networks [3, 12].

24 Unfortunately, computing the posterior distribution in these settings is often intractable [13]. Meth-
25 ods such as the Laplace approximation [14] and variational inference [15] may be used to approxi-
26 mate the posterior in these cases, however their accuracy struggles for complicated posteriors, such
27 as those of a Bayesian neural network [13]. Regardless, sampling accurately from such posteriors
28 requires enormous computing resources [13].

29 Computational bottlenecks in Bayesian inference motivate the need for novel hardware accelera-
30 tors. Physics-based sampling hardware has been proposed for this purpose, including Ising ma-
31 chines [16, 17, 18, 19, 20], probabilistic bit computers [21, 22, 23], and thermodynamic com-
32 puters [24, 25, 26, 27, 28, 29]. Continuous-variable hardware is particularly suited to Bayesian
33 inference since continuous distributions are typically used in probabilistic machine learning [30].
34 However, a rigorous treatment of how such hardware can perform Bayesian inference with scalable
35 circuits has not yet been given.

36 The most computationally tractable algorithms for exact Bayesian inference are Monte Carlo sam-
 37 pling algorithms. The Langevin sampling algorithm [31, 32] is an elegant example inspired by
 38 statistical physics, based on the dynamics of a damped system in contact with a heat bath. What we
 39 propose in this work is to build a physical realization of the system that is simulated by the Langevin
 40 algorithm. The system must be designed to have a potential energy such that the Gibbs distribution
 41 $p(x) \propto e^{-\beta U(x)}$ is the desired posterior distribution which is reached at thermodynamic equilib-
 42 rium. We present circuit schematics for electronic implementations of such devices for Bayesian
 43 inference for two special cases. The first is a Gaussian-Gaussian model (where the prior and the
 44 likelihood are both multivariate normal, as found in linear regression and Kalman filtering), and the
 45 second is logistic regression (where the prior is Gaussian and the likelihood is Bernoulli parameter-
 46 ized by a logistic function). In each case, the parameters of the prior and likelihood are encoded in
 47 the values of components of the circuit, and then voltages or currents are measured to sample the
 48 random variable.

49 While thermodynamic algorithms have been proposed for linear algebra [27] and neural network
 50 training [33], our work can be viewed as the first thermodynamic algorithm for sampling from
 51 Bayesian posteriors. Moreover, our work provides the first concrete proposal for non-Gaussian sam-
 52 pling with thermodynamic hardware. Overall, our work opens up a new field of rigorous Bayesian
 53 inference with thermodynamic computers and lays the groundwork for scalable CMOS-based chips
 54 for probabilistic machine learning.

55 We show that in theory the device proposed for sampling the Gaussian-Gaussian model posterior can
 56 obtain N samples in d dimensions in time scaling with $O(N \ln d)$. This is a significant speedup over
 57 typical methods used digitally for the same problem, which involve matrix inversions taking time
 58 scaling with $O(d^\omega)$ where $2 < \omega < 3$. This speedup is larger than the polynomial speedups found
 59 in previous work on thermodynamic algorithms for linear algebra primitives [27] (where speedups
 60 were found to scale linearly with dimension).

61 2 Results

62 Suppose that we have samples of a random vector y , and would like to estimate a random vector θ
 63 on which y depends somehow. The Bayesian approach is to assume a prior distribution on θ given
 64 by a density function $p_\theta(\theta)$, and a likelihood function $p_{y|\theta}(y|\theta)$. The posterior distribution for θ is
 65 then given by Bayes's theorem $p_{\theta|y}(\theta|y) = p_{y|\theta}(y|\theta)p_\theta(\theta)/p_y(y)$. To sample from the posterior
 66 using the Langevin algorithm, one first computes the score

$$\nabla_\theta \ln p_{\theta|y}(\theta|y) = \nabla_\theta \ln p_{y|\theta}(y|\theta) + \nabla_\theta \ln p_\theta(\theta). \quad (1)$$

67 Then the score is used as the drift term in the following stochastic differential equation (SDE)

$$d\theta = \nabla_\theta \ln p_{\theta|y}(\theta|y) dt + \mathcal{N}[0, 2 dt]. \quad (2)$$

68 After this SDE is evolved for a sufficient time T , the value of θ will be a sample from $p_{\theta|y}$. This
 69 algorithm is equivalent to the equilibration of an overdamped system, as we will now describe. First
 70 let r be a vector of the same dimension as θ describing the state of a physical system, and satisfying
 71 $r = \theta \tilde{r}$ for some constant \tilde{r} (this factor is necessary because θ is unitless while the physical quantity r
 72 has units). Now we define the potential energy function $\beta U(r) = -\ln p_{\theta|y}(r/\tilde{r} | y)$. The dynamics
 73 of an overdamped system with potential energy U in contact with a heat bath at inverse temperature
 74 β can be modeled by the overdamped Langevin equation

$$dr = -\gamma^{-1} \nabla_r U(r) dt + \mathcal{N}[0, 2\gamma^{-1}\beta^{-1} dt], \quad (3)$$

75 where γ is a damping constant. Note that this implies that γ has dimensions of energy \cdot time/ $[r]^2$. If
 76 we introduce a constant $\tau = \gamma\beta\tilde{r}^2$, Eq. (3) can be written

$$d\theta = \nabla_\theta \ln p_{\theta|y}(\theta|y)\tau^{-1} dt + \mathcal{N}[0, 2\tau^{-1} dt], \quad (4)$$

77 which has the same form as Eq. (2), except with the time constant τ . It is clear that if Eq. (2)
 78 must be run for a dimensionless duration T to achieve convergence, then the physical system must
 79 be allowed to evolve for a physical time duration τT to achieve the same result. While we have
 80 addressed the case of conditioning on a single sample y above, the generalization of these ideas to
 81 the case of conditioning on multiple I.I.D. samples is given in Appendix D. In what follows we will
 82 present designs for circuits whose potential energy results in an overdamped Langevin equation that
 83 yields samples from Bayesian posteriors.

84 **2.1 Gaussian-Gaussian model**

85 A particularly simple special case of Bayesian inference is a when both the prior and the likelihood
 86 are multivariate normal, and we address this simple model first in order to illustrate our approach
 87 more clearly. Specifically, let $\theta \in \mathbb{R}^d$ have prior distribution $p_\theta(\theta) = \mathcal{N}[\mu, \Sigma]$, and let the likelihood
 88 be $p_{y|\theta}(y|\theta) = \mathcal{N}[\theta, \Sigma_{y|\theta}]$, where $y \in \mathbb{R}^d$ is an observed sample. In this case the posterior $p_{\theta|y}$ is
 89 also multivariate normal, with parameters [12]

$$\mu_{\theta|y} = \mu + \Sigma (\Sigma + \Sigma_{y|\theta})^{-1} (y - \mu), \quad (5)$$

90

$$\Sigma_{\theta|y} = \Sigma - \Sigma(\Sigma + \Sigma_{y|\theta})^{-1}\Sigma. \quad (6)$$

91 For this model, the posterior is tractable and can be computed on digital computers relatively effi-
 92 ciently, however for very large dimensions the necessary matrix inversion and matrix-matrix mul-
 93 tiplications can still create a costly computational bottleneck. As we will see, the thermodynamic
 94 approach provides a means to avoid the costly inversion and matrix products in the computation,
 95 and therefore to accelerate Bayesian inference for this model.

96 We begin by deriving the Langevin equation for sampling this posterior. For this prior and likelihood,
 97 the score of the posterior Eq. (1) is

$$\nabla_\theta \ln p_{\theta|y}(\theta|y) = -\Sigma^{-1}(\theta - \mu) - \Sigma_{y|\theta}^{-1}(\theta - y), \quad (7)$$

98 and so Eq. (4) becomes

$$d\theta = -\Sigma^{-1}(\theta - \mu)\tau^{-1}dt - \Sigma_{y|\theta}^{-1}(\theta - y)\tau^{-1}dt + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt]. \quad (8)$$

99 In fact, this SDE can be implemented by a circuit consisting of two resistor networks coupled by
 100 inductors, shown in Fig. 1 for the two-dimensional case.

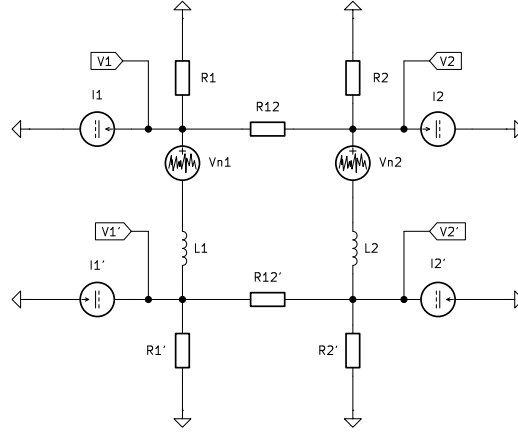


Figure 1: Circuit schematic for the Gaussian-Gaussian model posterior sampling device.

101 The full analysis of the circuit in Fig. 1 is given in Appendix A, but a few remarks are made here to
 102 explain its operation. First, we define the conductance matrices \mathcal{G} as

$$\mathcal{G} = \begin{pmatrix} R_{11}^{-1} + R_{12}^{-1} & -R_{12}^{-1} \\ -R_{12}^{-1} & R_{22} + R_{12}^{-1} \end{pmatrix}, \quad (9)$$

103 and \mathcal{G}' is defined in the same way for the primed resistors R'_1 , R'_2 , and R'_{12} . By applying Kirchoff's
 104 current law (KCL), the voltages across the resistors can be eliminated. Then the equation $V = L\dot{I}$
 105 is used to derive the following stochastic differential equation for the currents through the inductors

$$dI_L = -L^{-1}\mathcal{G}^{-1}(I_L - I) dt - L^{-1}\mathcal{G}'^{-1}(I_L - I') dt + L^{-1}\sqrt{S}\mathcal{N}[0, \mathbb{I} dt], \quad (10)$$

106 where $I_L = (I_{L1} \ I_{L2})^\top$ and S is the power spectral density of each noise source. This equation
 107 has the same form as Eq. (8), so it is only necessary to determine an appropriate mapping of
 108 distributional parameters to physical properties of the circuit's components (see Appendix A). By

109 including more inductors and coupling resistors (as well as current and voltage sources), the design
 110 can be generalized to arbitrary dimension.

111 To verify that the proposed circuit does indeed evolve according to the correct SDE, we ran SPICE
 112 circuit simulations. Figure 2 shows the results of such a simulation where a 2-dimensional Gaussian
 113 prior and a 2-dimensional Gaussian likelihood are encoded into the conductances while the current
 114 in each inductor is measured to determine the resulting posterior.

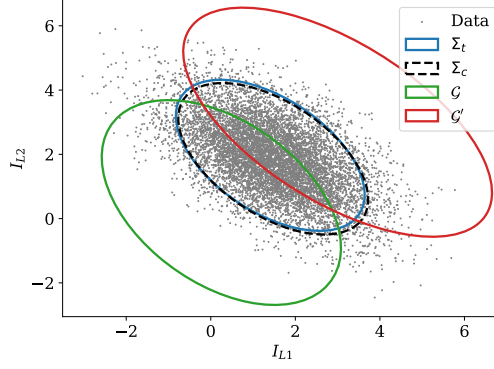


Figure 2: SPICE simulations of proposed Gaussian-Gaussian circuit in Fig. 1. The grey points represent the simulated circuit’s induct currents. The dashed black and solid blue ellipses represent the empirical sample covariance and the target posterior covariance from a Gaussian Bayesian update, respectively. The red and green ellipses represent the prior and likelihood.

115 As shown in Appendix C, the asymptotic runtime complexity for this algorithm is

$$t = O(N\kappa\tau \ln(\kappa^{3/2}d^{1/2}W_0^{-1})), \quad (11)$$

116 where κ is the condition number of the posterior covariance, $\tau = L/\tilde{R}$, d is the dimension, and W_0 is
 117 the Wasserstein distance between the true posterior and the distribution sampled by the device. The
 118 assumptions used to derive this result can also be found in Appendix C. Remarkably, the required
 119 time is sublinear in dimension, a large improvement over digital algorithms where complexity of
 120 constructing and sampling from the Gaussian-Gaussian posterior (67 - 68) is $O(d^\omega)$ where ω is the
 121 matrix multiplication constant (or more practically $O(d^3)$ via common implementations of Cholesky
 122 factorization). In Figure 3(a), we report the convergence of simulated thermodynamic samples for
 123 the Gaussian-Gaussian model with zero prior mean and covariances $\Sigma, \Sigma_{y|\theta}$ randomly sampled from
 124 a Wishart distribution with $2d$ degrees of freedom. We see fast convergence in Wasserstein distance
 125 to the true posterior, supporting our theoretical claims.

126 2.2 Bayesian linear regression and Kalman filtering

127 A generalization of the Gaussian-Gaussian model is that of Bayesian linear regression [8] (or equiv-
 128 alently a Kalman filter update step [9, 12]). In full generality we have

$$p_\theta(\theta) = \mathcal{N}[\mu, \Sigma], \quad (12)$$

$$p_{y|\theta}(y | \theta) = \mathcal{N}[H\theta, \Sigma_{y|\theta}], \quad (13)$$

129 Then the overdamped Langevin SDE becomes

$$\begin{aligned} d\theta &= -\Sigma^{-1}(\theta - \mu)\tau^{-1} dt - H^\top \Sigma_{y|\theta}^{-1}(y - H\theta)\tau^{-1} dt + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt], \\ &= -(A\theta - b)\tau^{-1} + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt], \text{ for } A = \Sigma^{-1} + H^\top \Sigma_{y|\theta}^{-1}H \text{ and } b = \mu + H^\top \Sigma_{y|\theta}^{-1}y. \end{aligned} \quad (14)$$

130 The form of the SDE (Ornstein-Uhlenbeck process) in 14 is exactly that of the thermodynamic
 131 device in [27] which if given input A and b above will produce samples from the Gaussian Bayesian
 132 posterior $p_{\theta|y}(\theta | y)$. Compared to the simpler Gaussian-Gaussian model above, a disadvantage of
 133 this approach is that the covariances Σ and $\Sigma_{y|\theta}$ have to be inverted prior to input as A . However,
 134 for linear regression, these matrices are often assumed to be diagonal and otherwise they can be

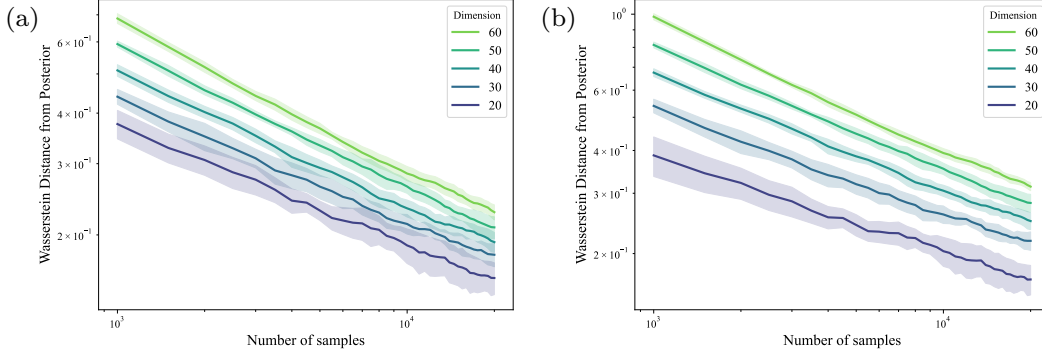


Figure 3: Convergence in Wasserstein distance between simulated thermodynamic samples and the true Gaussian posterior as a function of the number of samples (sampling time). All results are simulated exactly with `thermox` [35] and averaged over 10 random seeds with one standard deviation shown. Panel (a): Gaussian-Gaussian model with zero prior mean and covariances sampled from a Wishart distribution. Panel (b): Bayesian linear regression on the diabetes dataset [36] with dimension (number of features) varied by including higher-order cross terms of the 10 input data features.

135 efficiently inverted using the thermodynamic procedures in [27] as preprocessing. Additionally, the
 136 formulation of A requires matrix-matrix multiplications which can be costly (even in the case of
 137 diagonal covariances). Although, this can be accelerated with parallelization.

138 On the other hand, the generality of (12-13) makes the approach highly practical. Encompassing
 139 Bayesian linear regression [34] and the update step of the Kalman filter [9]. Moreover in the setting
 140 of Kalman filtering, the matrices Σ and $\Sigma_{y|\theta}$ are typically shared across time points and thus only
 141 need to be inverted once in comparison to the Bayesian posterior update which is applied at every
 142 time step (and typically represents the computation bottleneck due to the required matrix inversion).

143 In Figure 3(b), we simulate the evaluation of the thermodynamic linear algebra device [27] for a
 144 Bayesian linear regression task. We use the diabetes dataset [36] which has $N = 442$ continuous
 145 response variables y and 10 input features. We vary the number of features and therefore posterior
 146 dimension for the linear regression by extending to include the first d cross terms in the Taylor
 147 expansion over the input features. These input features are loaded as rows in the matrix $H \in \mathbb{R}^{N \times d}$.
 148 Both covariances are set to diagonal, $\Sigma = \mathbb{I}$ and $\Sigma_{y|\theta} = 0.1\mathbb{I}$. We observe that the Wasserstein
 149 distance converges quickly as more samples are collected and scales reasonably with dimension,
 150 indicating a sublinear scaling similar to the Gaussian-Gaussian model.

151 2.3 Bayesian logistic regression

152 Logistic regression is a method for classification tasks (both binary and multiclass) that models the
 153 dependence of class probabilities on independent variables using a logistic function. In the Bayesian
 154 setting, a prior can be assumed on the parameters of a logistic regression model, for example it is
 155 common to assume a Gaussian prior. However, after conditioning on observed data a posterior dis-
 156 tribution is produced that has no analytical closed form, making Bayesian logistic regression far
 157 less efficient than obtaining a point estimate of the parameters. In this section we present a ther-
 158 modynamic hardware architecture capable of sampling the posterior for binary logistic regression,
 159 and show some preliminary evidence that this architecture can do so more efficiently than existing
 160 methods.

161 Given a parameter vector $\theta \in \mathbb{R}^d$ and an independent variable vector $x \in \mathbb{R}^d$, binary logis-
 162 tic regression outputs a class probability $p_{y|\theta,x}(y|\theta,x)$, where $y \in \{-1, 1\}$ (often $y \in \{0, 1\}$)
 163 is written instead but we choose this notation to simplify the presentation). The likelihood is
 164 $p_{y|\theta,x}(y|\theta,x) = L(y\theta^\top x)$ where $L(z) = 1/(1 + e^{-z})$ is the standard logistic function [37]. Note
 165 that we will first consider the case of conditioning on a single sample, and in this case the likelihood
 166 will be denoted $p_{y|\theta}(y|\theta)$ as x is constant. Additionally, a multivariate normal prior is assumed for
 167 the parameters $\theta \sim \mathcal{N}[\mu, \Sigma]$. The Langevin equation for sampling the posterior is therefore:

$$d\theta = -\Sigma^{-1}(\theta - \mu)\tau^{-1}dt + L(-y\theta^\top x)yx\tau^{-1}dt + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt]. \quad (15)$$

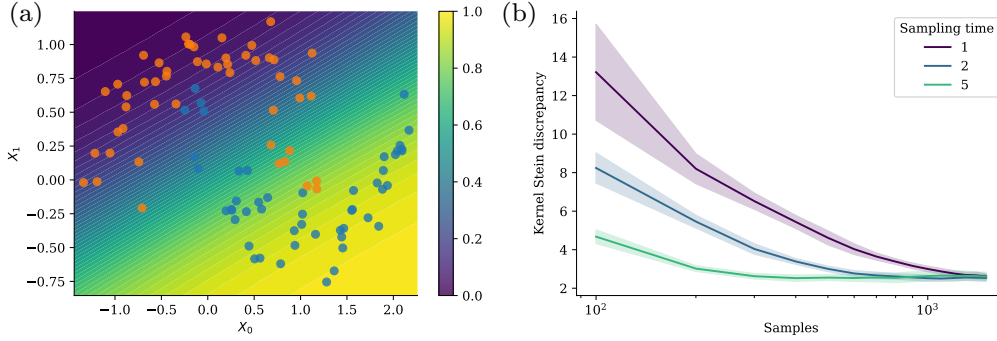


Figure 4: Panel (a): Probability surface to belong to class 1 (blue points). The dataset is also shown, where class 0 (blue points) and class 2 (orange points) are arranged in two intersecting moons. Panel (b): Kernel Stein discrepancy (KSD) of the collected samples with an ideal thermodynamic sampler, for varying sampling times. The sampling time is given in units of $10^{-3}\tau$. The KSD is averaged over five sets of random samples and $\tau = 1$.

168 A circuit implementing Eq. (15) is shown in Fig. 7, and the detailed analysis of this circuit is given
 169 in Appendix B. Equation (15) is valid for a single data sample, however, as mentioned, in practice
 170 we generally take gradients over a larger number of examples such that the gradients are less noisy.
 171 This can be done by enlarging the hardware, resulting in the second term of Eq. (15) being replaced
 172 by a sum $\sum_{i=1}^N L(-y_i \theta^T x_i) y_i x_i dt$, with N the number of data points. One may also consider
 173 minibatches, and the sum is only over a batch of size b . This is achievable by summing currents,
 174 which is detailed in the circuit implementation in Appendix B. At a high-level, implementing this
 175 protocol in hardware is very simple in the case of a full batch, since the data only needs to be sent
 176 once onto the hardware. The following steps are taken to collect the samples: (1) Map the data
 177 labels to $\{+1, -1\}$. (2) Map the data (X, Y) onto the hardware (full batch setting). (3) Initialize the
 178 state of the system, set the mean and the covariance matrix of the prior. (4) At every interval t_s (the
 179 sampling time), measure the state of the system $\theta(t)$ to collect samples.

180 In Fig. 4, we present results for a Bayesian logistic regression on a two-moons dataset, made of
 181 points separated in two classes that are arranged in intersecting moons in the 2D planes, as shown in
 182 Fig. 4(a). These results are obtained by running the SDE of Eq. (15), hence corresponds to an ideal
 183 simulation of the thermodynamic hardware. In this scenario, there are 3 parameters to sample, and
 184 $N = 100$ points are considered. In Fig. 4(a), we see that even for such a simple model, only a few
 185 points are misclassified. As mentioned, previously, this setting also gives access to better methods
 186 to estimate uncertainty in predictions. In Fig. 4(b), the Kernel Stein discrepancy (KSD) [38] is
 187 shown as a function of the number of collected samples for varying sampling rates. These results
 188 indicate that the number of samples to reach a low KSD (close to convergence) can be reduced by
 189 increasing the sampling time, indicating correlated samples, as is often the case.

190 3 Conclusion

191 In this work, we proposed the first thermodynamic algorithms for sampling from Bayesian poste-
 192 riors. We provided explicit constructions of CMOS-compatible analog circuits to implement these
 193 algorithms with scalable silicon chips. Our circuit for performing logistic regression represents the
 194 first concrete proposal for non-Gaussian sampling with a thermodynamic computer. In the case of
 195 Gaussian Bayesian inference (Gaussian prior, Gaussian likelihood), our analysis showed a sublinear
 196 complexity in d , leading to a speedup over standard digital methods that is greater than linear. This
 197 is an even larger speedup than those previously observed for thermodynamic linear algebra [27],
 198 suggesting that Bayesian inference is an ideal application for thermodynamic computers. Our work
 199 lays the foundation for accelerating Bayesian inference, a key component of probabilistic machine
 200 learning, with physics-based hardware.

References

- 201 [1] Andrew Gelman, Aki Vehtari, Daniel Simpson, Charles C Margossian, Bob Carpenter, Yuling
202 Yao, Lauren Kennedy, Jonah Gabry, Paul-Christian Bürkner, and Martin Modrák. Bayesian
203 workflow. *arXiv preprint arXiv:2011.01808*, 2020.
- 205 [2] Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, Zheng Wen, et al. A tutorial
206 on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.
- 207 [3] Andrew G Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective
208 of generalization. *Advances in neural information processing systems*, 33:4697–4708, 2020.
- 209 [4] Samuel Duffield, Kaelan Donatella, Johnathan Chiu, Phoebe Klett, and Daniel Simpson. Scal-
210 able bayesian learning with posteriors. *arXiv preprint arXiv:2406.00104*, 2024.
- 211 [5] Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad
212 Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A
213 review of uncertainty quantification in deep learning: Techniques, applications and challenges.
214 *Information fusion*, 76:243–297, 2021.
- 215 [6] Theodore Papamarkou, Maria Skoularidou, Konstantina Palla, Laurence Aitchison, Julyan
216 Arbel, David Dunson, Maurizio Filippone, Vincent Fortuin, Philipp Hennig, José Miguel
217 Hernández-Lobato, et al. Position: Bayesian deep learning is needed in the age of large-scale
218 ai. In *Forty-first International Conference on Machine Learning*, 2024.
- 219 [7] Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. *Automated machine learning: methods,*
220 *systems, challenges*. Springer Nature, 2019.
- 221 [8] Christopher M Bishop and Michael E Tipping. Bayesian regression and classification. In
222 *Advances in learning theory: methods, models and applications*, pages 267–285. IOS Press,
223 2003.
- 224 [9] Simo Särkkä and Lennart Svensson. *Bayesian filtering and smoothing*, volume 17. Cambridge
225 university press, 2023.
- 226 [10] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, An-
227 dree A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al.
228 Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy*
229 *of sciences*, 114(13):3521–3526, 2017.
- 230 [11] Cuong V. Nguyen, Yingzhen Li, Thang D. Bui, and Richard E. Turner. Variational continual
231 learning. In *International Conference on Learning Representations*, 2018.
- 232 [12] Samuel Duffield and Sumeetpal S Singh. Ensemble kalman inversion for general likelihoods.
233 *Statistics & Probability Letters*, 187:109523, 2022.
- 234 [13] Pavel Izmailov, Sharad Vikram, Matthew D Hoffman, and Andrew Gordon Gordon Wilson.
235 What are bayesian neural network posteriors really like? In *International conference on ma-*
236 *chine learning*, pages 4629–4640. PMLR, 2021.
- 237 [14] Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer,
238 and Philipp Hennig. Laplace redux – effortless bayesian deep learning, 2022.
- 239 [15] David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for
240 statisticians. *Journal of the American Statistical Association*, 112(518):859–877, April 2017.
- 241 [16] Takahiro Inagaki, Yoshitaka Haribara, Koji Igarashi, Tomohiro Sonobe, Shuhei Tamate, Toshi-
242 mori Honjo, Alireza Marandi, Peter L. McMahon, Takeshi Umeki, Koji Enbutsu, Osamu
243 Tadanaga, Hirokazu Takenouchi, Kazuyuki Aihara, Ken-ichi Kawarabayashi, Kyo Inoue,
244 Shoko Utsunomiya, and Hiroki Takesue. A coherent ising machine for 2000-node optimization
245 problems. *Science*, 354(6312):603–606, 2016.
- 246 [17] Jeffrey Chou, Suraj Bramhavar, Siddhartha Ghosh, and William Herzog. Analog coupled
247 oscillator based weighted ising machine. *Scientific reports*, 9(1):14786, 2019.

- 248 [18] Naeimeh Mohseni, Peter L. McMahon, and Tim Byrnes. Ising machines as hardware solvers
249 of combinatorial optimization problems. *Nat. Rev. Phys.*, 4(6):363–379, 2022.
- 250 [19] Y Yamamoto, T Leleu, S Ganguli, and H Mabuchi. Coherent ising machines—quantum optics
251 and neural network perspectives. *Applied Physics Letters*, 117(16):160501, 2020.
- 252 [20] Tianshi Wang and Jaijeet Roychowdhury. Oim: Oscillator-based ising machines for solving
253 combinatorial optimisation problems. In *Unconventional Computation and Natural Computa-
254 tion: 18th International Conference, UCNC 2019, Tokyo, Japan, June 3–7, 2019, Proceedings
255 18*, pages 232–256. Springer, 2019.
- 256 [21] J. Kaiser, S. Datta, and B. Behin-Aein. Life is probabilistic—why should all our computers
257 be deterministic? computing with p-bits: Ising solvers and beyond. In *2022 International
258 Electron Devices Meeting (IEDM)*, pages 21–4. IEEE, 2022.
- 259 [22] Kerem Y. Camsari, Brian M. Sutton, and Supriyo Datta. p-bits for probabilistic spin logic.
260 *Appl. Phys. Rev.*, 6(1):011305, mar 2019.
- 261 [23] Navid Anjum Aadit, Andrea Grimaldi, Mario Carpentieri, Luke Theogarajan, John M. Marti-
262 nis, Giovanni Finocchio, and Kerem Y. Camsari. Massively parallel probabilistic computing
263 with sparse Ising machines. *Nat. Electron.*, 5(7):460–468, 2022.
- 264 [24] Tom Conte, Erik DeBenedictis, Natesh Ganesh, Todd Hylton, John Paul Strachan, R Stanley
265 Williams, Alexander Alemi, Lee Altenberg, Gavin E. Crooks, James Crutchfield, et al. Ther-
266 modynamic computing. *arXiv preprint arXiv:1911.01968*, 2019.
- 267 [25] Todd Hylton. Thermodynamic neural network. *Entropy*, 22(3):256, 2020.
- 268 [26] Denis Melanson, Mohammad Abu Khater, Maxwell Aifer, Kaelan Donatella, Max Hunter
269 Gordon, Thomas Ahle, Gavin Crooks, Antonio J Martinez, Faris Sbahi, and Patrick J Coles.
270 Thermodynamic computing system for AI applications. *arXiv preprint arXiv:2312.04836*,
271 2023.
- 272 [27] Maxwell Aifer, Kaelan Donatella, Max Hunter Gordon, Samuel Duffield, Thomas Ahle, Daniel
273 Simpson, Gavin E Crooks, and Patrick J Coles. Thermodynamic linear algebra. *arXiv preprint
274 arXiv:2308.05660*, 2023.
- 275 [28] Samuel Duffield, Maxwell Aifer, Gavin Crooks, Thomas Ahle, and Patrick J Coles. Thermody-
276 namic matrix exponentials and thermodynamic parallelism. *arXiv preprint arXiv:2311.12759*,
277 2023.
- 278 [29] Patryk Lipka-Bartosik, Martí Perarnau-Llobet, and Nicolas Brunner. Thermodynamic comput-
279 ing via autonomous quantum thermal machines. *arXiv preprint arXiv:2308.15905*, 2023.
- 280 [30] Patrick J Coles, Collin Szczepanski, Denis Melanson, Kaelan Donatella, Antonio J Martinez,
281 and Faris Sbahi. Thermodynamic ai and the fluctuation frontier. In *2023 IEEE International
282 Conference on Rebooting Computing (ICRC)*, pages 1–10. IEEE, 2023.
- 283 [31] Radford M Neal. Mcmc using hamiltonian dynamics. *arXiv preprint arXiv:1206.1901*, 2012.
- 284 [32] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics.
285 In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages
286 681–688. Citeseer, 2011.
- 287 [33] Kaelan Donatella, Samuel Duffield, Maxwell Aifer, Denis Melanson, Gavin Crooks, and
288 Patrick J Coles. Thermodynamic natural gradient descent. *arXiv preprint arXiv:2405.13817*,
289 2024.
- 290 [34] Thomas Minka. Bayesian linear regression. Technical report, Citeseer, 2000.
- 291 [35] Samuel Duffield, Kaelan Donatella, and Denis Melanson. thermox: Exact ou processes with
292 jax. <https://github.com/normal-computing/thermox>, 2024.
- 293 [36] Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression.
294 2004.

- 295 [37] Kevin P Murphy. *Probabilistic machine learning: an introduction*. MIT press, 2022.
- 296 [38] Jackson Gorham and Lester Mackey. Measuring sample quality with stein’s method. *Advances*
297 *in neural information processing systems*, 28, 2015.
- 298 [39] Crispin W Gardiner. *Handbook of stochastic methods for physics, chemistry and the natural*
299 *sciences. Springer series in synergetics*, 1985.
- 300 [40] DC Dowson and BV666017 Landau. The fréchet distance between multivariate normal distri-
301 *butions. Journal of multivariate analysis*, 12(3):450–455, 1982.
- 302 [41] J Leo van Hemmen and Tsuneya Ando. An inequality for trace ideals. *Communications in*
303 *Mathematical Physics*, 76:143–148, 1980.

304 **A Analysis of Gaussian Bayesian Inference Circuit**

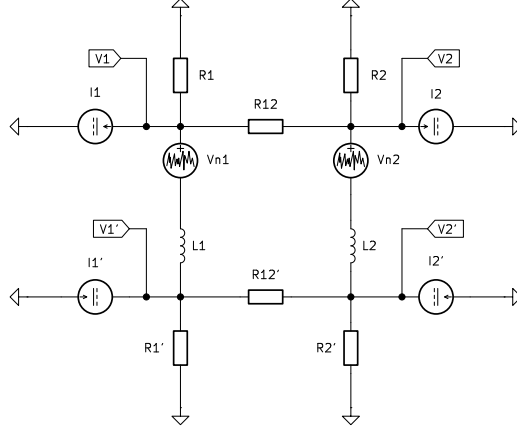


Figure 5: Circuit schematic for the Gaussian-Gaussian model posterior sampling device.

305 In Figure 5, positive current goes up through the two inductors, left to right through R_{12} and R'_{12} ,
 306 and towards ground in the other resistors. The two inductors have the same inductance L . KCL
 307 gives

$$308 \quad I_{L1} - I_1 = I_{R1} + I_{12} \quad (16)$$

$$309 \quad I_{L2} - I_2 = I_{R2} - I_{12} \quad (17)$$

$$310 \quad -I_{L1} + I'_1 = I'_{R1} + I'_{12} \quad (18)$$

$$311 \quad -I_{L2} + I'_2 = I'_{R2} - I'_{12}. \quad (19)$$

311 Using Ohm's law,

$$312 \quad I_{L1} - I_1 = R_1^{-1}V_1 + R_{12}^{-1}(V_1 - V_2) = (R_1^{-1} + R_{12}^{-1})V_1 - R_{12}^{-1}V_2 \quad (20)$$

$$313 \quad I_{L2} - I_2 = R_2^{-1}V_2 - R_{12}^{-1}(V_1 - V_2) = (R_2^{-1} + R_{12}^{-1})V_2 - R_{12}^{-1}V_1. \quad (21)$$

313 These can be written as a single vector equation as follows

$$314 \quad I_L - I = \mathcal{G}V, \quad (22)$$

314 where $I_L = (I_{L1} \ I_{L2})^\top$, $I = (I_1 \ I_2)^\top$, and

$$315 \quad \mathcal{G} = \begin{pmatrix} R_1^{-1} + R_{12}^{-1} & -R_{12}^{-1} \\ -R_{12}^{-1} & R_2^{-1} + R_{12}^{-1} \end{pmatrix}. \quad (23)$$

315 Similarly, for the lower subcircuit we have

$$316 \quad -I_L + I' = \mathcal{G}'V'. \quad (24)$$

316 The inductors obey the equations

$$317 \quad L_1 \dot{I}_{L1} = V'_1 - (V_1 - V_{n1}) \quad (25)$$

$$318 \quad L_2 \dot{I}_{L2} = V'_2 - (V_2 - V_{n2}), \quad (26)$$

318 or in vector notation

$$319 \quad L \dot{I}_L = V' - V + V_n. \quad (27)$$

319 Substituting in the expressions for V and V' derived before, we have

$$320 \quad L \dot{I}_L = \mathcal{G}'^{-1}(I' - I_L) - \mathcal{G}^{-1}(I_L - I) + V_n, \quad (28)$$

320 or

$$321 \quad \dot{I}_L = -L^{-1}\mathcal{G}^{-1}(I_L - I) - L^{-1}\mathcal{G}'^{-1}(I_L - I') + L^{-1}V_n. \quad (29)$$

$$322 \quad dI_L = -L^{-1}\mathcal{G}^{-1}(I_L - I) dt - L^{-1}\mathcal{G}'^{-1}(I_L - I') dt + L^{-1}\sqrt{SN}[0, \mathbb{I} dt]. \quad (30)$$

321 We now proceed to non-dimensionalize the above equation. $\mathcal{G} = \tilde{R}^{-1}A$.

$$\tilde{I}d\theta = -\tilde{I}\tilde{R}L^{-1}A^{-1}(\theta - \mu)dt - \tilde{I}\tilde{R}L^{-1}A^{-1}(\theta - \mu')dt + L^{-1}\sqrt{S}\mathcal{N}[0, \mathbb{I}dt]. \quad (31)$$

322 Define $\tau = L/\tilde{R}$, giving

$$d\theta = -A^{-1}(\theta - \mu)\tau^{-1}dt - A'^{-1}(\theta - \mu')\tau^{-1}dt + \tilde{I}^{-1}L^{-1}\sqrt{S}\mathcal{N}[0, \mathbb{I}dt]. \quad (32)$$

323 If we set $S = 2\tilde{I}^2L\tilde{R}$, then we have

$$d\theta = -A^{-1}(\theta - \mu)\tau^{-1}dt - A'^{-1}(\theta - \mu')\tau^{-1}dt + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt]. \quad (33)$$

324 B Analysis of Bayesian Logistic Regression Circuit

325 We now analyze the circuit in Figure 7. The boxes labeled Diff. Pair represent differential pairs of
 326 NPN bipolar junction transistors (BJTs), as shown in Fig. 6. To achieve a working implementation,
 327 additional circuitry is needed to support the differential pair and assure that it is appropriately biased,
 328 including a power source and possibly current mirrors.

329 The following conventions for current flow will be used

- 330 • I_C is the current *into* the collector of a transistor. I_B is the current *into* the base of a
- 331 transistor. I_E is the current *out of* the emitter of a transistor.
- 332 • The output current I_o of a differential pair is the current that flows *into* the collector of the
- 333 BJT labeled Q_a .
- 334 • Positive current flows in the direction of the arrow through all current sources.
- 335 • Positive current flows downwards through C_1 and C_2 and from left to right through R_{12} .
- 336 • Through resistors R_{A11} , R_{B11} , etc. positive current always flows towards the base of the
- 337 transistor.

338 B.1 Analysis of the BJT differential pair

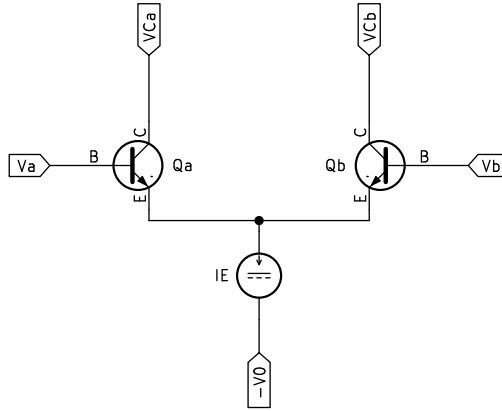


Figure 6: Circuit schematic for the BJT differential pair.

339 We first consider the behavior of differential pair subcircuit, which can be explained using the Ebers-
 340 Moll model. The Ebers-Moll model describes the BJT in active mode, meaning when $V_E < V_B <$
 341 V_C , and the circuit must be appropriately biased at all times to ensure the device is always in active
 342 mode. According to this model, in active mode the following relations are satisfied

$$I_C = I_S \left(e^{(V_B - V_E)/V_T} - 1 \right), \quad (34)$$

343

$$I_C = \alpha I_E, \quad (35)$$

344 where I_S is the saturation current, V_T the thermal voltage, and α is the common-base current gain.
 345 I_S is typically on the order of 10^{-15} to 10^{-12} Amps, and at room temperature $V_T = 25.3\text{mV}$.
 346 The parameter α is between 0.98 and 1. It follows from Kirchoff's current law (KCL) that $I_B =$
 347 $(1 - \alpha)I_E$. For these typical values of the parameters appearing in Eq. (34) the subtraction of unity
 348 in parentheses can safely be ignored, which we will do in what follows. In order for the Ebers-Moll
 349 model to be valid, the voltage V_0 should be determined such that $V_C > V_B > V_E$ for all transistors
 350 at all times, but the value of V_0 is otherwise unimportant.

351 To analyze the differential pair of transistors Q_a and Q_b , observe that (by KCL)

$$I_{Ea} + I_{Eb} = I_E. \quad (36)$$

352 We must distinguish between the two base voltages V_a and V_b , but the two emitter voltages are the
 353 same, so we write $V_E = V_{Ea} = V_{Eb}$. Using Eqs. (34) and (35) then,

$$I_E = \frac{I_S}{\alpha} e^{-V_E/V_T} \left(e^{V_a/V_T} + e^{V_b/V_T} \right), \quad (37)$$

354 where we have dropped the -1 as explained earlier. Now the emitter current I_{Ea} can be written as

$$I_{Ea} = \frac{I_S}{\alpha} e^{(V_a - V_E)/V_T} \quad (38)$$

$$= \frac{I_E e^{V_a/V_T}}{e^{V_a/V_T} + e^{V_b/V_T}} \quad (39)$$

$$= \frac{I_E}{1 + e^{-(V_a - V_b)/V_T}}, \quad (40)$$

355 and similarly

$$I_{Eb} = \frac{I_E}{1 + e^{(V_a - V_b)/V_T}}. \quad (41)$$

356 Equation (35) is then used to find the collector currents

$$I_{Ca} = \frac{\alpha I_E}{1 + e^{-(V_a - V_b)/V_T}}, \quad (42)$$

357

$$I_{Cb} = \frac{\alpha I_E}{1 + e^{(V_a - V_b)/V_T}}. \quad (43)$$

358 The base voltages V_a and V_b are still undetermined. However, we will assume the limit $\alpha \rightarrow 1$,
 359 where the base current goes to zero. In this limit, the two transistor bases may be connected to
 360 nodes in an external circuit to set their voltages. As there is no base current, these connections do
 361 not affect the voltages in the external circuit. In what follows, we will consider I_{Ca} the output of the
 362 differential pair, and label this current I_o . Again taking the limit $\alpha \rightarrow 1$, we have

$$I_o = \frac{I_E}{1 + e^{(V_a - V_b)/V_T}} = I_E L(-(V_a - V_b)/V_T), \quad (44)$$

363 where $L(z) = 1/(1 + e^{-z})$ is the standard logistic function. Note that the support circuitry may
 364 include a current mirror that inverts the sign of the output current. As this formally has the same
 365 effect as a negative value of I_E , we will allow I_E to be negative in what follows.

366 B.2 Analysis of the logistic regression circuit

367 As the BJT bases draw negligible current, the voltages V_{a1} , V_{b1} , V_{a2} , and V_{b2} in the circuit can be
 368 determined by considering the circuit in the absence of the differential pairs. In this case, we see
 369 that (by KCL)

$$R_{a11}^{-1}(V_{C1} - V_{a1}) + R_{a12}^{-1}(V_{C2} - V_{a1}) - R_{a10}^{-1}V_{a1} = 0, \quad (45)$$

370 and solving for V_{a1} gives

$$V_{a1} = \frac{R_{a11}V_{C1} + R_{a12}V_{C2}}{R_{a11} + R_{a12} + R_{a11}R_{a12}R_{a10}^{-1}} = \frac{R_{a11}^{-1}V_{C1} + R_{a12}^{-1}V_{C2}}{R_{a10}^{-1} + R_{a11}^{-1} + R_{a12}^{-1}}. \quad (46)$$

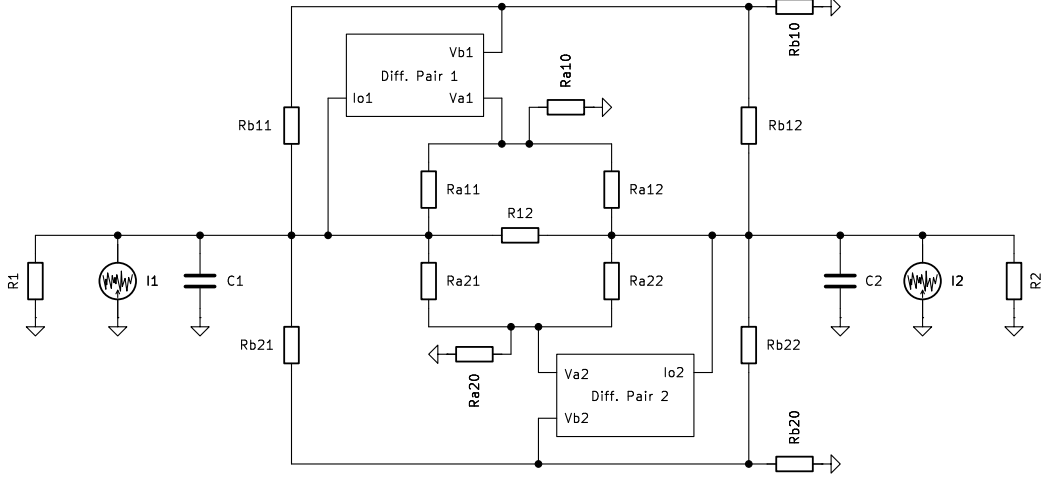


Figure 7: Circuit schematic for the logistic regression posterior sampling device.

371 The same reasoning applies for V_{b1} , resulting in

$$V_{b1} = \frac{R_{b11}^{-1}V_{C1} + R_{b12}^{-1}V_{C2}}{R_{b10}^{-1} + R_{b11}^{-1} + R_{b12}^{-1}}, \quad (47)$$

372 SO

$$V_{a1} - V_{b1} = \frac{g_{a11}V_{C1} + g_{a12}V_{C2}}{g_{a10} + g_{a11} + g_{a12}} - \frac{g_{b11}V_{C1} + g_{b12}V_{C2}}{g_{b10} + g_{b11} + g_{b12}}, \quad (48)$$

373 where we have written the previous results in terms of the conductance $g = R^{-1}$. The above can be
 374 written more conveniently by defining the vectors $\hat{g}_{a1} = (g_{a10} + g_{a11} + g_{a12})^{-1}(g_{a11}, g_{a12})^\top$ and
 375 $\hat{g}_{b1} = (g_{b10} + g_{b11} + g_{b12})^{-1}(g_{b11}, g_{b12})^\top$, in terms of which we have

$$V_{a1} - V_{b1} = (\hat{g}_a - \hat{g}_b)^\top V_C. \quad (49)$$

376 or, defining $\hat{g}_1 = \hat{g}_{a1} - \hat{g}_{b1}$, we simply have

$$V_{a1} - V_{b1} = \hat{g}_1^\top V_C. \quad (50)$$

377 The latter result can be plugged into Eq. (44) to get I_{o1} ,

$$I_{o1} = I_{E1}L(-\hat{g}_1^\top V_C/V_T), \quad (51)$$

378 where, as before, $L(z) = 1/(1 + e^{-z})$ is the standard logistic function. By an identical derivation
 379 to the one above, a similar relation holds for the lower subcircuit

$$I_{o2} = I_{E2}L(-\hat{g}_2^\top V_C/V_T). \quad (52)$$

380 We also assume that all resistors R_{aij} , R_{bij} are very large compared to R_{12} so the current flowing
 381 through these resistors can be treated as negligible. This assumption does not affect the function
 382 of resistors R_{aij} , R_{bij} because only the ratios of these resistances determine the voltages V_{ai} , V_{bi} .
 383 Next, we apply KCL to the nodes at the top of capacitors C_1 and C_2

$$-I_{C1} + I_1 - R_1^{-1}V_{C1} + R_{12}^{-1}(V_{C2} - V_{C1}) - I_{o1} = 0, \quad (53)$$

384 Similarly, KCL for the node above capacitor C_2 reads

$$-I_{C2} + I_2 - R_2^{-1}V_{C2} + R_{12}^{-1}(V_{C1} - V_{C2}) - I_{o2} = 0. \quad (54)$$

385 Substituting in the expressions derived for the collector currents, we then have

$$-I_{C1} + I_1 - R_1^{-1}V_{C1} + R_{12}^{-1}(V_{C2} - V_{C1}) - I_{E1}L(-\hat{g}_1^\top V_C/V_T) = 0, \quad (55)$$

386

$$-I_{C2} + I_2 - R_2^{-1}V_{C2} + R_{12}^{-1}(V_{C1} - V_{C2}) - I_{E2}L(-\hat{g}_2^\top V_C/V_T) = 0. \quad (56)$$

387 Next we define the conductance matrix

$$\mathcal{G} = \begin{pmatrix} R_1^{-1} + R_{12}^{-1} & -R_{12}^{-1} \\ -R_{12}^{-1} & R_2^{-1} + R_{12}^{-1} \end{pmatrix}, \quad (57)$$

388 allowing us to write a single vector equation

$$-I_C + I - \mathcal{G}V_C - I_E L(-\hat{g}^\top V_C/V_T) = 0, \quad (58)$$

389 where we have also set $\hat{g}_1 = \hat{g}_2$. Now using the fact that $dV_C/dt = C^{-1}I_C$, we have the following
390 vector differential equation

$$C \frac{dV_C}{dt} = -\mathcal{G}V_C - I_E L(-\hat{g}^\top V_C/V_T) + I. \quad (59)$$

391 We assume the current vector I has a DC component I_{DC} and a noise component I_{noise} . The noise
392 component is assumed to be an ideal white noise process of infinite bandwidth and power spectral
393 density S , which we write $I_{\text{noise}} = \sqrt{S}\xi(t)$. Altogether, we get the stochastic differential equation

$$dV_C = -C^{-1}\mathcal{G}V_C dt + C^{-1}I_{DC} dt - C^{-1}\alpha I_E L(\hat{g}^\top V_C/V_T) dt + C^{-1}\sqrt{S}\xi(t) dt. \quad (60)$$

394 Using the identity $\xi(t) dt = \mathcal{N}[0, dt]$, this becomes

$$dV_C = -C^{-1}\mathcal{G}V_C dt + C^{-1}I_{DC} dt - C^{-1}\alpha I_E L(\hat{g}^\top V_C/V_T) dt + C^{-1}\sqrt{S}\mathcal{N}[0, dt]. \quad (61)$$

395 At this point it is convenient to define dimensionless quantities which are mapped to the physical
396 parameters of the circuit. Define $\theta = V_C/\tilde{V}$, $\Sigma^{-1} = \tilde{R}\mathcal{G}$, $\Sigma^{-1}\mu = I_{DC}/\tilde{I}$, and $yx = -I_E/\tilde{I}$. Our
397 equation now takes the form

$$\tilde{V} dx = -\tilde{V}\tilde{R}^{-1}C^{-1}\Sigma^{-1}x dt + C^{-1}\tilde{I}\Sigma^{-1}\mu dt + C^{-1}\tilde{I}L(-\hat{g}^\top V_C/V_T)yx dt + C^{-1}\sqrt{S}\mathcal{N}[0, \mathbb{I} dt]. \quad (62)$$

398 Next, let $\tau = \tilde{R}C$, and set $\tilde{I} = \tilde{V}C/\tau$ and $S = 2\tilde{V}^2C^2/\tau$. In this case,

$$d\theta = -\Sigma^{-1}\theta\tau^{-1}dt + \Sigma^{-1}\mu\tau^{-1}dt + L(-\hat{g}^\top V_C/V_T)yx\tau^{-1}dt + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt]. \quad (63)$$

399 Finally, we set $\hat{g} = yxV_T/\tilde{V}$ to obtain

$$d\theta = -\Sigma^{-1}(\theta - \mu)\tau^{-1}dt + L(-y\theta^\top x)yx\tau^{-1}dt + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt], \quad (64)$$

400 which is identical to Eq. (15)

401 C Analysis of complexity of Gaussian-Gaussian posterior sampling

402 In this section we analyze the time-complexity of sampling the Bayesian posterior of the Gaussian-
403 Gaussian model using the device in Fig. 1. As shown in Appendix A, the SDE for this circuit
404 is

$$d\theta = -\Sigma^{-1}(\theta - \mu)\tau^{-1}dt - \Sigma_{y|\theta}^{-1}(\theta - y)\tau^{-1}dt + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt], \quad (65)$$

405 where $\tau = L/\tilde{R}$. This SDE may also be written in terms of the posterior parameters,

$$d\theta = -\Sigma_{\theta|y}^{-1}(\theta - \mu_{\theta|y})\tau^{-1}dt + \mathcal{N}[0, 2\mathbb{I}\tau^{-1}dt], \quad (66)$$

406 where

$$\mu_{\theta|y} = \mu + \Sigma(\Sigma + \Sigma_{y|\theta})^{-1}(y - \mu), \quad (67)$$

407

$$\Sigma_{\theta|y} = \Sigma - \Sigma(\Sigma + \Sigma_{y|\theta})^{-1}\Sigma. \quad (68)$$

408 The above equation is in the form of a multivariate Ornstein-Uhlenbeck (OU) process [39].

409 The squared Wasserstein distance between the distribution at time t and the target posterior distri-
410 bution is [40]

$$W(t)^2 = \|\mu(t) - \mu_{\theta|y}\|_2^2 + D(\Sigma(t), \Sigma_{\theta|y}), \quad (69)$$

411 where

$$D(\Sigma_1, \Sigma_2) = \text{tr} \left\{ \Sigma_1 + \Sigma_2 - 2 \left(\Sigma_2^{1/2} \Sigma_1 \Sigma_2^{1/2} \right)^{1/2} \right\}. \quad (70)$$

412 Note that the second term in the squared Wasserstein distance is bounded above as [41]

$$D(\Sigma_1, \Sigma_2) \leq \frac{\|\Sigma_1 - \Sigma_2\|_F^2}{(\sqrt{\lambda_{\min}(\Sigma_1)} + \sqrt{\lambda_{\min}(\Sigma_2)})^2} \leq \frac{d\|\Sigma_1 - \Sigma_2\|^2}{(\sqrt{\lambda_{\min}(\Sigma_1)} + \sqrt{\lambda_{\min}(\Sigma_2)})^2} \quad (71)$$

413 Let α_{\min} be the smallest eigenvalue of $\Sigma_{\theta|y}$. At all times $D(\Sigma(t), \Sigma_{\theta|y})$ is bounded above as

$$D(\Sigma(t), \Sigma_{\theta|y}) \leq \alpha_{\min}^{-1} d \|\Sigma_1 - \Sigma_2\|^2 \quad (72)$$

414 For an OU process the mean behaves as $\mu(t) - \mu_{\theta|y} = e^{-\Sigma_{\theta|y}t/\tau}(\mu(0) - \mu_{\theta|y})$ [39], so the distance
415 between the mean and the target mean is bounded as

$$\|\mu(t) - \mu_{\theta|y}\| = \left\| e^{-\Sigma_{\theta|y}t/\tau}(\mu(0) - \mu_{\theta|y}) \right\| \leq e^{-\alpha_{\min}t/\tau} \|\mu(0) - \mu_{\theta|y}\|. \quad (73)$$

416 We now bound the distance between the covariance $\Sigma(t)$ and the target covariance. We use the
417 formula [39]

$$\Sigma(t) = \mathcal{P}_t(\Sigma(0)) + \int_0^t dt' \mathcal{P}_{t'}(2\tau^{-1}\mathbb{I}). \quad (74)$$

418 where $\mathcal{P}_t(X) = e^{-\Sigma_{\theta|y}t/\tau} X e^{-\Sigma_{\theta|y}t/\tau}$. As $\lim_{t \rightarrow \infty} \Sigma(t) = \Sigma_{\theta|y}$, we have

$$\Sigma(t) - \Sigma_{\theta|y} = \mathcal{P}_t(\Sigma(0)) + \int_0^t dt' \mathcal{P}_{t'}(2\tau^{-1}\mathbb{I}) - \int_0^\infty dt' \mathcal{P}_{t'}(2\tau^{-1}\mathbb{I}) \quad (75)$$

$$= \mathcal{P}_t(\Sigma(0)) - \int_t^\infty dt' \mathcal{P}_{t'}(2\tau^{-1}\mathbb{I}). \quad (76)$$

419 Taking the norm and using the triangle inequality gives

$$\|\Sigma(t) - \Sigma_{\theta|y}\| \leq \|\mathcal{P}_t(\Sigma(0))\| + \int_t^\infty dt' \|\mathcal{P}_{t'}(2\tau^{-1}\mathbb{I})\|. \quad (77)$$

420 As the norm is submultiplicative $\|\mathcal{P}(X)\| \leq e^{-2\alpha_{\min}t/\tau} \|X\|$, so

$$\|\Sigma(t) - \Sigma_{\theta|y}\| \leq e^{-2\alpha_{\min}t/\tau} \|\Sigma(0)\| + 2\tau^{-1} \int_t^\infty dt' e^{-2\alpha_{\min}t'/\tau} \quad (78)$$

$$= e^{-2\alpha_{\min}t/\tau} (\|\Sigma(0)\| + \alpha_{\min}^{-1}). \quad (79)$$

421 If we would like to have $W(t) \leq W_0$, then we can demand that $\|\mu(t) - \mu_{\theta|y}\|^2 \leq \frac{1}{2}W_0^2$ and
422 $D(\Sigma(t), \Sigma_{\theta|y}) \leq \frac{1}{2}W_0^2$. The inequality involving the mean is satisfied when

$$t \geq \alpha_{\min}^{-1} \tau \ln(\sqrt{2} \|\mu(0) - \mu_{\theta|y}\| W_0^{-1}). \quad (80)$$

423 The inequality involving the covariance is satisfied when

$$\|\Sigma(t) - \Sigma_{\theta|y}\| \leq \frac{1}{\sqrt{2}} W_0 \alpha_{\min}^{1/2} d^{-1/2} \quad (81)$$

424 This, in turn, is satisfied if

$$e^{-2\alpha_{\min}t/\tau} (\|\Sigma(0)\| + \alpha_{\min}^{-1}) \leq \frac{1}{\sqrt{2}} W_0 \alpha_{\min}^{1/2} d^{-1/2}. \quad (82)$$

425 The matrix $\Sigma(0)$ can always be chosen to be zero, so the above inequality becomes

$$e^{-2\alpha_{\min}t/\tau} \leq \frac{1}{\sqrt{2}} W_0 \alpha_{\min}^{3/2} d^{-1/2}. \quad (83)$$

426 We arrive at the requirement

$$t \geq \frac{1}{2} \alpha_{\min}^{-1} \tau \ln \left[\sqrt{2} \alpha_{\min}^{-3/2} d^{1/2} W_0^{-1} \right]. \quad (84)$$

427 Therefore $W(t) \leq W_0$ if the following unified bound is satisfied

$$t \geq \alpha_{\min}^{-1} \tau \max \left(\ln \left[\sqrt{2} \|\mu(0) - \mu_{\theta|y}\| W_0^{-1} \right], \frac{1}{2} \ln \left[\sqrt{2} \alpha_{\min}^{-3/2} d^{1/2} W_0^{-1} \right] \right). \quad (85)$$

428 The quantity $\|\mu(0) - \mu_{\theta|y}\|$ may hide some dependence on dimension, which is discussed presently.
 429 It is assumed that the quantity $c = \|\mu_{\theta|y}\|/\sqrt{\alpha_{\max}}$ has an upper bound c_{\max} which is independent
 430 of dimension, where α_{\max} is the largest eigenvalue of $\Sigma_{\theta|y}$. That is, the mean of the posterior may
 431 be at most c_{\max} standard deviations away from the origin, independent of dimension. This choice
 432 represents a particular scaling regime, which we feel is a realistic representation of the accuracy
 433 requirements for many applications. We may also choose $\mu(0) = 0$, and this leads to the requirement

$$t \geq \max \alpha_{\min}^{-1} \tau \left(\ln \left[\sqrt{2} c \sqrt{\alpha_{\max}} W_0^{-1} \right], \frac{1}{2} \ln \left[\sqrt{2} \alpha_{\min}^{-3/2} d^{1/2} W_0^{-1} \right] \right). \quad (86)$$

434 In general, the problem may be rescaled in such a way that $\alpha_{\max} \leq 1$, and some rescaling of this
 435 kind is realistic given that a particular device will have a specific signal range (that is, the range over
 436 which voltages and currents may vary). Redefining the problem this way will also cause the smallest
 437 eigenvalue of $\Sigma_{\theta|y}$ to be reduced by a factor of α_{\max} , and in this case the bound would be

$$t \geq \max \kappa \tau \left(\ln \left[\sqrt{2} c W_0^{-1} \right], \frac{1}{2} \ln \left[\sqrt{2} \kappa^{3/2} d^{1/2} W_0^{-1} \right] \right), \quad (87)$$

438 where $\kappa = \alpha_{\max}/\alpha_{\min}$ is the condition number. Subject to these assumptions, we may express the
 439 asymptotic time complexity as

$$t = O(\kappa \tau \ln(\kappa^{3/2} d^{1/2} W_0^{-1})) \quad (88)$$

440 In order to collect N samples the same process is run N times, resulting in complexity

$$t = O(N \kappa \tau \ln(\kappa^{3/2} d^{1/2} W_0^{-1})). \quad (89)$$

441 D Conditioning on multiple I.I.D. samples

442 When conditioning on a single sample y , the energy U can be separated into two terms, one mapping
 443 to the prior and the other to the likelihood:

$$U(r) = U_{\pi}(r) + U_{\ell}(r), \quad (90)$$

444 where $\beta U_{\pi}(r) = -\ln p_{\theta}(r/\tilde{r})$ and $\beta U_{\ell}(r) = -\ln p_{y|\theta}(y|r/\tilde{r})$. In general we may have a number
 445 of I.I.D. samples $Y = (y_1, \dots, y_N)$, and would like to sample from $p_{\theta|Y}(\theta|Y)$. Because the samples
 446 of y are I.I.D., we have

$$p_{Y|\theta}(Y|\theta) = \prod_{i=1}^N p_{y_i|\theta}(y_i|\theta). \quad (91)$$

447 In this case the likelihood part of the potential energy takes the form

$$\beta U_{\ell}(r) = -\sum_{i=1}^N \ln p_{y_i|\theta}(y_i|r/\tilde{r}), \quad (92)$$

448 while the prior part is the same as in the single-sample case, $\beta U_{\pi}(r) = -\ln p_{\theta}(r/\tilde{r})$. This form
 449 of the potential energy has a convenient physical interpretation: the function U_{π} can be interpreted
 450 as the self-energy of the system in state r (that is when it is decoupled from an external system),
 451 while the function $U_{\ell}(\theta)$ can be viewed as an interaction energy between the state r and the state
 452 y of an external system. When there are multiple I.I.D. samples, this is analogous to the state r
 453 interacting with a collection of external systems in states $Y = (y_1 \dots y_N)$, and each such interaction
 454 contributes its own term to the interaction energy. This provides a framework for building a physical
 455 device to sample from the posterior conditioned on multiple I.I.D. samples; one must simply couple
 456 a collection of external systems in states $Y = (y_1 \dots y_N)$ to the system in such a way that each
 457 interaction contributes an energy of $-\ln p_{y_i|\theta}(y_i|r/\tilde{r})$.

458 We will now describe another approach to building a physical device that samples from the posterior
 459 conditioned on multiple I.I.D. samples of y . We first observe that the Langevin equation for the
 460 device in this case must be

$$d\theta = \nabla_{\theta} \ln p_{\theta}(\theta) \tau^{-1} dt + \sum_{i=1}^N \nabla_{\theta} \ln p_{y_i|\theta}(y_i|\theta) dt + \mathcal{N}[0, 2\tau^{-1} dt], \quad (93)$$

461 As discussed above, if we have a device that can implement the N likelihood drift terms simultane-
 462 ously then the problem is solved. However, suppose that we have a device that is only capable of
 463 implementing a single likelihood term at a time, but y may be varied as a function of time. Addi-
 464 tionally, we make the interaction energy for this device larger by a factor of N for reasons that will
 465 become clear. That is, we have a device that implements an SDE of the form

$$d\theta = \nabla_{\theta} \ln p_{\theta}(\theta) \tau^{-1} dt + N \nabla_{\theta} \ln p_{y|\theta}(y(t)|\theta) dt + \mathcal{N}[0, 2\tau^{-1} dt]. \quad (94)$$

466 We may choose a short time duration Δt , and set

$$y(t) = y_{\lfloor t/\Delta t \rfloor \bmod N+1}. \quad (95)$$

467 So for $0 \leq t \leq \Delta t$ we set $y(t) = y_1$, for $\Delta t < t \leq 2\Delta t$ we set $y(t) = y_2$, and so on. Once
 468 $t > N\Delta t$ we start over at y_1 and continue cycling over all of the I.I.D. samples. Suppose that Δt is
 469 short enough that all of the samples are cycled over before the state θ changes significantly. We may
 470 then average drift term $N \nabla_{\theta} \ln p_{y|\theta}$ over a period of time $N\Delta t$ and consider θ constant within this
 471 average. Carrying out this time average, we find

$$\frac{1}{N\Delta t} \sum_{i=1}^N \Delta t N \nabla_{\theta} \ln p_{y|\theta}(y_i|\theta) = \sum_{i=1}^N \nabla_{\theta} p_{y|\theta}(y_i|\theta), \quad (96)$$

472 resulting in the correct form of the Langevin equation.

473 E Computational complexity of logistic regression

474 The runtime complexity of digital Langevin sampling of a logistic regression model is $O(n_{\delta t} d N)$,
 475 with $n_{\delta t}$ the number of time steps, K the number of trainable parameters, and N the number of data
 476 points (in the case of minibatching b replaces N). Added to this, there can be some discretization
 477 error if the step size is chosen too large, which generally means the number of time steps is made
 478 quite large to avoid this (meaning that $n_{\delta t} \gg n_s$). The memory complexity is that of storing the data
 479 and the samples, hence is $O(d n_s + N)$. In contrast, running the thermodynamic logistic regression
 480 algorithm only includes two digital steps: i) pre-processing and sending over the data to the hardware
 481 and ii) initializing the system, which involves setting the prior distribution and the initial state. The
 482 gradient evaluations are all done in analog, which incurs a cost of $O(t)$, with t the analog dynamics
 483 time. In the best case scenario, where we do not oversample correlated samples, we have $t = n_s \tau_c$,
 484 with τ_c the correlation time. The runtime complexity of the thermodynamic solver is therefore
 485 $O(d + N + n_s \tau_c)$, which is a large improvement over the digital case since there is no discretization
 486 factor and less multiplicative factors. In addition, note that t can be made extremely small (of the
 487 order of the microsecond) in practice thanks to the value of the physical time constants of electronic
 488 systems.¹

¹Since the system is nonlinear, similar bounds to those presented in [27] cannot be obtained.