
Fast and accurate training and sampling of Restricted Boltzmann Machines

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Thanks to their simple architecture, Restricted Boltzmann Machines (RBMs) are
2 powerful tools for modeling complex systems and extracting interpretable insights
3 from data. However, training RBMs, as other energy-based models, on highly
4 structured data poses a major challenge, as effective training relies on mixing the
5 Markov chain Monte Carlo simulations used to estimate the gradient. This process
6 is often hindered by multiple second-order phase transitions and the associated
7 critical slowdown. In this paper, we present an innovative method in which the
8 principal directions of the dataset are integrated into a low-rank RBM through a
9 convex optimization procedure. This approach enables efficient sampling of the
10 equilibrium measure via a static Monte Carlo process. By starting the standard
11 training process with a model that already accurately represents the main modes of
12 the data, we bypass the initial phase transitions. Our results show that this strategy
13 successfully trains RBMs to capture the full diversity of data in datasets where
14 previous methods fail. Furthermore, we use the training trajectories to propose a
15 new sampling method, *parallel trajectory tempering*, which allows us to sample
16 the equilibrium measure of the trained model much faster than previous optimized
17 MCMC approaches and a better estimation of the log-likelihood. We illustrate the
18 success of the training method on several highly structured datasets.

19 1 Introduction

20 Energy-based models (EBMs) are a classic approach to generative modeling that has been studied for
21 decades. They were introduced using the Restricted Boltzmann Machine formulation by Smolen-
22 sky [1] and later further developed by Sejnowski et al. [2]. They provide a straightforward method for
23 modeling effective interactions within complex data distributions and for sufficiently simple energy
24 functions, such as the Boltzmann machine (BM) [3], it is also possible to interpret and infer the
25 underlying constituent rules from the observed data. This inference strategy is often associated with
26 the inverse Ising problem and pairwise interaction models [4], and it has found a great variety of
27 applications in fields such as neuroscience [5] or computational biology [6]. A recent work has
28 proposed replacing the use of pairwise models with the Restricted Boltzmann Machine (RBM) [7],
29 as it allows the same direct interpretation of its energy function as an explicit many-body interaction
30 model while greatly extending the expressive power of the model. RBMs are also very useful for
31 grouping data into hierarchical families [8]. On the diametrically opposite side (on interpretability)
32 are generative ConvNets [9, 10], where the energy function is formulated as a deep neural network,
33 which are capable of synthesizing photorealistic images but are almost impossible to interpret as a
34 physical model.

35 The applications of simple EBMs in science are very diverse. For example, they are often used today
36 to encode the Hamiltonian of physical many-body systems, such as Quantum wave functions [11]
37 or the accurate determination of ground state wave functions of strongly interacting and entangled

38 quantum spins [12] or they have proven to be suitable for the representation of the AdS/CFT
39 correspondence in theories of quantum gravity [13, 14]. Simple EBMs are also very common to
40 encode the evolutionary constraints in protein families [6, 15], and to predict mutations [16], or to
41 generate realistic synthetic sequences, such as fake human genomes [17, 18]. These examples show
42 that, despite their somewhat old-fashioned architecture, shallow EBMs are increasingly seen as useful
43 tools for better understanding modern physics/biology, as they allow for a certain level of analytical
44 description.

45 Despite the appealing modeling properties of RBMs, they are notoriously difficult to train, a challenge
46 common to EBMs in general. The main difficulty arises from the computation of the log-likelihood
47 gradient, which requires an ergodic exploration of a dynamically evolving and potentially complex
48 free energy landscape using Markov Chain Monte Carlo (MCMC) processes. Recent studies have
49 shown that models trained with non-convergent MCMC processes suffer from out-of-equilibrium
50 dynamic memory effects [19, 20, 21]. This dynamical behavior can be explained analytically using
51 moment-matching arguments [19, 22]. While exploiting these effects can yield fast and accurate
52 generative models, even for highly structured data [23] or high-quality images with RBMs [24], this
53 approach results in a sharp separation between the model’s Gibbs-Boltzmann distribution and the
54 dataset distribution, thereby undermining the interpretability of the model parameters [22, 7]. Thus,
55 to extract meaningful information from datasets using RBMs, it is essential to ensure proper mixing
56 of the chains during training, in short, one needs *equilibrium models*.

57 Both the ability to train an RBM in equilibrium and to generate convincing new samples from its
58 equilibrium measure strongly depend on the dataset in question. For typical image datasets such
59 as MNIST or CIFAR-10, good RBMs can be obtained by increasing the number of MCMC steps.
60 However, this approach is no longer feasible for highly structured datasets [25]. Datasets from
61 which one seeks scientific insights are often highly structured, such as genomics/proteomics data
62 or low-temperature many-body physical systems. These datasets typically exhibit distinct clusters,
63 identifiable via principal component analysis (PCA) which form distant groups of similar entries.
64 We show an example of the PCA of 4 clustered dataset we will be studying in this work in Fig. 1;
65 details about these datasets are given in the caption and in the Supplemental Information (SI). During
66 training, the model must evolve from an initial normal distribution to an increasingly multimodal
67 distribution. Sampling from multimodal distributions is particularly challenging because the mixing
68 times are determined by the transition times between modes. But this is not the only difficulty. These
69 distant modes are encoded by second-order phase transitions during training [26, 27, 28], leading to
70 diverging mixing times in these regions — a phenomenon known as *critical slowdown* —, which
71 means that mixing times are expected to grow with a power of their system size. This sampling
72 challenge not only hinders the training process, but also limits the model’s ability to generate new
73 samples. Obtaining new and independent configurations would require an impractically large number
74 of sampling steps.

75 2 Related work

76 Training EBMs by maximizing log-likelihood has long been a challenge in the community [29, 10].
77 EBMs gained popularity with the introduction of the contrastive divergence algorithm [30], in which
78 a set of parallel chains is initialized on independent examples in the minibatch and the MCMC
79 process iterates for a few steps. Despite its widespread use, this algorithm yields models with poor
80 equilibrium properties that are ineffective as generative models [31, 32, 21]. An improvement is the
81 persistent contrastive divergence (PCD) algorithm [33], which maintains a permanent chain in which
82 the last configurations used to estimate the previous gradient update are reused. PCD acts like a
83 slow annealing process improving gradient estimation quality. However, it often fails on clustered
84 data as the statistical properties of the permanent chain quickly move away from the equilibrium
85 measure and degrade the model [25]. This problem, which is primarily related to phase coexistence,
86 can be addressed with constrained MCMC methods if appropriate order parameters are identified.
87 For RBMs, these order parameters are related to the singular value decomposition of the model
88 coupling matrix, which enables efficient reconstruction of multimodal distributions [25]. Although
89 this method is effective for evaluating model quality, it is too computationally intensive to be used
90 in training, even if it leads to models with good equilibrium properties. Other optimized MCMC
91 methods, such as the Parallel Tempering (PT) [34] algorithm, simulate multiple models at different
92 temperatures, facilitating mixing through temperature exchange [35, 32]. However, PT is costly and

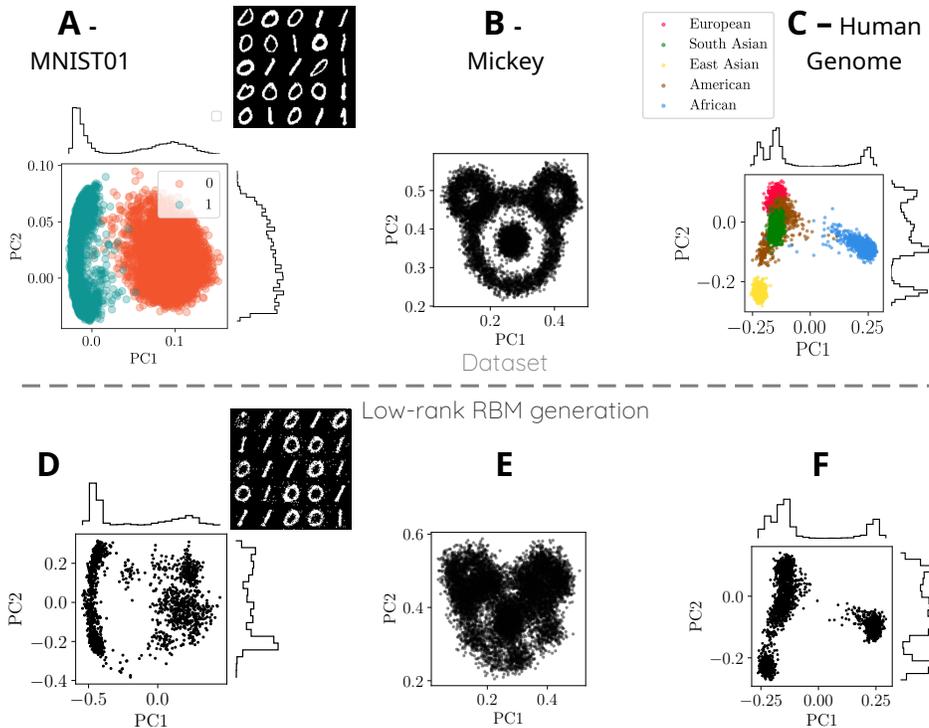


Figure 1: **Clustered datasets.** In A-C we show the 4 different clustered data sets that we will consider in this paper, projected onto their first two PCA components. In A we show the data of the MNIST 01 dataset (both projected and some instances), which contains only the 0-1 images of the complete MNIST dataset. In B, we show the Mickey dataset, an artificial dataset whose PCA forms a “Mickey” face shape. In C, we show data from the Human Genome Dataset (HGD), which contains binary vectors each corresponding to a human individual and whose sites correspond to selected genes. A value of 1 at a particular position means that a mutation was observed there compared to an individual reference sequence. Details of these data sets can be found in the SI. In D-F we show the samples we generate with the low-rank RBMs that are used as initial point of a standard training.

93 often ineffective, especially because EBMs undergo first-order phase transitions at the temperature
 94 where PT typically fails because one needs too many temperatures to make the moves accepted. We
 95 will see below that a more appropriate approach exchanges the models at different training times,
 96 which only implies crossing second-order phase transitions.

97 The population annealing algorithm, which reweights parallel chains during learning based on their
 98 relative weight changes during parameter updates, was proposed as an alternative [36]. Similarly,
 99 reweighting chains using non-equilibrium physics concepts such as the Jarzynski equality has been
 100 proposed [37]. Both approaches struggle with highly structured data sets. To prevent the different
 101 chains to get too correlated around the training phase transitions, one must either increase the
 102 number of sampling steps or decrease the learning rate, which in practice means very long training
 103 processes to ensure a proper equilibrium training. Another strategy is to use EBMs as corrections
 104 for straightforward-to-sample flow-based models [38]. This simplifies sampling and learning, but
 105 sacrifices the interpretability of the energy function, which was our goal. An evolving flow model
 106 can be used as a fast sampling moves proposer for the EBM [39] objective. This method requires the
 107 training of two different networks in parallel and may result in the drop of the move acceptancy as
 108 the EBM becomes specialized.

109 For RBMs, a recent method called “stacked tempering” [40] dramatically speeds up sampling by
 110 training smaller RBMs with latent variables from previous models, allowing fast updates to be
 111 proposed using a PT like algorithm. Authors also showed that this algorithm was much faster than
 112 the standard PT. While effective, it is too cumbersome for use in training. Also for RBMs, it has

113 recently been shown that it is possible to train a low-rank RBM that accurately reproduces the
 114 statistics of the data projected along the d first data principal directions through a convex and very
 115 fast optimization process (see [41] and the discussion below). This low-rank model can be seen as a
 116 good approximation to the correct RBM needed to describe the data, and has the nice property that it
 117 can be efficiently sampled via a static Monte Carlo process.

118 In this paper, we will show how to drastically reduce training times by starting the RBM training
 119 process at this low-rank RBM, as this means that the first and strongest dynamic effects associated
 120 with them are directly bypassed. We also show that one can exploit the training trajectory to develop
 121 an effective sampling method, the *parallel trajectory tempering* (PTT) that outperforms the “stacked
 122 tempering” [40] and only requires saving a reduced number of models during the training. This
 123 strategy also allows to obtain reliable estimations for the log-likelihood in well-trained models, much
 124 better than those obtained with the standard Annealing Important Sampling (AIS) techniques [42].
 125 Using both strategies, we show that we are able to train and evaluate methods that accurately represent
 126 the different modes in the dataset, where standard methods lead to mode collapse effects.

127 3 The Restricted Boltzmann Machine

128 The RBM is composed by N_v visible nodes and N_h hidden nodes. In our study, we primarily use
 129 binary variables $\{0, 1\}$ or ± 1 for both layers. The two layers (visible and hidden) interact via a
 130 weight matrix \mathbf{w} , with no direct couplings within a given layer. Variables are also adjusted by visible
 131 and hidden local biases, $\boldsymbol{\theta}$ and $\boldsymbol{\eta}$, respectively. The Gibbs-Boltzmann distribution for this model is
 132 expressed as

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp[-\mathcal{H}(\mathbf{v}, \mathbf{h})] \text{ where } \mathcal{H}(\mathbf{v}, \mathbf{h}) = - \sum_{ia} v_i w_{ia} h_a - \sum_i \theta_i v_i - \sum_a \eta_a h_a, \quad (1)$$

133 where Z is the partition function of the system. As with other models containing hidden variables, the
 134 training objective is to minimize the distance between the empirical distribution of the data, $p_{\mathcal{D}}(\mathbf{v})$,
 135 and the model’s marginal distribution over the visible variables, $p(\mathbf{v}) = \sum_{\mathbf{h}} \exp[-\mathcal{H}(\mathbf{v}, \mathbf{h})] / Z =$
 136 $\exp[-H(\mathbf{v})] / Z$. Minimizing the Kullback-Leibler divergence is equivalent to maximizing the
 137 likelihood of observing the dataset in the model. Thus, the log-likelihood $\mathcal{L} = \langle -H(\mathbf{v}) \rangle_{\mathcal{D}} -$
 138 $\log Z$ can be maximized using the classical stochastic gradient ascent. For a training dataset $\mathcal{D} =$
 139 $\{\mathbf{v}^{(m)}\}_{m=1, \dots, M}$, the log-likelihood gradient is given by

$$\frac{\partial \mathcal{L}}{\partial w_{ia}} = \langle v_i h_a \rangle_{\mathcal{D}} - \langle v_i h_a \rangle_{\text{RBM}}, \quad \frac{\partial \mathcal{L}}{\partial \theta_i} = \langle v_i \rangle_{\mathcal{D}} - \langle v_i \rangle_{\text{RBM}}, \quad \frac{\partial \mathcal{L}}{\partial \eta_a} = \langle h_a \rangle_{\mathcal{D}} - \langle h_a \rangle_{\text{RBM}}, \quad (2)$$

140 where $\langle \cdot \rangle_{\mathcal{D}}$ denotes the average with respect to the entries in the dataset, and $\langle \cdot \rangle_{\text{RBM}}$ with respect
 141 to $p(\mathbf{v}, \mathbf{h})$. Since Z is intractable, the model averages in the gradient are typically estimated using
 142 N_s independent MCMC processes, and observable averages $\langle o(\mathbf{v}, \mathbf{h}) \rangle_{\text{RBM}}$ are replaced by
 143 $\sum_{r=1}^R o(\mathbf{v}^{(r)}, \mathbf{h}^{(r)}) / R$, with $(\mathbf{v}^{(r)}, \mathbf{h}^{(r)})$ being the last configurations reached with each of the
 144 $r = 1, \dots, R$ parallel chains. To obtain reliable estimates, it should be ensured that each of the
 145 Markov chains mix well before each parameter update. However, ensuring equilibrium at each update
 146 is impractical, slow and tedious. The common use of non-convergent MCMC processes is the cause
 147 of most difficulties and weird dynamical behaviors encountered in training RBMs [21].

148 Typical MCMC mixing times in RBMs are very small at the beginning of the training and grow as
 149 it progresses [21], suffering with sharp increases every-time the training trajectory crosses each of
 150 the critical transitions that give birth to new modes [28]. In order to minimize out-of-equilibrium
 151 effects, it is often useful to keep R *permanent (or persistent)* chains, which means that the last
 152 configurations reached with the MCMC process used to estimate the gradient at a given parameter
 153 update t , $\mathbf{P}_t \equiv \{(\mathbf{v}_t^{(r)}, \mathbf{h}_t^{(r)})\}_{r=1}^R$, are used to initialize the chains of the subsequent update $t + 1$.
 154 This algorithm is typically referred as PCD. In this scheme, the process of training can be mimicked
 155 to a slow cooling process, only that instead of varying a single parameter, the temperature, a whole
 156 set of parameters $\boldsymbol{\Theta}_t = (\mathbf{w}_t, \boldsymbol{\theta}_t, \boldsymbol{\eta}_t)$ are updated at every step to $\boldsymbol{\Theta}_{t+1} = \boldsymbol{\Theta}_t + \gamma \nabla \mathcal{L}_t$ with $\nabla \mathcal{L}_t$
 157 being the gradient in Eq. (2) estimated using the configurations in \mathbf{P}_t , and γ being the learning rate.

158 **4 The low-rank RBM pretrained**

159 In Ref. [41], it was shown that it is possible to train exactly (i.e. by direct numerical integration
 160 instead of MCMC sampling) an RBM containing a reduced number of modes in the weight matrix
 161 \mathbf{W} by exploiting a mapping between the RBM and a Restricted Coulomb Machine and solving a
 162 convex optimization problem, see the SI. In other words, it is possible to train a RBM with a coupling
 163 matrix of this simplified form

$$\mathbf{W} = \sum_{\alpha=1}^d w_{\alpha} \bar{\mathbf{u}}_{\alpha} \mathbf{u}_{\alpha}^{\top}, \quad \text{with} \quad (\mathbf{u}_{\alpha}, \bar{\mathbf{u}}_{\alpha}) \in \mathbb{R}^{N_v} \times \mathbb{R}^{N_h}, \quad (3)$$

164 and where the right singular vectors $\{\mathbf{u}_{\alpha}\}_{\alpha=1}^d$ correspond exactly to the first d principal directions
 165 of the data set. Under this assumption, it is possible to write $p(\mathbf{v})$ only as a function of d order
 166 parameters given by the *magnetizations* along each of the u_{α} components, $m_{\alpha}(\mathbf{v}) = \mathbf{u}_{\alpha} \cdot \mathbf{v} / \sqrt{N_v}$,
 167 and in particular,

$$H(\mathbf{v}) = - \sum_a \log \cosh \left(\sqrt{N_v} \bar{u}_a \sum_{\alpha=1}^d w_{\alpha} m_{\alpha} + \eta_a \right) = \mathcal{H}(\mathbf{m}(\mathbf{v})), \quad (4)$$

168 where $\mathbf{m} = (m_1, \dots, m_{\alpha})$. As proposed in [41], the optimal parameters of such a model can
 169 basically be determined by solving a regression problem. We describe this method in details in the
 170 SI. This means that once the model is trained, we obtain a probability $p(\mathbf{m})$ defined on a much
 171 lower dimension than the original $p(\mathbf{v})$. Such a probability can be straightforwardly sampled using
 172 *inverse transform sampling*. Since this method requires a discretization of the \mathbf{m} -space both for
 173 training and generation, we cannot consider intrinsic space dimension $d > 4$ dimensions in practice.
 174 These low-rank RBMs are then trained to reproduce the statistics of the dataset projected in its first d
 175 principal components. Despite their simplicity, the low-rank models are already able to generate an
 176 approximate version of the dataset, as shown in Fig. 1–D–F for the 4 datasets previously presented.

177 In the initial stage of the standard learning process, the model encodes the strongest PCA components
 178 of the data through multiple critical transitions [26, 27, 28]. Pre-training with the low-rank construc-
 179 tion allows us to bypass these transitions and avoid out-of-equilibrium effects caused by critical
 180 slowing down associated to these transitions. Once the main directions are incorporated, training can
 181 efficiently continue with standard algorithms like PCD, as the mixing times of pre-trained machines
 182 tend to be much shorter. In particular, in the PCD-100 training with MNIST01, relaxation times for
 183 the visible variables’ time correlation reach $5 \cdot 10^5$ MCMC steps at the first three transitions, coincid-
 184 ing with the growth of singular values in the model weight matrix W . In contrast, the pre-trained
 185 machine has a much shorter relaxation time of $\sim 10^3$, allowing us to safely restart the PCD process
 186 from a set of equilibrium samples generated by static sampling of the low-rank RBM.

187 Overcoming these transitions has dramatic implications for the quality of the models we can train
 188 and how accurately they reproduce the statistics of the data. In Fig. 2, we show for 3 datasets the
 189 equilibrium samples drawn from 3 RBMs trained with identical number of samples, minibatch size,
 190 $k = 100$ Gibbs steps, and learning rate $\gamma = 0.01$, but different training strategies. In particular,
 191 we consider 2 RBMs trained from scratch with the standard PCD [33] and the recently proposed
 192 Jarzynski reweighing method [23] (see SI for our specific implementation in the RBM), and a final
 193 machine trained with PCD and pre-trained with a low-rank RBM. In all cases, the quality of the
 194 generated samples is significantly better when pre-training is used. For the Mickey dataset, neither
 195 JarRBM nor normal PCD are able to generate convincing data. For the MNIST01 dataset, all 3
 196 methods are able to generate convincing data, but only Pretrain+PCD is able to correctly balance all
 197 modes, as can be seen in Fig. 3, where we compare the histograms of the generated data projected
 198 onto the first 3 PCA directions with those of the dataset and a random selection of the generated
 199 samples. We see that the pre-training+PCD training perfectly balances the different modes (here we
 200 show the first 3 directions, but it goes much further), unlike the other 2 methods, and also generates
 201 more diverse images. We can also compare the log-likelihood of all 3 models and find that the
 202 pre-trained RBM achieves higher values. At this point, it is important to emphasize that in order to
 203 properly quantify the increase in log-likelihood, we need to use the PTT algorithm (see section 5) to
 204 correctly thermalize in these well-trained machines. For comparison, we show our PTT measure in
 205 dark and solid lines, while the standard AIS [42] estimate is shown in light dashed lines.

206 Already from the scatter plots we see that the pre-training has a dramatic effect in obtaining models
 207 where all modes are properly balanced, but also has important effects in the maximum test-likelihood

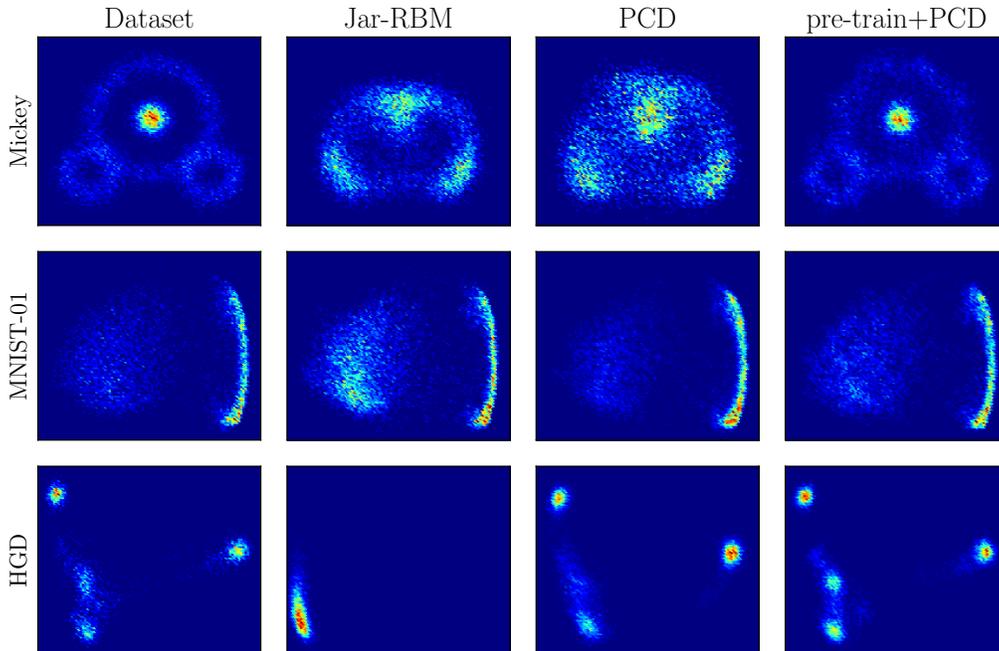


Figure 2: We compare the equilibrium samples generated by RBMs trained on the Mickey, MNIST01, and HGD datasets using three different training schemes: Jarzynski (JarRBM), PCD, and PCD initialized on low-rank RBMs (used to generate the samples in Fig. 1–D–F). To assess the fitting of the modes, we show a density plot of the projections of the data in the first two principal directions of each dataset. We compare these results with the density plot of the original datasets in the first column.

208 we can achieve. In all cases, these equilibrium samples are drawn using the trajectory PT algorithm
 209 that will be explained in the next section, and the log-likelihood obtained using the equilibrium
 210 configurations obtained at different epochs as a result of the trajectory PT flow.

211 5 Standard Gibbs sampling vs. Parallel Trajectory Tempering (PTT)

212 One major challenge with structured datasets is quantifying the model’s quality, since sampling the
 213 equilibrium measure of a well-trained model is often too time-consuming. This affects the reliability of
 214 generated samples and indirect measures as log-likelihood’s estimation through Annealing Importance
 215 Sampling (AIS) [42], making them inaccurate and meaningless.

216 To illustrate this problem, let us consider the MNIST01 and the HGD datasets. MNIST01 dataset is
 217 bimodal and the HGD highly multimodal as shown in their PCA in Figs. 1–A and C. Let us consider
 218 that we want to sample the equilibrium measure of the RBMs trained using low-rank RBM pretraining.
 219 In order to draw new samples from these models, one would typically run MCMC processes from
 220 random initialization and iterate them until convergence. The mixing time is controlled by the
 221 jumping time between clusters. To accurately estimate the relative weight between modes, the
 222 MCMC processes must be ergodic, requiring many back-and-forth jumps. However, as shown in
 223 Figs. 4–A and C for the MNIST01 and HGD datasets, Gibbs sampling dynamics are extremely slow,
 224 rarely producing jumps even after 10^4 MCMC steps. The yellow curves in Figs. 4–B and D show the
 225 mean number of jumps over 100 independent chains as a function of MCMC steps, indicating that a
 226 proper equilibrium generation would require at least $10^6 - 10^7$ MCMC steps.

227 One effective way to accelerate the dynamics is to exploit the training trajectory, where the model
 228 progressively specializes through second-order phase transitions. To achieve this, we save RBMs
 229 trained at various epochs and propose swaps between configurations of similarly trained models. We

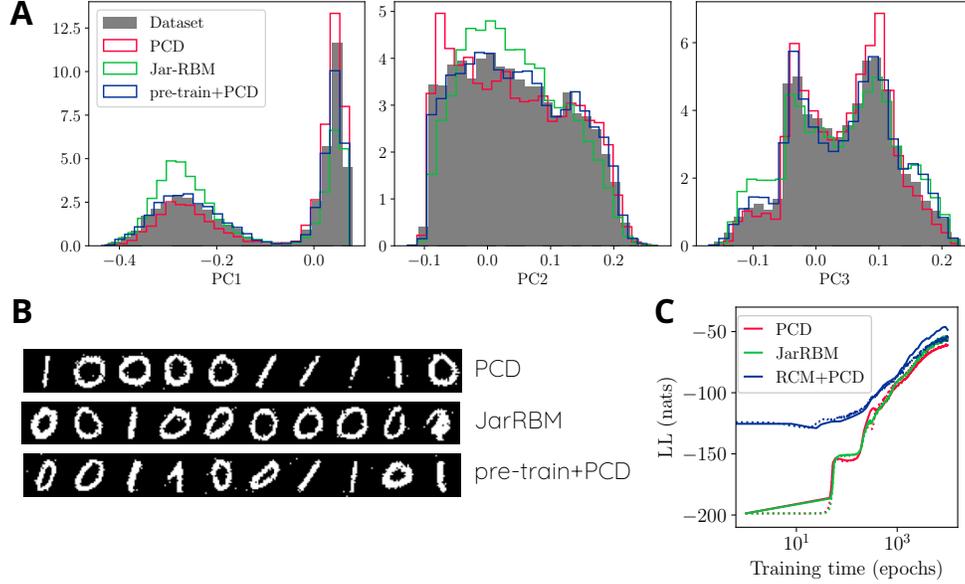


Figure 3: We compare the samples generated by the 3 RBMs (JarRBM, PCD, pretrain+PCD) trained with MNIST01 data. In A, we show the histograms of the generated data projected on the first, second and third principal directions with those of the dataset. We see that only the pretrain+PCD correctly balances the different modes. In B we show 10 images generated by each machine. In C, we compare the log-likelihood of each model’s dataset as a function of training time. The dark and full curves were obtained using the PTT algorithm discussed in section 5, and the lighter and dashed curves using the AIS method [42].

230 call this the *Parallel Trajectory Tempering* (PTT) algorithm. Unlike the standard Parallel Tempering
 231 (PT) algorithm, which attempts swaps configurations between different temperatures, the PTT swaps
 232 between model parameters with different degrees of specialization. This approach is more natural
 233 for this problem because it involves crossing only second-order transitions, unlike the first-order
 234 transitions occurring in temperature annealing. And in fact, we show in Figs. 4–A and C, that this
 235 approach allows us to sharply accelerate the dynamics, as opposed to the standard PT algorithm
 236 (studied in detail for the MNIST dataset in [40]).

237 In the PTT algorithm, the configurations $\mathbf{x} = (\mathbf{v}, \mathbf{h})$ of neighboring machines indexed by t and $t - 1$
 238 are interchanged with the probability

$$p_{\text{acc}}(\mathbf{x}^t \leftrightarrow \mathbf{x}^{t-1}) = \min(1, \exp(\Delta\mathcal{H}^t(\mathbf{x}^t) - \Delta\mathcal{H}^t(\mathbf{x}^{t-1}))).$$

239 This move satisfies detailed balance with our target equilibrium distribution $p(\mathbf{x}) = \exp(-\mathcal{H}(\mathbf{x}))/Z$,
 240 ensuring that the moves lead to the same equilibrium measure. As “nonspecialized” models mix
 241 very quickly, either because the distribution is essentially Gaussian at the initialisation of a standard
 242 training, or because the low-rank RBM can be sampled with a static Monte Carlo process (yielding
 243 independent configurations each time), the trajectory flow significantly accelerates convergence
 244 to equilibrium. The time interval between successive machines is selected in such a way that the
 245 probability of accepting interchanges between neighboring machines remains around 0.3. Pre-trained
 246 machines require a significant fewer number of models to be effective, because most selected models
 247 are positioned at the most prominent phase transitions. We give the number of machines used for
 248 each sampling process in the SI. We also provide there a specific and detailed description of the
 249 algorithm used.

250 In the red curves in Fig. (4)–B and D, we show the number of jumps between clusters as a function
 251 of the number of elementary MCMC steps, which in the PTT scheme refer to 1 Gibbs sampling step
 252 + one swap proposal. For the DNA dataset, we have two measures corresponding to jumps along the
 253 two principal component directions. We observe at 10^4 MCMC steps an increase of the number of
 254 jumps by a factor of 80 for MNIST01 and by a factor of 1350 for the HGD in this machine, although

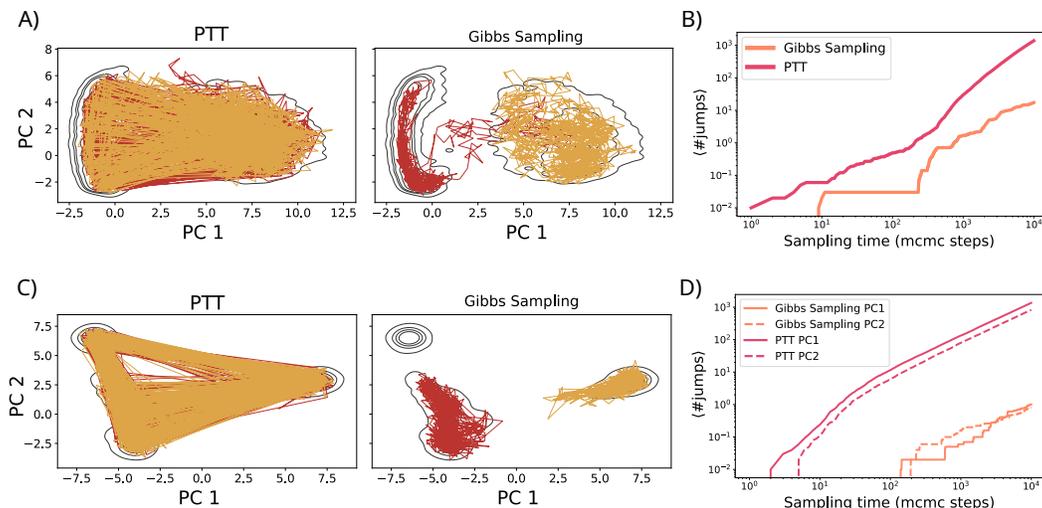


Figure 4: Comparison between PTT and classical Gibbs sampling for the MNIST01 dataset (A and B, respectively) and the human genome dataset (C and D, respectively). In A and C, we show the trajectory of two independent chains (red and orange) projected onto the PCA along the sampling process of the pretraining+PCD model for 10^4 MCMC steps. The black contour represents the density profile of the dataset and the position of the chains is plotted every 10 steps. In B and D we show the average number of jumps from one cluster to another as a function of the MCMC steps performed. The average is calculated over a population of 100 chains. In D, we show the average jump time between clusters along the first (solid line) and second (dashed line) principal components of the data.

255 we achieve higher factors in other machines, as we show in the SI. The sampling of RBMs training
 256 on the MNIST01 dataset was the subject of the study of the “stacked tempering” algorithm in [40].
 257 If we compare the numbers with their work, we see that we achieve a 3-4 times higher speedup factor,
 258 where our model has the advantage that it does not need additional training, but simply uses the stored
 259 machines correctly.

260 Another desirable advantage of our PTT algorithm is that we can easily use it to compute an improved
 261 estimate of the AIS log-likelihood, except that in our case we consider the training trajectory instead
 262 of a cooling process and use the equilibrium samples obtained for each of the models to compute the
 263 model averages. In Figs. 3–C 5–A we compare the log-likelihood estimates obtained with our method
 264 (AIS-PTT) in full and dark lines and in light and dashed lines the AIS estimate (AIS). We see that
 265 both measures coincide for most parts of the training and that they split when the sampling becomes
 266 too long to thermalize along the temperature annealing curve in AIS. This effect is particularly evident
 267 for the JarRBM run in 5–A, where AIS takes a long time to recognize that the model suffers from a
 268 strong mode-collapse effect.

269 6 Overfitting and privacy loss as quality indicators

270 In this section, we examine the quality of the samples generated, regarding overfitting and privacy
 271 criteria which have been defined for genomic data in particular. We look at this on the models trained
 272 with PCD with and without pre-training. We do not include the Jarzynski method here, as this method
 273 fails to obtain a reliable model as clearly shown in the evolution of the Log-likelihood in Fig. 5. We
 274 focus on the human genome dataset, as shown in Fig. 1–C, to evaluate the ability of various state-
 275 of-the-art generative models to generate realistic fake genomes while minimizing privacy concerns
 276 (i.e., reducing overfitting). Recent studies [17, 18] have thoroughly investigated this for a variety of
 277 generative models. Both studies concluded that the RBM was the most accurate method for generating
 278 high-quality and private synthetic genomes. The comparison between models relies primarily on the
 279 Nearest Neighbor Adversarial Accuracy (AA_{TS}) and privacy loss indicators, introduced in Ref. [43],
 280 which quantify the similarity and the level of “privacy” of the data generated by a model w.r.t. the
 281 training set. We have $AA_{TS} = \frac{1}{2}(AA_{True} + AA_{Synth})$ where AA_{True} [resp. AA_{Synth}] are two

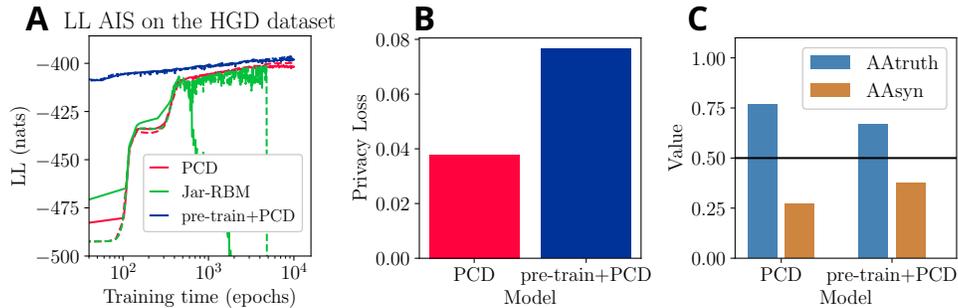


Figure 5: We compare the quality of the RBMs trained with the human genome data (HGD). In A, we show the log-likelihood as a function of the training epochs for the 3 training procedures. Solid lines correspond to AIS-PTT and dashed lines to AIS. The JarRBM falls down because the training breaks eventually. In B and C we compare privacy and overfitting based on the AA_{TS} indicator.

282 quantities in $[0, 1]$ obtained by merging two sets of real and synthetic data of equal size N_s and
 283 measuring respectively the frequency that a real [rep. synthetic] has a synthetic [resp. real] as
 284 nearest neighbor. If the generated samples are statistically indistinguishable from real samples, both
 285 frequencies AA_{True} and AA_{Synth} should converge to 0.5 at large N_s . AA_{TS} can be evaluated both
 286 with train or test samples and the privacy loss indicator is defined as $Privacy\ loss = AA_{TS}^{test} - AA_{TS}^{train}$
 287 and is expected to be strictly positive. Fig. 5 shows the comparison of AA_{TS} and privacy loss values
 288 obtained with our two models, demonstrating that the pre-trained RBM clearly outperforms the
 289 other model, and even achieves better results (AA_{TS} values much closer to 0.5) than those discussed
 290 in [17, 18].

291 7 Conclusions

292 We have shown that the strategy of initiating the training on a pre-trained low-rank RBM is an
 293 extremely effective strategy to obtain high quality models for structured datasets that accurately
 294 represent all the modes in the datasets and with significantly higher log-likelihoods. We have also
 295 shown that the models obtained in that way are: (i) better generative models than those obtained
 296 with standard trainings, both, in the sense that they over-fit less at the same time they are more
 297 indistinguishable from the test samples, (ii) they display faster relaxational dynamics.

298 We have also proposed a new fast sampling method that exploits the progressive learning of features
 299 in the training of RBMs to design an efficient trajectory PT strategy that allows accelerating the
 300 parallel Gibbs sampling dynamics by many orders of magnitude and overcome the performance
 301 of recent efficient sampling methods without adding any extra cost than saving models during the
 302 training.

303 Both strategies for training and sampling are very general, and could be generalized to more complex
 304 EBMs. In this sense, the low-rank RBM model could be used as a more efficient pre-initialisation
 305 in deeper structures, and the trajectory PT algorithm is suitable to be directly used in any EBM no
 306 matter how complex it is.

307 8 Code availability

308 The code and datasets are available at <https://github.com/nboreux/fast-RBM>.

309 **References**

- 310 [1] Paul Smolensky. In *Parallel Distributed Processing: Volume 1 by D. Rumelhart and J. McClelland*, chapter 6: Information Processing in Dynamical Systems: Foundations of Harmony
311 Theory. 194-281. MIT Press, 1986.
- 313 [2] David H Ackley, Geoffrey E Hinton, and Terrence J Sejnowski. A learning algorithm for
314 Boltzmann machines. *Cognitive science*, 9(1):147–169, 1985.
- 315 [3] Geoffrey E Hinton and Terrence J Sejnowski. Optimal perceptual inference. In *Proceedings of*
316 *the IEEE conference on Computer Vision and Pattern Recognition*, volume 448, pages 448–453.
317 Citeseer, 1983.
- 318 [4] H Chau Nguyen, Riccardo Zecchina, and Johannes Berg. Inverse statistical problems: from the
319 inverse Ising problem to data science. *Advances in Physics*, 66(3):197–261, 2017.
- 320 [5] John Hertz, Yasser Roudi, and Joanna Tyrcha. Ising models for inferring network structure from
321 spike data. 2011.
- 322 [6] Simona Cocco, Christoph Feinauer, Matteo Figliuzzi, Rémi Monasson, and Martin Weigt.
323 Inverse statistical physics of protein sequences: a key issues review. *Reports on Progress in*
324 *Physics*, 81(3):032601, 2018.
- 325 [7] Aurélien Decelle, Cyril Furtlehner, Alfonso de Jesús Navas Gómez, and Beatriz Seoane. In-
326 ferring effective couplings with restricted Boltzmann machines. *SciPost Physics*, 16(4):095,
327 2024.
- 328 [8] Aurélien Decelle, Beatriz Seoane, and Lorenzo Rosset. Unsupervised hierarchical clustering
329 using the learning dynamics of restricted Boltzmann machines. *Phys. Rev. E*, 108:014110, Jul
330 2023.
- 331 [9] Jianwen Xie, Yang Lu, Song-Chun Zhu, and Yingnian Wu. A theory of generative convnet. In
332 Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International*
333 *Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*,
334 pages 2635–2644, New York, New York, USA, 20–22 Jun 2016. PMLR.
- 335 [10] Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint*
336 *arXiv:2101.03288*, 2021.
- 337 [11] Giuseppe Carleo and Matthias Troyer. Solving the quantum many-body problem with artificial
338 neural networks. *Science*, 355(6325):602–606, 2017.
- 339 [12] Roger G Melko, Giuseppe Carleo, Juan Carrasquilla, and J Ignacio Cirac. Restricted Boltzmann
340 machines in quantum physics. *Nature Physics*, 15(9):887–892, 2019.
- 341 [13] Koji Hashimoto, Sotaro Sugishita, Akinori Tanaka, and Akio Tomiya. Deep learning and the
342 ads/cft correspondence. *Physical Review D*, 98(4):046019, 2018.
- 343 [14] Koji Hashimoto. Ads/cft correspondence as a deep Boltzmann machine. *Physical Review D*,
344 99(10):106017, 2019.
- 345 [15] Jérôme Tubiana, Simona Cocco, and Rémi Monasson. Learning protein constitutive motifs
346 from sequence data. *Elife*, 8:e39397, 2019.
- 347 [16] Juan Rodriguez-Rivas, Giancarlo Croce, Maureen Muscat, and Martin Weigt. Epistatic models
348 predict mutable sites in SARS-CoV-2 proteins and epitopes. *Proceedings of the National Academy*
349 *of Sciences*, 119(4):e2113118119, 2022.
- 350 [17] Burak Yelmen, Aurélien Decelle, Linda Ongaro, Davide Marnetto, Corentin Tallec, Francesco
351 Montinaro, Cyril Furtlehner, Luca Pagani, and Flora Jay. Creating artificial human genomes
352 using generative neural networks. *PLoS genetics*, 17(2):e1009303, 2021.
- 353 [18] Burak Yelmen, Aurélien Decelle, Leila Lea Boulos, Antoine Szatkownik, Cyril Furtlehner,
354 Guillaume Charpiat, and Flora Jay. Deep convolutional and conditional neural networks for
355 large-scale genomic data generation. *PLoS Computational Biology*, 19(10):e1011584, 2023.

- 356 [19] Erik Nijkamp, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Learning non-convergent
357 non-persistent short-run mcmc toward energy-based model. *Advances in Neural Information*
358 *Processing Systems*, 32, 2019.
- 359 [20] Erik Nijkamp, Mitch Hill, Tian Han, Song-Chun Zhu, and Ying Nian Wu. On the anatomy
360 of mcmc-based maximum likelihood learning of energy-based models. In *Proceedings of the*
361 *AAAI Conference on Artificial Intelligence*, volume 34, pages 5272–5280, 2020.
- 362 [21] Aurélien Decelle, Cyril Furtlehner, and Beatriz Seoane. Equilibrium and non-equilibrium
363 regimes in the learning of restricted boltzmann machines. *Advances in Neural Information*
364 *Processing Systems*, 34:5345–5359, 2021.
- 365 [22] Elisabeth Agoritsas, Giovanni Catania, Aurélien Decelle, and Beatriz Seoane. Explaining the
366 effects of non-convergent sampling in the training of energy-based models. *arXiv preprint*
367 *arXiv:2301.09428*, 2023.
- 368 [23] Alessandra Carbone, Aurélien Decelle, Lorenzo Rosset, and Beatriz Seoane. Fast and
369 functional structured data generators rooted in out-of-equilibrium physics. *arXiv preprint*
370 *arXiv:2307.06797*, 2023.
- 371 [24] Renjie Liao, Simon Kornblith, Mengye Ren, David J Fleet, and Geoffrey Hinton. Gaussian-
372 bernoulli rbms without tears. *arXiv preprint arXiv:2210.10318*, 2022.
- 373 [25] Nicolas Béreux, Aurélien Decelle, Cyril Furtlehner, and Beatriz Seoane. Learning a restricted
374 boltzmann machine using biased monte carlo sampling. *SciPost Physics*, 14(3):032, 2023.
- 375 [26] A. Decelle, G. Fissore, and C. Furtlehner. Spectral dynamics of learning in restricted boltzmann
376 machines. *Europhysics Letters*, 119(6):60001, nov 2017.
- 377 [27] Aurélien Decelle, Giancarlo Fissore, and Cyril Furtlehner. Thermodynamics of restricted
378 boltzmann machines and related learning dynamics. *Journal of Statistical Physics*, 172:1576–
379 1608, 2018.
- 380 [28] A. Anonymous. Cascade of phase transitions in the training of energy-based models. *arXiv*,
381 2024.
- 382 [29] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and Fugie Huang. A tutorial on energy-
383 based learning. *Predicting structured data*, 1(0), 2006.
- 384 [30] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural*
385 *computation*, 14(8):1771–1800, 2002.
- 386 [31] Ruslan Salakhutdinov and Iain Murray. On the quantitative analysis of deep belief networks. In
387 *Proceedings of the 25th international conference on Machine learning*, pages 872–879, 2008.
- 388 [32] Guillaume Desjardins, Aaron Courville, Yoshua Bengio, Pascal Vincent, and Olivier Delalleau.
389 Tempered markov chain monte carlo for training of restricted boltzmann machines. In *Proceed-*
390 *ings of the thirteenth international conference on artificial intelligence and statistics*, pages
391 145–152. JMLR Workshop and Conference Proceedings, 2010.
- 392 [33] Tijmen Tieleman. Training restricted boltzmann machines using approximations to the likeli-
393 hood gradient. In *Proceedings of the 25th international conference on Machine learning*, pages
394 1064–1071, 2008.
- 395 [34] Enzo Marinari and Giorgio Parisi. Simulated tempering: a new monte carlo scheme. *Europhysics*
396 *letters*, 19(6):451, 1992.
- 397 [35] Russ R Salakhutdinov. Learning in markov random fields using tempered transitions. *Advances*
398 *in neural information processing systems*, 22, 2009.
- 399 [36] Oswin Krause, Asja Fischer, and Christian Igel. Population-contrastive-divergence: Does
400 consistency help with rbm training? *Pattern Recognition Letters*, 102:1–7, 2018.

- 401 [37] Davide Carbone, Mengjian Hua, Simon Coste, and Eric Vanden-Eijnden. Efficient training of
402 energy-based models using jarzynski equality. *Advances in Neural Information Processing*
403 *Systems*, 36, 2024.
- 404 [38] E Nijkamp, R Gao, P Sountsov, S Vasudevan, B Pang, S-C Zhu, and YN Wu. Mcmc should
405 mix: learning energy-based model with neural transport latent space mcmc. In *International*
406 *Conference on Learning Representations (ICLR 2022)*., 2022.
- 407 [39] Louis Grenioux, Éric Moulines, and Marylou Gabrié. Balanced training of energy-based models
408 with adaptive flow sampling. In *ICML 2023 Workshop on Structured Probabilistic Inference*
409 *and Generative Modeling*, 2023.
- 410 [40] Clément Roussel, Jorge Fernandez-de Cossio-Diaz, Simona Cocco, and Remi Monasson.
411 Accelerated sampling with stacked restricted boltzmann machines. In *The Twelfth International*
412 *Conference on Learning Representations*, 2023.
- 413 [41] Aurélien Decelle and Cyril Furtlehner. Exact training of restricted boltzmann machines on
414 intrinsically low dimensional data. *Physical Review Letters*, 127(15):158303, 2021.
- 415 [42] Oswin Krause, Asja Fischer, and Christian Igel. Algorithms for estimating the partition function
416 of restricted boltzmann machines. *Artificial Intelligence*, 278:103195, 2020.
- 417 [43] Andrew Yale, Saloni Dash, Ritik Dutta, Isabelle Guyon, Adrien Pavao, and Kristin P Ben-
418 nett. Generation and evaluation of privacy preserving synthetic health data. *Neurocomputing*,
419 416:244–255, 2020.

420 A Details of the pre-training of a low rank RBM

421 A.1 The low-rank RBM and its sampling procedure

422 Our goal is to pre-train an RBM to directly encode the first d principal modes of the dataset in the
423 model’s coupling matrix. This approach avoids the standard procedure of progressively encoding
424 these modes through a series of second-order phase transitions, which negatively impact the quality
425 of gradient estimates during standard training. It also helps prevent critical relaxation slowdown of
426 MCMC dynamics in the presence of many separated clusters.

427 Given a dataset, we want to find a good set of model parameters (\mathbf{w} , $\boldsymbol{\theta}$ and $\boldsymbol{\eta}$) for which the statistics
428 of the generated samples exactly match the statistics of the data projected onto the first d directions
429 of the PCA decomposition of the training set. Let us call each of these $\alpha = 1, \dots, d$ projections
430 $m_\alpha = \mathbf{u}_\alpha \cdot \mathbf{v} / \sqrt{N_v}$ the *magnetizations* along the mode α , where \mathbf{u}_α is the α -th mode of the PCA
431 decomposition of the dataset. A simple way to encode these d -modes is to parameterize the \mathbf{w} -matrix
432 as:

$$433 \mathbf{w} = \sum_{\alpha=1}^d w_\alpha \bar{\mathbf{u}}_\alpha \mathbf{u}_\alpha^\top, \quad \text{with} \quad (\mathbf{u}_\alpha, \bar{\mathbf{u}}_\alpha) \in \mathbb{R}^{N_v} \times \mathbb{R}^{N_h}, \quad (5)$$

433 where \mathbf{u} and $\hat{\mathbf{u}}$ are respectively the right-hand and left-hand singular vectors of \mathbf{w} , the former being
434 directly given by the PCA, while w_α are the singular values of \mathbf{w} . Using this decomposition, the
435 marginal energy on the visible variables, $\mathcal{H}(\mathbf{v}) = \log \sum_{\mathbf{h}} \exp \mathcal{H}(\mathbf{v}, \mathbf{h})$ can be rewritten in terms of
436 these magnetizations $\mathbf{m} \equiv (m_1, \dots, m_d)$

$$437 \mathcal{H}(\mathbf{v}) = - \sum_a \log \cosh \left(\sqrt{N_v} \bar{u}_a \sum_{\alpha=1}^d w_\alpha m_\alpha + \eta_a \right) = \mathcal{H}(\mathbf{m}(\mathbf{v})). \quad (6)$$

437 Now, the goal of our pre-training is not to match the entire statistics of the data set, but only the
438 marginal probability of these magnetizations. In other words, we want to model the marginal
439 distribution

$$439 p_{\text{emp}}(\mathbf{m}) \equiv \sum_{\mathbf{v}} p_{\text{emp}}(\mathbf{v}) \prod_{\alpha=1}^d \delta \left(m_\alpha - \frac{1}{\sqrt{N_v}} \mathbf{u}_\alpha^\top \mathbf{v} \right), \quad (7)$$

440 where δ is the Dirac δ -distribution. In this formulation, the distribution of the model over the
 441 magnetization \mathbf{m} can be easily characterized

$$p(\mathbf{m}) = \frac{1}{Z} \sum_{\mathbf{v}} e^{-\mathcal{H}(\mathbf{v})} \prod_{\alpha=1}^d \delta \left(m_{\alpha} - \frac{1}{\sqrt{N_v}} \mathbf{u}_{\alpha}^T \mathbf{v} \right) \quad (8)$$

$$= \frac{1}{Z} \mathcal{N}(\mathbf{m}) \exp \sum_a \log \cosh \left(\bar{u}_a \sum_{\alpha=1}^d w_{\alpha} m_{\alpha} + \eta_a \right) \quad (9)$$

$$= \frac{1}{Z} e^{-\mathcal{H}(\mathbf{m}) + N_v s(\mathbf{m})} = \frac{1}{Z} e^{-N_v f(\mathbf{m})} \quad (10)$$

442 where $\mathcal{N}(\mathbf{m}) = \sum_{\mathbf{v}} \prod_{\alpha=1}^d \delta \left(m_{\alpha} - \frac{1}{\sqrt{N_v}} \mathbf{u}_{\alpha}^T \mathbf{v} \right)$ is the number of configurations with magnetiza-
 443 tions \mathbf{m} , and thus $S(\mathbf{m}) = \log \mathcal{N}(\mathbf{m}) / N_v$ is the associated entropy. Now, for large N_v the entropic
 444 term can be determined using large deviation theory, and in particular the Gärtner-Ellis theorem:

$$p_{\text{prior}}(\mathbf{m}) = \frac{e^{N_v s(\mathbf{m})}}{2^{N_v}} \approx \exp(-N_v \mathcal{I}(\mathbf{m})), \quad (11)$$

445 with the rate function

$$\mathcal{I}(\mathbf{m}) = \sup_{\boldsymbol{\mu}} [\mathbf{m}^T \boldsymbol{\mu} - \phi(\boldsymbol{\mu})] = \mathbf{m}^T \boldsymbol{\mu}^* - \phi(\boldsymbol{\mu}^*), \quad (12)$$

446 and

$$\phi(\boldsymbol{\mu}) = \lim_{N_v \rightarrow \infty} \frac{1}{N_v} \log \langle e^{N_v \mathbf{m}^T \boldsymbol{\mu}} \rangle = \lim_{N_v \rightarrow \infty} \frac{1}{N_v} \log \frac{1}{2^{N_v}} \sum_{\mathbf{v}} e^{\sqrt{N_v} \sum_{\alpha=1}^d \mu_{\alpha} \sum_i u_{\alpha,i} v_i} \quad (13)$$

$$= \lim_{N_v \rightarrow \infty} \frac{1}{N_v} \sum_{i=1}^{N_v} \log \cosh \left(\sqrt{N_v} \sum_{\alpha=1}^d \mu_{\alpha} u_{\alpha,i} \right). \quad (14)$$

447 Then, given a magnetization \mathbf{m} , we can compute the minimizer $\boldsymbol{\mu}^*(\mathbf{m})$ of $\phi(\boldsymbol{\mu}) - \mathbf{m}^T \boldsymbol{\mu}$ which is
 448 convex, using e.g. Newton method which converge really fast since we are in small dimension. Note
 449 that in practice we will obviously use finite estimates of ϕ , assuming N_v is large enough. As a result
 450 we get $\boldsymbol{\mu}^*(\mathbf{m})$ satisfying implicit equations given by the constraints given at given N_v :

$$m_{\alpha} = \frac{1}{\sqrt{N_v}} \sum_{i=1}^{N_v} u_i^{\alpha} \tanh \left(\sqrt{N_v} \sum_{\beta=1}^d u_i^{\beta} \mu_{\beta}^* \right). \quad (15)$$

451 It is then straightforward to check that spins distributed as

$$p_{\text{prior}}(\mathbf{v} | \mathbf{m}) \propto e^{N_v \boldsymbol{\mu}^{*T} \mathbf{m}(\mathbf{v})} \quad (16)$$

452 fulfill well the requirement, as $\langle \mathbf{u}_{\alpha}^T \mathbf{v} / \sqrt{N_v} \rangle_{p_{\text{prior}}} = m_{\alpha}$. In other words, we can generate samples
 453 having mean magnetization m_{α} just by choosing v_i as

$$p_{\text{prior}}(v_i = 1 | \mathbf{m}) = \text{sigmoid} \left(2 \sqrt{N_v} \sum_{\alpha=1}^d u_{\alpha,i} \mu_{\alpha}^*(\mathbf{m}) \right) \quad (17)$$

454 The training can therefore be done directly in the subspace of dimension d . In Ref. [41], it has been
 455 shown that such RBM can be trained by mean of the Restricted Coulomb Machine, where the gradient
 456 is actually convex in the parameter's space. It is then possible to do a mapping from the RCM to
 457 the RBM to recover the RBM's parameters. In brief, the training of the low-dimensional RBM is
 458 performed by the RCM, and then the parameters are obtained via a direct relation between the RCM
 459 and the RBM's parameters. The detail of the definition and of the training of the RCM is detailed in
 460 the appendix A.2.

461 **A.2 The Restricted Coulomb Machine**

462 As introduced in [41], it is possible to exactly train a surrogate model for the RBM, called the
 463 Restricted Coulomb Machine (RCM), on a low dimensional dataset without explicitly sampling the
 464 machine allowing to learn even heavily clustered datasets. We will briefly outline the main steps to
 465 train the RCM. A more detailed explanation can be found in Appendix A.2.

466 The RCM is an approximation of the marginal distribution of the RBM with $\{-1, 1\}$ binary variables:

$$\mathcal{H}(\mathbf{v}) = - \sum_i v_i \theta_i - \sum_a \log \cosh \left(\sum_i w_{ia} v_i + \eta_a \right). \quad (18)$$

467 We then project both the parameters and variables of the RBM on the first d principal components of
 468 the dataset:

$$m_\alpha := \frac{1}{\sqrt{N_v}} \sum_{i=1}^{N_v} s_i u_{i\alpha}, \quad w_{\alpha a} := \sum_{i=1}^{N_v} w_{ia} u_{i\alpha}, \quad \theta_\alpha := \frac{1}{\sqrt{N_v}} \sum_{i=1}^{N_v} \theta_i u_{i\alpha} \quad (19)$$

469 with $\alpha \in \{1, \dots, d\}$ and \mathbf{v} the projection matrix of the PCA. The projected distribution of the model
 470 is then given by

$$p_{\text{RBM}}(\mathbf{m}) = \frac{\exp \left(N_v \left[\mathcal{S}(\mathbf{m}) + \sum_{\alpha=1}^d \theta_\alpha m_\alpha + \frac{1}{N_v} \sum_{a=1}^{N_h} \log \cosh \left(\sqrt{N_v} \sum_{\alpha=1}^d m_\alpha w_{\alpha a} + \eta_a \right) \right] \right)}{Z} \quad (20)$$

471 where we ignore the fluctuations related to the transverse directions and $\mathcal{S}[\mathbf{m}]$ accounts for the
 472 non-uniform prior on \mathbf{m} due to the projection of the uniform prior on \mathbf{s} for the way to compute it.

473 The RCM is then built by approximating

$$\log \cosh(x) \simeq |x| - \log 2, \quad (21)$$

474 which is valid for x large enough. The probability of the RCM is thus given by:

$$p_{\text{RCM}}(\mathbf{m}) = \frac{\exp \left(N_v \left[\mathcal{S}(\mathbf{m}) + \sum_{\alpha=1}^d \theta_\alpha m_\alpha + \sum_{a=1}^{N_h} q_a \left| \sum_{\alpha=1}^d n_\alpha m_\alpha + z_a \right| \right] \right)}{Z} \quad (22)$$

475 where

$$q_a = \sqrt{N_v \sum_{\alpha=1}^d w_{\alpha a}^2}, \quad n_a = \frac{w_{\alpha a}}{\sqrt{\sum_{\alpha=1}^d w_{\alpha a}^2}}, \quad z_a = \frac{\eta_a}{\sqrt{N_v \sum_{\alpha=1}^d w_{\alpha a}^2}}. \quad (23)$$

476 This can be easily inverted as

$$w_{\alpha a} = \frac{1}{\sqrt{N_v}} q_a n_a \quad \text{and} \quad \eta_a = q_a z_a,$$

477 in order to obtain the RBM from the RCM. The model is then trained through log-likelihood
 478 maximization over its parameters. However, this objective is non-convex if all the parameters are
 479 trained through gradient ascent. To relax the problem, since we're in low dimension, we can define a
 480 family of hyperplanes (\mathbf{n}, z) covering the space and let the model only learn the weights of each to
 481 the hyperplane. We can then discard the ones with a weight low enough for the approximation (21) to
 482 be bad.

483 The gradients are given by

$$\frac{\partial J(\Theta)}{\partial q_a} = \mathbb{E}_{\mathbf{m} \sim p_{\mathcal{D}}(\mathbf{m})} \left[|\mathbf{n}_a^T \mathbf{m} + z_a| \right] - \mathbb{E}_{\mathbf{m} \sim p_{\text{RCM}}(\mathbf{m})} \left[|\mathbf{n}_a^T \mathbf{m} + z_a| \right], \quad (24)$$

484

$$\frac{\partial J(\Theta)}{\partial \theta_\alpha} = \mathbb{E}_{\mathbf{m} \sim p_{\mathcal{D}}(\mathbf{m})} [m_\alpha] - \mathbb{E}_{\mathbf{m} \sim p_{\text{RCM}}(\mathbf{m})} [m_\alpha]. \quad (25)$$

485 The positive term is straightforward to compute. For the negative term, we rely on a discretization of
 486 the longitudinal space to estimate the probability density of the model and compute the averages.

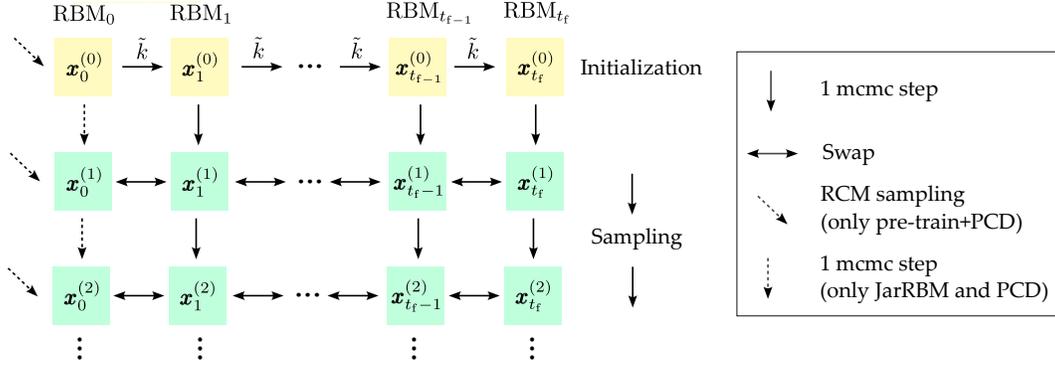


Figure 6: Scheme of PTT. We Initialize the chains of the models by starting from a configuration $\mathbf{x}_0^{(0)}$ and passing it through the machines along the training trajectory, each time performing \tilde{k} mcmc steps. For pre-train+PCD, $\mathbf{x}_0^{(0)}$ is a sampling from the RCM, otherwise it is a uniform random initialization. The sampling consists of alternating one mcmc step for each model with a swap attempt between adjacent machines. For pre-train+PCD, at each step we sample a new independent configuration for RBM_0 using the RCM.

Model	Dataset	# of machines	Alg. # of steps	acc. factor @ 10^4 steps
pre-train+PCD	MNIST01	6 (+1)	10000	80
JarJar	MNIST01	28	10000	50
PCD	MNIST01	13	10000	30
pre-train+PCD	Human Genome	6 (+1)	10000	1350
PCD	Human Genome	13	10000	7100

Table 1: Performance comparison of different models on various datasets for the sampling using PTT versus Gibbs sampling for 10^4 mcmc steps. The acceleration factor is defined as the ratio of the average number of jumps obtained until 10^4 steps between PTT and Gibbs sampling. For pre-train+PCD, the RCM machine has not to be counted among the list of models (hence the +1) because it is very fast to sample from.

487 B Sampling via Parallel Tempering using the learning trajectory

488 Assuming we have successfully trained a robust equilibrium model, there remains the challenge of
 489 efficiently generating equilibrium configurations from this model. Although models trained at equi-
 490 librium exhibit faster and more ergodic dynamics compared to poorly trained models, the sampling
 491 time can still be excessively long when navigating a highly rugged data landscape. Consequently,
 492 we devised a novel method for sampling equilibrium configurations that draws inspiration from
 493 the well-established parallel tempering approach. In this traditional method, multiple simulations
 494 are conducted in parallel at various temperatures, and configurations are exchanged among them
 495 using the Metropolis rule. Unlike this conventional technique, our method involves simultaneously
 496 simulating different models that are selected from various points along the training trajectory. This
 497 approach is motivated by the perspective that learning represents an annealing process for the model,
 498 encountering second-order type phase transitions during training. In contrast, annealing related to
 499 temperature changes involves first-order phase transitions, making traditional parallel tempering less
 500 effective for sampling from clustered multimodal distributions.

501 A sketch of the Parallel Trajectory Tempering (PTT) is represented in fig. 6. Specifically, we save
 502 t_f models at checkpoints $t = 1, \dots, t_f$ along the training trajectory. We denote the Hamiltonian of
 503 the model at checkpoint t as \mathcal{H}_t , and refer to the Hamiltonian of the RCM model as \mathcal{H}_0 . We define
 504 $\text{GibbsSampling}(\mathcal{H}, \mathbf{x}, k)$ as the operation of performing k Gibbs sampling updates using the model
 505 \mathcal{H} starting from the state \mathbf{x} . In all our sampling simulations we used $k = 1$.

506 The first step is to initialize the models’ configurations efficiently. This involves sampling N
507 chains from the RCM model, $\mathbf{x}_0^{(0)} \sim \text{RCMSampling}(\mathcal{H}^0)$, and then passing the chains through
508 all the models from $t = 1$ to $t = t_f$, performing k Gibbs steps at each stage: $\mathbf{x}_t^{(0)} \sim$
509 $\text{GibbsSampling}(\mathcal{H}_t, \mathbf{x}_{t-1}^{(0)}, k)$.

510 The sampling process proceeds in steps where we update the configuration of each model except \mathcal{H}_0
511 with k Gibbs steps, and sample a completely new configuration for the RCM model \mathcal{H}_0 . Following
512 this update step, we propose swapping chains between adjacent models with an acceptance probability
513 given by:

$$p_{\text{acc}}(\mathbf{x}_t \leftrightarrow \mathbf{x}_{t-1}) = \min(1, \exp(\Delta\mathcal{H}_t(\mathbf{x}_t) - \Delta\mathcal{H}_t(\mathbf{x}_{t-1}))), \quad (26)$$

514 where $\Delta\mathcal{H}_t(\mathbf{x}) = \mathcal{H}_t(\mathbf{x}) - \mathcal{H}_{t-1}(\mathbf{x})$.

515 We continue alternating between the update step and the swap step until a total of N_{mcmc} steps is
516 reached. The sampling procedure is illustrated in the following pseudo-code:

517 **Input:** Set of models $\{\mathcal{H}_t\}$, $t = 0, \dots, t_f$, Number of Gibbs steps k , Number of MCMC steps
518 N_{mcmc}

519 **Output:** Configurations \mathbf{x}_t for $t = 1, \dots, t_f$

520 **Initialize:** Sample N chains from the RCM model $\mathbf{x}_0^{(0)} \sim \text{RCMSampling}(\mathcal{H}_0)$

521 **for** $t = 1$ to t_f **do**

522 $\mathbf{x}_t^{(0)} \sim \text{GibbsSampling}(\mathcal{H}_t, \mathbf{x}_{t-1}^{(0)}, \tilde{k})$

523 **end for**

524 **for** $n = 1$ to $\lfloor N_{\text{mcmc}}/k \rfloor$ **do**

525 **for** $t = 1$ to t_f **do**

526 $\mathbf{x}_t^{(n)} \sim \text{GibbsSampling}(\mathcal{H}_t, \mathbf{x}_t^{(n-1)}, k)$

527 **end for**

528 Resample $\mathbf{x}_0^{(n)} \sim \text{RCMSampling}(\mathcal{H}_0)$

529 **for** $t = 1$ to t_f **do**

530 Compute acceptance probability

$$p_{\text{acc}}(\mathbf{x}_t^{(n)} \leftrightarrow \mathbf{x}_{t-1}^{(n)}) = \min\left(1, \exp\left(\Delta\mathcal{H}_t(\mathbf{x}_t^{(n)}) - \Delta\mathcal{H}_t(\mathbf{x}_{t-1}^{(n)})\right)\right)$$

531 Swap $\mathbf{x}_t^{(n)}$ and $\mathbf{x}_{t-1}^{(n)}$ with probability $p_{\text{acc}}(\mathbf{x}_t^{(n)} \leftrightarrow \mathbf{x}_{t-1}^{(n)})$

532 **end for**

533 **end for**

534 A comparison of performances between PTT and standard Gibbs sampling is reported in Tab. 1.

535 C Training details

536 We describe in Tables 2 and 3 the datasets and hyperparameters used during training. The test set
was used to evaluate the metrics. All experiments were run on a RTX 4090 with an AMD Ryzen 9

Table 2: Details of the datasets used during training.

Name	#Samples	#Dimensions	Train size	Test size
Human Genome Dataset (HGD)	4500	805	60%	40%
MNIST-01	10 610	784	60%	40%
Mickey	16 000	1000	60%	40%

537
538 5950X.

539 D Training of the RBM using the Jarzynski equation

540 In this section, we describe a procedure similar to the one introduced in [37] for training the RBM
541 by leveraging the Jarzynski equation. In one of its formulations, the Jarzynski equation states that
542 we can relate the ensemble average of an observable \mathcal{O} with the average obtained through many

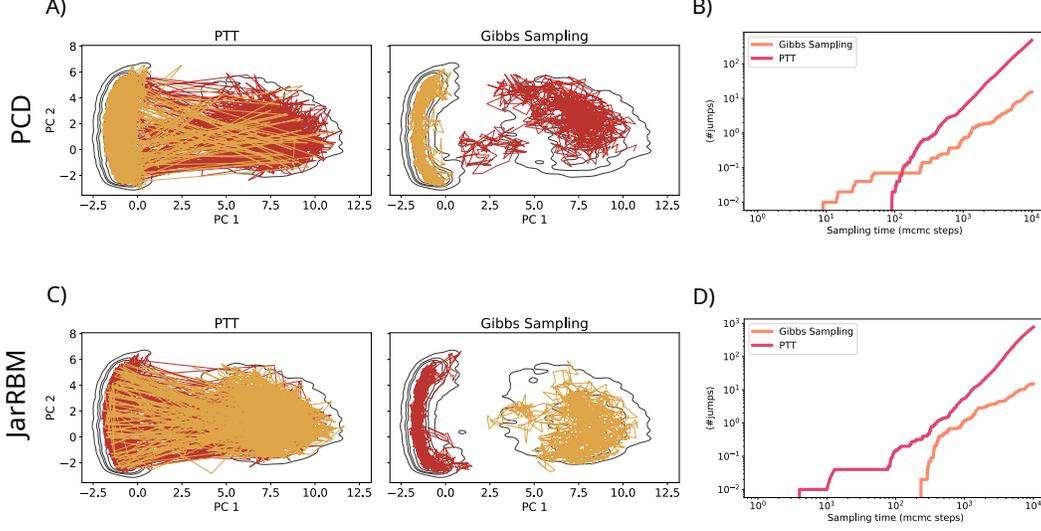


Figure 7: Comparison between PTT and standard Gibbs Sampling for RBMs trained using PCD (A and B) and JarRBM (C and D) on the MNIST01 dataset. A and C show the sampling trajectory of two chains recorded every 10 steps for a total of 10^4 mcmc steps. B and D show the average number of jumps of a population of 100 chains as a function of the sampling time.

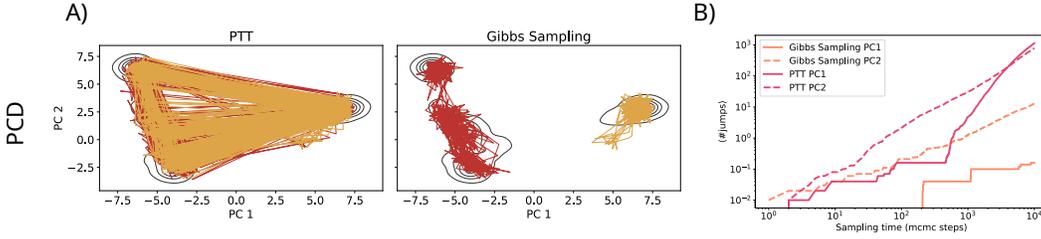


Figure 8: Comparison between PTT and standard Gibbs Sampling for RBMs trained using PCD on the Human Genome dataset. The sampling has been performed under the same conditions of fig. 7.

543 repetitions of an out-of-equilibrium dynamical process. If we consider the training trajectory of an
 544 RBM, $p_0 \rightarrow p_1 \rightarrow \dots \rightarrow p_{t-1} \rightarrow p_t$, we can write

$$\langle \mathcal{O} \rangle_t = \frac{\langle \mathcal{O} e^{-W_t} \rangle_{\text{traj}}}{\langle e^{-W_t} \rangle_{\text{traj}}}, \quad (27)$$

545 where the average on the lhs is done over the last model p_t , the averages on the rhs are taken
 546 across many different trajectory realizations and W_t is a trajectory-dependent importance factor. By
 547 all practical means, under the assumption of having quasi-adiabatic parameters updates, namely
 548 $p(\Theta_{t-1} \rightarrow \Theta_t) = p(\Theta_t \rightarrow \Theta_{t-1})$, this means that we can assign to each Markov chain of the
 549 simulation $\mathbf{x}^{(r)}$, $r = 1, \dots, R$, an importance weight given by:

$$W_t^{(r)} = \sum_{\tau=1}^t [\mathcal{H}_\tau(\mathbf{x}_{\tau-1}^{(r)}) - \mathcal{H}_{\tau-1}(\mathbf{x}_{\tau-1}^{(r)})] \quad (28)$$

550 and then compute the gradient of the log-likelihood by means of a weighted average over the chains:

$$\langle \mathcal{O} \rangle_t \simeq \frac{\sum_{r=1}^R \mathcal{O}(\mathbf{x}^{(r)}) e^{-W_t^{(r)}}}{\sum_{r=1}^R e^{-W_t^{(r)}}}. \quad (29)$$

Table 3: Hyperparameters used for the training of RBMs.

Name	Batch size	#Chains	#Epochs	Learning rate	#MCMC steps	#Hidden nodes
HGD						
PCD	2000	2000	10 000	0.01	100	185
Jar-RBM	2000	10 000	10 000	0.01	100	185
Pre-train+PCD	2000	2000	10 000	0.01	100	185
MNIST-01						
PCD	2000	2000	10 000	0.01	100	200
Jar-RBM	2000	10 000	10 000	0.01	100	200
Pre-train+PCD	2000	2000	10 000	0.01	100	200
Mickey						
PCD	2000	2000	10 000	0.01	100	100
Jar-RBM	2000	10 000	10 000	0.01	100	100
Pre-train+PCD	2000	2000	10 000	0.01	100	100

551 Notice that, since Eq. (27) is an exact result, the importance weights should, in principle, eliminate
552 the bias brought by the non-convergent chains used for approximating the log-likelihood gradient in
553 the classical PCD scheme. However, after many updates of the importance weights, one finds that
554 only a few chains carry almost all the importance mass. In other words, the vast majority of the chains
555 we are simulating are statistically irrelevant, and we expect to get large fluctuations in the estimate of
556 the gradient because of the small effective number of chains contributing to the statistical average. A
557 good observable for monitoring this effect is the Effective Sample Size (ESS), defined as [37]

$$\text{ESS} = \frac{\left(R^{-1} \sum_{r=1}^R e^{-W^{(r)}} \right)^2}{R^{-1} \sum_{r=1}^R e^{-2W^{(r)}}} \in [0, 1], \quad (30)$$

558 which measures the relative dispersion of the weights distribution. A way of circumventing the weight
559 concentration on a few chains, then, is to resample the chain population according to the importance
560 weights every time the ESS drops below a certain threshold, for instance 0.5. After this resampling,
561 all the chain weights have to be set to 1 ($W^{(r)} = 0 \forall r = 1, \dots, R$).

562 **NeurIPS Paper Checklist**

563 The checklist is designed to encourage best practices for responsible machine learning research,
564 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
565 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
566 follow the references and precede the (optional) supplemental material. The checklist does NOT
567 count towards the page limit.

568 Please read the checklist guidelines carefully for information on how to answer these questions. For
569 each question in the checklist:

- 570 • You should answer [Yes] , [No] , or [NA] .
- 571 • [NA] means either that the question is Not Applicable for that particular paper or the
572 relevant information is Not Available.
- 573 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

574 **The checklist answers are an integral part of your paper submission.** They are visible to the
575 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
576 (after eventual revisions) with the final version of your paper, and its final version will be published
577 with the paper.

578 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
579 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
580 proper justification is given (e.g., "error bars are not reported because it would be too computationally
581 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
582 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
583 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
584 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
585 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
586 please point to the section(s) where related material for the question can be found.

587 **1. Claims**

588 Question: Do the main claims made in the abstract and introduction accurately reflect the
589 paper's contributions and scope?

590 Answer: [Yes]

591 Justification:

592 Guidelines:

- 593 • The answer NA means that the abstract and introduction do not include the claims
594 made in the paper.
- 595 • The abstract and/or introduction should clearly state the claims made, including the
596 contributions made in the paper and important assumptions and limitations. A No or
597 NA answer to this question will not be perceived well by the reviewers.
- 598 • The claims made should match theoretical and experimental results, and reflect how
599 much the results can be expected to generalize to other settings.
- 600 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
601 are not attained by the paper.

602 **2. Limitations**

603 Question: Does the paper discuss the limitations of the work performed by the authors?

604 Answer: [Yes]

605 Justification:

606 Guidelines:

- 607 • The answer NA means that the paper has no limitation while the answer No means that
608 the paper has limitations, but those are not discussed in the paper.
- 609 • The authors are encouraged to create a separate "Limitations" section in their paper.

- 610 • The paper should point out any strong assumptions and how robust the results are to
611 violations of these assumptions (e.g., independence assumptions, noiseless settings,
612 model well-specification, asymptotic approximations only holding locally). The authors
613 should reflect on how these assumptions might be violated in practice and what the
614 implications would be.
- 615 • The authors should reflect on the scope of the claims made, e.g., if the approach was
616 only tested on a few datasets or with a few runs. In general, empirical results often
617 depend on implicit assumptions, which should be articulated.
- 618 • The authors should reflect on the factors that influence the performance of the approach.
619 For example, a facial recognition algorithm may perform poorly when image resolution
620 is low or images are taken in low lighting. Or a speech-to-text system might not be
621 used reliably to provide closed captions for online lectures because it fails to handle
622 technical jargon.
- 623 • The authors should discuss the computational efficiency of the proposed algorithms
624 and how they scale with dataset size.
- 625 • If applicable, the authors should discuss possible limitations of their approach to
626 address problems of privacy and fairness.
- 627 • While the authors might fear that complete honesty about limitations might be used by
628 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
629 limitations that aren't acknowledged in the paper. The authors should use their best
630 judgment and recognize that individual actions in favor of transparency play an impor-
631 tant role in developing norms that preserve the integrity of the community. Reviewers
632 will be specifically instructed to not penalize honesty concerning limitations.

633 3. Theory Assumptions and Proofs

634 Question: For each theoretical result, does the paper provide the full set of assumptions and
635 a complete (and correct) proof?

636 Answer: [NA]

637 Justification:

638 Guidelines:

- 639 • The answer NA means that the paper does not include theoretical results.
- 640 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
641 referenced.
- 642 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 643 • The proofs can either appear in the main paper or the supplemental material, but if
644 they appear in the supplemental material, the authors are encouraged to provide a short
645 proof sketch to provide intuition.
- 646 • Inversely, any informal proof provided in the core of the paper should be complemented
647 by formal proofs provided in appendix or supplemental material.
- 648 • Theorems and Lemmas that the proof relies upon should be properly referenced.

649 4. Experimental Result Reproducibility

650 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
651 perimental results of the paper to the extent that it affects the main claims and/or conclusions
652 of the paper (regardless of whether the code and data are provided or not)?

653 Answer: [Yes]

654 Justification:

655 Guidelines:

- 656 • The answer NA means that the paper does not include experiments.
- 657 • If the paper includes experiments, a No answer to this question will not be perceived
658 well by the reviewers: Making the paper reproducible is important, regardless of
659 whether the code and data are provided or not.
- 660 • If the contribution is a dataset and/or model, the authors should describe the steps taken
661 to make their results reproducible or verifiable.

- 662
- 663
- 664
- 665
- 666
- 667
- 668
- 669
- 670
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
 - While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

687 5. Open access to data and code

688 Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

691 Answer: [Yes]

692 Justification:

693 Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

713 6. Experimental Setting/Details

714 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

715

716

717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766

Answer: [Yes]

Justification:

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification:

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification:

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

767 Question: Does the research conducted in the paper conform, in every respect, with the
768 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

769 Answer: [Yes]

770 Justification:

771 Guidelines:

- 772 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 773 • If the authors answer No, they should explain the special circumstances that require a
774 deviation from the Code of Ethics.
- 775 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
776 eration due to laws or regulations in their jurisdiction).

777 10. Broader Impacts

778 Question: Does the paper discuss both potential positive societal impacts and negative
779 societal impacts of the work performed?

780 Answer: [NA]

781 Justification:

782 Guidelines:

- 783 • The answer NA means that there is no societal impact of the work performed.
- 784 • If the authors answer NA or No, they should explain why their work has no societal
785 impact or why the paper does not address societal impact.
- 786 • Examples of negative societal impacts include potential malicious or unintended uses
787 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
788 (e.g., deployment of technologies that could make decisions that unfairly impact specific
789 groups), privacy considerations, and security considerations.
- 790 • The conference expects that many papers will be foundational research and not tied
791 to particular applications, let alone deployments. However, if there is a direct path to
792 any negative applications, the authors should point it out. For example, it is legitimate
793 to point out that an improvement in the quality of generative models could be used to
794 generate deepfakes for disinformation. On the other hand, it is not needed to point out
795 that a generic algorithm for optimizing neural networks could enable people to train
796 models that generate Deepfakes faster.
- 797 • The authors should consider possible harms that could arise when the technology is
798 being used as intended and functioning correctly, harms that could arise when the
799 technology is being used as intended but gives incorrect results, and harms following
800 from (intentional or unintentional) misuse of the technology.
- 801 • If there are negative societal impacts, the authors could also discuss possible mitigation
802 strategies (e.g., gated release of models, providing defenses in addition to attacks,
803 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
804 feedback over time, improving the efficiency and accessibility of ML).

805 11. Safeguards

806 Question: Does the paper describe safeguards that have been put in place for responsible
807 release of data or models that have a high risk for misuse (e.g., pretrained language models,
808 image generators, or scraped datasets)?

809 Answer: [NA]

810 Justification:

811 Guidelines:

- 812 • The answer NA means that the paper poses no such risks.
- 813 • Released models that have a high risk for misuse or dual-use should be released with
814 necessary safeguards to allow for controlled use of the model, for example by requiring
815 that users adhere to usage guidelines or restrictions to access the model or implementing
816 safety filters.
- 817 • Datasets that have been scraped from the Internet could pose safety risks. The authors
818 should describe how they avoided releasing unsafe images.

819 • We recognize that providing effective safeguards is challenging, and many papers do
820 not require this, but we encourage authors to take this into account and make a best
821 faith effort.

822 12. Licenses for existing assets

823 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
824 the paper, properly credited and are the license and terms of use explicitly mentioned and
825 properly respected?

826 Answer: [NA]

827 Justification:

828 Guidelines:

- 829 • The answer NA means that the paper does not use existing assets.
- 830 • The authors should cite the original paper that produced the code package or dataset.
- 831 • The authors should state which version of the asset is used and, if possible, include a
832 URL.
- 833 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 834 • For scraped data from a particular source (e.g., website), the copyright and terms of
835 service of that source should be provided.
- 836 • If assets are released, the license, copyright information, and terms of use in the
837 package should be provided. For popular datasets, `paperswithcode.com/datasets`
838 has curated licenses for some datasets. Their licensing guide can help determine the
839 license of a dataset.
- 840 • For existing datasets that are re-packaged, both the original license and the license of
841 the derived asset (if it has changed) should be provided.
- 842 • If this information is not available online, the authors are encouraged to reach out to
843 the asset's creators.

844 13. New Assets

845 Question: Are new assets introduced in the paper well documented and is the documentation
846 provided alongside the assets?

847 Answer: [NA]

848 Justification:

849 Guidelines:

- 850 • The answer NA means that the paper does not release new assets.
- 851 • Researchers should communicate the details of the dataset/code/model as part of their
852 submissions via structured templates. This includes details about training, license,
853 limitations, etc.
- 854 • The paper should discuss whether and how consent was obtained from people whose
855 asset is used.
- 856 • At submission time, remember to anonymize your assets (if applicable). You can either
857 create an anonymized URL or include an anonymized zip file.

858 14. Crowdsourcing and Research with Human Subjects

859 Question: For crowdsourcing experiments and research with human subjects, does the paper
860 include the full text of instructions given to participants and screenshots, if applicable, as
861 well as details about compensation (if any)?

862 Answer: [No]

863 Justification:

864 Guidelines:

- 865 • The answer NA means that the paper does not involve crowdsourcing nor research with
866 human subjects.
- 867 • Including this information in the supplemental material is fine, but if the main contribu-
868 tion of the paper involves human subjects, then as much detail as possible should be
869 included in the main paper.

870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891

- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.