

Decelle, Furtlehner, Seoane (NeurIPS 2021)

Agoritsas, Catania, Decelle, Seoane (ICML 2023)

On the effect of MCMC in the training of EBMs

Beatriz Seoane

LISN Paris-Saclay University

EBMs... Yes! ... But!

EBMs

Pros: Appealing for modeling and interpretability applications

Cons: Very hard to train :

- The quality of the training is hard to control
- Sampling are unstable

Generating samples : **what one expects**

Empirical

$$p_{\text{data}}(\mathbf{x}) \sim \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}}$$

Model

Dominated minimum
free-energy
configurations

$$\{\mathbf{x}\}_{\text{eq}, \boldsymbol{\theta}} \sim \mathcal{D}$$



Markov Chain Monte Carlo
Langevin dynamics

Generate new samples

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\text{data}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\boldsymbol{\theta}}} \quad \forall \theta_i$$

Generating samples : what one expects

MCMC sampling steps



(some time to equilibrate/converge)

10⁴ steps



Trained model $E_{\theta}(\mathbf{x})$

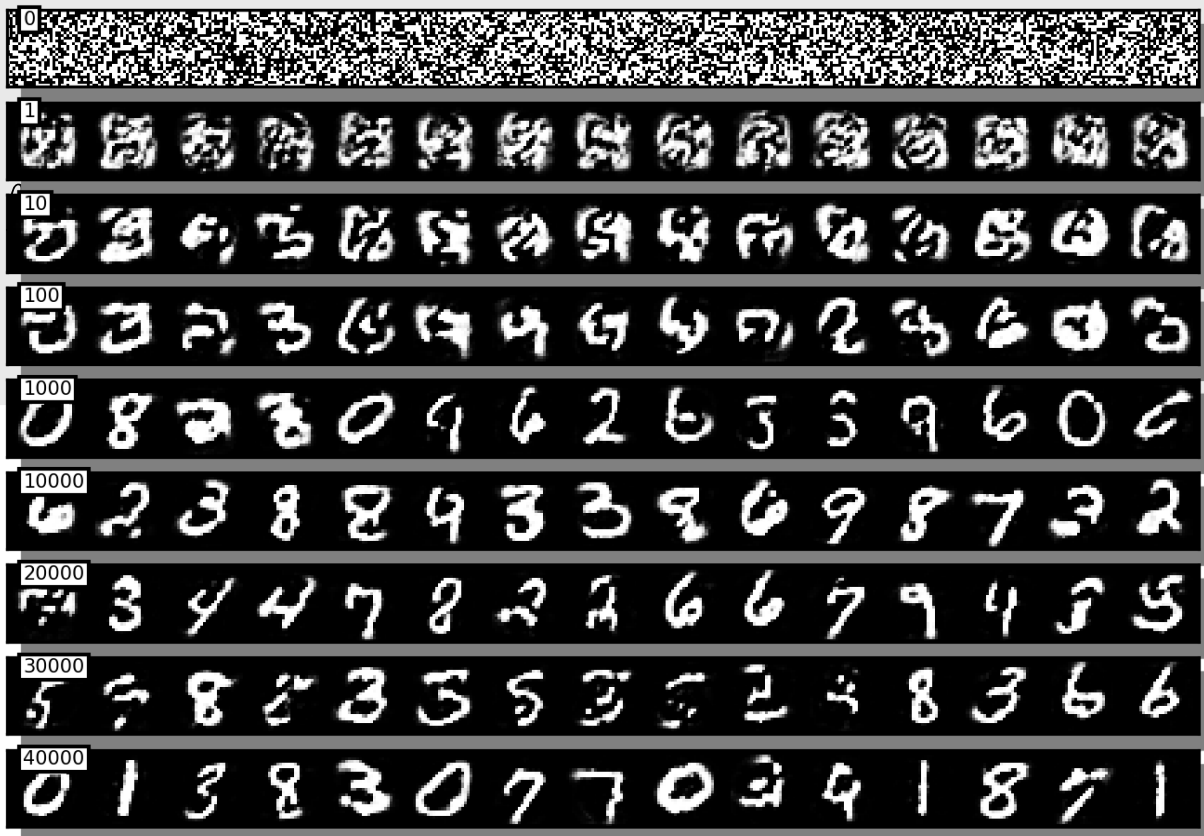


Generate samples with $p_{\theta}(\mathbf{x})$

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\text{data}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\theta}} \quad \forall \theta_i$$

Generating samples : **what one expects...**

MCMC sampling steps



Markov Chain Monte Carlo
Langevin dynamics

Generate new samples

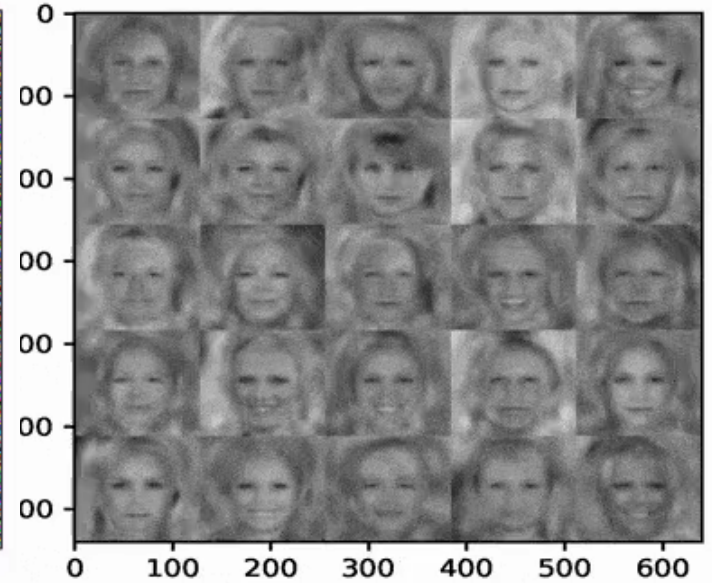
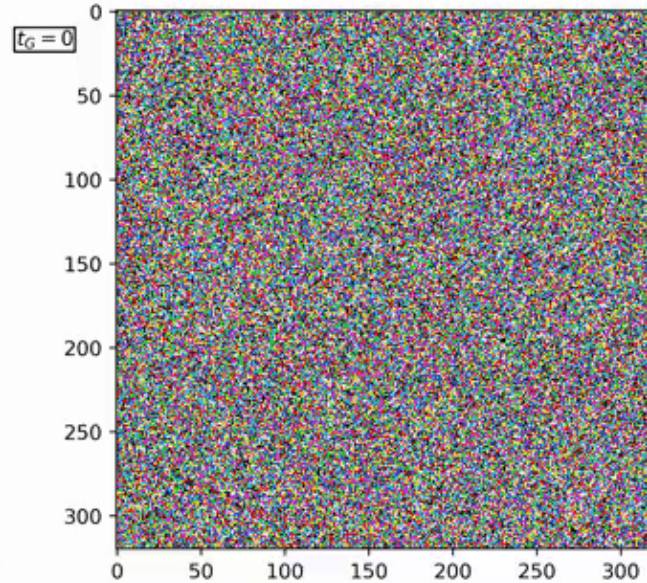
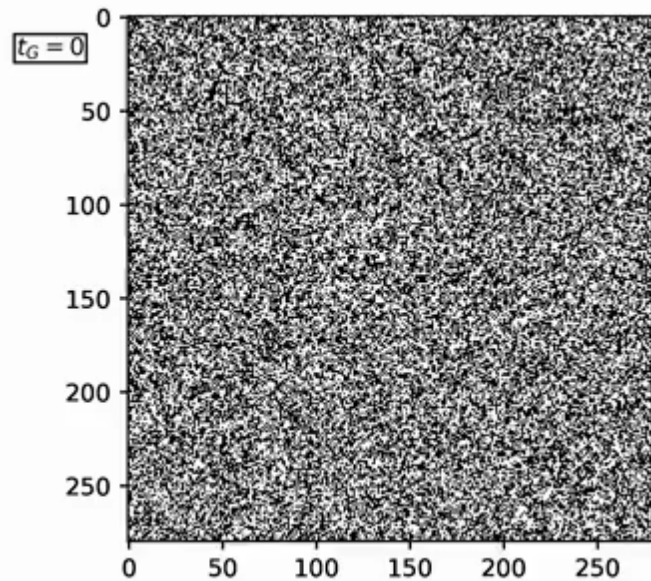
$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\text{data}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\theta}} \quad \forall \theta_i$$

+10⁴ steps

+10⁴ steps

+10⁴ steps

But this is **not** what one **typically** observes...





What's going wrong?

The sampling problem

$$\nabla \mathcal{L} = \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\text{data}}}}_{\text{Easy}} - \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\theta}}}_{\text{Hard} \Rightarrow \text{Markov Chain MC sampling}}$$

The sampling problem

$$\nabla \mathcal{L} = \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\text{data}}}}_{\text{Easy}} - \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\theta}}}_{\text{Hard} \Rightarrow \text{Markov}}$$

RBM

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_{ia}} &= \langle x_i h_a \rangle_{p_{\mathcal{D}}} - \langle x_i h_a \rangle_{\mathcal{E}} \\ \frac{\partial \mathcal{L}}{\partial \eta_a} &= \langle h_a \rangle_{p_{\mathcal{D}}} - \langle h_a \rangle_{\mathcal{E}} \\ \frac{\partial \mathcal{L}}{\partial \zeta_i} &= \langle x_i \rangle_{p_{\mathcal{D}}} - \langle x_i \rangle_{\mathcal{E}} \end{aligned}$$

The sampling problem

$$\nabla \mathcal{L} = \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\text{data}}}}_{\text{Easy}} - \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\theta}}}_{\text{Hard} \Rightarrow \text{Markov}}$$

RBM

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_{ia}} &= \langle x_i h_a \rangle_{p_{\mathcal{D}}} - \langle x_i h_a \rangle_{\mathcal{E}} \\ \frac{\partial \mathcal{L}}{\partial \eta_a} &= \langle h_a \rangle_{p_{\mathcal{D}}} - \langle h_a \rangle_{\mathcal{E}} \\ \frac{\partial \mathcal{L}}{\partial \zeta_i} &= \langle x_i \rangle_{p_{\mathcal{D}}} - \langle x_i \rangle_{\mathcal{E}} \end{aligned}$$

$$\mathbf{x}_{\text{gen}}^{(m)} \quad m = 1, \dots, n_{\text{chains}}$$

$\mathbf{X}_{\text{gen}} \sim P_{\theta}$ Via a Markov Chain Monte Carlo process

$$\langle -\nabla E_{\theta} \rangle_{p_{\theta}} \approx \frac{1}{n_{\text{chains}}} \sum_{m=1}^{n_{\text{chains}}} \nabla E(\mathbf{x}_{\text{gen}}^{(m)})$$

The sampling problem

$$\nabla \mathcal{L} = \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\text{data}}}}_{\text{Easy}} - \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\theta}}}_{\text{Hard} \Rightarrow \text{Markov}}$$

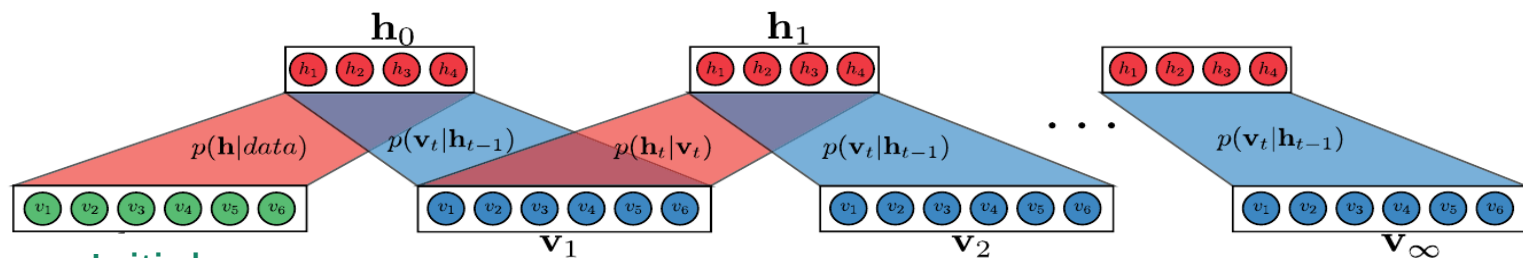
RBM

$$\frac{\partial \mathcal{L}}{\partial W_{ia}} = \langle x_i h_a \rangle_{p_{\mathcal{D}}} - \langle x_i h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \eta_a} = \langle h_a \rangle_{p_{\mathcal{D}}} - \langle h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \zeta_i} = \langle x_i \rangle_{p_{\mathcal{D}}} - \langle x_i \rangle_{\mathcal{E}}$$

Alternating/Block Gibbs sampling



Initial configuration

Glauber dynamics

$$p(h_{\nu} = 1 | \mathbf{v}) = \sum_{\mathbf{h}_{-\nu}} \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})} = \text{sigmoid} \left(b_{\nu} + \sum_i v_i w_{i\nu} \right)$$

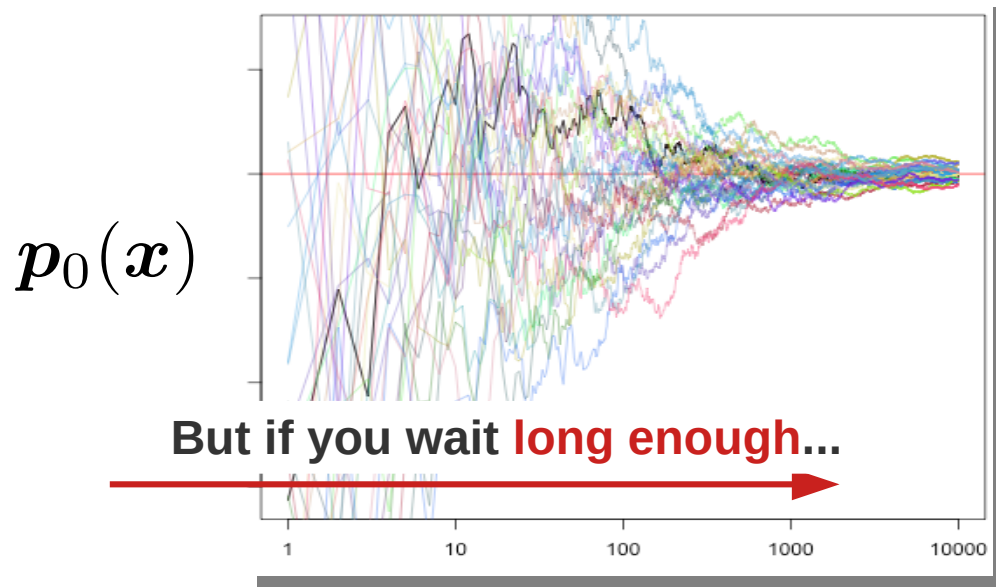
$$p(v_i = 1 | \mathbf{h}) = \sum_{\mathbf{v}_{-i}} \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{h})} = \text{sigmoid} \left(a_i + \sum_{\mu} w_{i\mu} h_{\mu} \right)$$

The sampling problem

$$\nabla \mathcal{L} = \langle -\nabla E_{\theta} \rangle_{p_{\text{data}}} - \langle -\nabla E_{\theta} \rangle_{p_{\theta}}$$

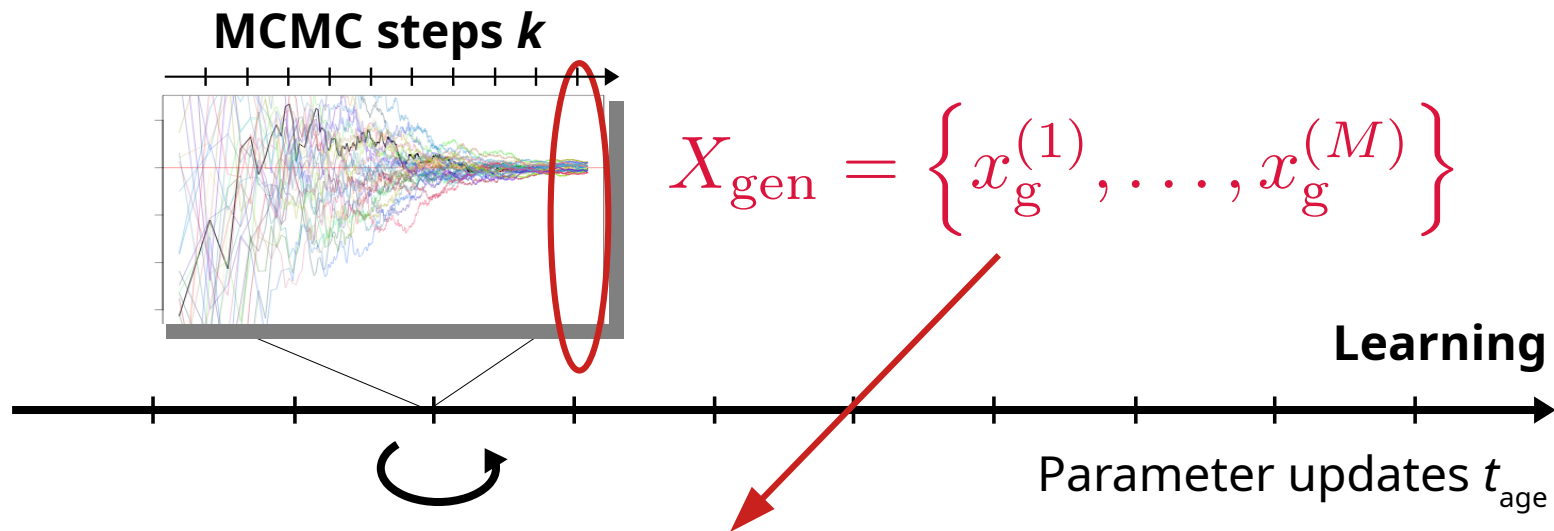
Easy

Hard \Rightarrow Markov Chain MC sampling



k MCMC steps $>$
thermalization/convergence time

Gibbs sampling + learning



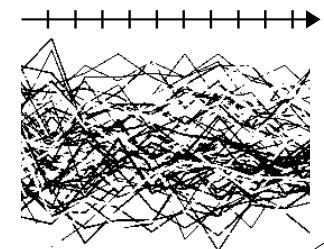
$$\theta_i^{(t+1)} \leftarrow \theta_i^t + \gamma \left. \frac{\partial \mathcal{L}}{\partial \theta_i} \right|_{\theta = \theta_i^{(t)}}$$

Repeat this process
 $\sim 10^5 - 10^6$ times

Gibbs sampling + learning

Standard approach

MCMC steps k



- 1) Use alternate/block sampling (Glauber) dynamics
- 2) Use some few MCMC steps $k \sim \mathcal{O}(1)$
- 3) Choose “good” **initializations**

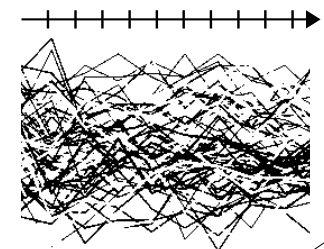
Parameter updates t_{age}



Gibbs sampling + learning



Standard approach

MCMC steps k

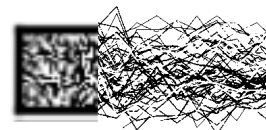


- 1) Use alternate/block sampling (Glauber) dynamics
- 2) Use some few MCMC steps $k \sim \mathcal{O}(1)$
- 3) Choose “good” initializations

Parameter updates t_{age}

- Contrastive divergence (CD) [Hinton (2002)] **Init – dataset** 
- Persistence CD (PCD) [Tieleman (ICML 2008)] **Init – previous end point** 
- **Other solutions** : TAP fixed points [Gabrié, Tramel, and Krzakala (NeurIPS 2015)], optimized sampling techniques (Parallel tempering, Simulated annealing, the Tethered MC method]

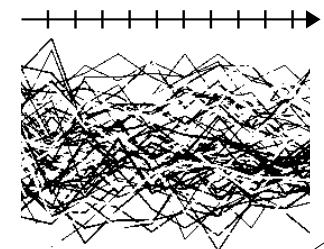
- **Random initialization**
 - [Nijkamp, Hill, Han, Wu, Zhu., (NeurIPS 2019)
 - Decelle, Furtlehner, Seoane (NeurIPS 2021)]



Gibbs sampling + learning

Standard approach

MCMC steps k



- 1) Use alternate/block sampling (Glauber) dynamics
- 2) Use some few MCMC steps $k \sim O(1)$
- 3) Choose "good" initializations

Parameter updates t_{age}

- Contrastive divergence (CD) [Hinton (2002)] **Init – dataset**
- Persistence CD (PCD) [Tieleman (ICML 2008)] **Init – previous end point**
- **Other solutions** : TAP fixed points [Gabrié, Tramel, and Krzakala (NeurIPS 2015)].

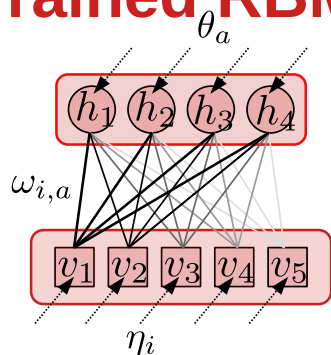


Total disinterest in controlling whether k was long enough to ensure a proper sampling

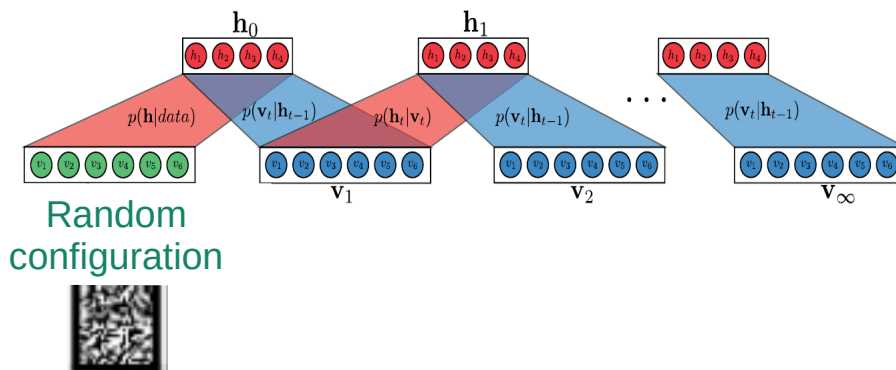


Generating new samples

Trained RBM



Alternating/Block Gibbs sampling

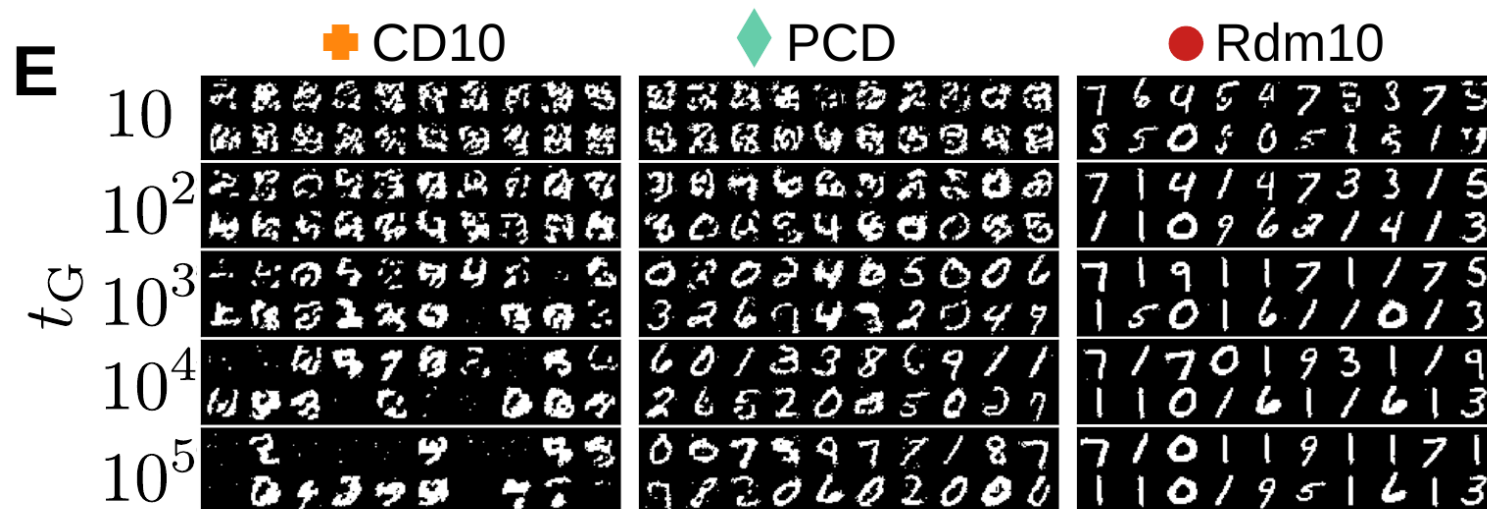
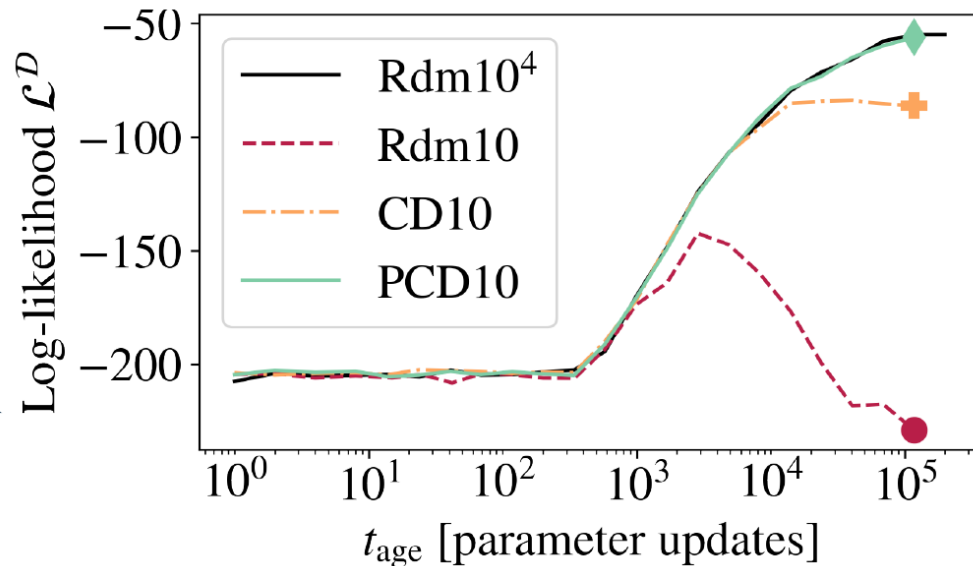


Converge to equilibrium

169547375

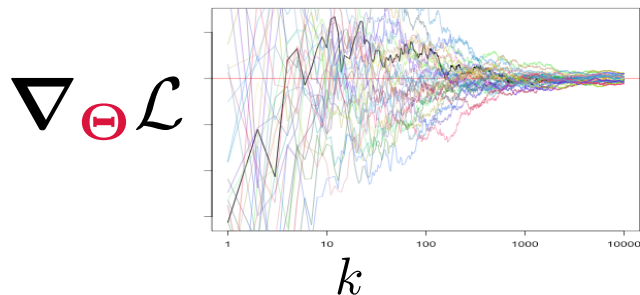
Generating samples : typical situation

The typical measures used to control the learning are not a good estimator of the quality of the generator



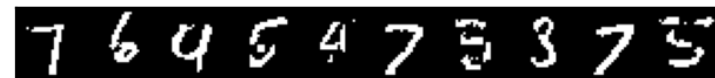
Equilibrium vs. Non-eq. regimes

Training



- **Equilibrium** $k > t_{\text{therm}}$
(Convergent MCMC)
- **Non-equilibrium** $k < t_{\text{therm}}$
(Non-convergent MCMC)

Sampling



$$p_{\text{data}}(x) \sim p_{\theta}(x) = \frac{e^{-E_{\theta}(x)}}{Z}$$

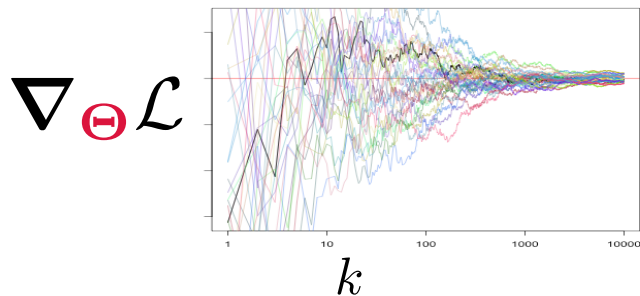
- **Learns a good model** for the data

$$p_{\text{data}}(x) \sim p(k, \mathbf{p}_0, x)$$

- **“Learns the dynamics”**
Memory effects

Equilibrium vs. Non-eq. regimes

Training



- **Equilibrium** $k > t_{\text{therm}}$
(Convergent MCMC)
- **Non-equilibrium** $k < t_{\text{therm}}$
(Non-convergent MCMC)

Sampling



Good for modeling

$$p_{\text{data}}(x) \sim p_{\theta}(x) = \frac{e^{-E_{\theta}(x)}}{Z}$$

- **Learns a good model for the data**

Optimal for generation

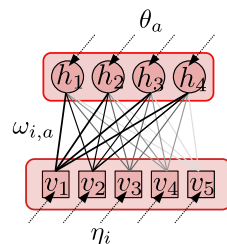
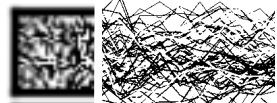
$$p_{\text{data}}(x) \sim p(k, \mathbf{p}_0, x)$$

- **"Learns the dynamics"**
Memory effects

Out-of-equilibrium regime

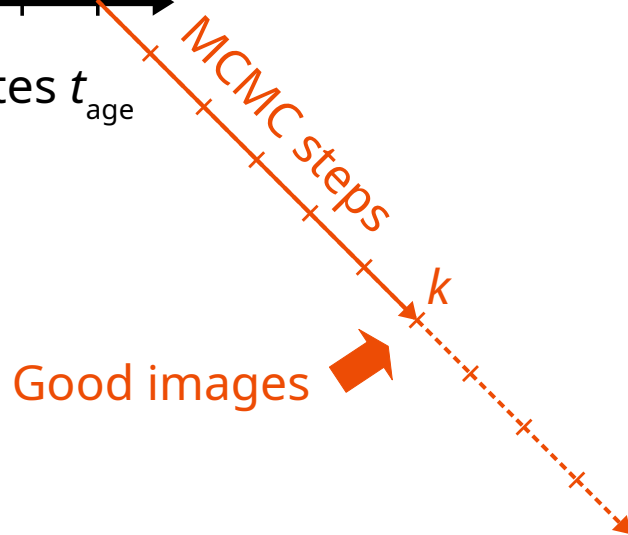
Gibbs sampling

MCMC steps k



Learning

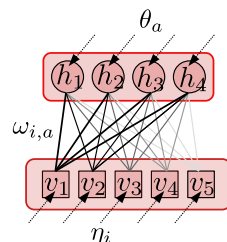
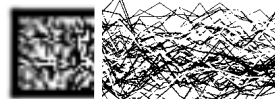
Parameter updates t_{age}



Out-of-equilibrium regime

Gibbs sampling

MCMC steps k



$k = 100$

Learning

memory

Parameter updates t_{age}

MCMC steps k

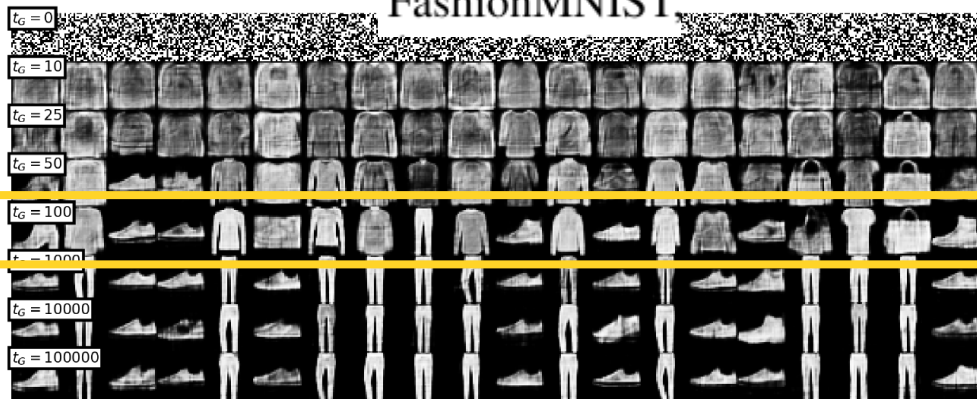
Good images

equilibrium



Out-of-equilibrium regime

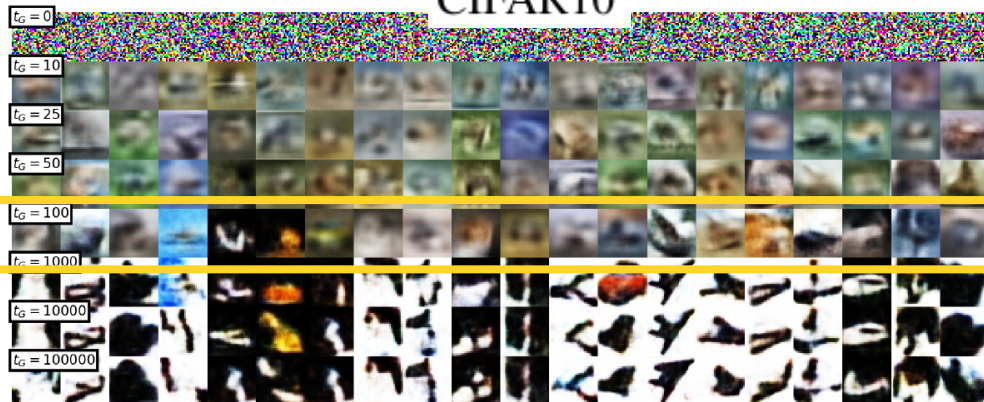
FashionMNIST



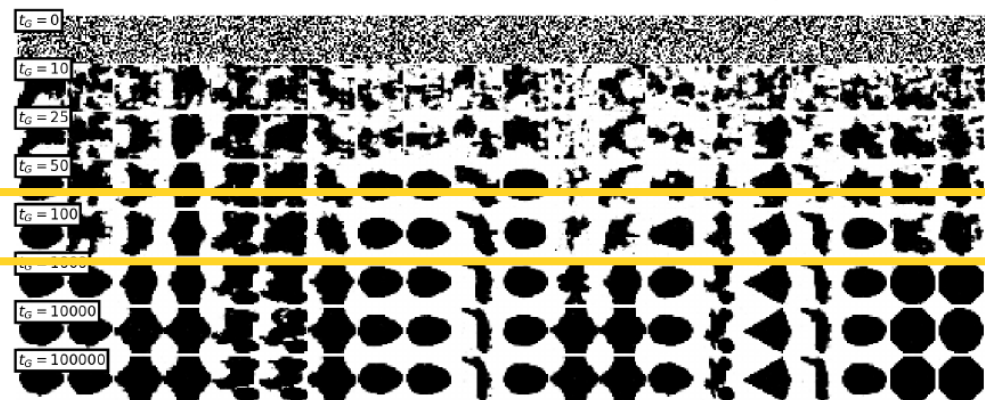
CELEBA



CIFAR10

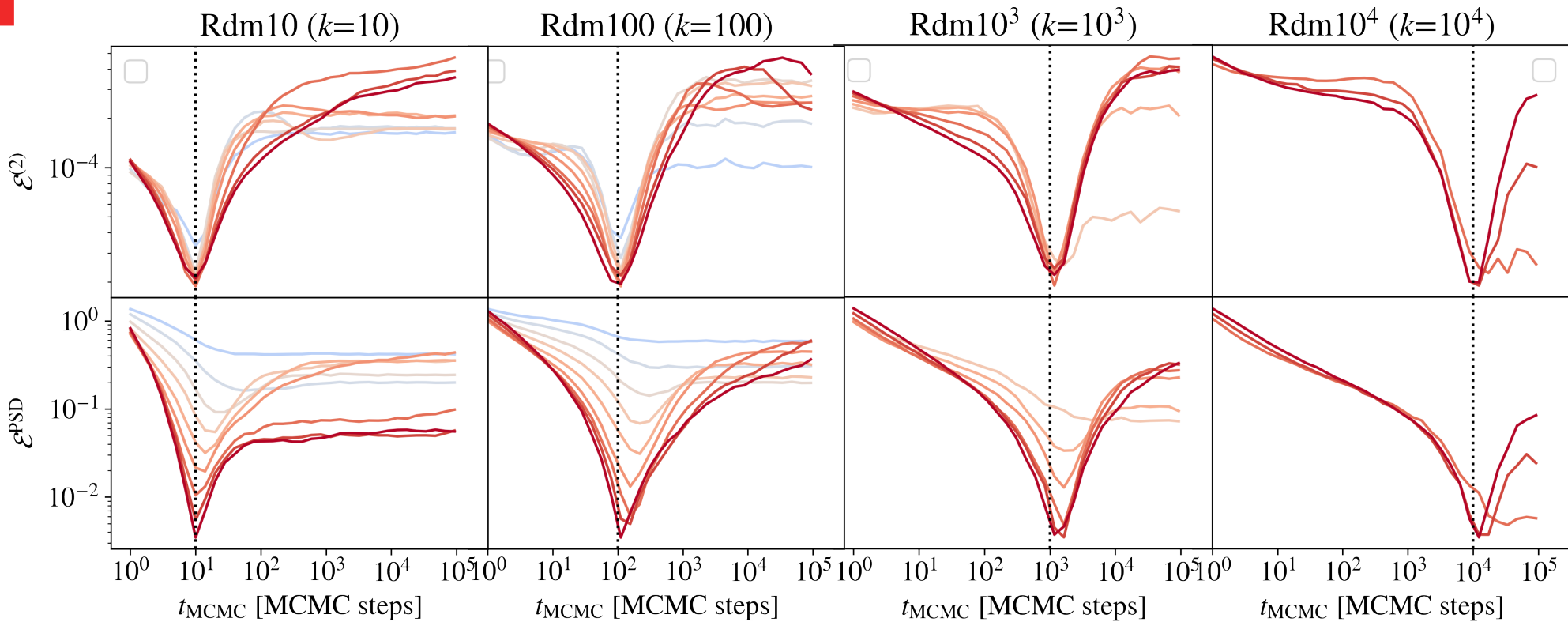


CALTECH101 Silhouettes,



Non-equilibrium regime : generation

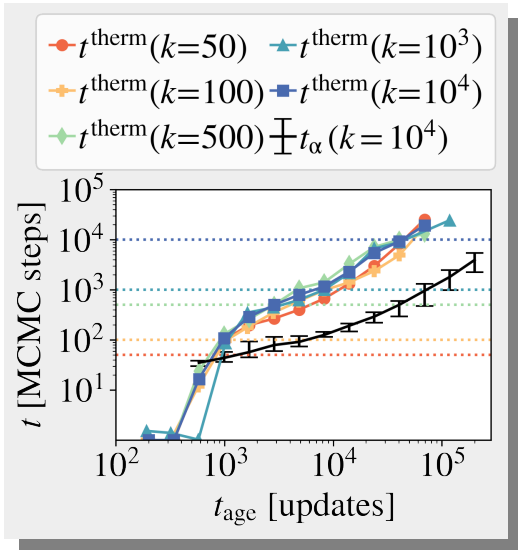
Quality of the generated samples



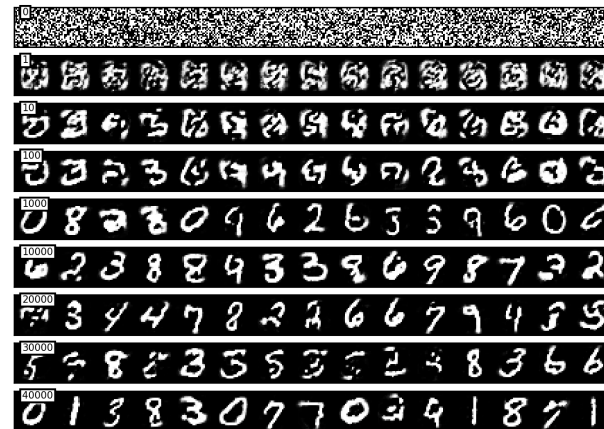
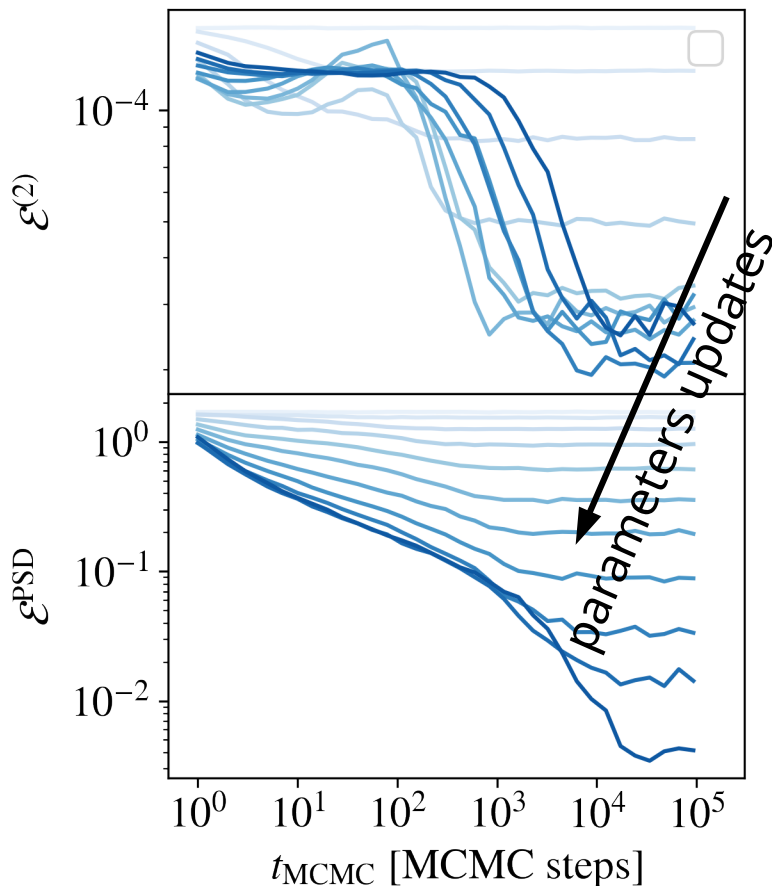
Best quality samples are obtained at $t_{\text{MCMC}} \sim k$

Equilibrium regime

Memory effects are caused by the **lack of convergence** of the Markov chains



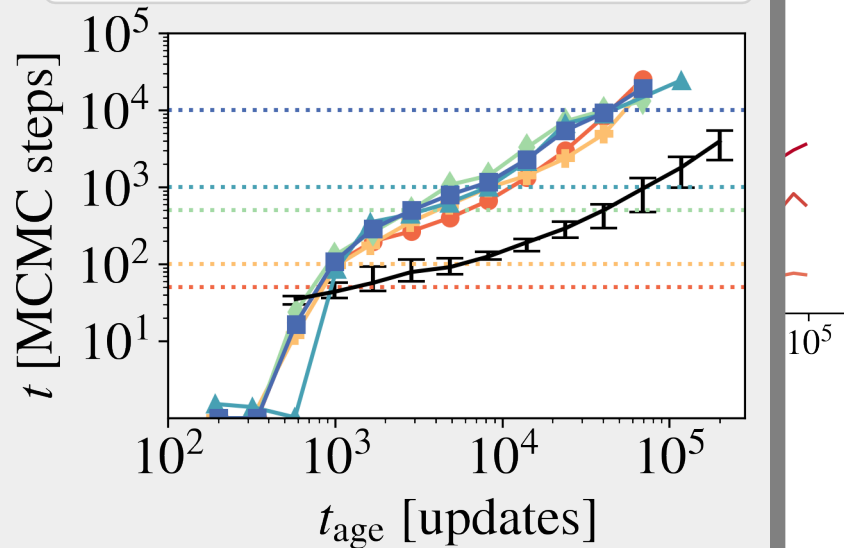
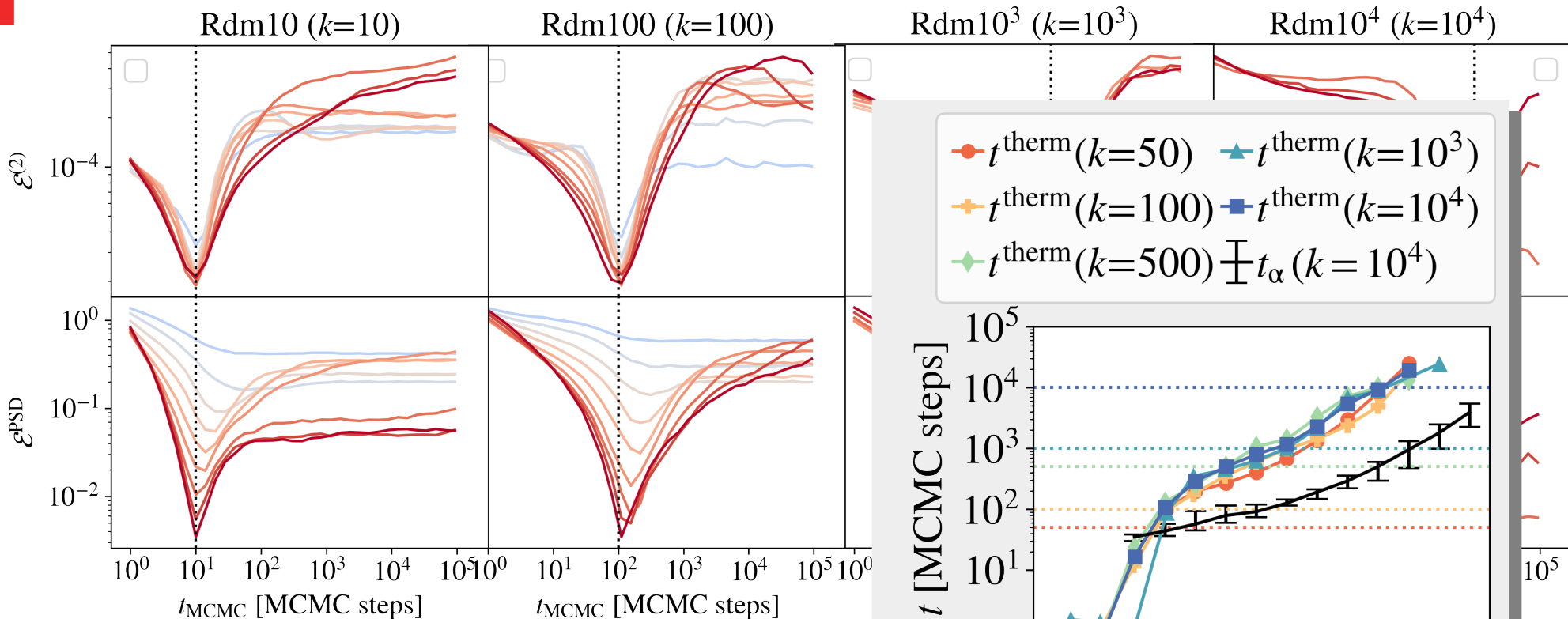
Rdm-10⁴



Dynamics are much faster

Non-equilibrium regime : generation

Quality of the generated samples



Best quality samples are reached at $t^{\text{therm}}(k=50)$, $t^{\text{therm}}(k=100)$, $t^{\text{therm}}(k=500)$, $t^{\text{therm}}(k=10^3)$, $t^{\text{therm}}(k=10^4)$, and $t_{\alpha}(k=10^4)$.

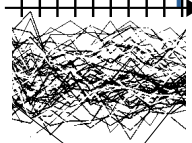
Dangers of Contrastive Divergence (CD)

Hinton, "A practical guide to training restricted Boltzmann machines" (2012)

$k=1$

Gibbs sampling

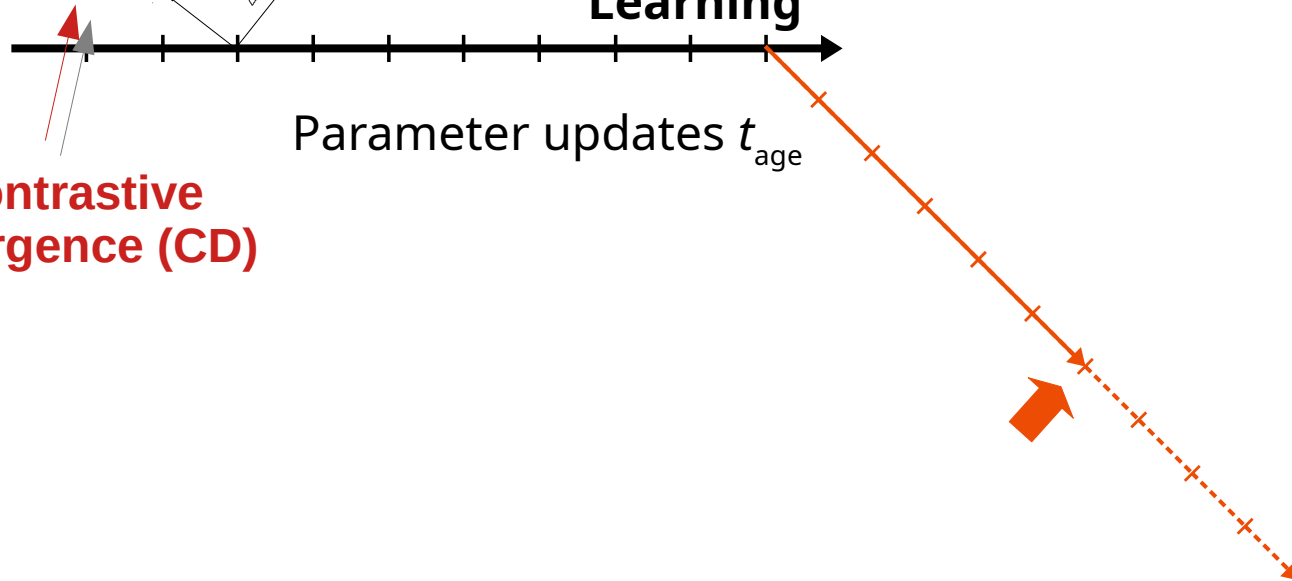
k MCMC steps



Learning

Parameter updates t_{age}

Contrastive Divergence (CD)



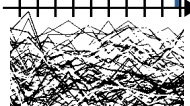
Dangers of Contrastive Divergence (CD)

Hinton, "A practical guide to training restricted Boltzmann machines" (2012)

$k=1$

Gibbs sampling

k MCMC steps

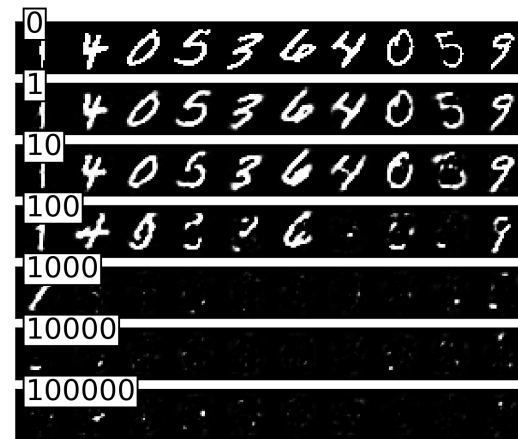
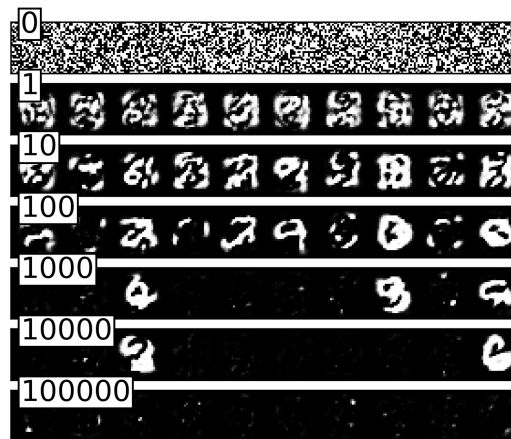


Learning

Parameter updates t_{age}

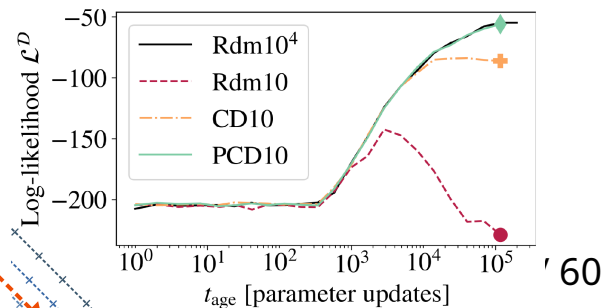
Contrastive Divergence (CD)

CD-10



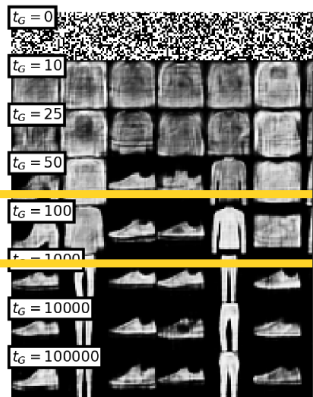
CD is a very bad recipe !

Salakhutdinov, R., & Murray, I. ICML (2008)
Desjardins, Courville, Bengio, Vincent, Delalleau, AISTATS, 2010



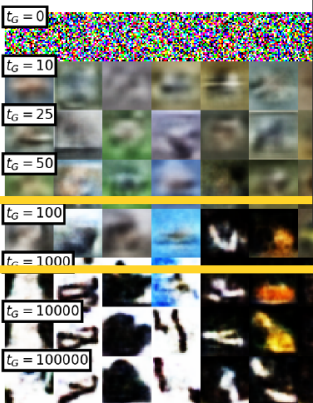
Out-of-equilibrium regime

FashionMNIST

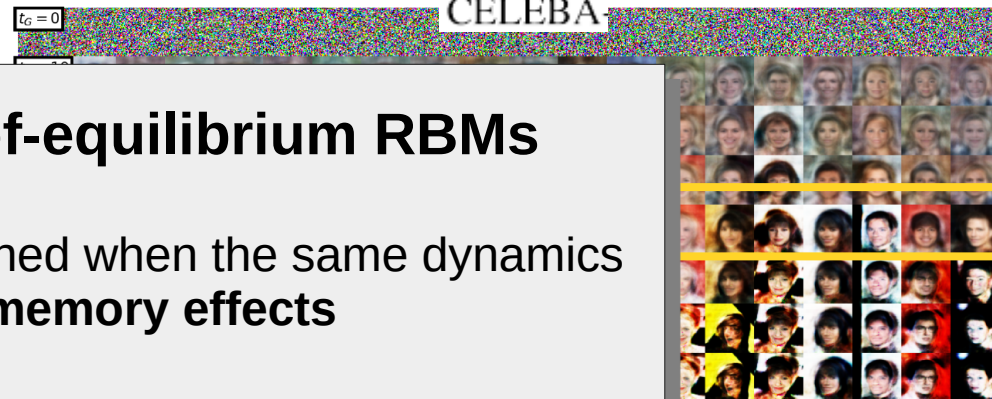


Characteristics out-of-equilibrium RBMs

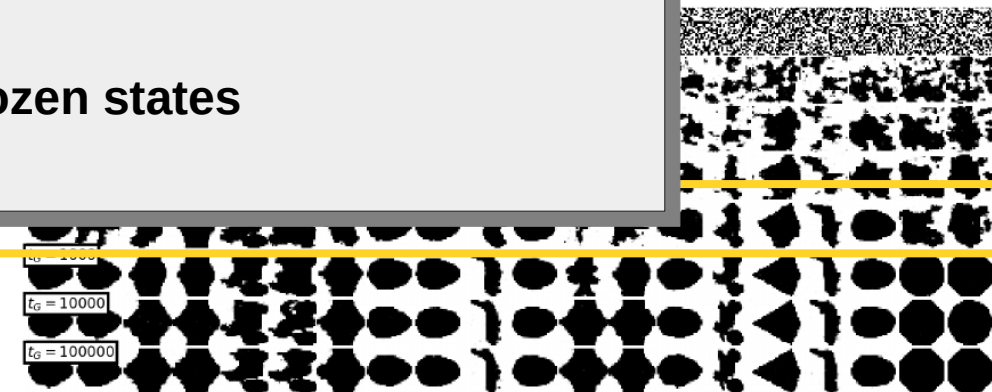
- Dataset-like samples are obtained when the same dynamics of the training are repeated → **memory effects**
- Long-time configurations are **biased** and **lack diversity**
- Extremely slow dynamics → **frozen states**



CELEBA



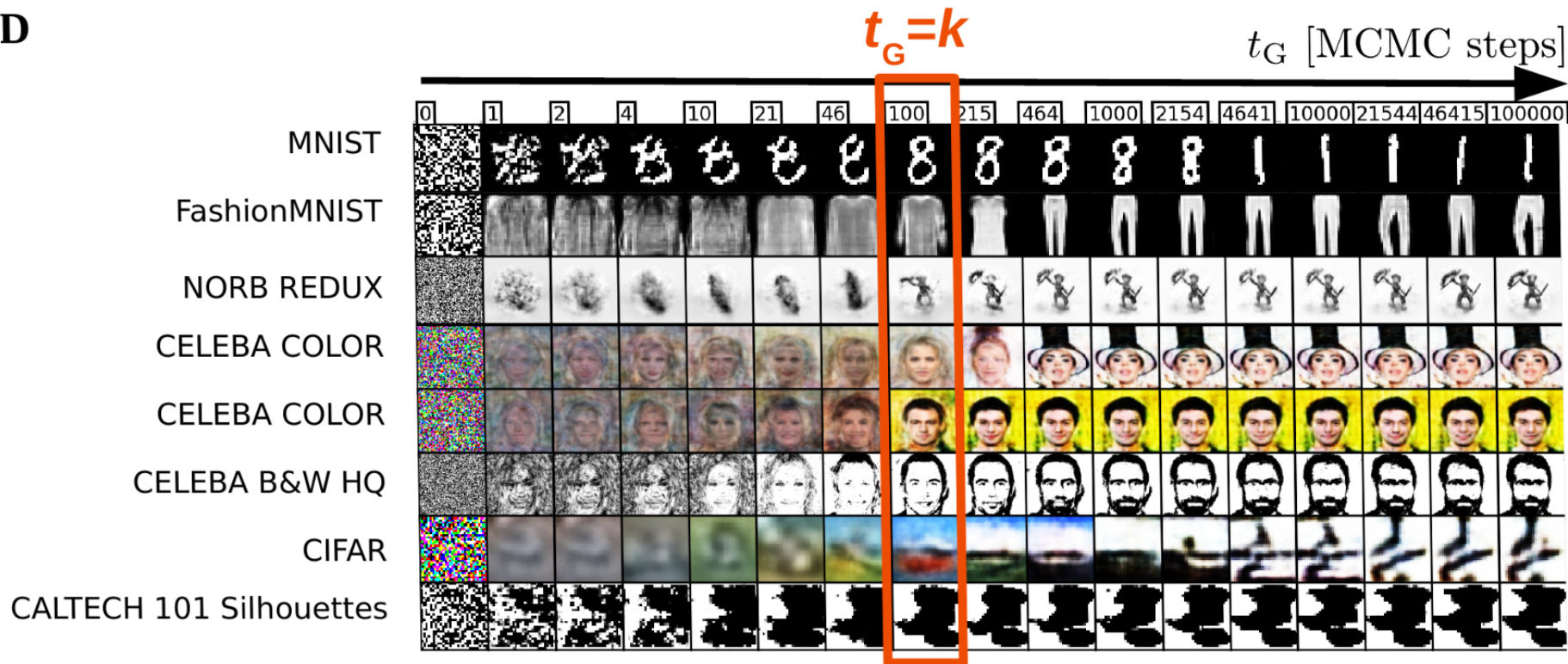
ettes,



Out-of-equilibrium regime : best for sample generation

[Decelle, Furtlehner, Seoane
NeurIPS (2021)]

D



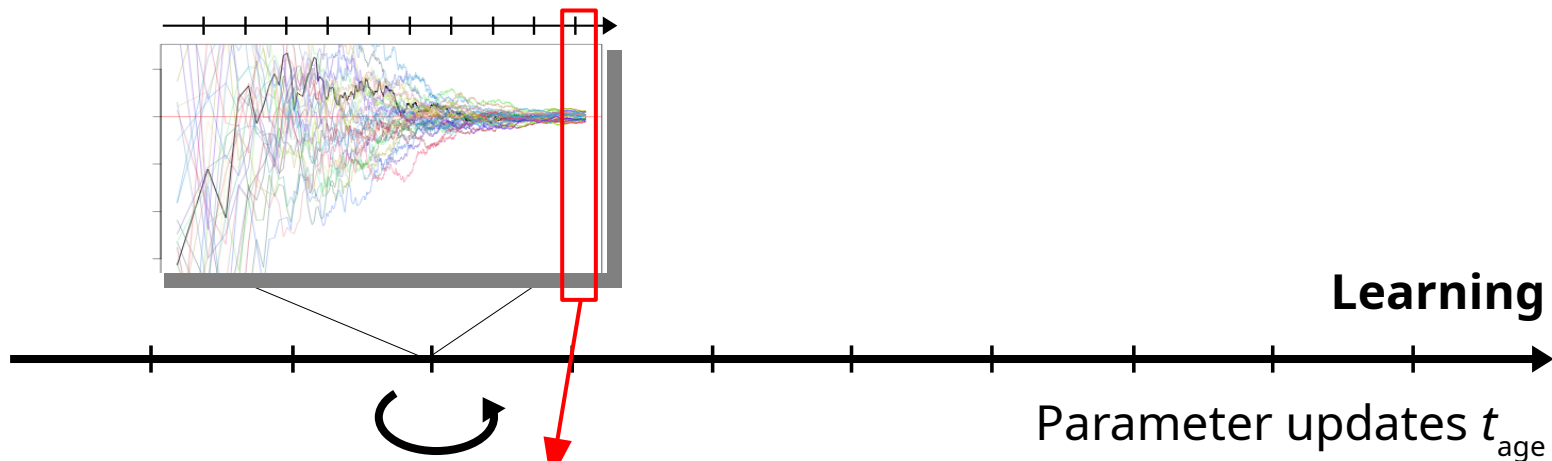


Explanation effects

Agoritsas, Catania, Decelle, Seoane ICML (2023)

Gibbs sampling + learning

Convergent MCMC



$$\theta(t + t) \leftarrow \theta(t) + \gamma \nabla \mathcal{L}(t)$$

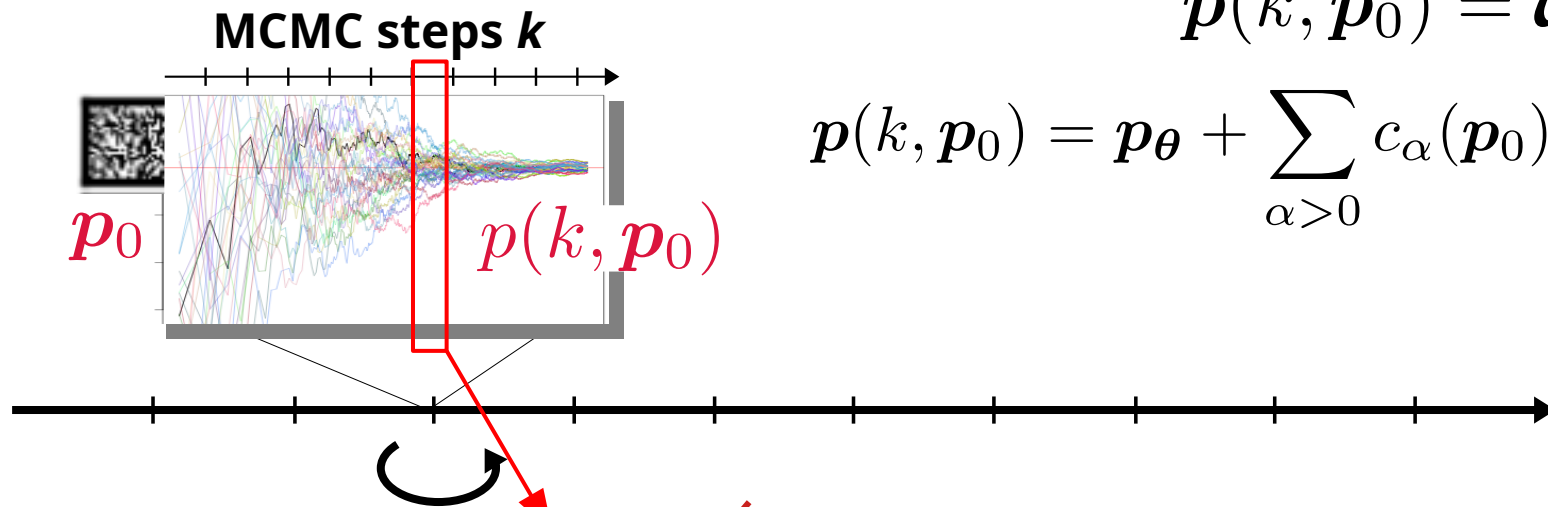
$$\nabla \mathcal{L} = \langle -\nabla E \rangle_{p_{\text{data}}} - \langle -\nabla E \rangle_{p_{\theta}}$$

Gibbs sampling + learning

\mathcal{U}_θ : stochastic matrix

$$p(k, p_0) = \mathcal{U}_\theta^k p_0$$

$$p(k, p_0) = p_\theta + \sum_{\alpha > 0} c_\alpha(p_0) e^{-k/\kappa_\alpha} \mathbf{u}_\alpha$$



$$\theta(t + t) \leftarrow \theta(t) + \gamma \nabla \mathcal{L}(t)$$

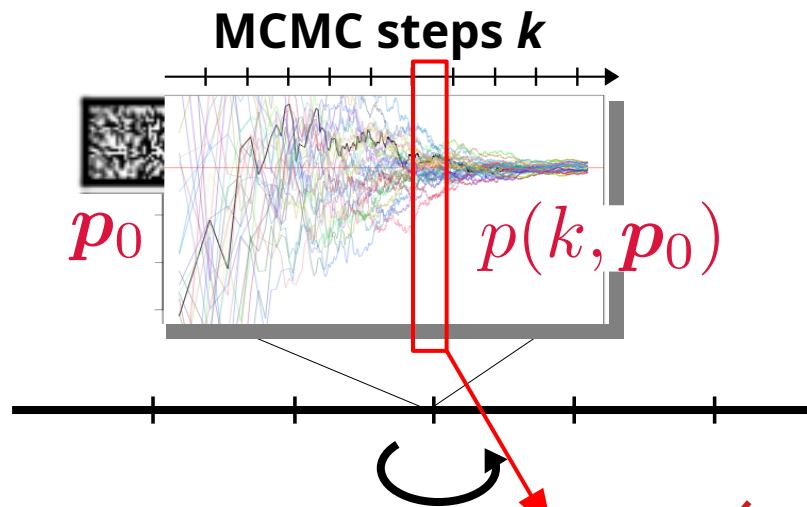
$$\nabla \mathcal{L} = \langle -\nabla E \rangle_{p_{\text{data}}} - \langle -\nabla E \rangle_{p_\theta}$$

Gibbs sampling + learning

\mathcal{U}_θ : stochastic matrix

$$p(k, p_0) = \mathcal{U}_\theta^k p_0$$

$$p(k, p_0) = p_\theta + \sum_{\alpha > 0} c_\alpha(p_0) e^{-k/\kappa_\alpha} \mathbf{u}_\alpha$$



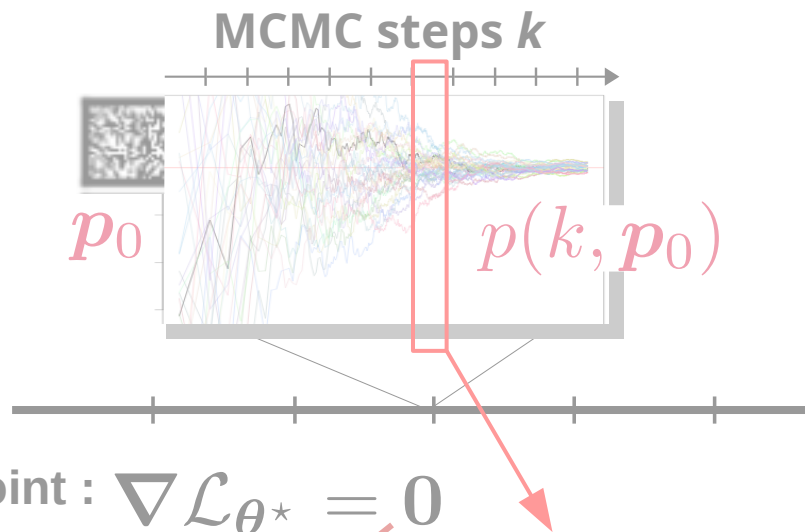
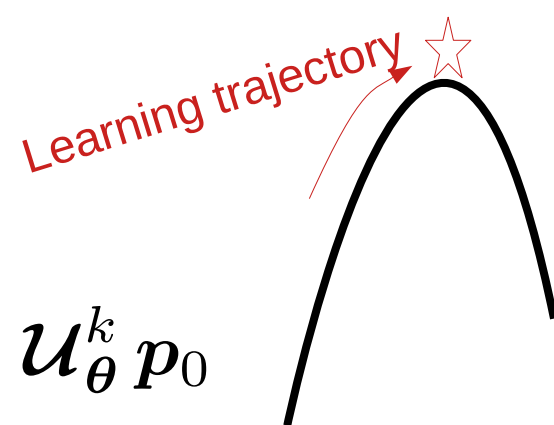
$$\langle -\nabla E_\theta \rangle_{p_\theta} \neq \langle -\nabla E_\theta \rangle_{p(k, p_0)}$$

$$\tilde{\nabla} \mathcal{L} = \langle -\nabla E_\theta \rangle_{p_{\text{data}}} - \langle -\nabla E_\theta \rangle_{p(k, p_0)}$$

$$\theta(t + t) \leftarrow \theta(t) + \gamma \nabla \mathcal{L}(t)$$

$$\nabla \mathcal{L} = \langle -\nabla E \rangle_{p_{\text{data}}} - \langle -\nabla E \rangle_{p_\theta}$$

Gibbs sampling + learning



$$p(k, p_0) = \mathcal{U}_{\theta}^k p_0$$

$$\langle -\nabla E_{\theta} \rangle_{p_{\theta}} \neq \langle -\nabla E_{\theta} \rangle_{p(k, p_0)}$$

$$\tilde{\nabla} \mathcal{L} = \langle -\nabla E_{\theta} \rangle_{p_{\text{data}}} - \langle -\nabla E_{\theta} \rangle_{p(k, p_0)}$$

~~$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p_{\theta^*}}$$~~

$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p(k, p_0)}$$

Gibbs sampling + learning

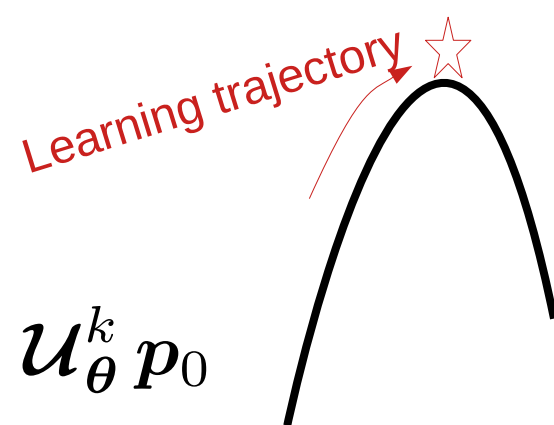
$$p_{\theta^*}(x) = \frac{e^{-E_{\theta^*}(x)}}{Z_{\theta^*}} \not\approx p_{\text{data}}(x)$$

Trainings with non-convergent MCMC fit **bad models** for the data

~~$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p_{\theta^*}}$$~~

$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p(k, p_0)}$$

$$p(k, p_0) = \mathcal{U}_{\theta}^k p_0$$



$$\langle -\nabla E_{\theta} \rangle_{p_{\theta}} \neq \langle -\nabla E_{\theta} \rangle_{p(k, p_0)}$$

$$\tilde{\nabla} \mathcal{L} = \langle -\nabla E_{\theta} \rangle_{p_{\text{data}}} - \langle -\nabla E_{\theta} \rangle_{p(k, p_0)}$$

Gibbs sampling + learning

Similarly:

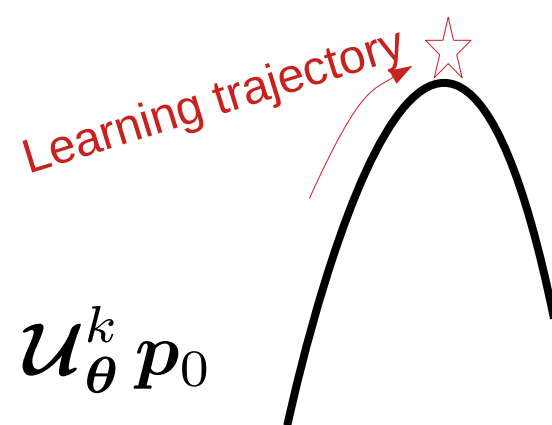
$$p_{\theta^*}(x) = \frac{e^{-E_{\theta^*}(x)}}{Z_{\theta^*}} \not\approx p_{\text{data}}(x)$$

Trainings with non-convergent MCMC fit **bad models** for the data

But they still can be used as **very good generative models**

$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p(k, \mathbf{p}_0)}$$

$$p(k, \mathbf{p}_0) = \mathcal{U}_{\theta}^k \mathbf{p}_0$$



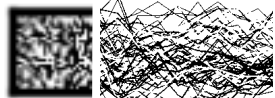
$$\langle -\nabla E_{\theta} \rangle_{p_{\theta}} \neq \langle -\nabla E_{\theta} \rangle_{p(k, \mathbf{p}_0)}$$

$$\tilde{\nabla} \mathcal{L} = \langle -\nabla E_{\theta} \rangle_{p_{\text{data}}} - \langle -\nabla E_{\theta} \rangle_{p(k, \mathbf{p}_0)}$$

Non-equilibrium regime

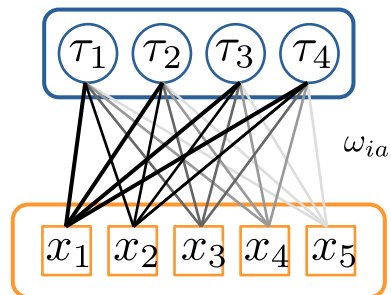
Gibbs sampling

MCMC steps $k=100$



Learning

Parameter updates t_{age}



$k = 100$

memory

MCMC steps

k

Good images

$$\langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p(k, \mathbf{p}_0)}$$

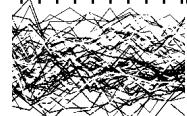
equilibrium



Change initial conditions

Gibbs sampling

MCMC steps $k=100$



Initialized at the dataset



Learning

Parameter updates t_{age}

MCMC steps

k

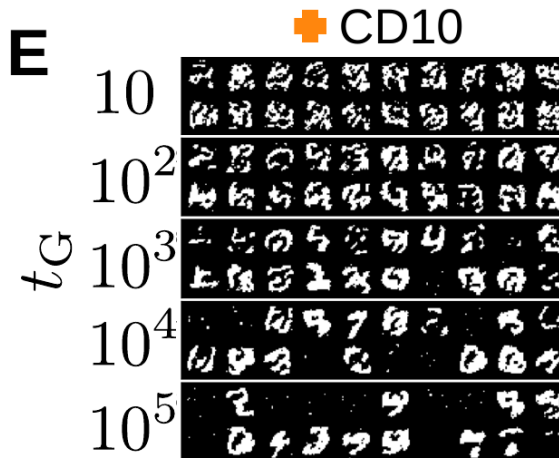
Good images

$$\langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p(k, p_0)}$$



Changing the sampling process...

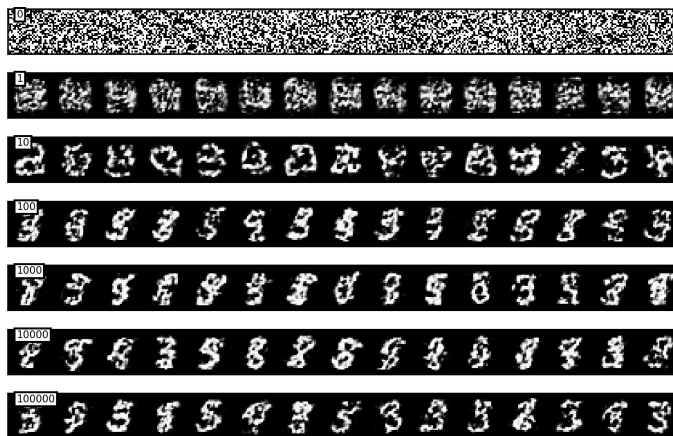
Change the initial conditions :



Gibbs sampling
MCMC steps $k=10$



Change the update rules:



RBM trained using

Gabrié, M., Tramel, E. W., &
Krzakala, F. (2015). NeurIPS

Out-of-equilibrium regime : best for sample generation



Generative Convnets:
Nijkamp, Hill, Han, Wu, Zhu.
NeurIPS 2019, AACL 2020.

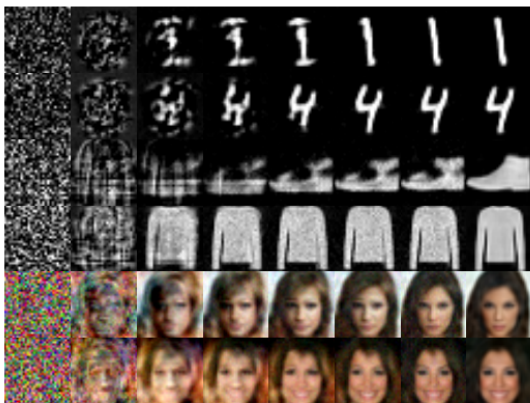


Figure 2: Intermediate samples from Gibbs-Langevin sampling.

GAUSSIAN-BERNOULLI RBMs WITHOUT TEARS

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

University of British Columbia¹, Google Research, Brain Team²,

New York University³, University of Toronto⁴, Vector Institute⁵

[ArXiv: 2210.10318 \(2022\)](https://arxiv.org/abs/2210.10318)

rjliao@ece.ubc.ca, mengye@cs.nyu.edu

{skornblith, davidfleet, geoffhinton}@google.com

Published as a conference paper at ICLR 2022

A TALE OF TWO FLOWS: COOPERATIVE LEARNING
OF LANGEVIN FLOW AND NORMALIZING FLOW
TOWARD ENERGY-BASED MODEL

Jianwen Xie, Yaxuan Zhu, Jun Li, Ping Li

Consequences (theorem)

Fixed point

$$\langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p(k, p_0)}$$

[Agoritsas, Catania, Decelle, Seoane ICML (2023)]

1) Very good quality generation

Fixed point

$$\text{EQ training } \langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p_{\theta}}$$

$$\text{OOE training } \langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p(k, \mathbf{p}_0)}$$



We reproduce
the same set of
statistics!

1) Very good quality generation

Fixed point

$$\text{EQ training } \langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p_{\theta}}$$

$$\text{OOE training } \langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p(k, \mathbf{p}_0)}$$



We reproduce
the same set of
statistics!

Samples generated using the OOE strategy should be as good as if the EBM was trained in EQ !

2) Fast generators

Fixed point

EQ training $\langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p_{\theta}}$

We need to thermalize first !

Typically very long...

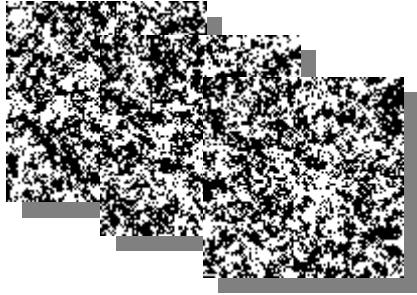
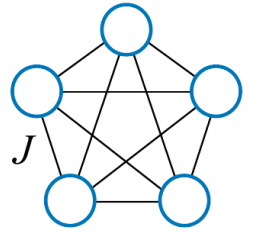
OOE training $\langle \nabla E_{\theta^*} \rangle_{p_D} = \langle \nabla E_{\theta^*} \rangle_{p(k, p_0)}$

We just need k steps !

The quality does not depend on

- k ,
- p_0 ,
- the dynamic rules used for sampling

2) Fast generators : **Inverse Ising**



Training set : Equilibrium samples of the Ising model

$$\{s_i^{(m)}\}_{m=1,\dots,M} \sim p(\mathbf{s}) = \frac{1}{Z(\beta)} e^{\beta \sum_{i<j} J_{ij} s_i s_j + \beta \sum_i h_i s_i}$$

Infer J by training a Boltzmann Machine (Ising model)

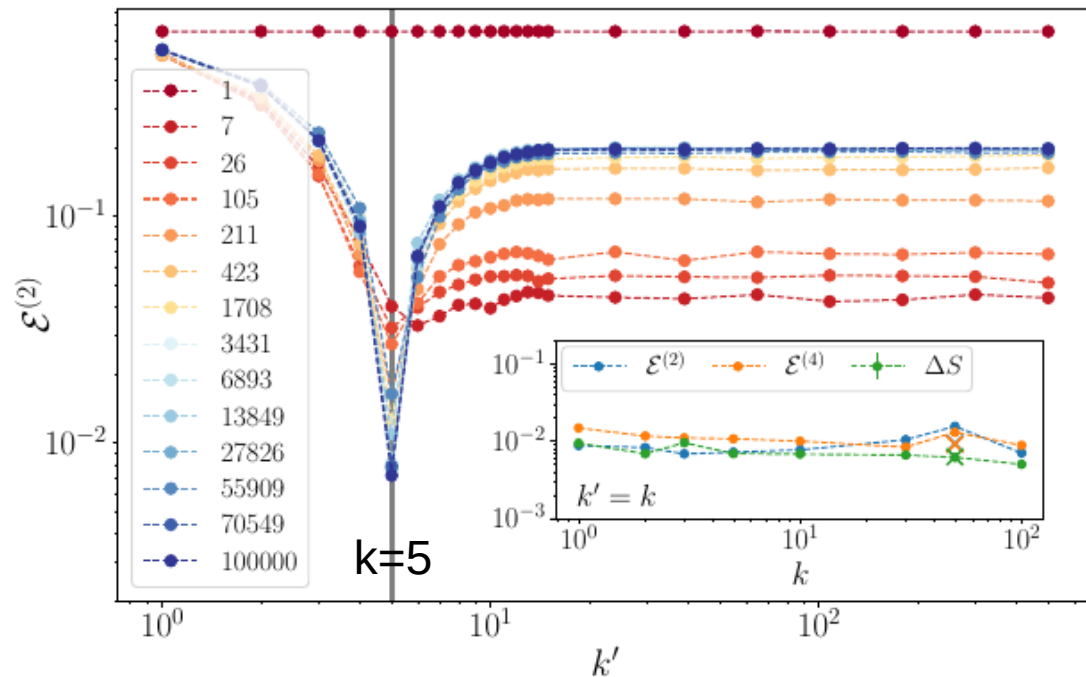
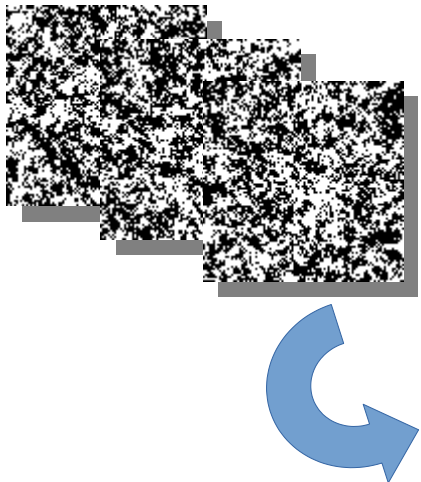
$$E(\mathbf{x}) = - \sum_{ij} J_{ij} x_i x_j - \sum_i h_i x_i$$

$$\langle S_i S_j \rangle_{p_{\text{data}}} = \langle S_i S_j \rangle_{p(k, \mathbf{p}_0)}$$

$$\langle S_i \rangle_{p_{\text{data}}} = \langle S_i \rangle_{p(k, \mathbf{p}_0)}$$

$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p(k, \mathbf{p}_0)}$$

2) Fast generators : **Inverse Ising**



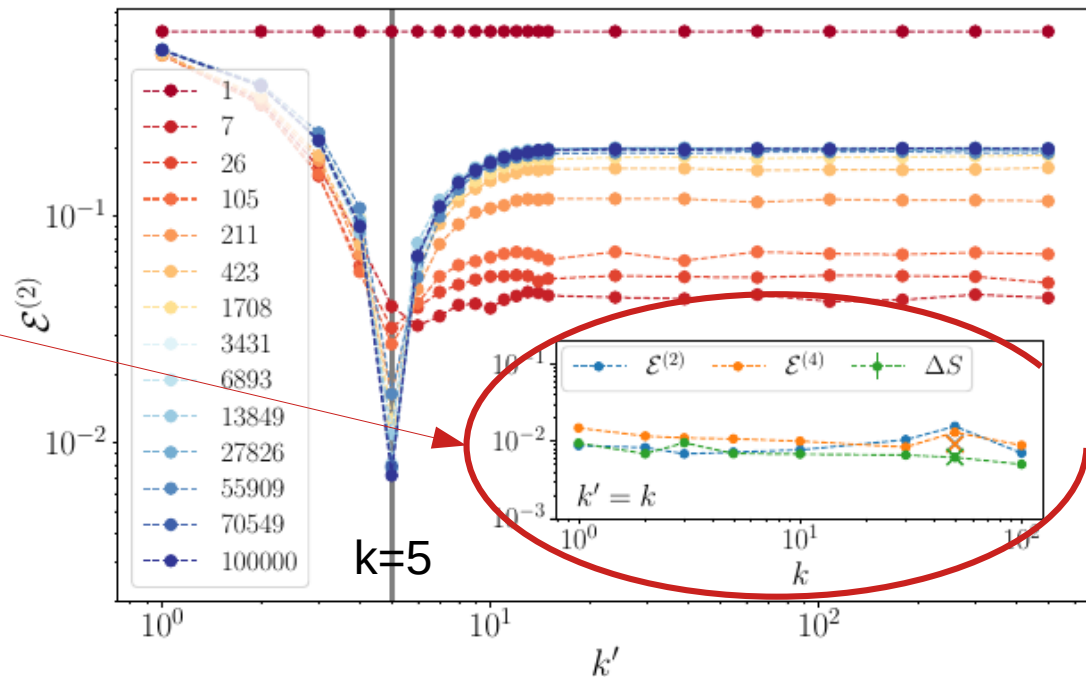
$$\langle S_i S_j \rangle_{p_{\text{data}}} = \langle S_i S_j \rangle_{p(k, \mathbf{p}_0)}$$

$$\langle S_i \rangle_{p_{\text{data}}} = \langle S_i \rangle_{p(k, \mathbf{p}_0)}$$

$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p(k, \mathbf{p}_0)}$$

2) Fast generators : **Inverse Ising**

Even at
k=1 !



$$\langle S_i S_j \rangle_{p_{\text{data}}} = \langle S_i S_j \rangle_{p(k, \mathbf{p}_0)}$$

$$\langle S_i \rangle_{p_{\text{data}}} = \langle S_i \rangle_{p(k, \mathbf{p}_0)}$$

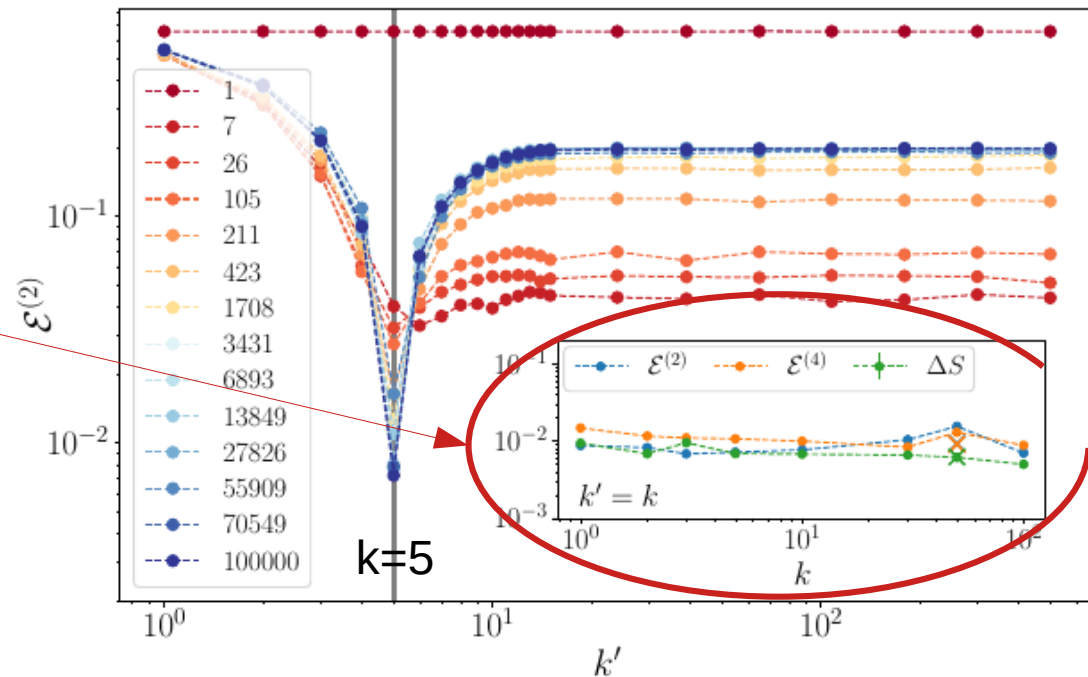
$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p(k, \mathbf{p}_0)}$$

2) Fast generators : Inverse Ising

Describe these curves analytically at high temperatures or in the Gaussian model...

-Agoritsas, Catania, Decelle, Seoane ICML (2023)

Even at $k=1$!

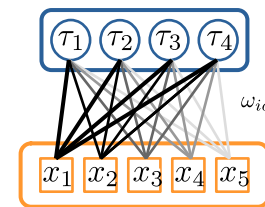


$$\langle S_i S_j \rangle_{p_{\text{data}}} = \langle S_i S_j \rangle_{p(k, \mathbf{p}_0)}$$

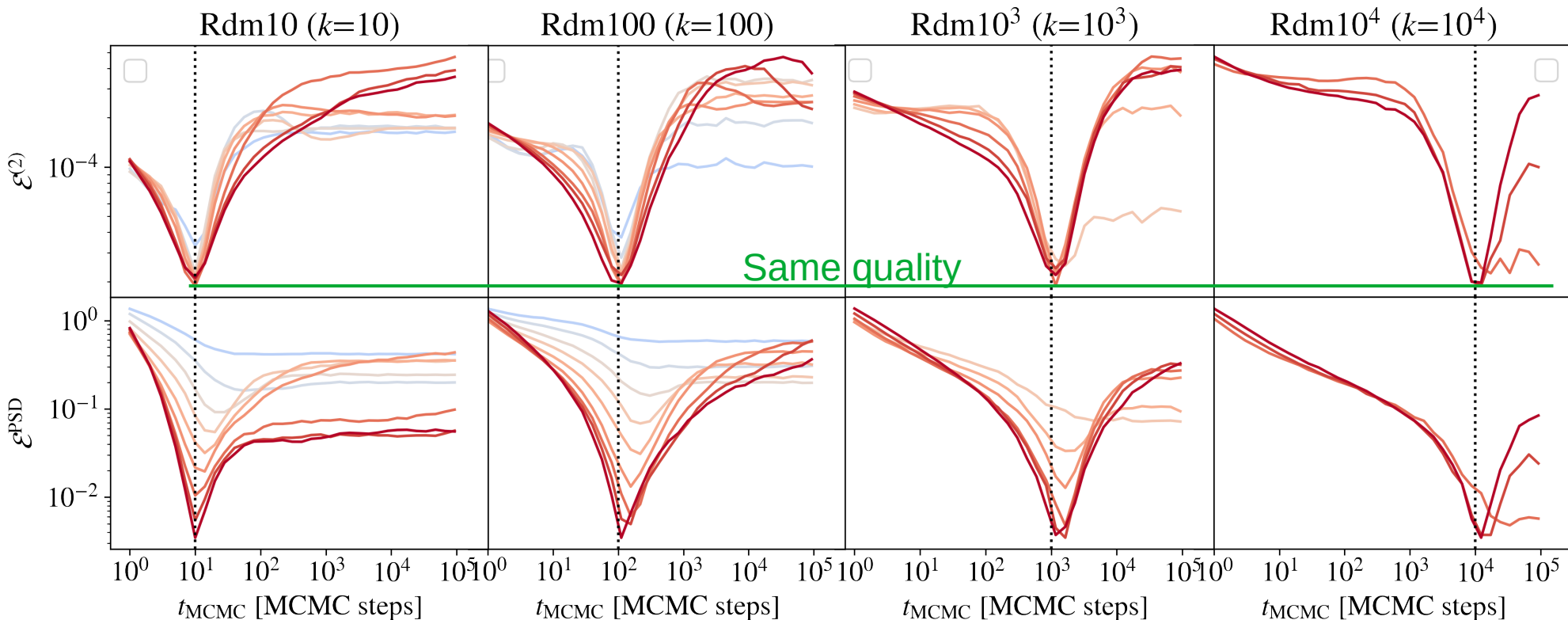
$$\langle S_i \rangle_{p_{\text{data}}} = \langle S_i \rangle_{p(k, \mathbf{p}_0)}$$

$$\langle \nabla E_{\theta^*} \rangle_{p_{\text{data}}} = \langle \nabla E_{\theta^*} \rangle_{p(k, \mathbf{p}_0)}$$

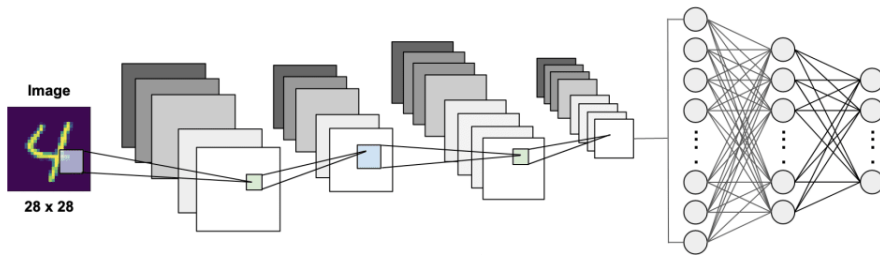
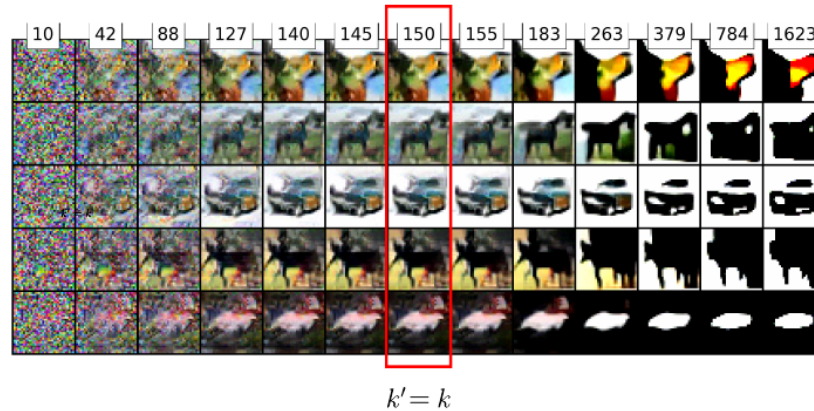
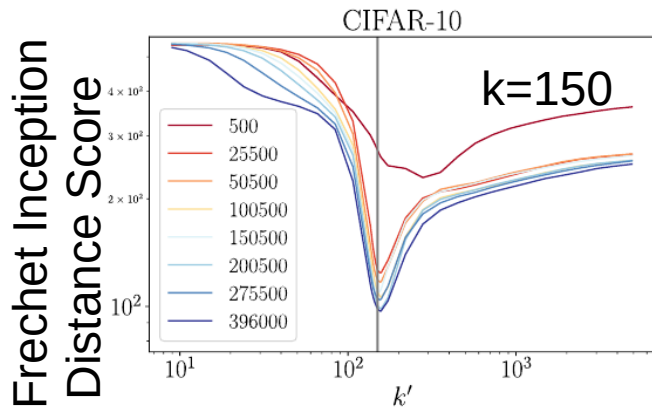
2) Fast generators : The RBM



Quality of the generated samples

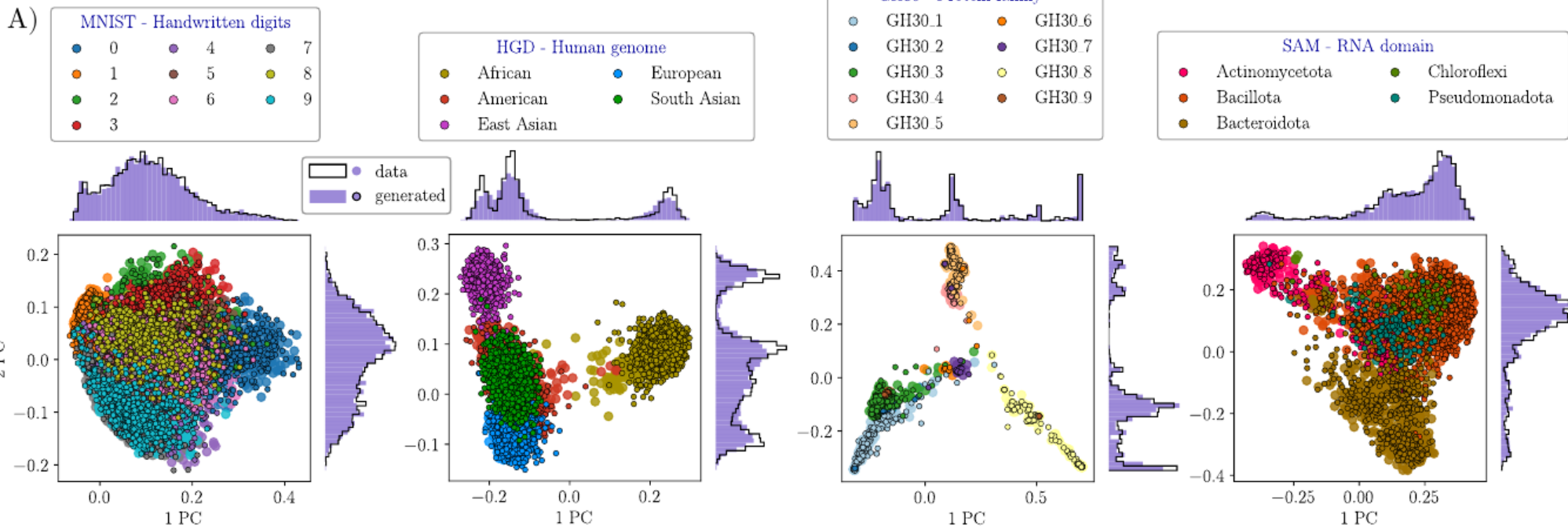


3) This effect is present in any EBM



CIFAR-10
ConvNet EBM

4) and any dataset: **Multimodal distributions**



Fast generation of multimodal data : $k=10$

Very similar to diffusion !

Sohl-Dickstein,
Weiss, Maheswaranathan,
Ganguli, PMLR (2015)

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

Encoder (Gaussian model)

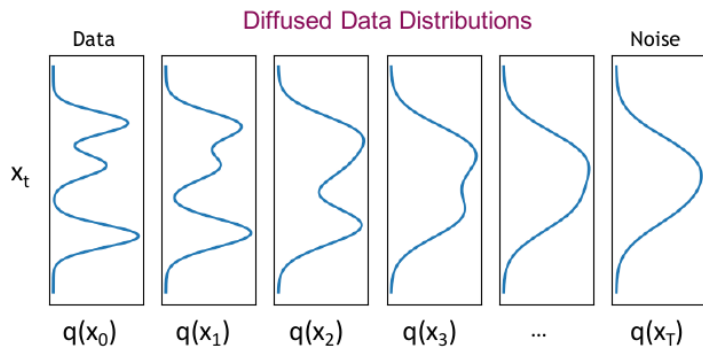
Data



Noise

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2\mathbf{I})$$

Decoder



- K. Kreis, R. Gao, and A. Vahdat. *Denoising diffusion-based generative modeling: foundations and applications*. CVPR Tutorial. 2022.
- H. Cao, C. Tan, Z. Gao, G. Chen, P.-A. Heng, and S. Z. Li. "A Survey on Generative Diffusion Model". ArXiv: 2209.02646 (2022).
- T. Karras, M. Aittala, T. Aila, and S. Laine. "Elucidating the Design Space of Diffusion-Based Generative Models". In: NeurIPS (2022)
- Gao, R., Song, Y., Poole, B., Wu, Y. N., & Kingma, D. P. **Learning energy-based models by diffusion recovery likelihood**. ICLR (2021)



Summary

- If the goal is to generate samples statistically similar to the data:

No point of trying to equilibrate (long training, long generation)

Go for out-of-equilibrium training !!

- But...

Summary

- If the goal is to generate samples statistically similar to the data:

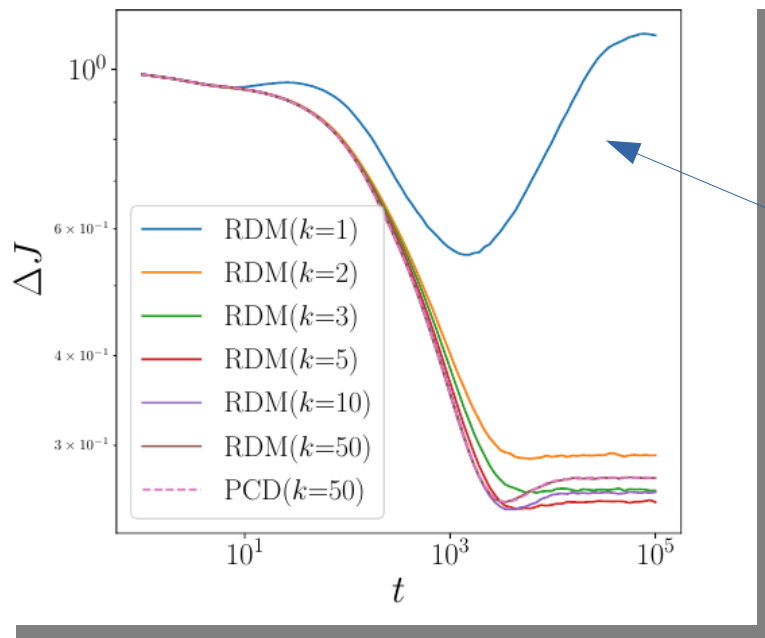
No point of trying to equilibrate (long training, long generation)

Go for out-of-equilibrium training !!

- But...

$$p_{\theta^*}(x) = \frac{e^{-E_{\theta^*}(x)}}{Z_{\theta^*}} \not\approx p_{\text{data}}(x)$$

Is this OOE energy function useful for interpretability ?

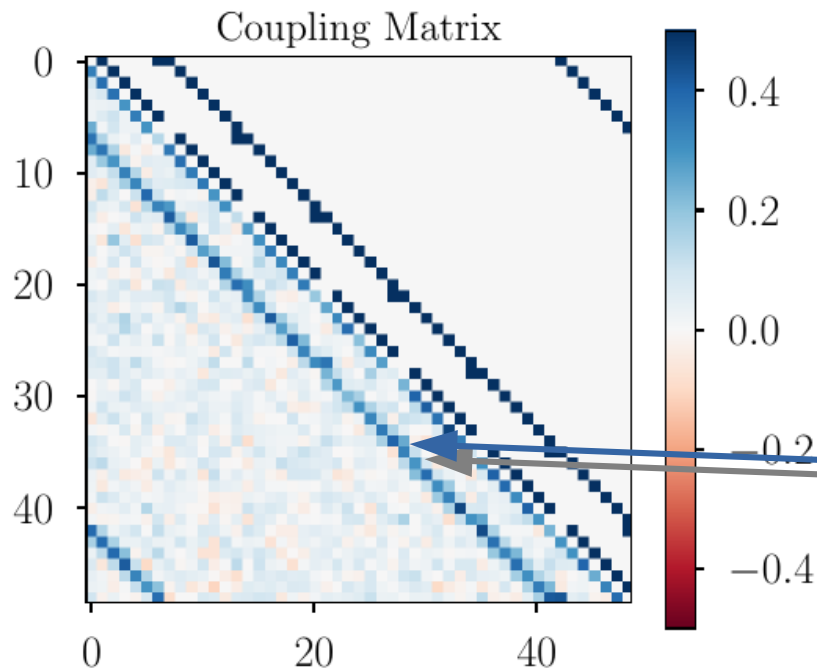


Inverse Ising case : $(J_{\text{learned}} - J_{\text{true}})^2$

Very bad model

$$p_{\theta^*}(x) = \frac{e^{-E_{\theta^*}(x)}}{Z_{\theta^*}} \not\approx p_{\text{data}}(x)$$

Is this OOE energy function useful for interpretability ?

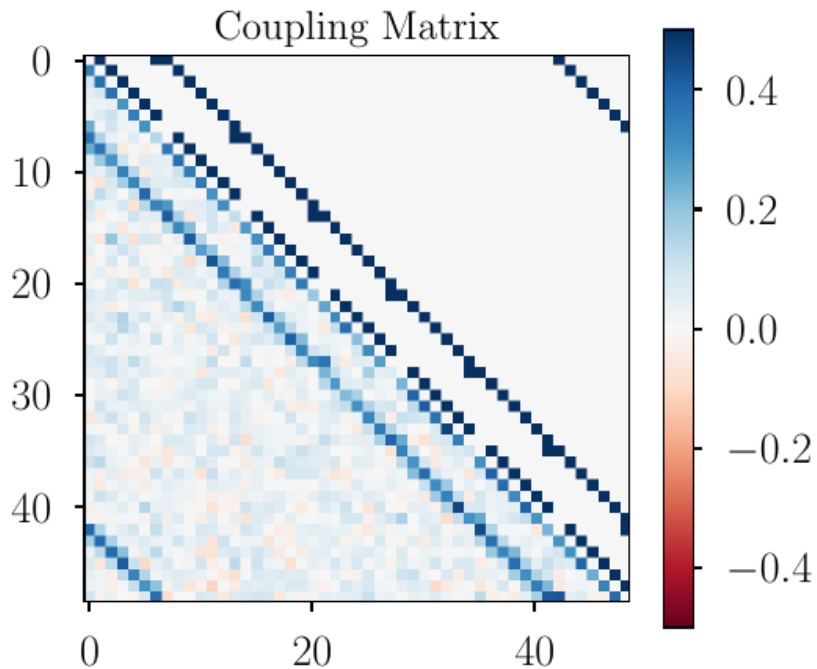


RBM Effective model : $(\mathbf{J}_{\text{learned}} - \mathbf{J}_{\text{true}})^2$

More interactions that what it should !

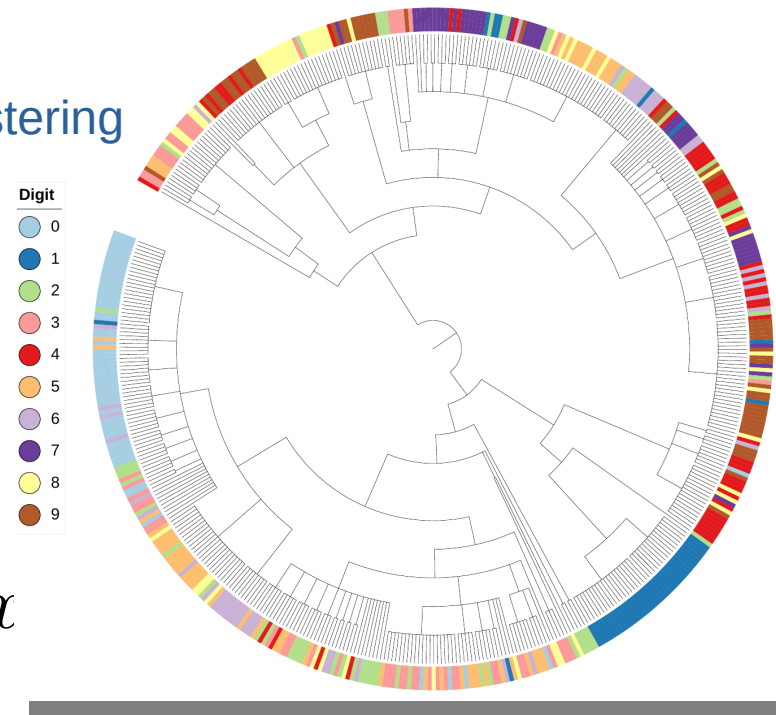
$$p_{\theta^*}(x) = \frac{e^{-E_{\theta^*}(x)}}{Z_{\theta^*}} \not\approx p_{\text{data}}(x)$$

Is this OOE energy function useful for interpretability ?



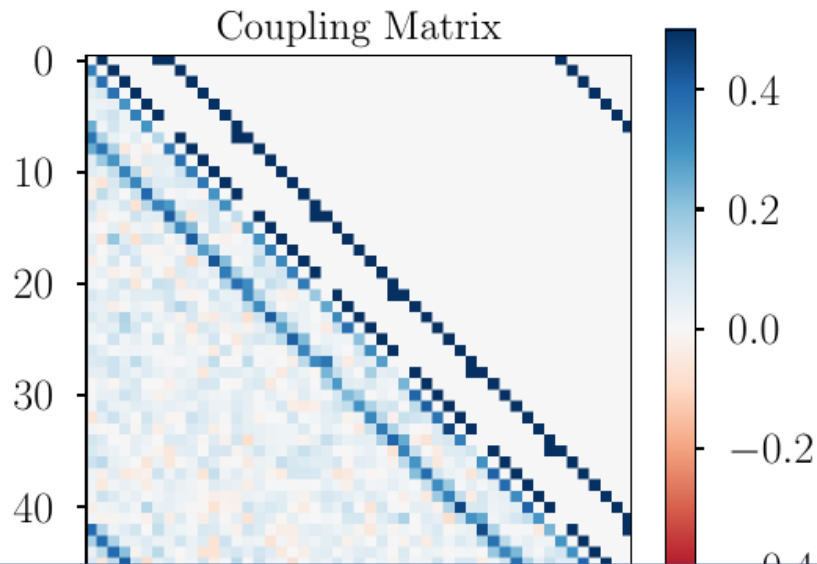
RBM Effective model : $(J_{\text{learned}} - J_{\text{true}})^2$

Poor clustering



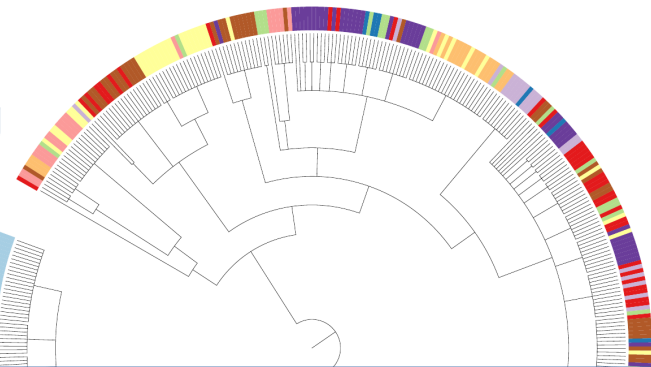
$$p_{\theta^*}(x) = \frac{e^{-E_{\theta^*}(x)}}{Z_{\theta^*}} \not\approx p_{\text{data}}(x)$$

Is this OOE energy function useful for interpretability ?



RBM Effective model : $(\mathbf{J}_{\text{learned}} - \mathbf{J}_{\text{true}})^2$

Poor clustering



If we want to get good models for our data...

we need to equilibrate the chains all along training

Very challenging with multimodal data

The good news is that we know how to do it ! (arriving soon)

