

# Assisting sampling (of physical states) with generative models

MCMC, Generative Models and Overlaps

ASC Summer School: Physics meet AI - September 12-16, 2022

**Marylou Gabrié** (CMAP, École Polytechnique, Paris)

# Manipulating high-dimensional probabilistic models: 1 motivations

## ▷ Statistical mechanics / Chemistry

$$\rho(x) = \frac{1}{Z_\beta} e^{-\beta U(x)}$$

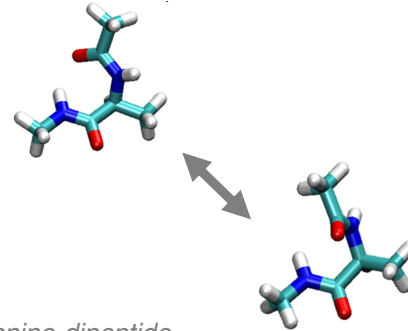
## ▷ Quantum mechanics (wave functions)

## ▷ Bayesian statistical modelling

$$\rho(\theta|D) = \frac{1}{Z_D} \rho(D|\theta)\rho(\theta)$$

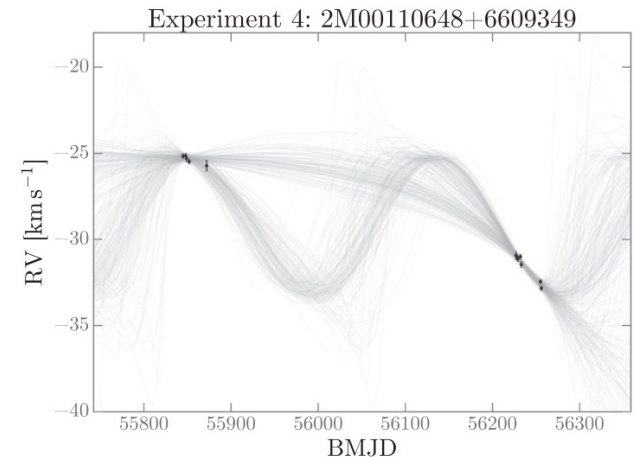
## ▷ Typically known up to normalization constant

ex: molecular configurations



Alanine-dipeptide  
*Jiang et al J. Phys. Chem. B 2019*

ex: Astrophysics data modelling



Price-Whelan et al. *The Astrophysical Journal* 2017

# Now that I know the Boltzmann distribution, what can I do?

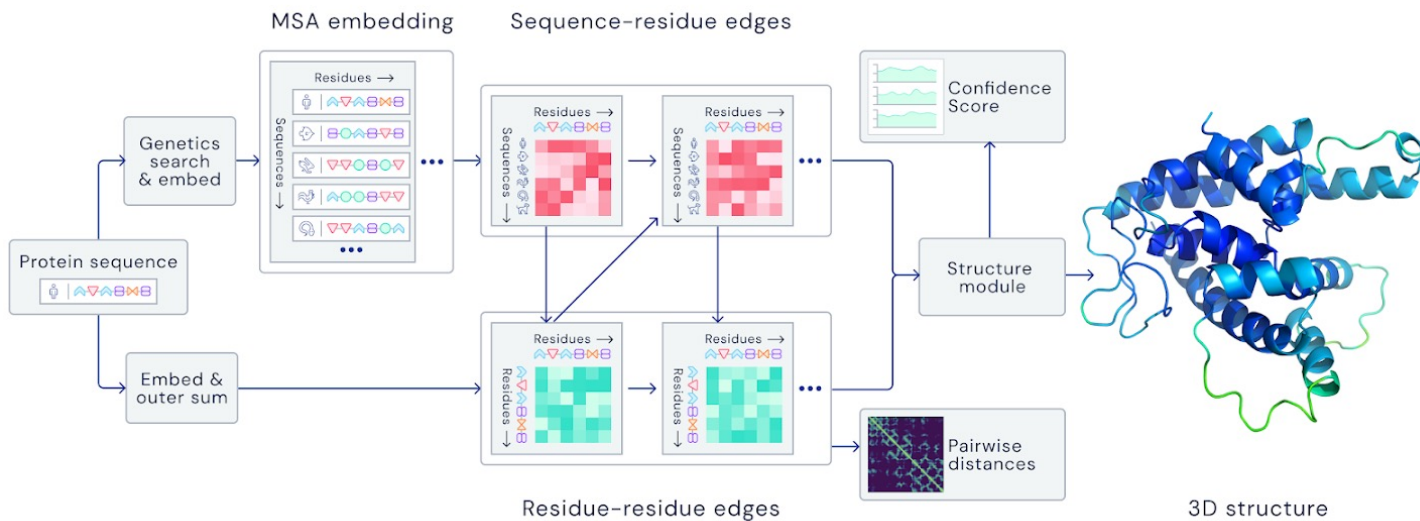
Distribution of physical state:

$$\rho(x) = \frac{1}{Z_\beta} e^{-\beta U(x)} \quad x \in \Omega \subset \mathbb{R}^D$$

e.g. particle positions, field values on a lattice etc ...

▷ Next?

- Look for ground states  $U_0 = U(x_0) = \min_x U(x)$



# Now that I know the Boltzmann distribution, what can I do?

Distribution of physical state:

$$\rho(x) = \frac{1}{Z_\beta} e^{-\beta U(x)} \quad x \in \Omega \subset \mathbb{R}^D$$

e.g. particle positions, field values  
on a lattice etc ...

▷ Next?

- Look for ground states  $U_0 = U(x_0) = \min_x U(x)$
- Compute equilibrium properties

- Approximate the partition function (and derivatives!)  $Z_\beta = \int_{\Omega} e^{-\beta U(x)} dx$  in some models/limits  
e.g. mean field limits
- Sample !

# Monte Carlo Methods

▷ Random variable  $x \in \Omega \subset \mathbb{R}^D$ , and density  $\rho(x) = \frac{1}{Z} e^{-U(x)}$  with unknown  $Z$

▷ Task: Compute expectations  $\mathbb{E}_\rho[f(x)] = \int_{\Omega} f(x)\rho(x)dx$

▷ Method: Monte Carlo approximations, generate  $x_1, \dots, x_N, \dots$

such that 
$$\mathbb{E}_\rho[f(x)] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(x_i)$$

Are we done already ??

In particular if  $x_1, \dots, x_N, \dots$  are i.i.d. draws from  $\rho(x)$

▷ Monte Carlo **Markov Chains** idea to obtain samples:

Design transition kernel  $\pi(x_{t+1}|x_t)$  such that

chain  $x_0, x_1, \dots, x_t =$  samples from  $\rho(x) \propto e^{-U(x)}$  for  $t$  large enough

e.g. Gibbs sampling, Metropolis-Hastings

▷ **Importance** sampling:

Reweight samples from an “easy” distribution  $\mathbb{E}_\rho[f(x)] \approx \frac{1}{N} \sum_{i=1}^N w_i f(x_i)$

## 1. A couple of important sampling methods

1.1 - Importance sampling

1.2 - Metropolis-Hasting

## 2. Unsupervised learning / generative models

2.1 - Latent deep generative models

2.2 - Normalizing flows

## 3. Combining traditional inference method and learning

3.1 - Variational Inference

3.2 - Adaptive algorithms

## 4. Will it scale?

4.1 - Local sampling in reparametrized space

4.2 - Local-global sampling

4.3 - Joining forces with annealing

4.4 - Leveraging physics

# 1.1 Importance Sampling

- ▷ Context:  $\rho_*(x) = \frac{1}{Z} e^{-U(x)}$  with unknown  $Z$
- ▷ Task: Compute expectations  $\mathbb{E}_\rho[f(x)] = \int_\Omega f(x)\rho_*(x)dx$
- ▷ Importance sampling

- Samples from proposal distribution  $x_i \sim \rho_p(x_i)$

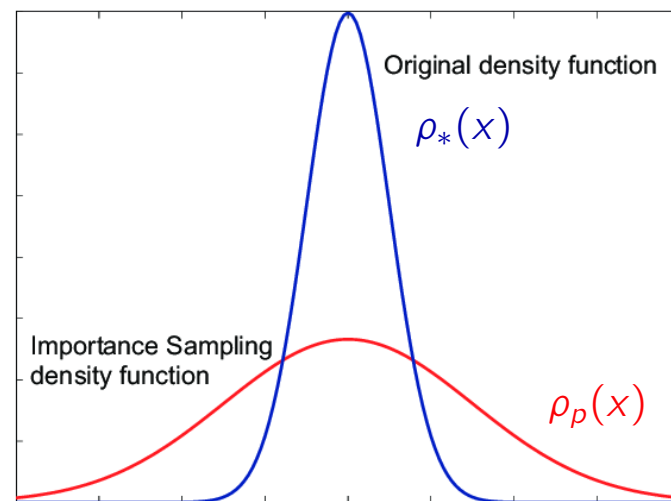
e.g. Gaussian, factorized, ...

- Self-normalized weights  $w_i = \frac{e^{-U(x_i)} / \rho_p(x_i)}{\sum_{i=1}^N e^{-U(x_i)} / \rho_p(x_i)}$

- Compute  $\mathbb{E}_{\rho_*}[f(x)] \approx \frac{1}{N} \sum_{i=1}^N w_i f(x_i)$

- Asymptotically “unbiased”  $N \rightarrow \infty$

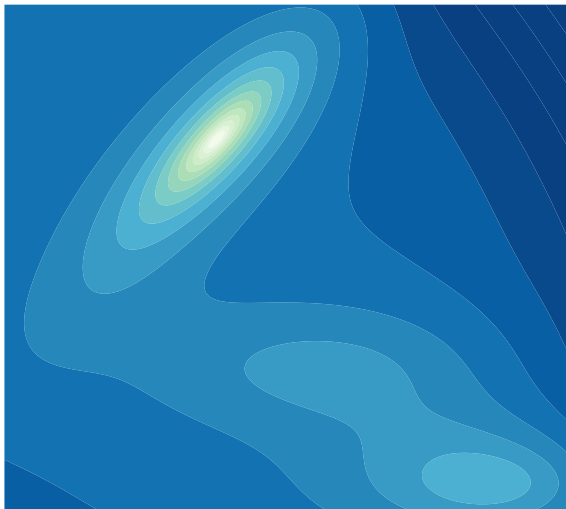
$$\mathbb{E}_\rho[f(x)] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N w_i f(x_i)$$



# 1.1 Importance Sampling: did it work?

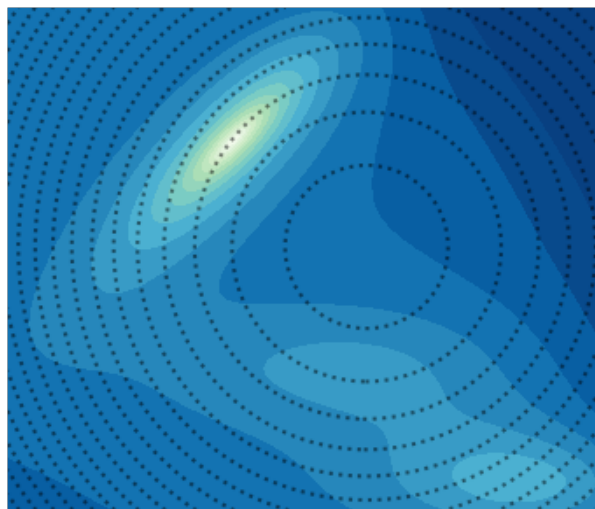
▷ Look at a concrete example: 2d Muller Brown potential

$$\rho_*(x) = e^{-U_*(x)} / Z$$



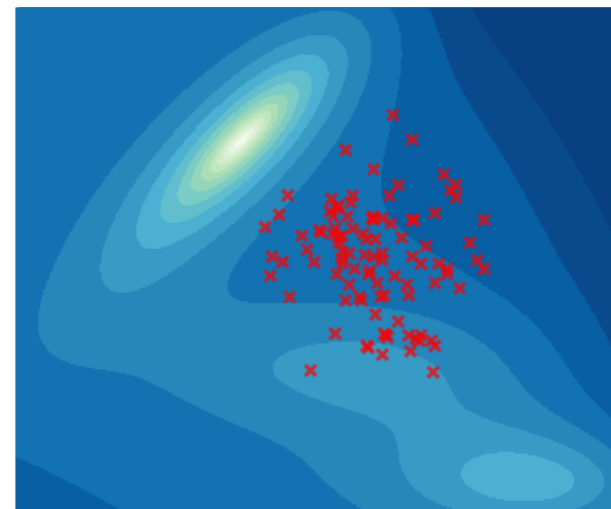
$$\rho_p(x)$$

Variance of proposal: 5.00e-02



$$x_i \sim \rho_p(x_i)$$

~~ESS: 1004~~



▷ What can go wrong?

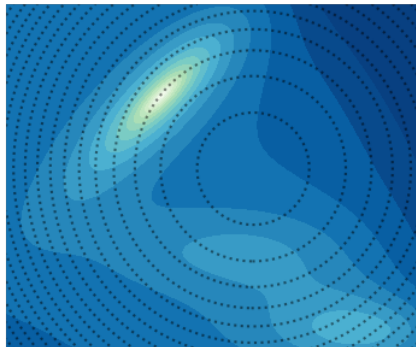
$$w_i = \frac{e^{-U(x_i)} / \rho_p(x_i)}{\sum_{i=1}^N e^{-U(x_i)} / \rho_p(x_i)}$$



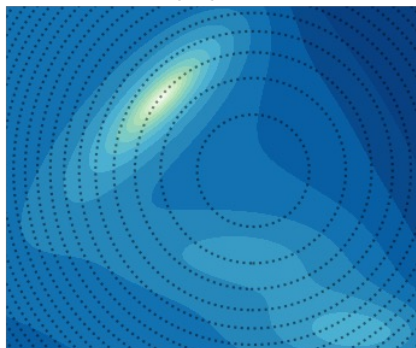
# 1.1 Importance Sampling: did it work?

$$\rho_p(x)$$

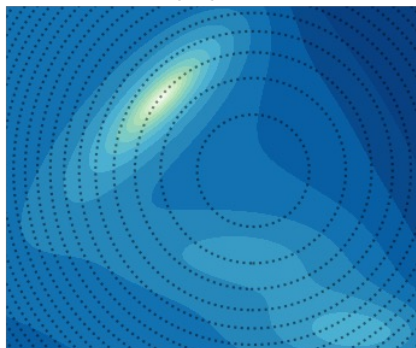
Variance of proposal: 5.00e-02



Variance of proposal: 5.00e-02

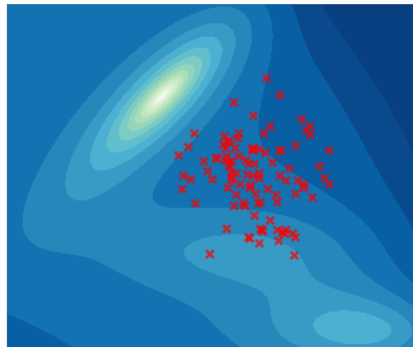


Variance of proposal: 5.00e-02

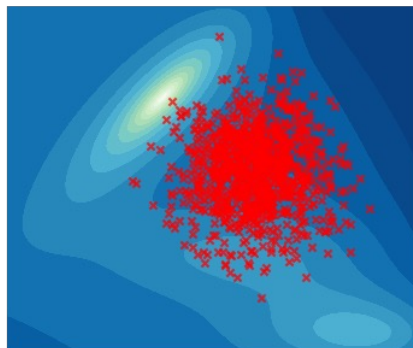


$$x_i \sim \rho_p(x_i)$$

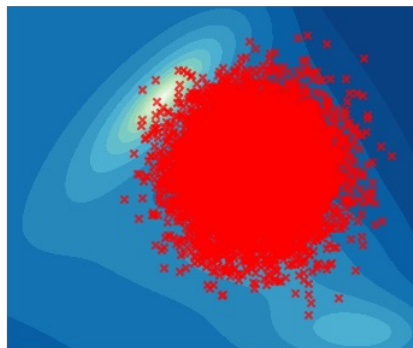
N = 100



N = 1000

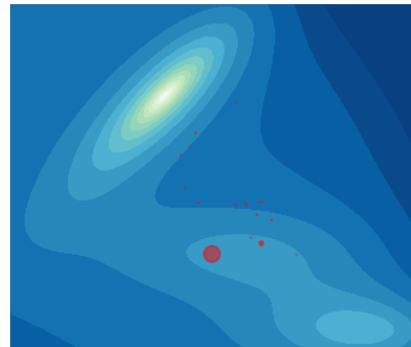


N = 10000

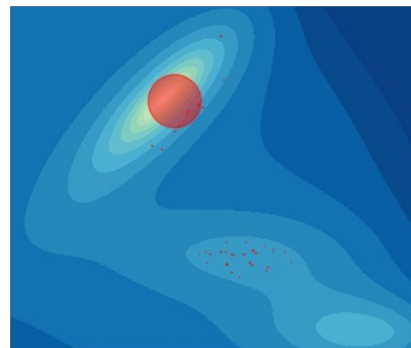


$$w_i = \frac{e^{-U(x_i)} / \rho_p(x_i)}{\sum_{i=1}^N e^{-U(x_i)} / \rho_p(x_i)}$$

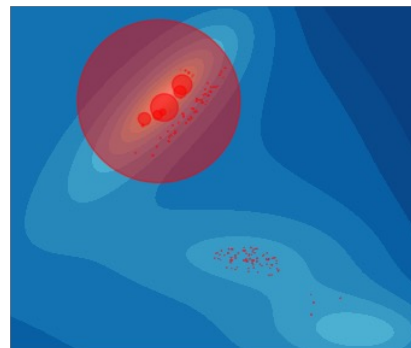
ESS: 1.14



ESS: 1.00

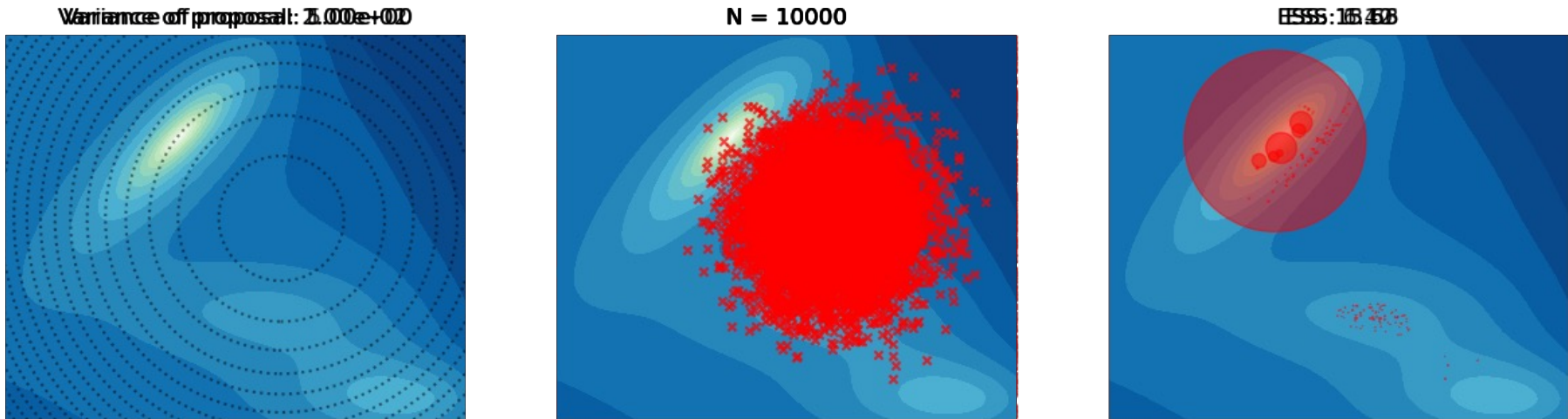


ESS: 1.12



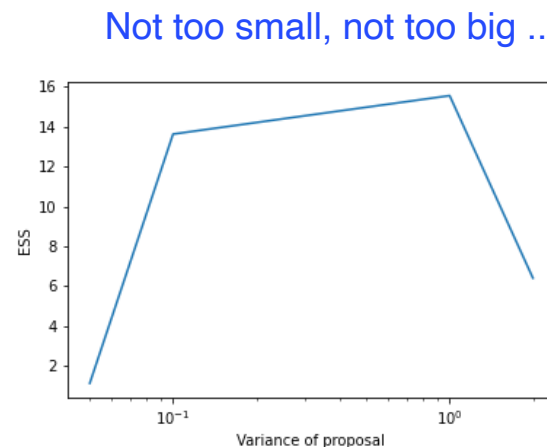
# 1.1 Importance Sampling: did it work?

## ▷ Changing the variance now



## ▷ How will I realize this is happening?

- Importance weights  $w_i = \frac{\rho_*(x_i)/\rho_p(x_i)}{\sum_{i=1}^N \rho_*(x_i)/\rho_p(x_i)}$
- Effective sample size  $ESS = \frac{(\sum_{i=1}^N w_i)^2}{\sum_{i=1}^N w_i^2}$
- All samples “participates”  $w_i = 1/N \Rightarrow ESS = N$
- Only one “participates”  $w_1 = 1, \Rightarrow ESS = 1$



**Proposal and target need to be well adapted!**

▷ Idea: design transition kernel  $\pi(x_{t+1}|x_t)$  such that chain  $x_0, x_1, \dots, x_t$  produces samples from  $\rho_*$  for  $t$  large enough

▷ Important example:

## Metropolis-Hastings sampler

Initialize:  $x_0$

Iterate:

○ Propose  $x_{t+1} \sim \rho_p(x_{t+1}|x_t)$

○ Accept/Reject with prob.

$$\text{acc}(x_{t+1}|x_t) = \min \left[ 1, \frac{\rho_*(x_{t+1})\rho_p(x_t|x_{t+1})}{\rho_*(x_t)\rho_p(x_{t+1}|x_t)} \right]$$

○ If reject stay  $x_{t+1} = x_t$

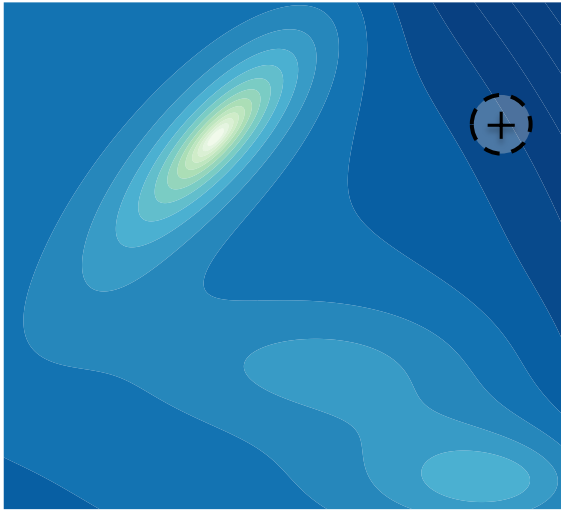
# Examples of Metropolis-Hastings MCMC

11

▷ Gaussian random walk  $\rho_p(x_{t+1}|x_t) = \mathcal{N}(x_t, \Sigma)$

e.g. 2d Müller-Brown potential

$$\rho_*(x) = e^{-U_*(x)} / Z$$



T = 100 steps

▷ (Metropolis Adjusted) Langevin algorithm (MALA)

$$\rho_p(x_{t+1}|x_t) = \mathcal{N}(x_t - dt\nabla U(x), \sqrt{2dt}I_d)$$

Metropolis-Hastings sampler

Initialize:  $x_0$

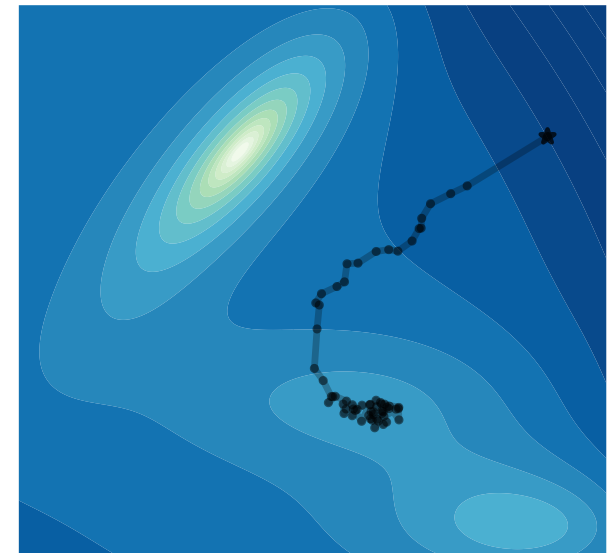
Iterate:

○ Propose  $x_{t+1} \sim \rho_p(x_{t+1}|x_t)$

○ Accept/Reject with prob.

$$\text{acc}(x_{t+1}|x_t) = \min \left[ 1, \frac{\rho_*(x_{t+1})\rho_p(x_t|x_{t+1})}{\rho_*(x_t)\rho_p(x_{t+1}|x_t)} \right]$$

○ If reject stay  $x_{t+1} = x_t$



T = 100 steps

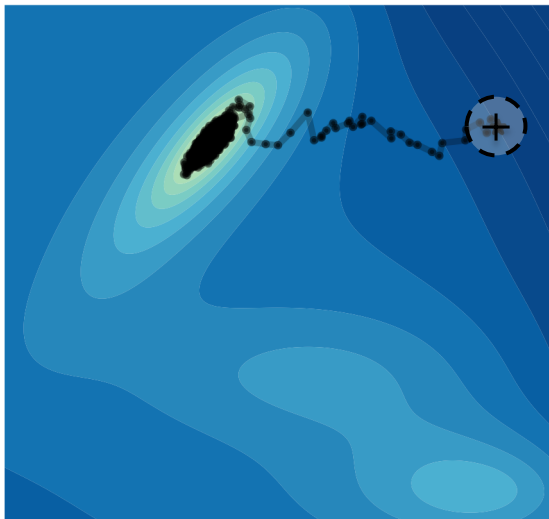
# Challenge: Decorrelation and convergence

12

▷ Gaussian random walk  $\rho_p(x_{t+1}|x_t) = \mathcal{N}(x_t, \Sigma)$

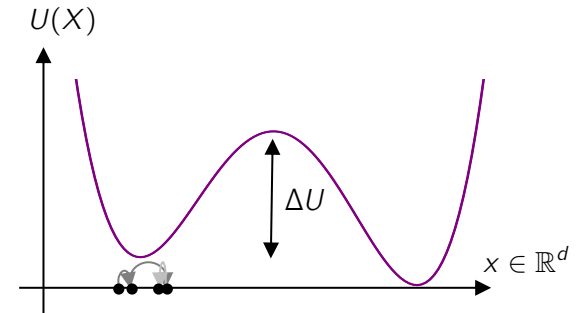
e.g. 2d Müller-Brown potential

$$\rho_*(x) = e^{-U_*(x)} / Z$$

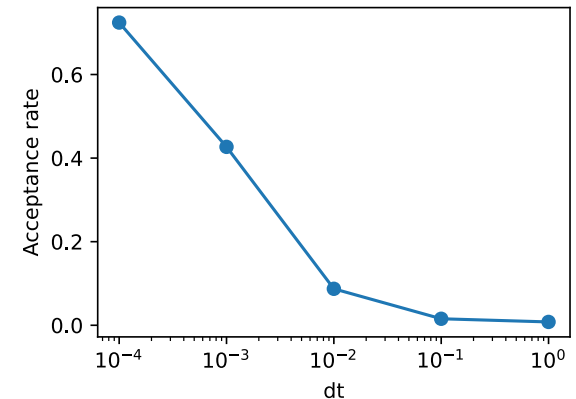


T = 10000 steps

$$\rho_*(x) = e^{-U(x)} / Z$$



▷ Trade-off size local moves / acceptance



▷ Many many proposition for faster “mixing”

○ Use gradient information: Langevin dynamics, Hamiltonian MC

still difficult to switch mode!

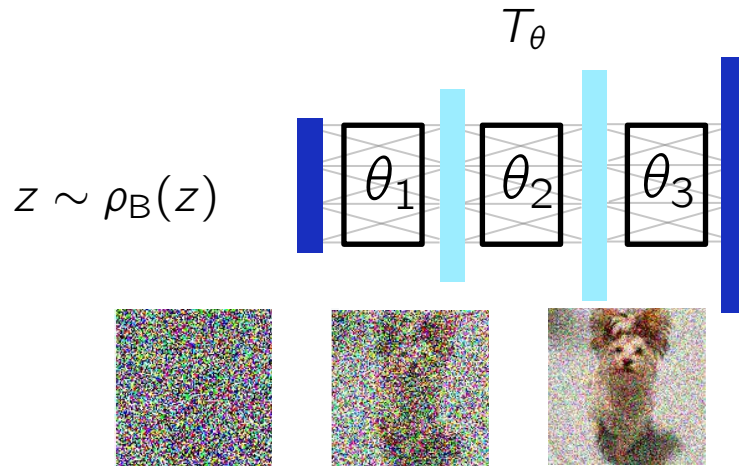
○ Gradually approach the target: Sequential Monte Carlo, Annealed Importance Sampling

powerful but computationally heavy!

1. A couple of important sampling methods
  - 1.1 - Importance sampling
  - 1.2 - Metropolis-Hasting
2. Unsupervised learning / generative models
  - 2.1 - Latent deep generative models
  - 2.2 - Normalizing flows
3. Combining traditional inference method and learning
  - 3.1 - Variational Inference
  - 3.2 - Adaptive algorithms
4. Will it scale ?
  - 4.1 - Local sampling in reparametrized space
  - 4.2 - Local-global sampling
  - 4.3 - Joining forces with annealing
  - 4.4 - Leveraging physics

## 2.1 Deep generative models

- ▷ Use transformation  $T_\theta$  (deep neural network) from simple base distribution  $\rho_B$  :



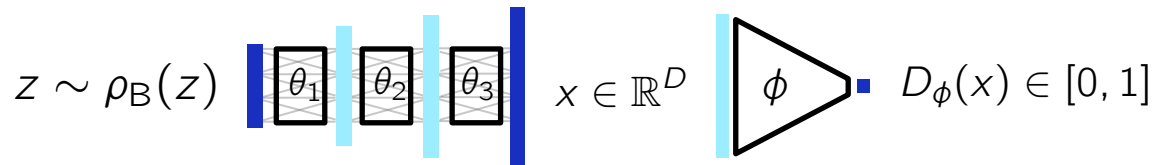
[“GANs” Goodfellow et al. *NeurIPS* 2014,  
“VAEs” Kingma & Welling *ICLR* 2014,  
“Normalizing flows” Papamakarios et al. *JMLR* 2021]

$x = T_\theta(z) \sim \rho_\theta(x)$  “push-forward”  
distribution

Song et al. *ICLR* 2021

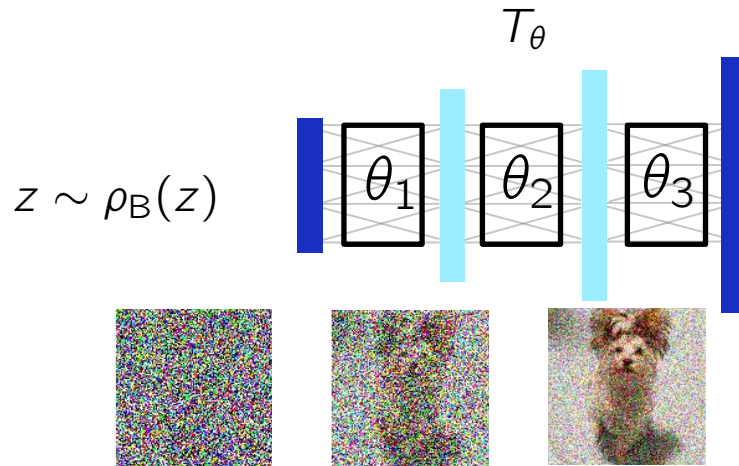
- ▷ Two main training methods of unsupervised learning:

- Maximum likelihood:  $L[\rho_\theta] = -\sum_{i=1}^N \log \rho_\theta(x_i)$  with  $x_i$  data samples + SGD!
- Adversarial training:  $\min_{\theta} \max_{\phi} [\mathbb{E}_{\rho_D} [\ln D_\phi(x)] + \mathbb{E}_{\rho_B} [\ln(1 - D_\phi(T_\theta(z)))]]$  with  $\rho_D$  data distribution



## 2.1 Deep generative models

- ▷ Use transformation  $T_\theta$  (deep neural network) from simple base distribution  $\rho_B$  :



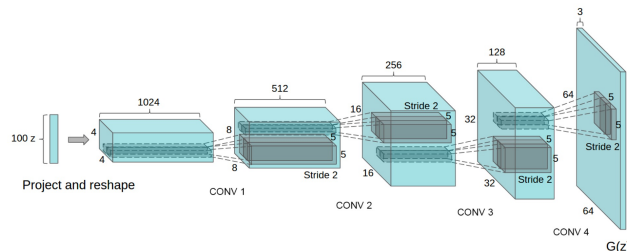
[“GANs” Goodfellow et al. *NeurIPS* 2014,  
 “VAEs” Kingma & Welling *ICLR* 2014,  
 “Normalizing flows” Papamakarios et al. *JMLR* 2021]

$$x = T_\theta(z) \sim \rho_\theta(x) \quad \text{“push-forward” distribution}$$

Song et al. *ICLR* 2021

- ▷ Two main training methods of unsupervised learning:

- Maximum likelihood:  $L[\rho_\theta] = -\sum_{i=1}^N \log \rho_\theta(x_i)$  with  $x_i$  data samples
- Adversarial training:  $\min_{\theta} \max_{\phi} [\mathbb{E}_{\rho_D} [\ln D_\phi(x)] + \mathbb{E}_{\rho_B} [\ln(1 - D_\phi(T_\theta(z)))]]$  with  $\rho_D$  data distribution

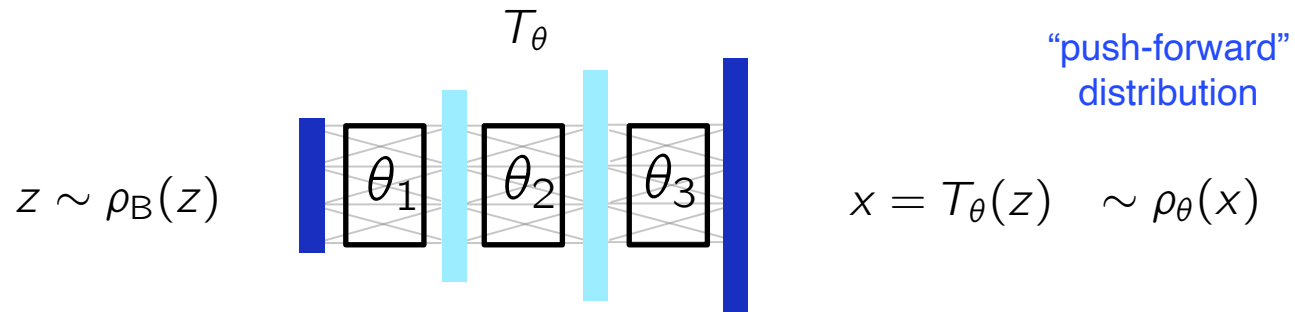


[Radford et al *ICLR* 2016; Karras et al *CVPR* 2019]



# Nota Bene: Intractability of the push-forward of many latent generative models

- ▷ In general latent dimension much smaller than data dimension



- Push-forward computation involves marginalization ...

$$\rho_\theta(x)dx = \int_{\mathbb{R}^d} dz \rho_B(z) \delta(T_\theta(z) - x)$$

- ▷ Hence difficult to do maximum likelihood:  
e.g. optimize ELBLO (evidence lower bound in VAE)

# 2.2 A special type of Deep Generative Models

## Normalizing Flows (NF): Invertible networks

▷ Parametrized invertible map  $T_\theta: \Omega \mapsto \Omega \quad \Omega \subset \mathbb{R}^d$

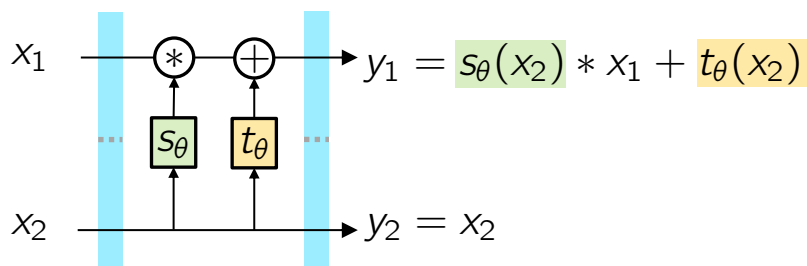
Most generative model are not invertible!  
Intractable push-forward.

○ Base distribution  $z \sim \rho_B(z)$

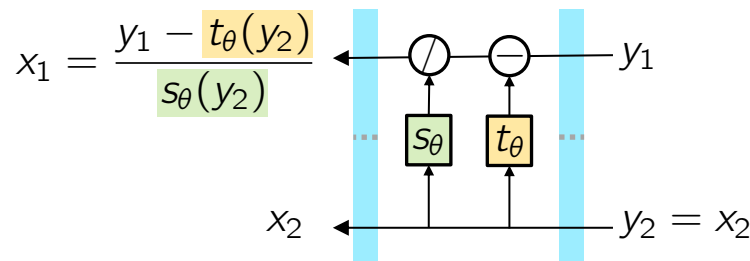
○ Push-forward distribution  $x = T_\theta(z) \sim \rho_\theta(x) = \rho_B(T_\theta^{-1}(x)) \det |\nabla_x T_\theta^{-1}|$

▷ e.g. “Coupling layers”: easy-to-compute inverse and Jacobian

Affine coupling layer  $T_\theta(x)$



Inverse layer  $T_\theta^{-1}(y)$



Block diagonal Jacobian:

$$\nabla_x T_\theta(x) = \begin{bmatrix} s_\theta(x_2) I_{d/2} & 0 \\ 0 & I_{d/2} \end{bmatrix}$$

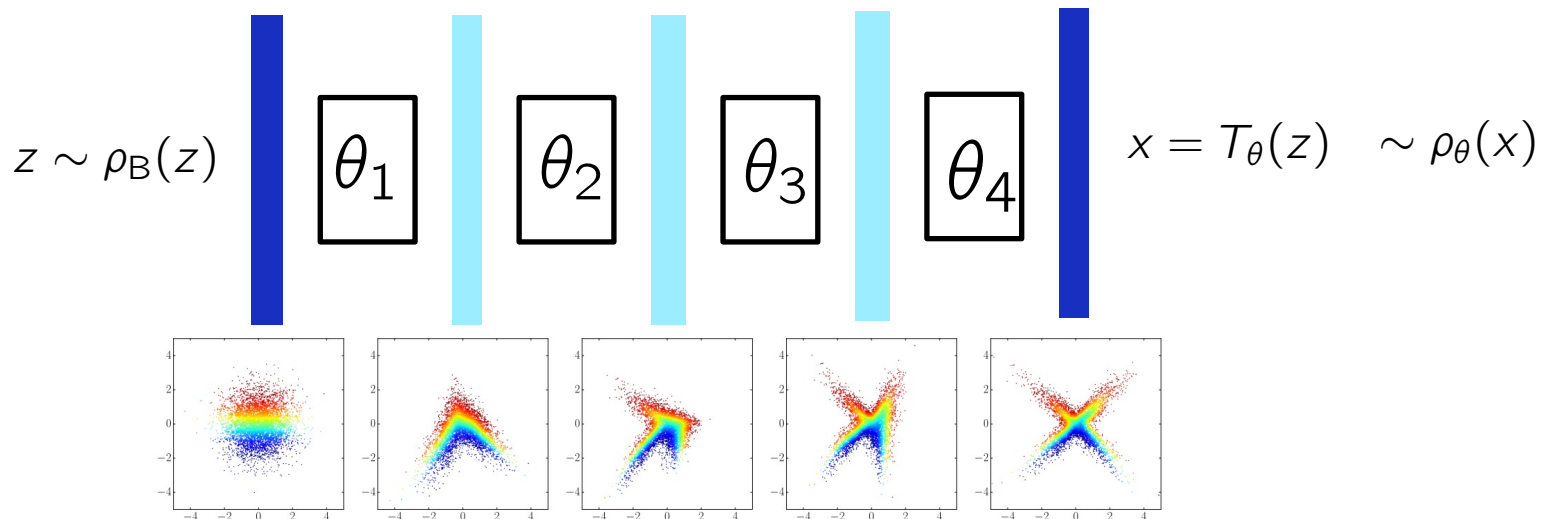
## 2.2 A special type of Deep Generative Models

### Normalizing Flows (NF): Invertible networks

- ▷ Parametrized invertible map  $T_\theta: \Omega \mapsto \Omega \quad \Omega \subset \mathbb{R}^d$ 
  - Base distribution  $z \sim \rho_B(z)$
  - Push-forward distribution  $x = T_\theta(z) \sim \rho_\theta(x) = \rho_B(T_\theta^{-1}(x)) \det |\nabla_x T_\theta^{-1}|$

- ▷ Composition to encode for sophisticated transformations

$$T_\theta = T_{\theta_4} \circ T_{\theta_3} \circ T_{\theta_2} \circ T_{\theta_1}$$

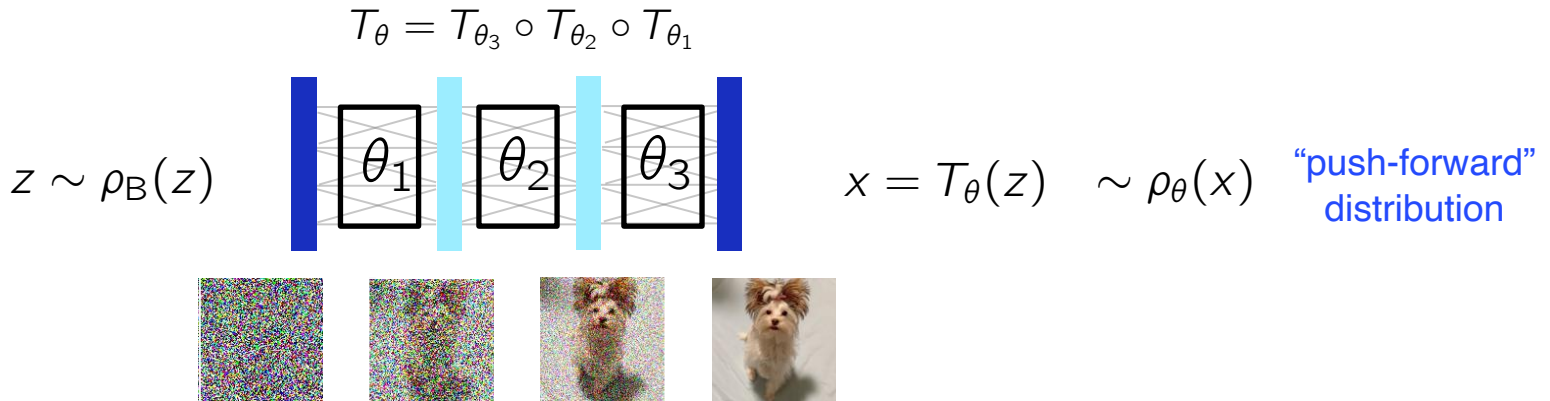


**Easy to sample and easy to evaluate density**

NADE networks also!

# Deep generative models for sampling target $\rho_*(x)$ ? 19

▷ Parametric model: Simple base random variable transformed by a deep neural network  $T_\theta$



▷ Opportunity alert! Sample complicated  $\rho_*(x)$  by modelling it with deep generative model?

- Need to learn  $T_\theta$  for which we need data -  $x_i \sim \rho_*(x)$  - do we? a.k.a. chicken-and-egg problem!
- Even if we get  $\rho_\theta(x) \approx \rho_*(x)$ , unlikely to learn perfect model  $\rho_\theta(x) = \rho_*(x)$ , right?

## 1. A couple of important sampling methods

1.1 - Importance sampling

1.2 - Metropolis-Hasting

## 2. Unsupervised learning / generative models

2.1 - Latent deep generative models

2.2 - Normalizing flows

## 3. Combining traditional inference method and learning

3.1 - Variational Inference

3.2 - Adaptive algorithms

## 4. Will it scale?

4.1 - Local sampling in reparametrized space

4.2 - Local-global sampling

4.3 - Joining forces with annealing

4.4 - Leveraging physics

# 3.1. Variational Inference or do we really need data? <sup>21</sup>

▷ Context: known  $\rho_*(x) = \frac{1}{Z} e^{-U(x)}$  up to normalizing factor  $Z$

▷ Task: Compute expectations  $\mathbb{E}_\rho[f(x)] = \int_{\Omega} f(x)\rho(x)dx$

▷ Variational inference (VI):

- Optimize surrogate tractable distribution to minimize Kullback-Leibler divergence

$$D_{\text{KL}}(\rho_\theta \parallel \rho_*) = \int \log \frac{\rho_\theta(x)}{\rho_*(x)} \rho_\theta(x) dx \quad \implies \quad L[\rho_\theta] = - \sum_{i=1}^N \log \frac{\rho_\theta(x_i)}{\rho_*(x_i)} \quad x_i \sim \rho_\theta(x)$$

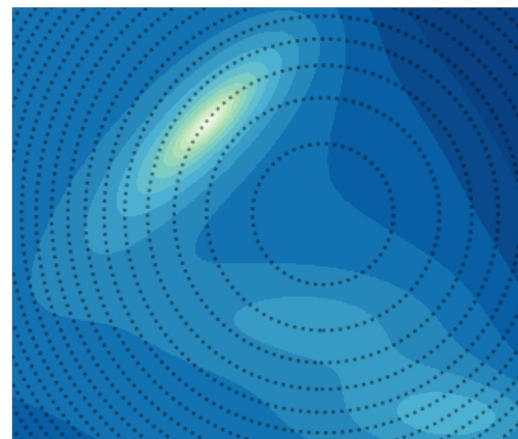
- Then use the proxy for all purpose  $\rho_\theta(x) \approx \rho_*(x)$

Quickly need expressive proxy!

▷ Questions:

1. Example of suitable  $\rho_\theta(x)$ ?
2. Which problems do you anticipate?

1. Factorized/mean-field, Gaussian ...
2. Quality of the approximation?



## ▷ Training without data?

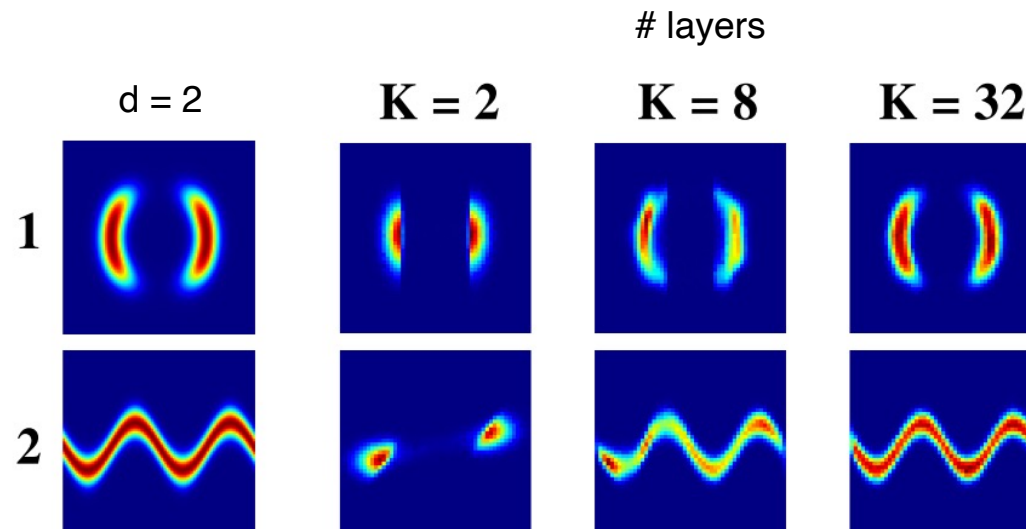
- minimize Kullback-Leibler  $D_{\text{KL}}(\rho_\theta \| \rho_*) =$  variational principle with expressive  $\rho_\theta(x)$  ansatz

$$D_{\text{KL}}(\rho_\theta \| \rho_*) = \int \log \frac{\rho_\theta(x)}{\rho_*(x)} \rho_\theta(x) dx \approx \sum_{i=1}^N \log \frac{\rho_\theta(x_i)}{\rho_*(x_i)} \quad x_i \sim \rho_\theta(x) \quad \text{easy to obtain!}$$

$$\rho_\theta(x) = \rho_B(T_\theta^{-1}(x)) \det |\nabla_x T_\theta^{-1}| \quad \text{explicit!}$$

D. Rezende: “good entry point for ML”

## ▷ First results: quality as a function of expressivity



Rezende & Mohamed, (2015). Variational inference with normalizing flows,

Albergo et al (2019). Flow-based generative models for Markov chain Monte Carlo in lattice field theory. *PRD* 2019.

Wu et al. (2019). Solving Statistical Mechanics Using Variational Autoregressive Networks.

# Correcting the samples with a MCMC

Target density:  $\rho_*(x) = e^{-U_*(x)} / Z$

Generative model parametrized density  $\rho_\theta(x)$  trained by Variational inference

▷ Algorithm: Metropolis-Hastings with generative model proposal

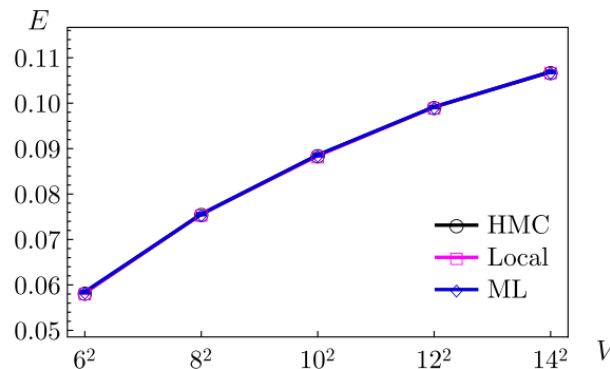
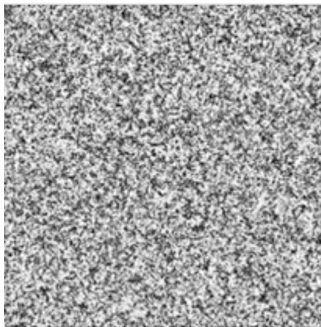
Initialize:  $x_0^i \quad i = 1 \dots N$

Loop:

○ Draw from generative model  $x_{t+1}^i \sim \rho_\theta(x)$  **NF proposal!**

○ Accept-reject  $\text{acc}(x_{t+1}^i | x_t^i) = \min \left[ 1, \frac{\rho_*(x_{t+1}^i) \rho_\theta(x_t^i)}{\rho_*(x_t^i) \rho_\theta(x_{t+1}^i)} \right]$

▷  $\phi^4$  model at  $T > T_c$



*"For simplicity in this initial work, all parameters were chosen to lie in the **symmetric phase**. In principle, the flow-based MCMC algorithm can be applied with identical methods to the broken-symmetry phase of the theory, but it **remains to be shown that models can be trained** for such choices of parameters."*



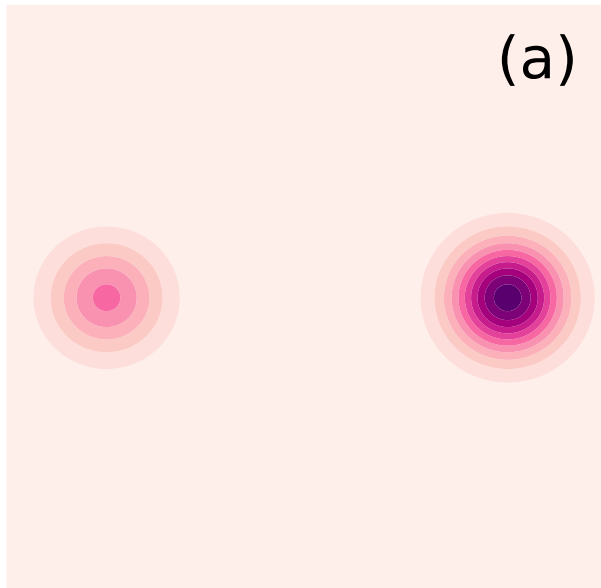
▷ No need for data?

- minimize Kullback-Leibler  $D_{\text{KL}}(\rho_\theta \| \rho_*) =$  variational principle with expressive  $\rho_\theta(x)$  ansatz

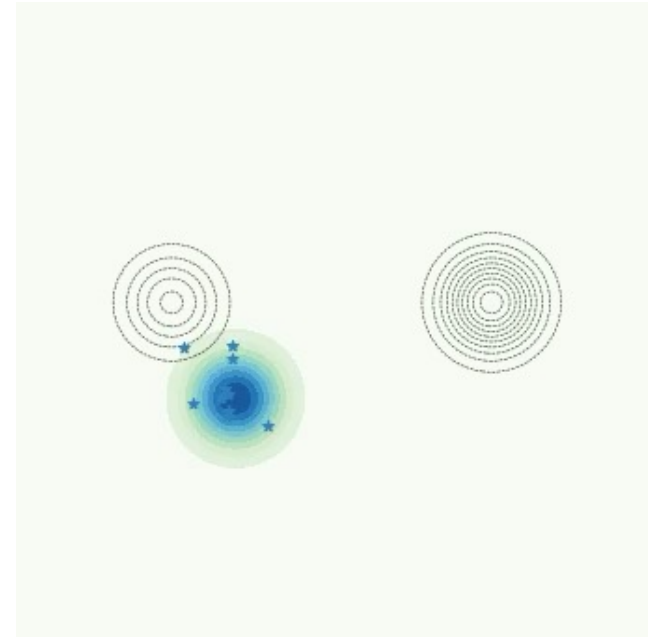
$$D_{\text{KL}}(\rho_\theta \| \rho_*) = \int \log \frac{\rho_\theta(x)}{\rho_*(x)} \rho_\theta(x) dx \approx \sum_{i=1}^N \log \frac{\rho_\theta(x_i)}{\rho_*(x_i)} \quad x_i \sim \rho_\theta(x) \quad \text{easy to obtain!}$$

▷ What can do wrong?

*example:*  
 $\rho_*(x)$  mixture of 2 Gaussians (2d)



$$\rho_{\theta_t}(x) = \rho_B(T_{\theta_t}^{-1}(x)) \det |\nabla_x T_{\theta_t}^{-1}|$$



prone to mode collapse !

# 3.1 Combining VI with a little dataset

“Boltzmann generator” Noé et al. (Science 2019)

▷ Minimize combined loss  $L[\rho_\theta] = L_{VI}[\rho_\theta] + L_{ML}[\rho_\theta]$

○ Training “by energy” (= Variational Inference)

$$L_{VI}[\rho_\theta] = - \sum_{i=1}^N \log \frac{\rho_\theta(x_i)}{\rho_*(x_i)} \quad x_i \sim \rho_\theta(x)$$

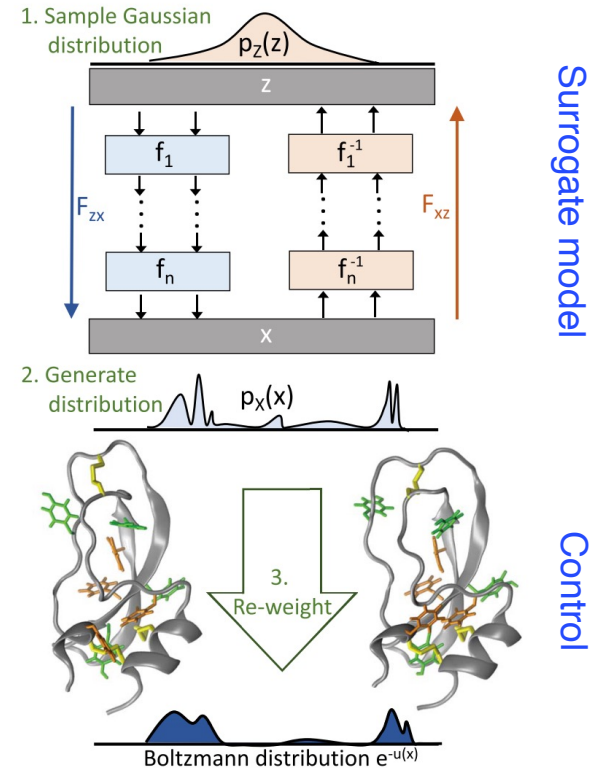
○ Training by data (= maximum likelihood (ML))

$$L_{ML}[\rho_\theta] = - \sum_{i=1}^N \log \rho_\theta(x_{d,i}) \quad x_{d,i} \text{ small data set (from MD)}$$

▷ Correct samples  $x_i \sim \rho_\theta(x)$  with importance weights

$$w_i = \frac{\rho_*(x_i)/\rho_\theta(x_i)}{\sum_{i=1}^N \rho_*(x_i)/\rho_\theta(x_i)}$$

$$\mathbb{E}_{\rho_*}[f(x)] \approx \frac{1}{N} \sum_{i=1}^N w_i f(x_i)$$



e.g. BPTI protein (58 amino acids)

How much data is enough data?

## 3.2 Adaptive MCMC with normalizing flow

Target density:  $\rho_*(x) = e^{-U_*(x)} / Z$

Generative model parametrized density:  $\rho_\theta(x)$

Another name could be:  
active learning!

### ▷ Algorithm: Metropolis-Hastings with generative model proposal

Initialize:  $x_0^i \quad i = 1 \dots N$

Loop:

Loop over parallel chains:  $i = 1 \dots N$

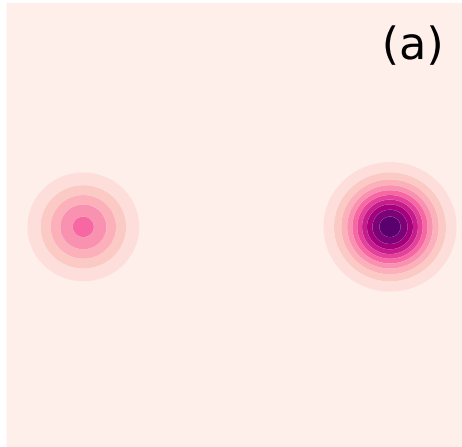
- Draw from generative model  $x_{t+1}^i \sim \rho_\theta(x)$
- Accept-reject  $\text{acc}(x_{t+1}^i | x_t^i) = \min \left[ 1, \frac{\rho_*(x_{t+1}^i) \rho_\theta(x_t^i)}{\rho_*(x_t^i) \rho_\theta(x_{t+1}^i)} \right]$
- Local resampling  $x_{t+1}^i \sim \pi_{\text{local}}(x_{t+1}^i | x_t^i)$
- Update NF paramters  $\theta \leftarrow \theta + \eta \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log \rho_\theta(x_{t+1}^i)$

Metropolis-Hastings  
with NF

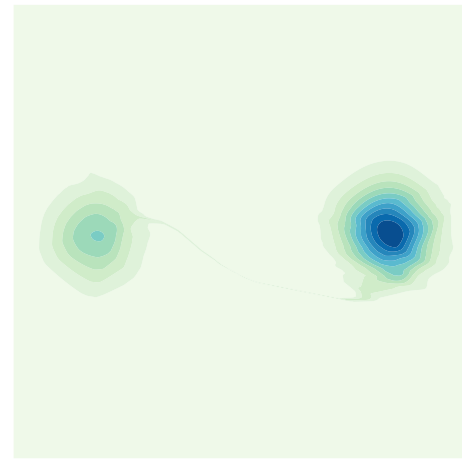
Maximum  
likelihood GD

# 3.4 Adaptive MCMC – 2d Mixture of two Gaussians 27

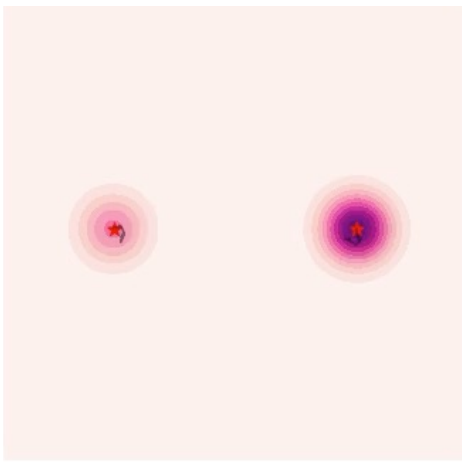
Target density:



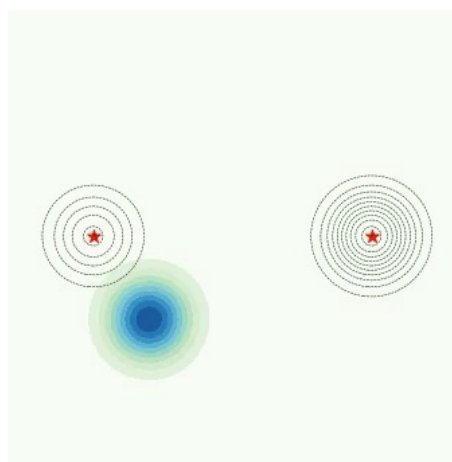
Final learned density:



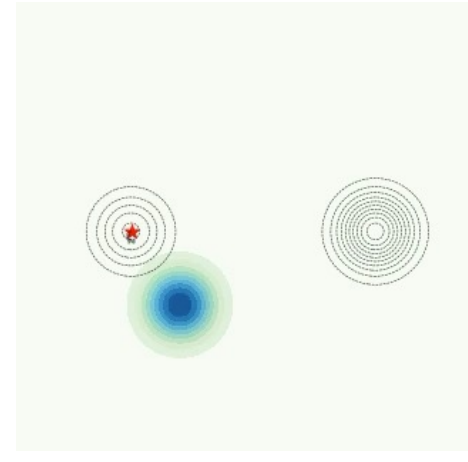
Local method only:



Concurrent:  
*careful initialization*

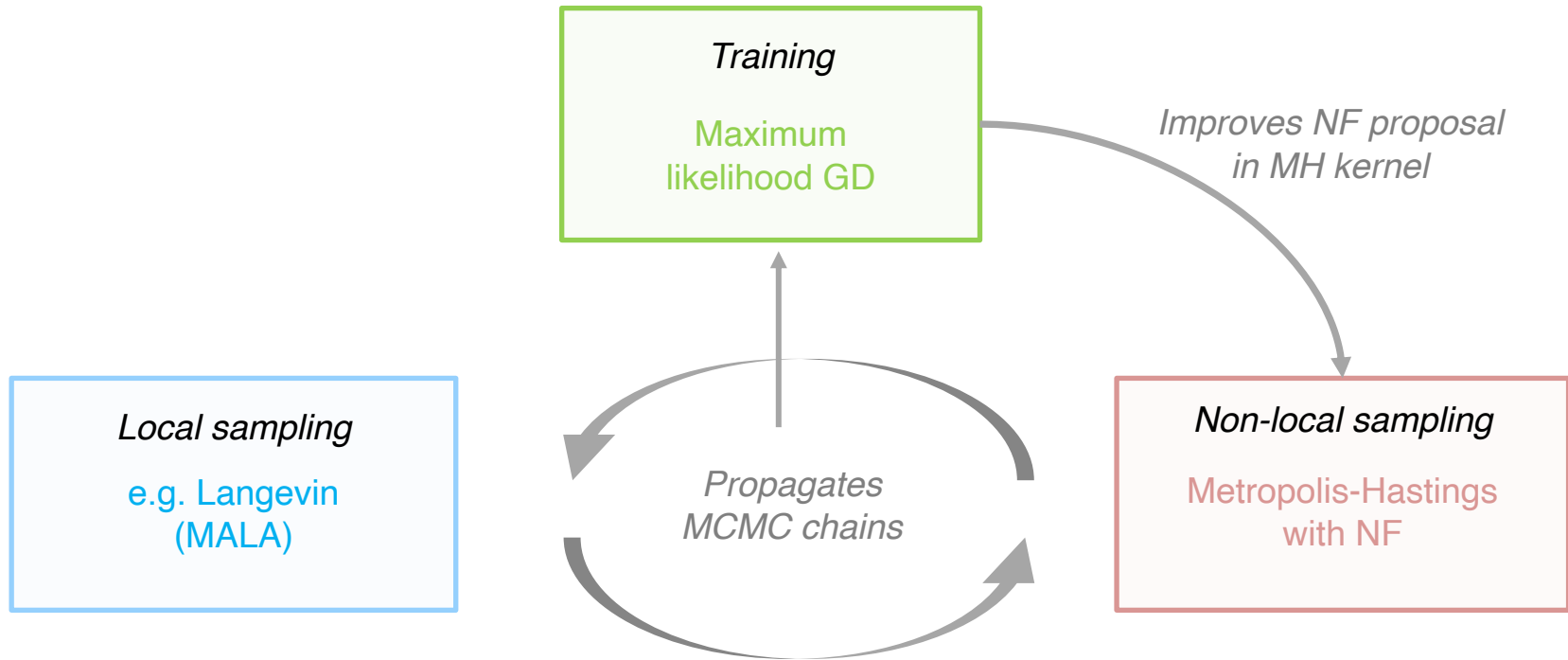


Concurrent:  
*starting with one walker*



No mode discovery!

# 3.4 Adaptive MCMC with normalizing flow



- Adaptive / “non-linear” Monte Carlo [Haario et al Bernoulli 2001, Jasra et al Statistics and Computing, 2007, Andrieu et al Bernoulli 2011, Sejdinovic et al ICML 2014, Parno & Marzouk 2018, Naesseth et al. Neurips 2020, Gabrié et al. PNAS 2022, ...]

○ Softwares:      pytorch:  
[marylou-gabrie / flonaco](#) Public

jax:  
[kazewong / NFSampler](#) Public

with Kaze Wong (Flatiron Institute)

## 1. A couple of important sampling methods

1.1 - Importance sampling

1.2 - Metropolis-Hasting

## 2. Unsupervised learning / generative models

2.1 - Latent deep generative models

2.2 - Normalizing flows

## 3. Combining traditional inference method and learning

3.1 - Variational Inference

3.2 - Adaptive algorithms

*“Looks like, we can find training procedures for flows in this context and speed up sampling”*

## 4. Will it scale?

4.1 - Local sampling in reparametrized space

4.2 - Local-global sampling

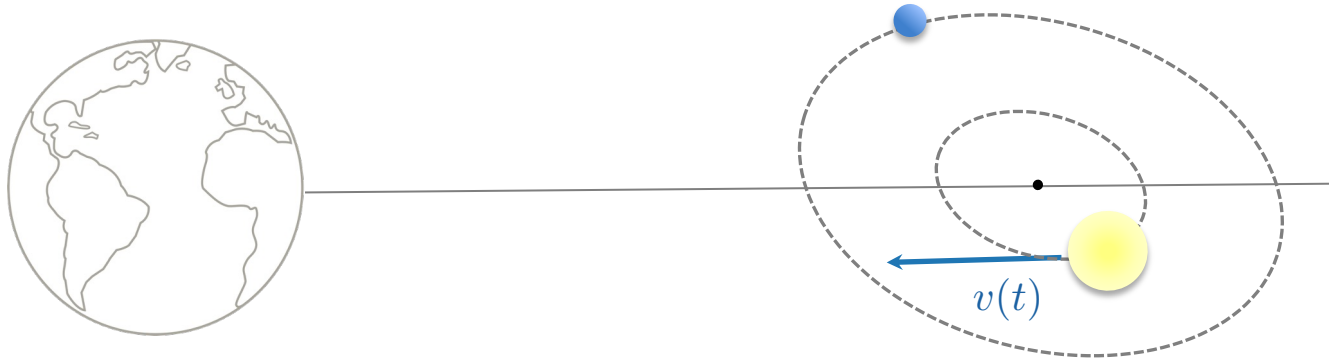
4.3 - Joining forces with annealing

4.4 - Leveraging physics

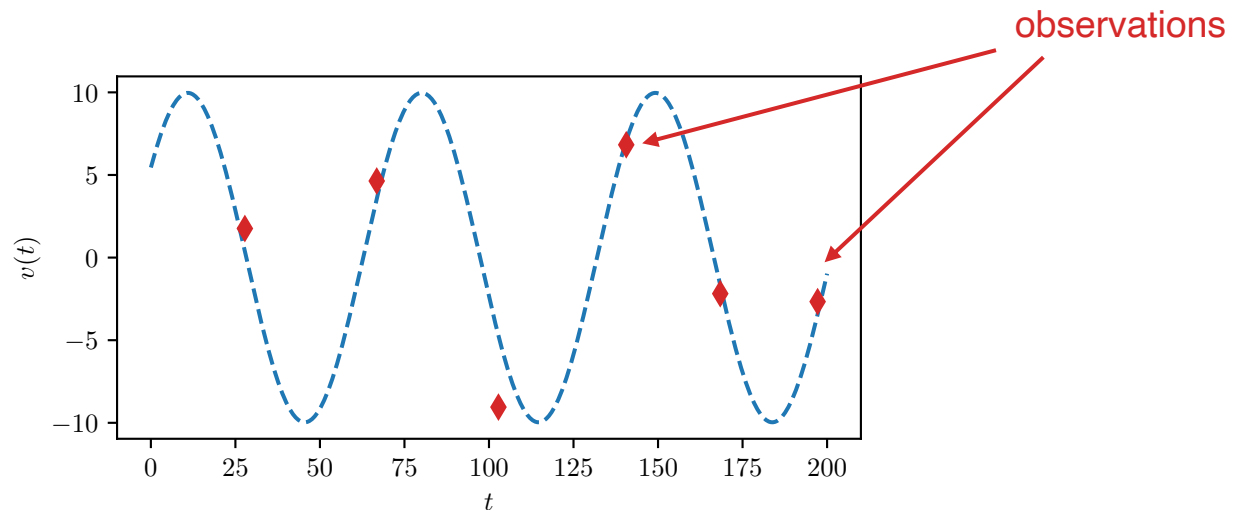
# EXAMPLES

# Bayesian inference: An example of model selection from astrophysics

- ▷ Star-exoplanet system orbiting center of mass



- ▷ Radial velocity along the orbit  $v(t; x) = v_0 + K \cos\left(\frac{2\pi}{P}t + \phi_0\right)$





# Bayesian model for velocity parameters

▷ **Radial velocity**  $v(t; x) = v_0 + K \cos\left(\frac{2\pi}{P}t + \phi_0\right)$

▷ **Parameters**  $x = (v_0, K, \phi_0, \ln P) \in \Omega \subset \mathbb{R}^4$

▷ **Likelihood from observations**  $L(x) = \mathcal{N}(v_k; v(t_k; x), \sigma_{\text{obs}}^2)$

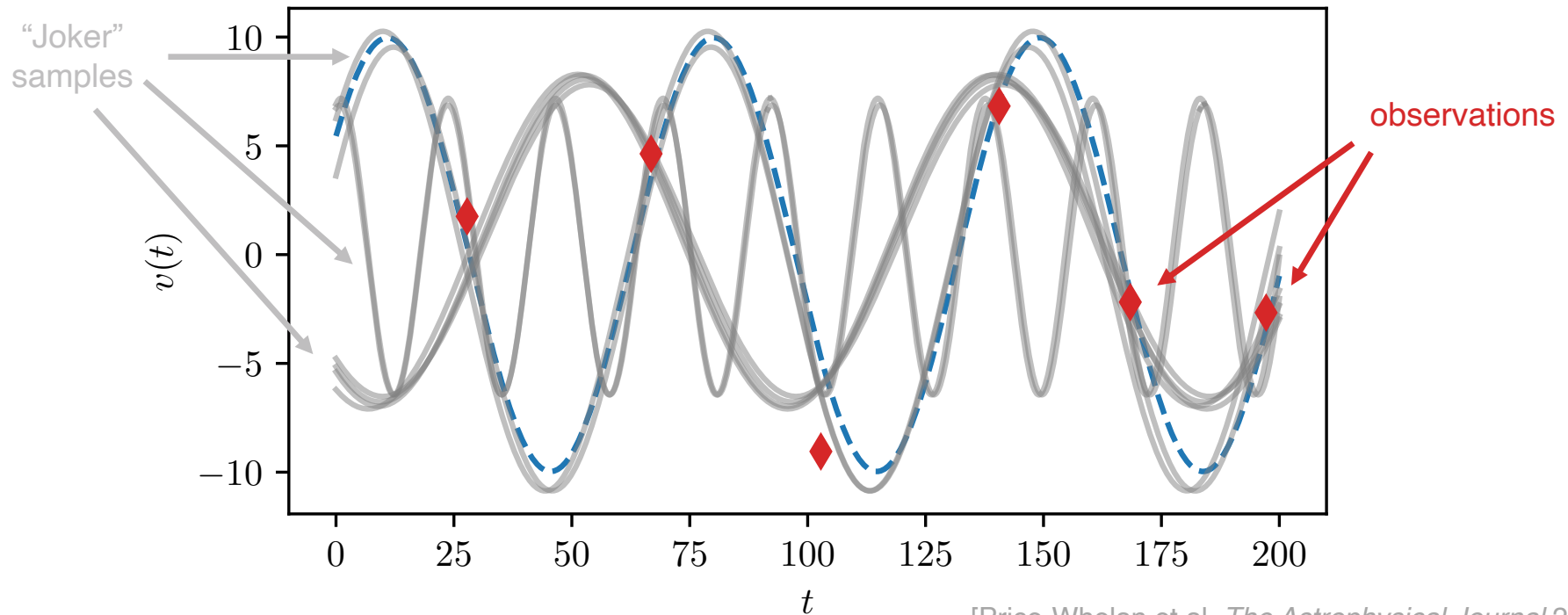
▷ **Priors**

$\ln P \sim \mathcal{U}(\ln P_{\text{min}}, \ln P_{\text{max}}),$

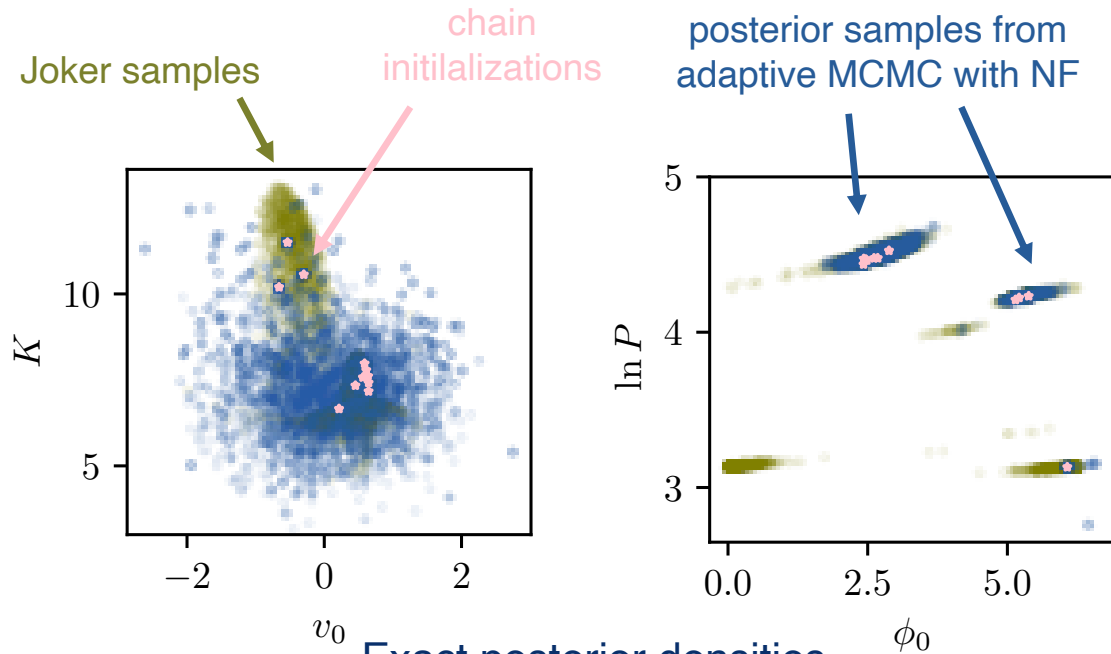
$\phi_0 \sim \mathcal{U}(0, 2\pi),$

$K \sim \mathcal{N}(\mu_K, \sigma_K^2),$

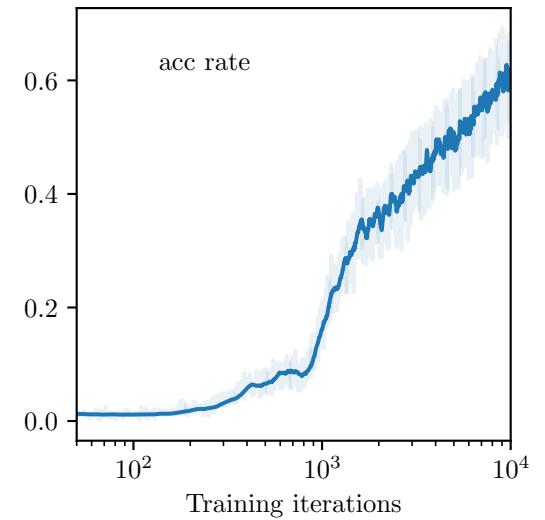
$v_0 \sim \mathcal{N}(0, \sigma_{v_0}^2).$



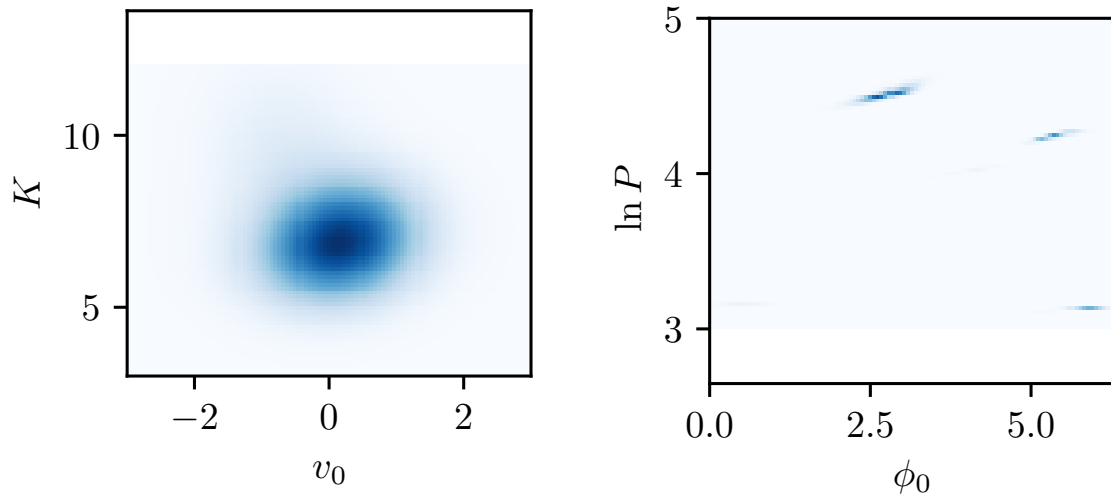
# Sampling from the posterior



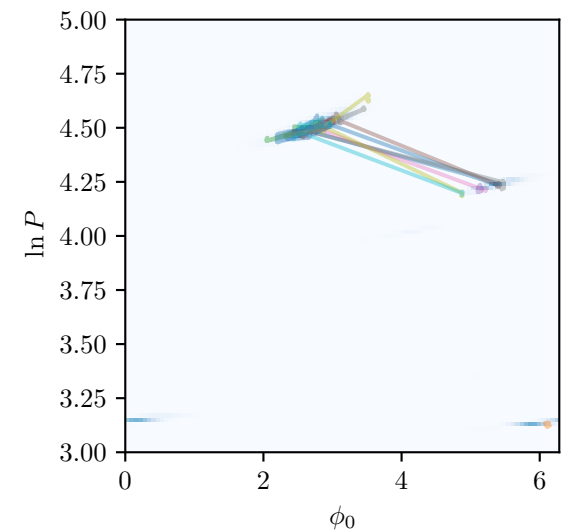
Acceptance  
along training



Exact posterior densities



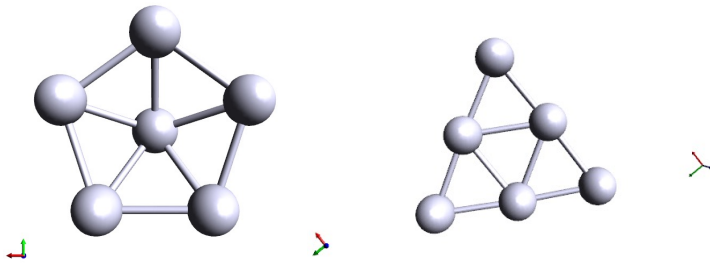
Fast mixing  
chains



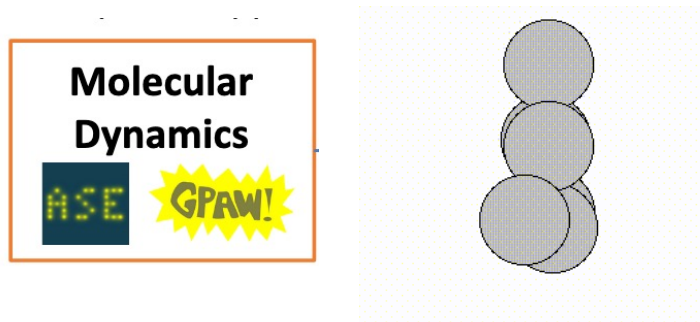
# Sampling metastable silver clusters

With Ana Molina Tarboda, Olga Lopez-Acevedo (Universidad de Antioquia), Pilar Cossio (Flatiron Institute)

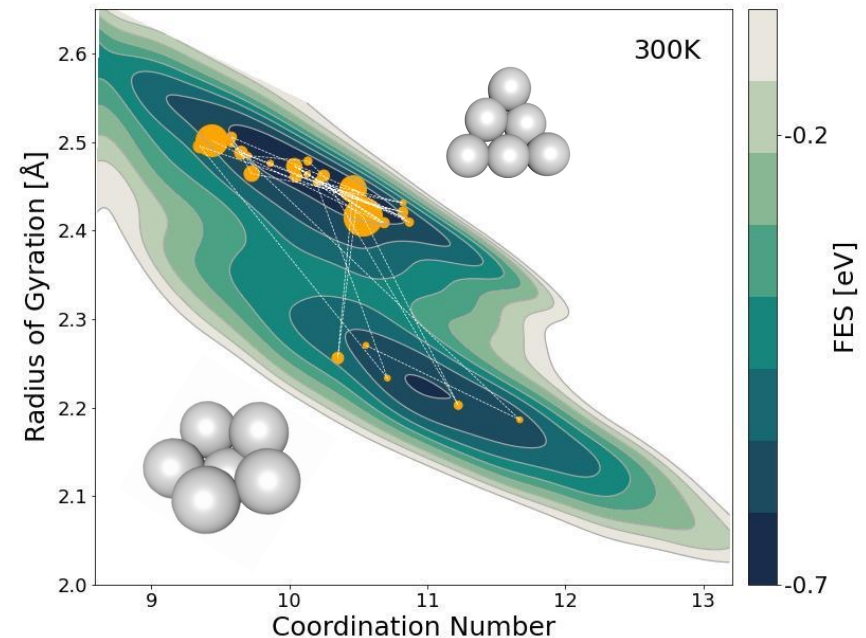
- ▷ Target density: Ground truth = Density Functional Theory: 2 metastable isomers



- ▷ Local sampler: Molecular Dynamics



2 dimension projection of the free energy surface



- ▷ Adaptive MCMC jumping between isomers

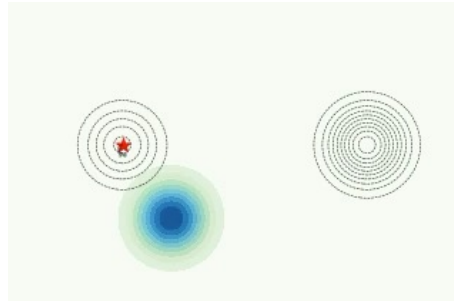
1. A couple of important sampling methods
  - 1.1 - Importance sampling
  - 1.2 - Metropolis-Hasting
2. Unsupervised learning / generative models
  - 2.1 - Latent deep generative models
  - 2.2 - Normalizing flows
3. Combining traditional inference method and learning
  - 3.1 - Variational Inference
  - 3.2 - Adaptive algorithms

4. Will it scale?
  - 4.1 - Local sampling in reparametrized space
  - 4.2 - Local-global sampling
  - 4.3 - Joining forces with annealing
  - 4.4 - Leveraging physics

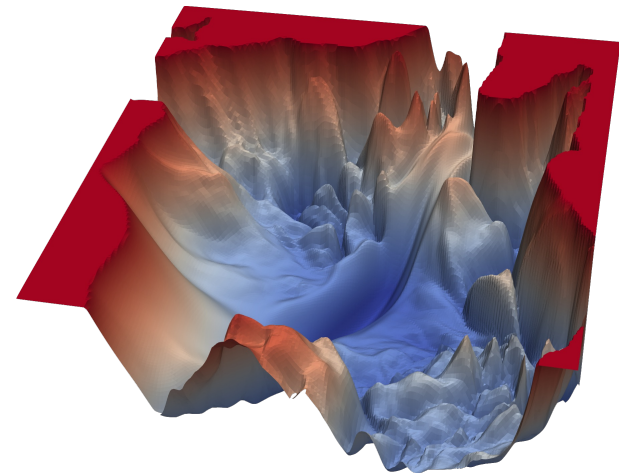
# Will it scale? A few hard problems

36

## ▷ Mode finding



## ▷ Disordered/Glassy landscapes



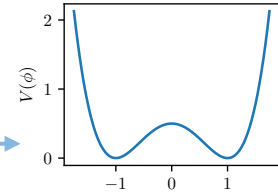
<https://www.cs.umd.edu/~tomg/projects/landscapes/>

## ▷ Probing bigger and bigger systems

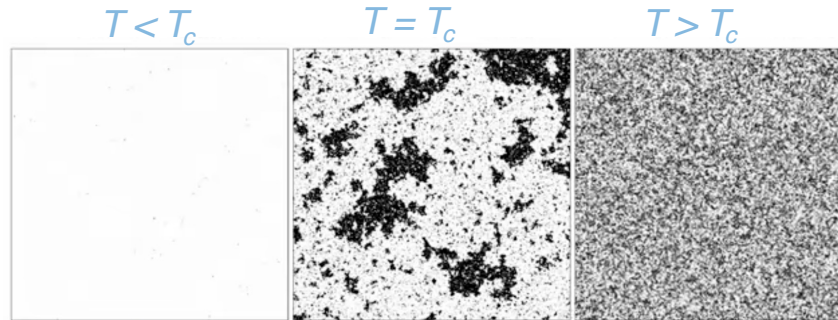
# Can the method scale to big systems?

▷ A systematic study on the 2d -  $\phi^4$  model (Del Debbio, *PRD* 2021)

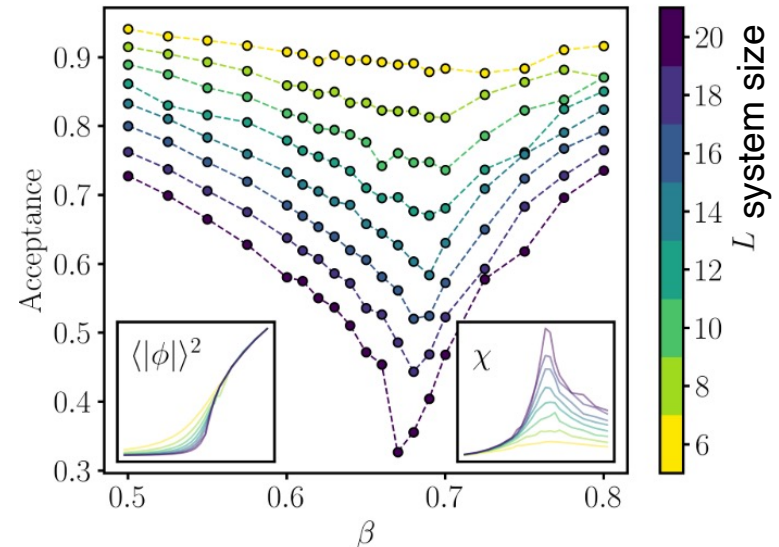
○ Action 
$$U(\Phi) = \sum_{i \in \Lambda} \left[ \underbrace{-\beta(\phi_{i+e_1}\phi_i + \phi_{i+e_2}\phi_i)}_{\text{coupling term}} + \underbrace{\phi_i^2 + \lambda(\phi_i^2 - 1)^2}_{\text{local potential}} \right]$$



○ Typical configurations



○ Mean acceptance probability after training



# Scaling to larger and larger systems

▷ Can surrogate probabilistic models scale?

▷ Metropolis acceptance scaling with dimension

$$\circ \text{acc}(x_{t+1}|x_t) = \min \left[ 1, e^{-(\Delta U_* - \Delta U_\theta)} \right] \sim \min \left[ 1, e^{-O(D)} \right]$$

Hm!

$$\Delta U_* = U_*(x_{t+1}) - U_*(x_t) \approx O(D) \text{ i.e. energy is typically extensive}$$

$$\Delta U_\theta = -\log \rho_\theta(x_{t+1}) + \log \rho_\theta(x_t) \approx O(D)$$

◦ Same story if importance sampling

▷ A first idea: be less ambitious and retain some locality in sampling

Loop over:

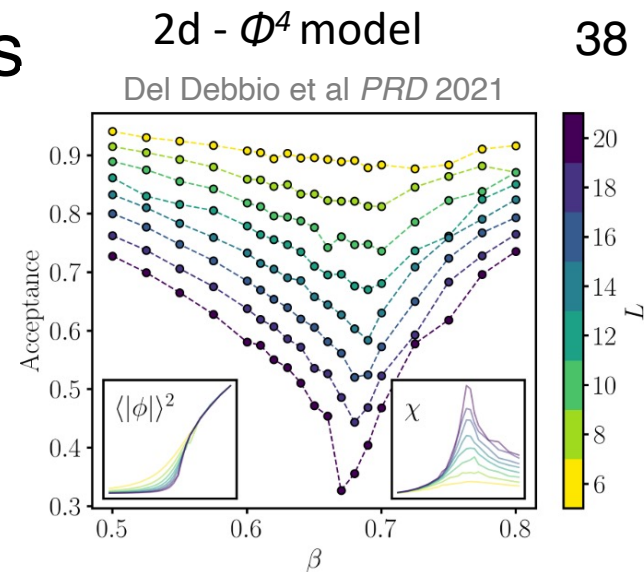
◦ Propose  $x_{t+1}^i \sim \rho_\theta(x)$

Independ flow proposal

$x_{t+1} \sim \pi_\theta(x_{t+1}|x_t)$

◦ Accept/reject  $\text{acc}(x_{t+1}^i|x_t^i) = \min \left[ 1, \frac{\rho_*(x_{t+1}^i)\rho_\theta(x_t^i)}{\rho_*(x_t^i)\rho_\theta(x_{t+1}^i)} \right]$

Conditional proposal less local than traditional kernel?

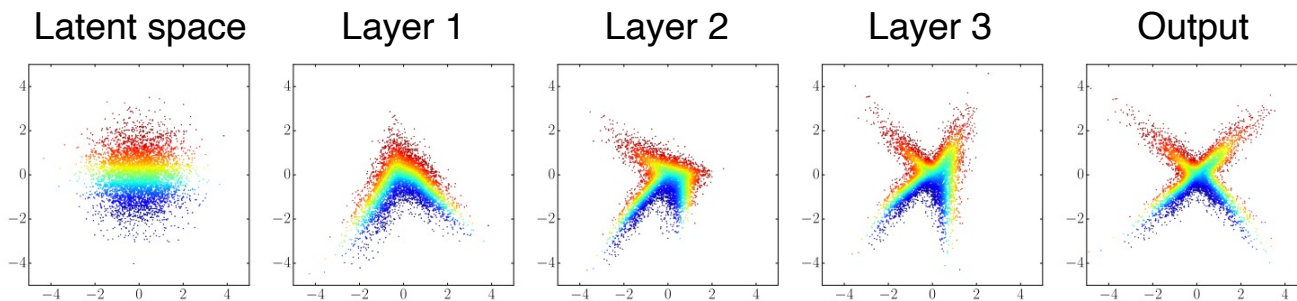


1. A couple of important sampling methods
  - 1.1 - Importance sampling
  - 1.2 - Metropolis-Hasting
2. Unsupervised learning / generative models
  - 2.1 - Latent deep generative models
  - 2.2 - Normalizing flows
3. Combining traditional inference method and learning
  - 3.1 - Variational Inference
  - 3.2 - Adaptive algorithms
4. Will it scale?
  - 4.1 - Local sampling in reparametrized space
  - 4.2 - Local-global sampling
  - 4.3 - Joining forces with annealing
  - 4.4 - Leveraging physics

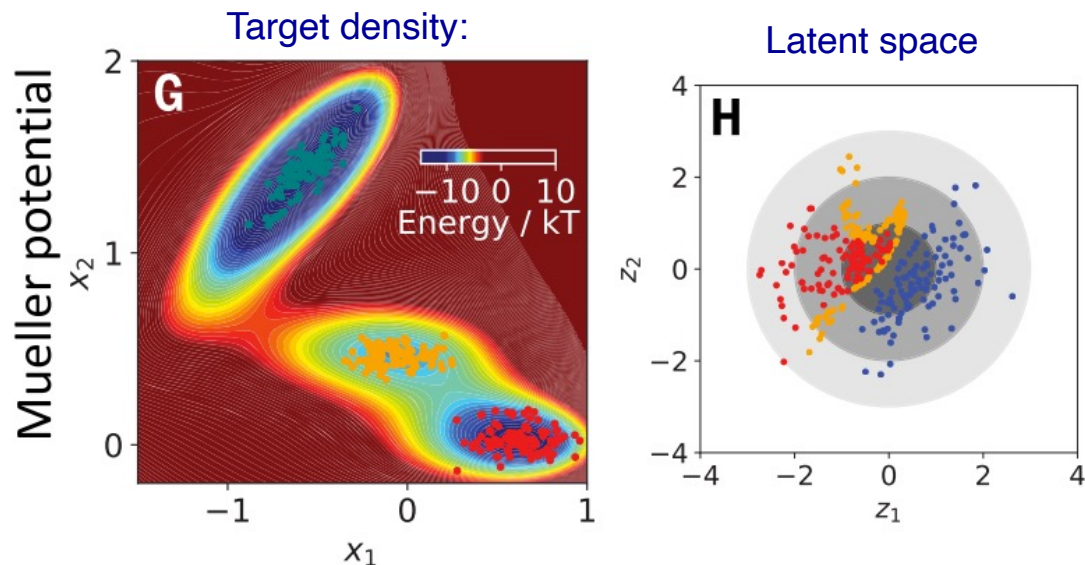


# 4.1 Reparametrization: reverse NF for MCMC

- ▷ Reverse transformation is normalizing = “Gaussianizing”



- ▷ Idea: train normalizing flow and use latent space to run traditional MCMC



# 4.1 Reparametrization: reverse NF for MCMC

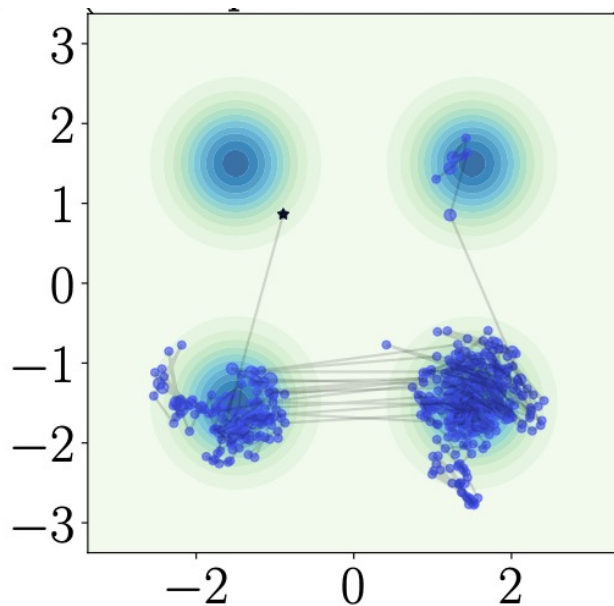
▷ NeuTra-lizing Bad Geometry in Hamiltonian Monte Carlo Using Neural Transport. (Hoffman et al 2019)

▷ Test case:

1. Train a flow on a mixture of Gaussian
2. Run MALA in the “latent space” on the push-backward

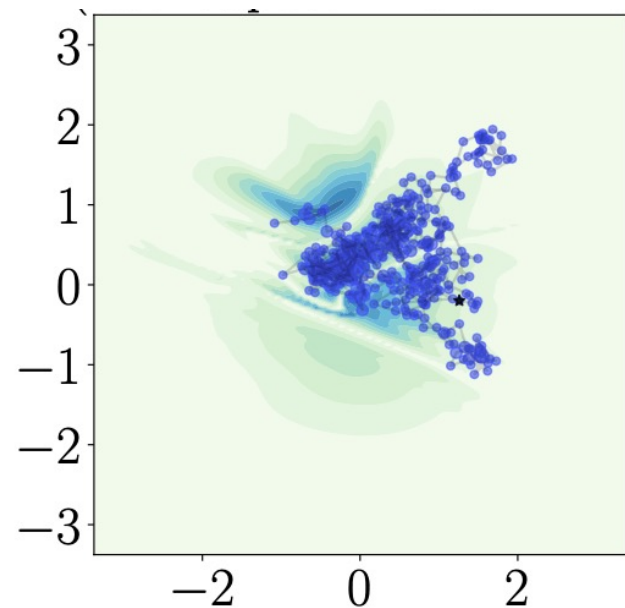
Target density:

$$\rho_*(x) = e^{-U(x)} / Z$$



Push backward:

$$\rho_*(T_\theta^{-1}(z)) \det |\nabla_z T_\theta| \approx \rho_B(z)$$



(Experiments by Louis Grenioux)

## 4.2 Global-Local samplers: the best of both worlds

With Eric Moulines, Sergey Samsonov and collaborators.

Target density:  $\rho_*(x) = e^{-U_*(x)} / Z$

Generative model parametrized density:  $\rho_\theta(x)$

### ▷ Algorithm: “Explore-Exploit MCMC”

Initialize:  $x_0^i \quad i = 1 \dots N \quad i = 1 \dots N$

Loop over parallel chains:

- Draw from generative model  $x_{t+1}^i \sim \rho_\theta(x)$
- Accept-reject  $\text{acc}(x_{t+1}^i | x_t^i) = \min \left[ 1, \frac{\rho_*(x_{t+1}^i) \rho_\theta(x_t^i)}{\rho_*(x_t^i) \rho_\theta(x_{t+1}^i)} \right]$
- Local resampling  $x_{t+1}^i \sim \pi_{\text{local}}(x_{t+1}^i | x_t^i)$

Metropolis-Hastings  
with NF

Local kernel

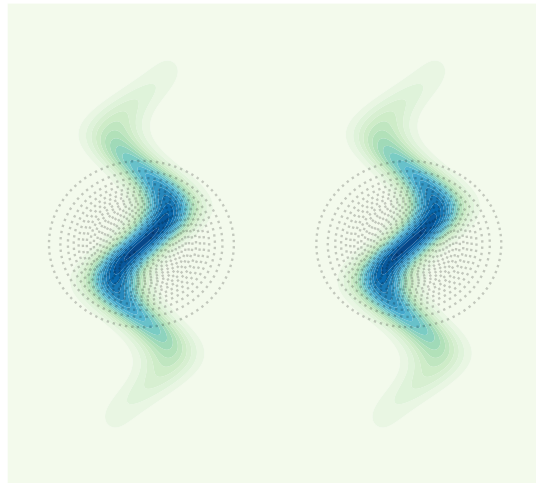
### ▷ Local + Mode jumping methods [Sminchisescu & Welling AISTAT 2017, Pompe et al. Ann. Stat 2020, Sbailò et al. J. Chem. Phys. 2021, ...]

## 4.2 Global-Local samplers: the best of both worlds

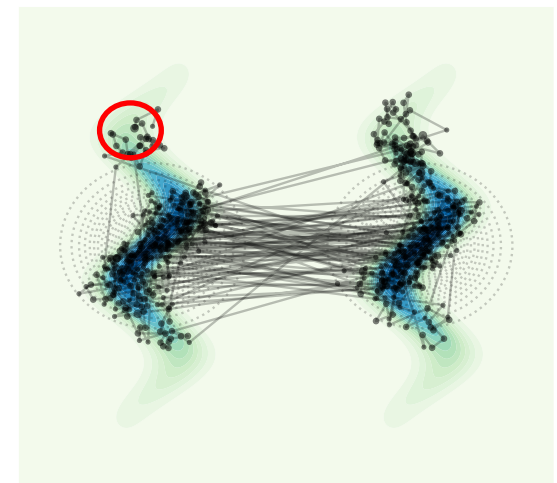
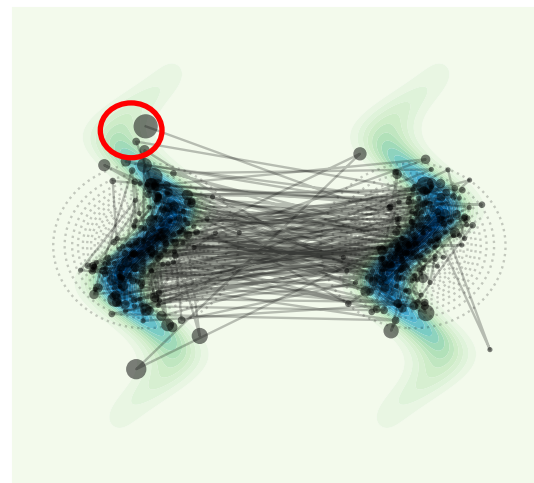
43

- ▷ In general tails of the distribution will be learned poorly

Global only



Local + global



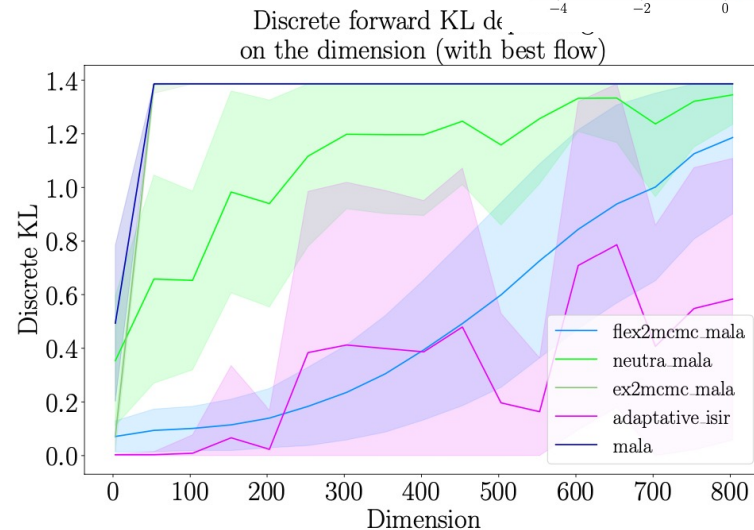
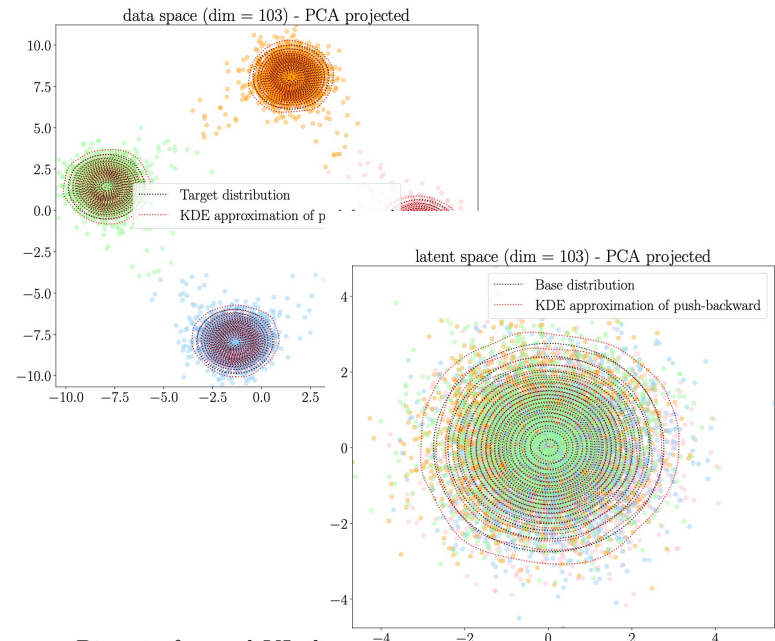
- ▷ Intuition: Local kernels allow to
  - Explore regions that are not (yet) properly learned
  - Drive learning there (if running an adaptive MCMC, a.k.a. adaptive learning)
- ▷ Possible to derive theoretical guarantees of the improvement they bring.

# 4.2 Global-Local samplers: the best of both worlds

▷ Does Global-local samplers scale better than NeuTraMCMC (transported MCMC)?

- Multimodal Gaussian mixture
- Train a flow to reproduce the mixture in higher dimension
- Measure of how many modes are visited by the different algorithms

Flows are better exploited with Explore-Exploit than the NeuTraMCMC !



1. A couple of important sampling methods
  - 1.1 - Importance sampling
  - 1.2 - Metropolis-Hasting
2. Unsupervised learning / generative models
  - 2.1 - Latent deep generative models
  - 2.2 - Normalizing flows
3. Combining traditional inference method and learning
  - 3.1 - Variational Inference
  - 3.2 - Adaptive algorithms
4. Will it scale?
  - 4.1 - Local sampling in reparametrized space
  - 4.2 - Local-global sampling
  - 4.3 - Joining forces with annealing
  - 4.4 - Leveraging physics

# 4.3 Joining forces with annealing

▷ Going back to Variational Inference:

- Minimize Kullback-Leibler  $D_{KL}(\rho_\theta || \rho_*) =$  variational principle with expressive  $\rho_\theta(x)$  ansatz

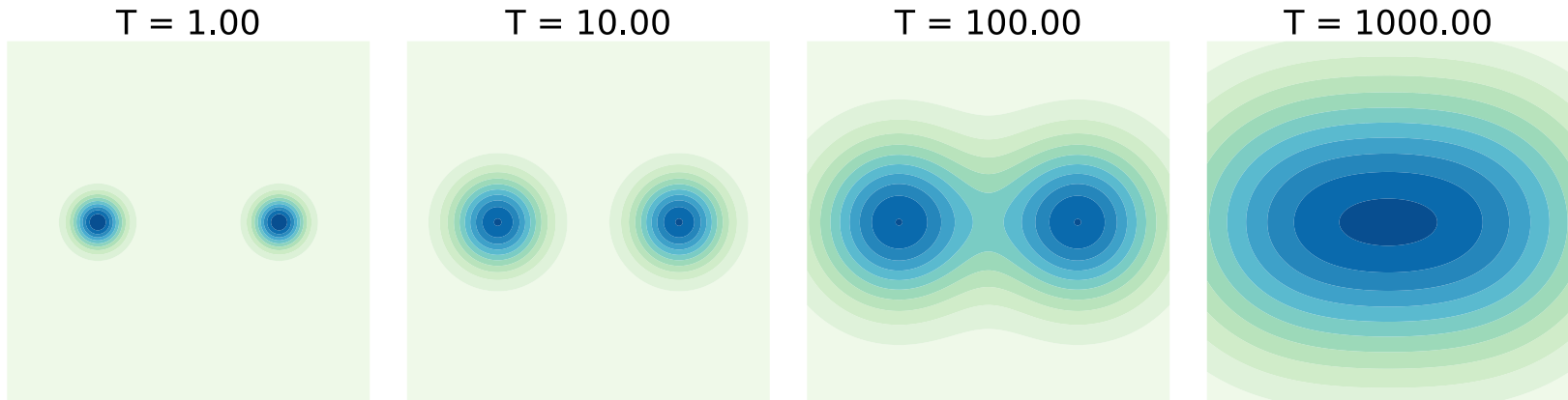
$$D_{KL}(\rho_\theta || \rho_*) = \int \log \frac{\rho_\theta(x)}{\rho_*(x)} \rho_\theta(x) dx \approx \sum_{i=1}^N \log \frac{\rho_\theta(x_i)}{\rho_*(x_i)} \quad x_i \sim \rho_\theta(x) \quad + \text{SGD!}$$

$$\rho_\theta(x) = \rho_B(T_\theta^{-1}(x)) \det |\nabla_x T_\theta^{-1}|$$

- Recall, issue is mode collapse.

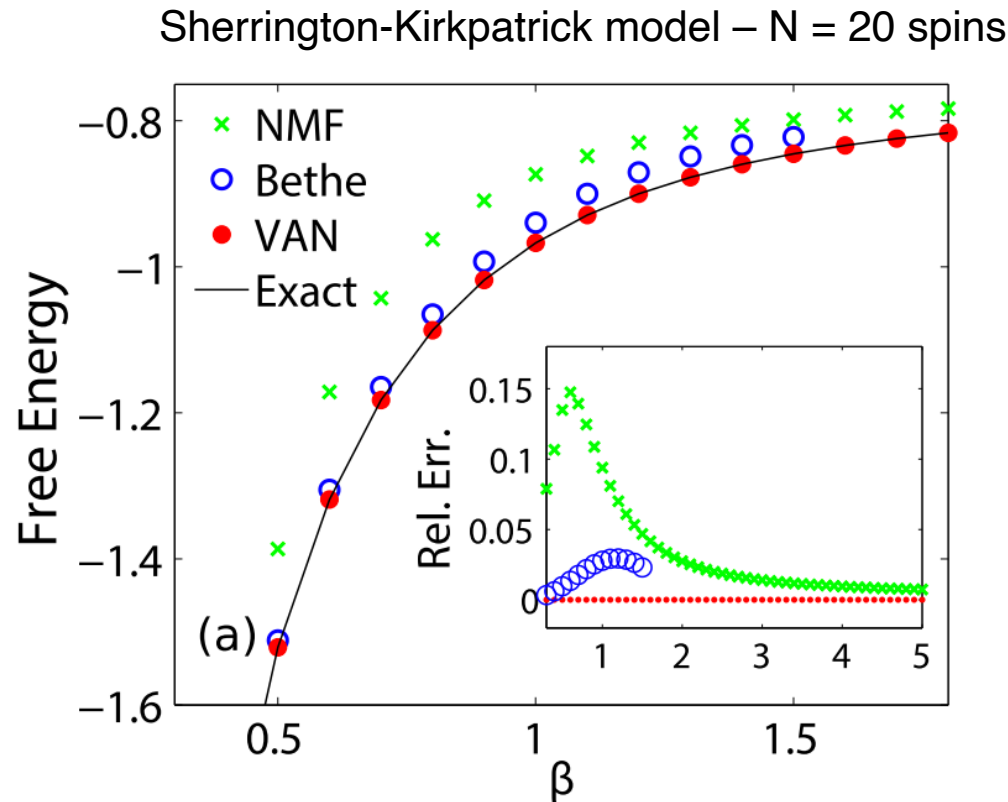
▷ Idea: start at high temperature and progressively decrease to desired model

$$\rho_{\beta_*}(x) = e^{-\beta U_*(x)} / Z$$



## 4.3 Joining forces with annealing

- ▷ “Solving Statistical Mechanics Using Variational Autoregressive Networks”  
Wu, Wang and Zhang (*PRL* 2019)
- ▷ Variational inference with annealed target + MCMC correction





## 4.3 Joining forces with annealing: Pushing towards more complicated models

### ▷ Annealing to create progressively dataset of training

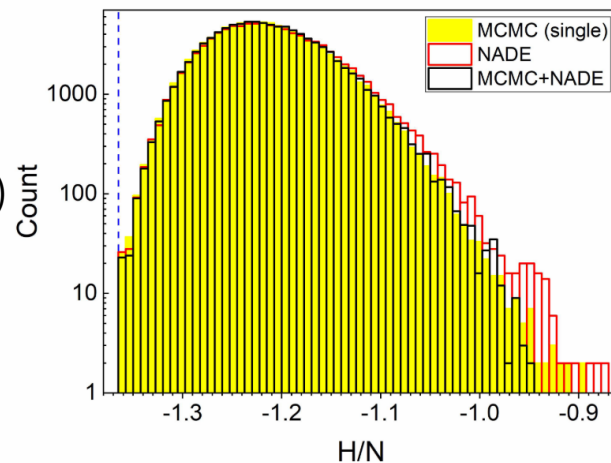
$$\rho_{\beta_*}(x) = e^{-\beta U_*(x)} / Z$$

From high-temperature repeat:

- Use  $\rho_{\theta_{k-1}}^{\beta_{k-1}}(x) = \rho_p(x)$  in MCMC to sample  $\rho_{\beta_*}^{\beta_k}(x)$
- Use  $x_i \sim \rho_{\beta_*}^{\beta_k}(x_i)$  as data to train  $\rho_{\theta_{k-1}}^{\beta_{k-1}}(x)$

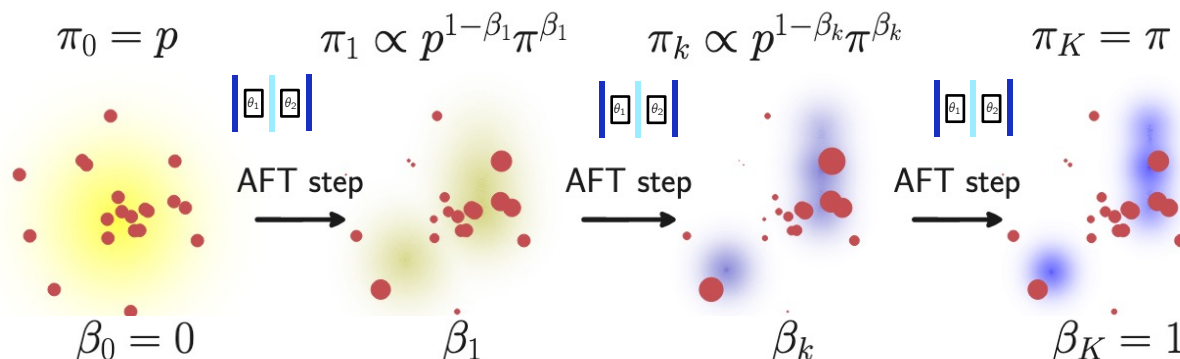
2d – Edwards Anderson

S. Pilati *PRE* 2020



### ▷ (Continuously Repeated) Annealed Flow Transport

- Add flow transport maps within steps of sequential Monte Carlo (SMC)



## ▷ Using symetries and invariance

- cf Danilo's talk
- cf cluster updates by Wu, Rossi & Carleo, PRR (2021)

## ▷ Using informed base measures

# 4.4 – Leveraging physics: informed base measures

## ▷ Example: 1d - $\phi^4$ model

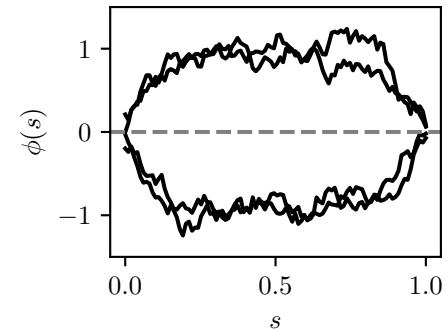
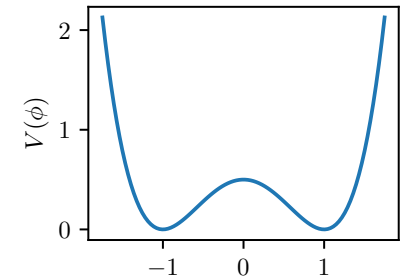
○ Random field  $\phi: [0, 1] \mapsto \mathbb{R} \in C([0, 1]; \mathbb{R})$  *local potential*

○ Energy functional  $U_*(\phi) = \int_{[0,1]} \left( \frac{a}{2} |\nabla_s \phi|^2 + V(\phi) \right) ds$

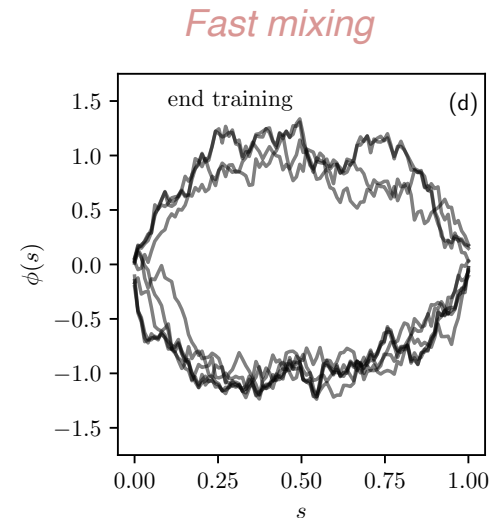
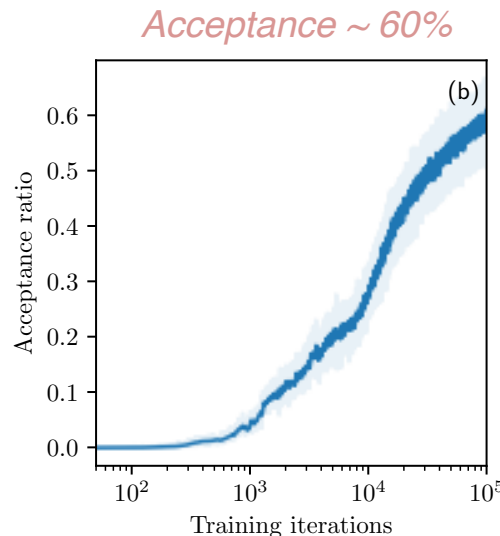
○ Local potential  $V(\phi) = \frac{1}{2}(\phi^2 - 1)^2$  *coupling term*

○ Dirichlet boundary conditions  $\phi(0) = 0, \phi(1) = 0$

○ Target distribution  $\rho(\phi) = \frac{1}{Z_\beta} e^{-\beta U(\phi)}$



## ▷ Discretized: N=100



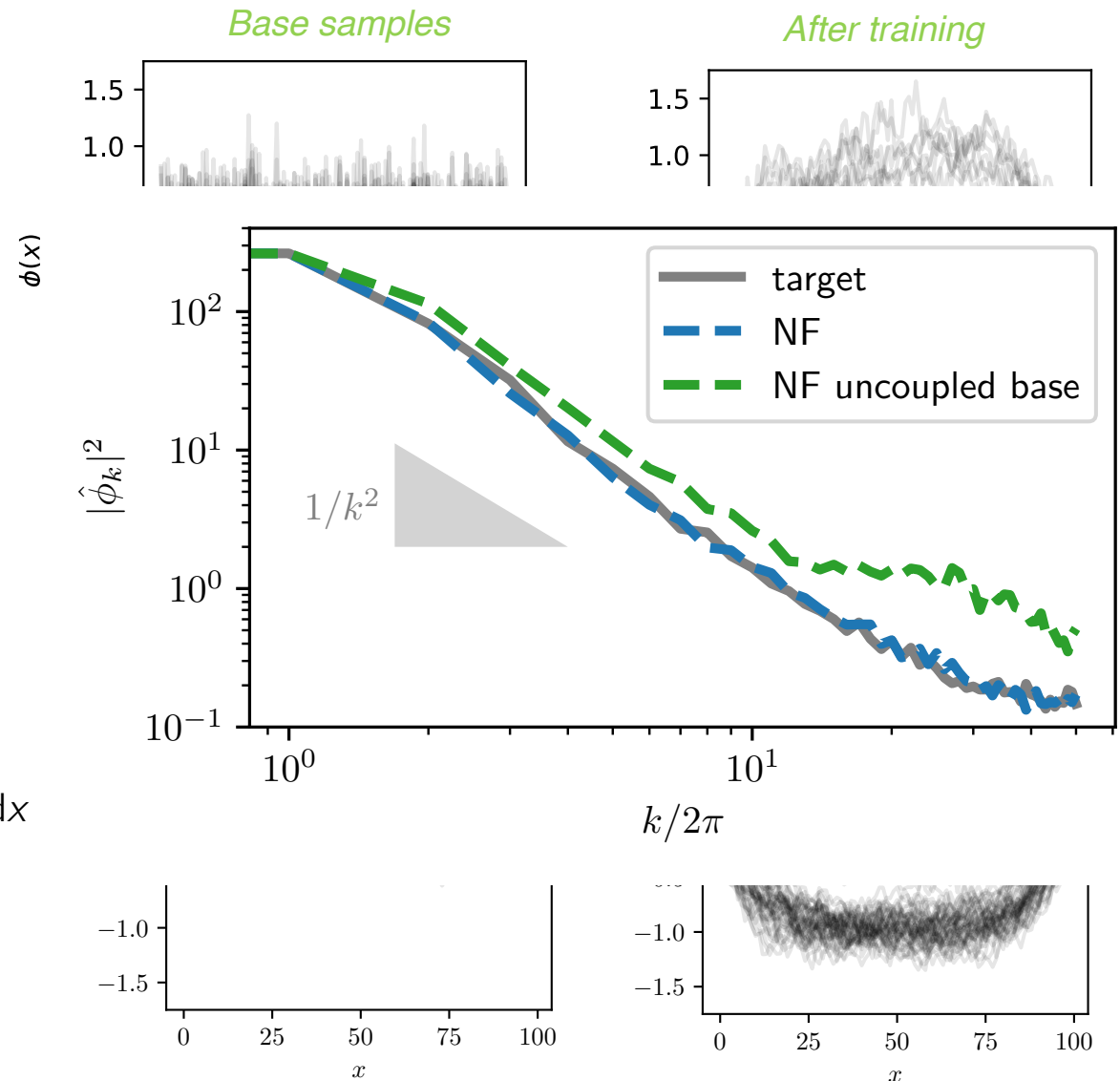
# Uncoupled vs coupled base distributions

Gaussian uninformed  
(uncoupled)

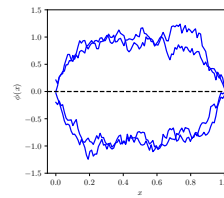
$$U_B(\phi) = \int \frac{1}{2\sigma^2} \phi^2 dx$$

Gaussian informed  
(coupled)

$$U_B(\phi) = \int \left( \frac{a}{2} |\nabla_x \phi|^2 + \frac{1}{2\sigma^2} \phi^2 \right) dx$$



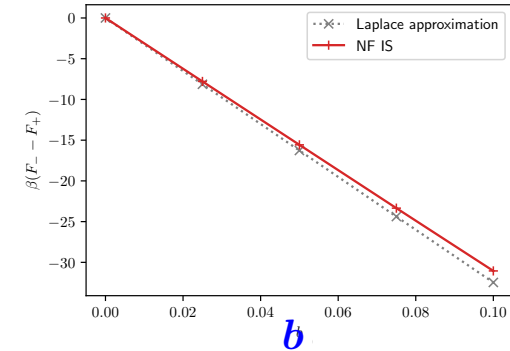
# More numerical checks



▷ Tilt distribution towards -1 configuration with local field

$$U_{*,b}(\phi) = \int \left( \frac{a}{2} |\nabla_x \phi|^2 + V(\phi) + \mathbf{b} \phi \right) dx$$

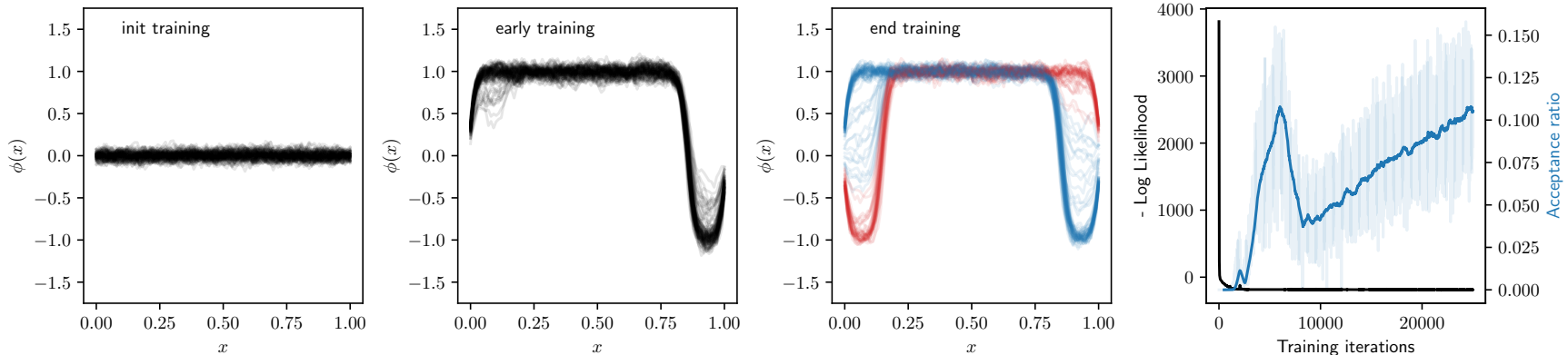
free energy difference



▷ Tilt distribution towards configuration with average value away from +1 or -1

$$U_{*,\lambda,\bar{\phi}}(\phi) = \int \left( \frac{a}{2} |\nabla_x \phi|^2 + V(\phi) + \lambda \left( \int \phi dx' - \bar{\phi} \right)^2 \right) dx$$

○ Train with  $\bar{\phi} = 0.7$



## ▷ Using symetries and invariance

- cf Danilo's talk
- cf cluster updates by Wu, Rossi & Carleo, PRR (2021)

## ▷ Using informed base measures

1. A couple of important sampling methods
  - 1.1 - Importance sampling
  - 1.2 - Metropolis-Hasting
2. Unsupervised learning / generative models
  - 2.1 - Latent deep generative models
  - 2.2 - Normalizing flows
3. Combining traditional inference method and learning
  - 3.1 - Variational Inference
  - 3.2 - Adaptive algorithms
4. Will it scale?
  - 4.1 - Local sampling in reparametrized space
  - 4.2 - Local-global sampling
  - 4.3 - Joining forces with annealing
  - 4.4 - Leveraging physics

## ▷ Opportunities

- VI, IS and MCMC can be powered by normalizing flows/NADE
- IS or MCMC allows to de-bias the training model

## ▷ Challenges for scaling things up and morals

- ML models should be joining forces with « traditional » sampling methods
- ML models should leverage the physical knowledge about systems

## ▷ Softwares

 PyTorch



### Adaptive MCMCs

 [marylou-gabrie / flonaco](#) Public

 [kazewong / flowMC](#) Public

### Sequential MCMCs

 [minaskar / pocomc](#) Public

 [deepmind / annealed\\_flow\\_transport](#) Public



# Workshop : Machine Learning-Assisted Sampling for Scientific Computing – Applications in Physics

## Dates of conference open to all :

- 3-4 October 2022
- Registration deadline : 16 September 2022

Conference registration [↗](#)

Join us in Paris  
or online!

## Organizers

- Valentin de Bortoli, Ecole Normale Supérieure, Paris
- Marylou Gabrié, Ecole Polytechnique, Paris
- Tony Lelièvre, Ecole Nationale des Ponts et Chaussées, Paris

## Invited Speakers

- David Aristoff, Colorado State University
- Peter Bolhuis, University of Amsterdam
- Freddy Bouchet, Ecole Normale Supérieure, Lyon
- Maria K. Cameron, University of Maryland
- Arnaud Doucet, University of Oxford
- Alain Durmus (ENS Paris Saclay & École Polytechnique)
- Stéphane Mallat, Ecole Normale Supérieure, Paris
- Pierre Monmarché, Pierre and Marie Curie University (UPMC)
- Jutta Rogal, New York University
- Phiala Shanahan, Massachusetts Institute of Technology
- Gabriel Stoltz, Ecole Nationale des Ponts et Chaussées, Paris
- Jonathan Weare, New York University
- Martin Weigt, Pierre and Marie Curie University (UPMC)
- Wei Zhang, Zuse Institute Berlin

Open opportunities for phds or postdocs!



École Polytechnique, (almost Paris)