

# Learning Coarse-Grained Models from Molecular Dynamics

Data-Driven and Embedded-Physics ML.

---

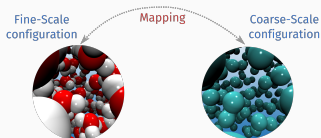
M. Schöberl<sup>1,2</sup>, N. Zabaras<sup>1</sup>, P.-S. Koutsourelakis<sup>2</sup>

SIAM CSE 2019, Spokane WA

February 27, 2019.

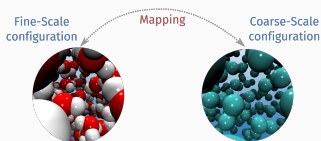
<sup>1</sup>Center for Informatics and Computational Science, University of Notre Dame, IN, USA.

<sup>2</sup>Continuum Mechanics Group, Technical University of Munich, Germany.



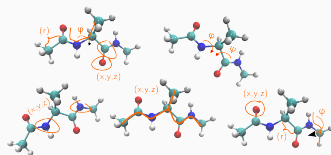
## Why coarse-graining?

- Overcome spatiotemporal limitations.
- Reveal physical insight from reduced representation.



## Why coarse-graining?

- Overcome spatiotemporal limitations.
- Reveal physical insight from reduced representation.



## Which set of collective variables (CVs) captures resilient and parsimonious features of atomistic systems?

- Relevant coordinates are highly clustered around a set of *lower dimensional collective variables*. [ban (2017), Chen and Ferguson (2017)]
- Vast combinatoric possibilities for *choosing CVs*. [Chakraborty et al. (2018)]

## Questions we address

- How to learn a *predictive* coarse-grained representation
  - in the *small data regime* and
  - *in absence of any data?*
- How to identify good CVs?
- What are *good* CVs?
- Are identified CVs physically interpretable?

## Atomistic model

$$p_{\text{target}}(\mathbf{x}) \propto e^{-\beta U(\mathbf{x})}$$

- $\mathbf{x} \in \mathcal{M}$ : atomistic coordinates
- $U(\mathbf{x})$ : atomistic potential
- Observables:

$$\mathbb{E}_{p(\mathbf{x})}[a] = \int a(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

## Coarse-grained representation

$$\mathbf{z} = \mathcal{R}(\mathbf{x}), \quad \dim(\mathbf{z}) \ll \dim(\mathbf{x})$$

- $\mathbf{z}$ : reduced CG / collective variables.
- $\mathcal{R}$ : mapping operator (mapping to CG variables).

## Data-driven

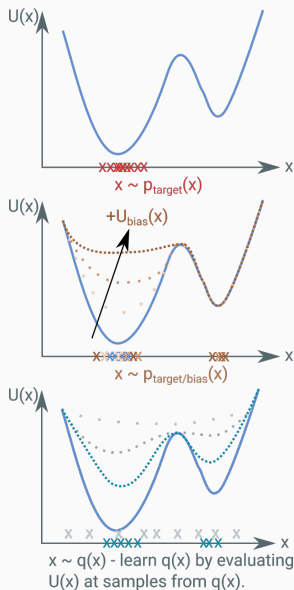
Simulate **sufficiently long** reference governing equation and obtain (limited) data approximating  $\rho_{\text{target}}(\mathbf{x})$ . [Shell (2008), Katsoulakis and Trushoras (2006), Trushoras and Tsagarogiannis (2010), Noid (2013), Brunton et al. (2016), Wehmeyer and Noé (2018)]

## Data-augmenting approach

E.g. enhanced sampling, learning based on consecutively gathered insight. [Laio and Parrinello (2002), Darve et al. (2008), Bilonis and Koutsourelakis (2012), Ferguson et al. (2011), Ferguson (2017), Chen and Tuckerman (2018)]

## Embedded-Physics approach

Instead of simulating the target (based on known  $U(\mathbf{x})$ ), directly incorporate physical constraints at our disposal, e.g. potential/force field. [Noé and Wu (2018)]



# Data-Driven Coarse-Graining

---

Introduce  $q(z)$  and  $q(x|z)$ :  $\underbrace{q(z)}_{\text{CV}} \xrightarrow{q(x|z)} q(x) = \int q(x|z) q(z) dz$  [Schöberl, Zabaras, and

Koutsourelakis (2017)]

Which latent CVs  $z$  give rise to observations  $x^{(i)}$ ?

CV

Latent -----  $z \sim q(z)$

Generate / Decode

$x \sim q(x|z)$

Observed -----

Data  $x^{(i)} \sim p_{\text{target}}(x)$



Introduce  $q(z)$  and  $q(x|z)$ :  $\underbrace{q(z)}_{\text{CV}} \xrightarrow{q(x|z)} q(x) = \int q(x|z) q(z) dz$  [Schöberl, Zabarás, and

Koutsourelakis (2017)]

Which latent CVs  $z$  give rise to observations  $x^{(i)}$ ?

CV

Latent -----  $z \sim q(z)$

Generate / Decode

$x \sim q(x|z)$

Observed -----

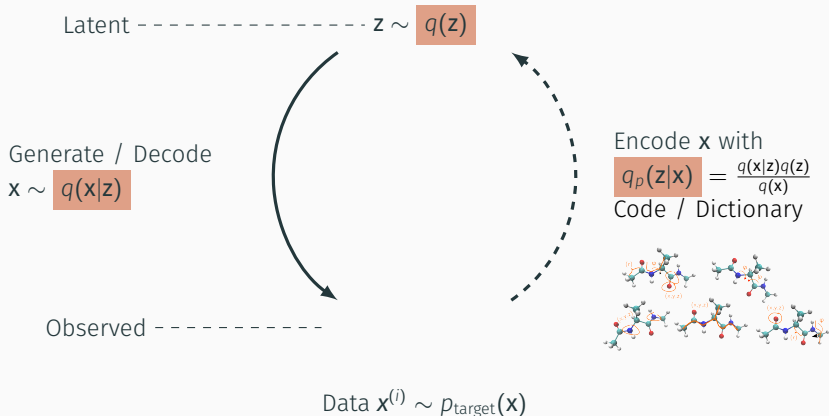
Data  $x^{(i)} \sim p_{\text{target}}(x)$

Introduce  $q(z)$  and  $q(x|z)$ :  $\underbrace{q(z)}_{\text{CV}} \xrightarrow{q(x|z)} q(x) = \int q(x|z) q(z) dz$  [Schöberl, Zabaras, and

Koutsourelakis (2017)]

Which latent CVs  $z$  give rise to observations  $x^{(i)}$ ?

CV



Parametrize

$$q(\mathbf{z}|\boldsymbol{\theta})$$

Model for CVs.

$$q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})$$

Probabilistic mapping given CVs to atomistic coordinates.

Parametrize  $\underbrace{q(\mathbf{z}|\boldsymbol{\theta})}$  Model for CVs.  $\underbrace{q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})}$  Probabilistic mapping given CVs to atomistic coordinates.

## Data-driven approach: Minimize

$$D_{\text{KL}}(p_{\text{target}}(\mathbf{x})||q(\mathbf{x})) = - \int p_{\text{target}}(\mathbf{x}) \log \frac{q(\mathbf{x})}{p_{\text{target}}(\mathbf{x})} d\mathbf{x}$$

equals maximizing marginal log-likelihood of the data set

$\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ :

$$\log q(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{i=1}^N \log q(\mathbf{x}^{(i)}).$$

## Stochastic Variational Bayesian approximation

The log-likelihood is decomposed into: [Beal and Ghahramani (2006); Kingma and Welling (2013); Rezende et al. (2014)]

$$\log q_{\theta}(\mathbf{x}^{(i)}) = \underbrace{\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})}_{\text{variational lower bound, since } D_{KL} \geq 0} + \underbrace{D_{KL}(r_{\phi}(\mathbf{z}^{(i)}|\mathbf{x}^{(i)})||q_{\theta}(\mathbf{z}^{(i)}|\mathbf{x}^{(i)}))}_{\geq 0}$$

$$\log q_{\theta}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = - \underbrace{D_{KL}(r_{\phi}(\mathbf{z}^{(i)}|\mathbf{x}^{(i)})||q_{\theta}(\mathbf{z}))}_{\text{Regularize } \phi, \text{ such that } r_{\phi}(\mathbf{z}^{(i)}|\mathbf{x}^{(i)}) \text{ in aggregation is close to } q_{\theta}(\mathbf{z}).} + \underbrace{\mathbb{E}_{r_{\phi}(\mathbf{z}^{(i)}|\mathbf{x}^{(i)})}[\log q_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i)})]}_{\text{Expected neg. reconstruction error.}}$$

with the approximate posterior  $r_{\phi}(\mathbf{z}^{(i)}|\mathbf{x}^{(i)})$ , e.g. a distribution of the exponential family and parameters  $\phi$ .

**Discovery of CVs as approximate Bayesian inference, i.e. identify the code/dictionary  $r_{\phi}(\mathbf{z}|\mathbf{x})$ .** [Schöberl, Zabarás, and Koutsourelakis (2019)]

- MLE estimate

$$\max_{\phi, \theta} \mathcal{L}(\theta, \phi; X)$$

- MAP estimate

$$\max_{\phi, \theta} \mathcal{L}(\theta, \phi; X) + \underbrace{\log p(\theta)}_{\text{log-prior}}$$

- MLE estimate

$$\max_{\phi, \theta} \mathcal{L}(\theta, \phi; X)$$

- MAP estimate

$$\max_{\phi, \theta} \mathcal{L}(\theta, \phi; X) + \underbrace{\log p(\theta)}_{\text{log-prior}}$$

- Approximate posterior of decoding parameters  $\theta$ ,  $p(\theta|X)$ , with Laplace approximation.

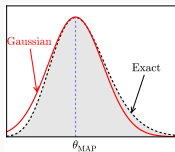
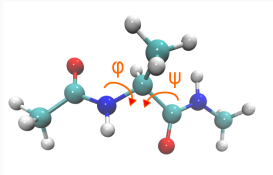


Figure 1: Laplace approximation:  
 $q(\theta|X) \approx \mathcal{N}(\mu, S)$

- $q(\theta|X) \approx \mathcal{N}(\mu, S)$
- $\mu = \theta_{MAP}$
- $S^{-1} = -\frac{\partial^2 \mathcal{L}(\theta, \phi; X)}{\partial \theta_r \partial \theta_l} - \frac{\partial^2 \log p(\theta)}{\partial \theta_r \partial \theta_l}$



## Simulation details



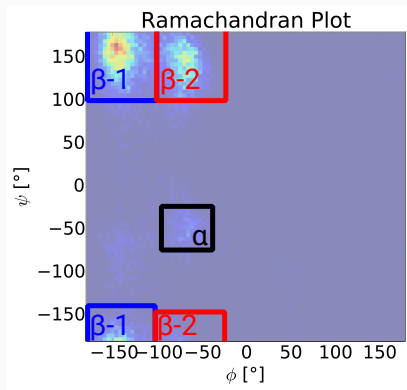
**Figure 2:** ALA-2 peptide with indicated dihedral angles.

- 22 atoms  $\rightarrow \dim(\mathbf{x}) = 66$  in implicit solvent [Still et al. (1990)] .
- AMBER-FF96 force field [Salomon-Ferrer et al.] with Andersen thermostat at  $T = 330\text{K}$  and  $\Delta t = 1 \text{ fs}$ .
- Equilibration for 50 ns. Snapshots taken every 10 ps.
- No pre-processing of data the data:  $\mathbf{x}^{(i)}$  are the Cartesian coordinates of *all* atoms in the system.





## Characteristic conformations of the ALA-2 peptide



**Figure 3:** Characteristic conformations ( $\alpha$ ,  $\beta$ -1,  $\beta$ -2) and their labelling as used in the sequel.



## Auto-Encoding Variational Bayes

According [Kingma and Welling (2013)].

$$q(\mathbf{x}|\boldsymbol{\theta}) = \int q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) q(\mathbf{z}) d\mathbf{z}$$

Mapping  $\mathbf{z} \rightarrow \mathbf{x}$  (decoder):

$$q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{z}; \boldsymbol{\theta}), S_{\boldsymbol{\theta}} = \text{diag}(\boldsymbol{\sigma}_{\boldsymbol{\theta}}^2)).$$

With,

- $\boldsymbol{\mu}(\mathbf{z}; \boldsymbol{\theta})$  output of fully connected decoding neural network.
- $\boldsymbol{\sigma}_{\boldsymbol{\theta}}^2$  independent of  $\mathbf{z}$  due to [Mattei and Frelsen (2018)].



## Auto-Encoding Variational Bayes

According [Kingma and Welling (2013)].

$$q(\mathbf{x}|\boldsymbol{\theta}) = \int q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) q(\mathbf{z}) d\mathbf{z}$$

Mapping  $\mathbf{z} \rightarrow \mathbf{x}$  (decoder):

$$q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{z}; \boldsymbol{\theta}), S_{\boldsymbol{\theta}} = \text{diag}(\boldsymbol{\sigma}_{\boldsymbol{\theta}}^2)).$$

With,

- $\boldsymbol{\mu}(\mathbf{z}; \boldsymbol{\theta})$  output of fully connected decoding neural network.
- $\boldsymbol{\sigma}_{\boldsymbol{\theta}}^2$  independent of  $\mathbf{z}$  due to [Mattei and Frelsen (2018)].
- Proposed model does not pre-assume any physical insight.
- No a priori assumption or data pre-processing needed.



## Auto-Encoding Variational Bayes

$$q(\mathbf{x}|\boldsymbol{\theta}) = \int q(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) q(\mathbf{z}|\boldsymbol{\theta}) dz$$

We assume CVs are normal distributed according,

$$q(\mathbf{z}) = \mathcal{N}(\mathbf{0}, I).$$

- Here:  $\dim(\mathbf{z}) = 2$ .
- Explore assigned meaning of CVs given  $\dim(\mathbf{z}) = 2$ .



## Auto-Encoding Variational Bayes

The approximate posterior is of the following form,

$$r(\mathbf{z}|\mathbf{x}; \phi) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}; \phi), \mathbf{S}_\phi = \text{diag}(\boldsymbol{\sigma}(\mathbf{x}; \phi)^2)).$$

With,

- $\boldsymbol{\mu}(\mathbf{x}; \phi)$  and  $\boldsymbol{\sigma}^2(\mathbf{x}; \phi)$  described by a fully connected encoding network.
- Last layers separate into  $\boldsymbol{\mu}(\mathbf{x}; \phi)$  and  $\boldsymbol{\sigma}^2(\mathbf{x}; \phi)$ .



Prediction of atomistic configurations for  $\{z|z_1 = [-4, 4], z_2 = 0\}$

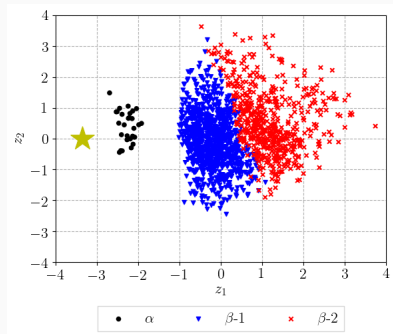


Figure 4: CVs are highly correlated with the dihedral angles ( $\phi, \psi$ ).

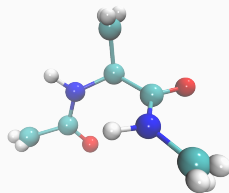


Figure 5: Mean prediction  $\mathbf{x} = \mu_{\theta}(\mathbf{z})$  of decoder  $q_{\theta}(\mathbf{x}|\mathbf{z})$ , given  $\{z|z_1 = [-4, 4], z_2 = 0\}$ .



*Prediction of atomistic configurations for  $\{z|z_1 = [-4, 4], z_2 = 0\}$*

**Figure 4:** CVs are highly correlated with the dihedral angles  $(\phi, \psi)$ .

**Figure 5:** Mean prediction  $\mathbf{x} = \mu_{\theta}(\mathbf{z})$  of decoder  $q_{\theta}(\mathbf{x}|\mathbf{z})$ , given  $\{z|z_1 = [-4, 4], z_2 = 0\}$ .



Identified CVs show high correlation to known optimal description, the dihedral angles ( $\phi, \psi$ )

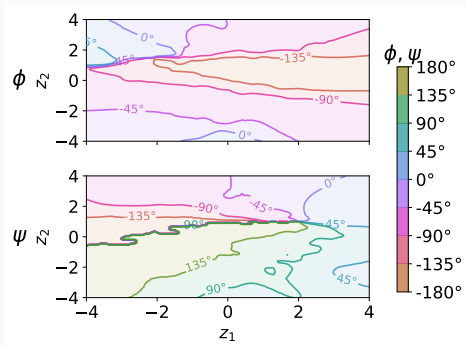


Figure 6: Predicted dihedral angles ( $\phi, \psi$ ) given the latent variables  $z \in [-4, 4]^2$ .

- Discovered CVs  $z$ : Map to the ( $\phi, \psi$ ) angles.
- Note, instead of having a distinct ( $\phi, \psi$ ) value, we obtain a distribution of CVs implied by  $r_\phi(z|x)$ .

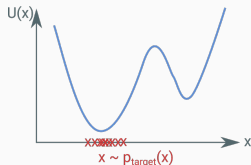


Objective utilizing data  $x^{(i)} \sim p_{\text{target}}(\mathbf{x})$

$$\min_{q(\mathbf{x})} D_{\text{KL}}(p_{\text{target}}(\mathbf{x}) || q(\mathbf{x})) \leq -\mathbb{E}_{r(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) - \log r(\mathbf{z}|\mathbf{x}; \phi)]$$

## Trade-off

- ✓ Revealing latent lower-dimensional embedding.
- ✓ Generative model, throughout Bayesian formulation feasible.
- ✗ Explorative capabilities are limited.
- ✗ Biased data-based approximation of  $p_{\text{target}}(\mathbf{x})$  (e.g. unseen modes) yields biased models.
- ✗ Biased predictions without being aware of it.



# Embedded-Physics Approach

---

## Incorporate available physics of $p_{\text{target}}(\mathbf{x})$

Instead of the forward KL-divergence,

$$\min_{q(\mathbf{x})} D_{KL}(p_{\text{target}}(\mathbf{x})||q(\mathbf{x})) = \min_{q(\mathbf{x})} \left[ - \int p_{\text{target}}(\mathbf{x}) \log \frac{q(\mathbf{x})}{p_{\text{target}}(\mathbf{x})} d\mathbf{x} \right],$$

## Incorporate available physics of $p_{\text{target}}(\mathbf{x})$

Instead of the forward KL-divergence,

$$\min_{q(\mathbf{x})} D_{\text{KL}}(p_{\text{target}}(\mathbf{x})||q(\mathbf{x})) = \min_{q(\mathbf{x})} \left[ - \int p_{\text{target}}(\mathbf{x}) \log \frac{q(\mathbf{x})}{p_{\text{target}}(\mathbf{x})} d\mathbf{x} \right],$$

“flip” the distance measure:

$$\min_{q(\mathbf{x})} D_{\text{KL}}(q(\mathbf{x})||p_{\text{target}}(\mathbf{x})) = \min_{q(\mathbf{x})} \left[ - \int q(\mathbf{x}) \log \frac{p_{\text{target}}(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \right].$$

## Incorporate available physics of $p_{\text{target}}(\mathbf{x})$

Instead of the forward KL-divergence,

$$\min_{q(\mathbf{x})} D_{\text{KL}}(p_{\text{target}}(\mathbf{x})||q(\mathbf{x})) = \min_{q(\mathbf{x})} \left[ - \int p_{\text{target}}(\mathbf{x}) \log \frac{q(\mathbf{x})}{p_{\text{target}}(\mathbf{x})} d\mathbf{x} \right],$$

“flip” the distance measure:

$$\min_{q(\mathbf{x})} D_{\text{KL}}(q(\mathbf{x})||p_{\text{target}}(\mathbf{x})) = \min_{q(\mathbf{x})} \left[ - \int q(\mathbf{x}) \log \frac{p_{\text{target}}(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \right].$$

- Forward KL-divergence involves expectations with respect to  $p_{\text{target}}(\mathbf{x}) \rightarrow$  Requires data.
- Reverse KL-divergence: **Evaluate**  $\log p_{\text{target}}(\mathbf{x})$  (i.e. evaluate  $U(\mathbf{x})$ ) at samples  $\mathbf{x}^{(i)} \sim q(\mathbf{x})$ .
- *Problem independent and explorative CG approach.*

## The objective

$$\min_{q(\mathbf{x})} D_{KL}(q(\mathbf{x}) || p_{\text{target}}(\mathbf{x})) = \min_{q(\mathbf{x})} \left[ - \int q(\mathbf{x}) \log \frac{p_{\text{target}}(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \right]$$

## How to model $q(\mathbf{x})$ ?

- Hierarchical Variational Models [Ranganath et al. (2016)]

$$q(\mathbf{x}) = \int q(\mathbf{x}|\mathbf{z})q(\mathbf{z})d\mathbf{z}.$$

with an expressive probabilistic mapping  $q(\mathbf{x}|\mathbf{z})$ .

- ✓ Hierarchical extension facilitates the construction of an expressive  $q(\mathbf{x})$ .
- ✓ Inference is feasible by bounding the objective.

## Remarks on the objective

$$\min_{q(\mathbf{x})} D_{KL}(q(\mathbf{x}) || p_{\text{target}}(\mathbf{x})) = \max_{q(\mathbf{x})} \underbrace{\left[ \int q(\mathbf{x}) \log \frac{p_{\text{target}}(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \right]}_{\mathcal{L}}$$

Consider a parametrization  $\theta$ :

$$\mathcal{L}(\theta) = \mathbb{E}_{q(\mathbf{x}; \theta)} [\log p_{\text{target}}(\mathbf{x}) - \log q(\mathbf{x}; \theta)]$$

## Remarks on the objective

$$\min_{q(\mathbf{x})} D_{KL}(q(\mathbf{x}) || p_{\text{target}}(\mathbf{x})) = \max_{q(\mathbf{x})} \underbrace{\left[ \int q(\mathbf{x}) \log \frac{p_{\text{target}}(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \right]}_{\mathcal{L}}$$

Consider a parametrization  $\theta$ :

$$\mathcal{L}(\theta) = \mathbb{E}_{q(\mathbf{x}; \theta)} \left[ \log p_{\text{target}}(\mathbf{x}) - \log q(\mathbf{x}; \theta) \right]$$

- $\mathbb{E}_{q(\mathbf{x}; \theta)} [\log p_{\text{target}}(\mathbf{x})]$  tractable if we can sample from  $q(\mathbf{x}; \theta)$ .



## Remarks on the objective

$$\min_{q(\mathbf{x})} D_{KL}(q(\mathbf{x}) || p_{\text{target}}(\mathbf{x})) = \max_{q(\mathbf{x})} \underbrace{\left[ \int q(\mathbf{x}) \log \frac{p_{\text{target}}(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \right]}_{\mathcal{L}}$$

Consider a parametrization  $\theta$ :

$$\mathcal{L}(\theta) = \mathbb{E}_{q(\mathbf{x}; \theta)} \left[ \log p_{\text{target}}(\mathbf{x}) - \log q(\mathbf{x}; \theta) \right]$$

- $\mathbb{E}_{q(\mathbf{x}; \theta)} [\log p_{\text{target}}(\mathbf{x})]$  tractable if we can sample from  $q(\mathbf{x}; \theta)$ .
- $-\mathbb{E}_{q(\mathbf{x}; \theta)} [\log q(\mathbf{x}; \theta)]$  is the entropy of  $q(\mathbf{x}; \theta)$ ,  $\mathbb{H}(q(\mathbf{x}; \theta))$ .
  - It comprises an integration with respect to  $\mathbf{z}$ :  $q(\mathbf{x}) = \int q(\mathbf{x}|\mathbf{z})q(\mathbf{z}) d\mathbf{z}$
  - In general not analytically tractable.

**Construct a tractable (lower-)bound on  $-\mathbb{E}_{q(\mathbf{x}; \theta)} [\log q(\mathbf{x}; \theta)]$ .**

Lower bound the entropy  $\mathbb{H}(q(\mathbf{x}; \boldsymbol{\theta}))$

$$-\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] = -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} \left[ \log q(\mathbf{x}; \boldsymbol{\theta}) + \underbrace{D_{\text{KL}}(q_{\text{P}}(\mathbf{z}|\mathbf{x}) || q_{\text{P}}(\mathbf{z}|\mathbf{x}))}_{=0} \right]$$

$$-\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] \geq -\mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [\log q(\mathbf{z}) + \log q(\mathbf{x}|\mathbf{z}) - \log r(\mathbf{z}|\mathbf{x})]$$

## Tractable objective

By employing the entropy bound derived before in

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{x};\boldsymbol{\theta})} [\log p_{\text{target}}(\mathbf{x}) - \log q(\mathbf{x}; \boldsymbol{\theta})],$$

we obtain a tractable lower-bound for the variational approach with,

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [\log p_{\text{target}}(\mathbf{x}) + \log r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi}) - \log q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}; \boldsymbol{\theta})].$$

**Proposed approach does not require any data and any physical insight - but learns from the potential energy and the force field**

How to connect the objective with the physics, e.g. the atomistic potential  $U(\mathbf{x})$ ?

With  $p_{\text{target}} \propto -\beta U(\mathbf{x})$ :

$$\begin{aligned}\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi}) &= -\beta \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [U(\mathbf{x})] + \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [\log r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi}) - \log q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}; \boldsymbol{\theta})] \\ &= -\beta \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} \left[ \underbrace{U(\mathbf{x}) - \frac{1}{\beta} \log r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi})}_{\tilde{U}(\mathbf{x}, \mathbf{z})} \right] + \mathbb{H}(q(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}))\end{aligned}$$

How to connect the objective with the physics, e.g. the atomistic potential  $U(\mathbf{x})$ ?

With  $p_{\text{target}} \propto -\beta U(\mathbf{x})$ :

$$\begin{aligned}\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi}) &= -\beta \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [U(\mathbf{x})] + \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [\log r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi}) - \log q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}; \boldsymbol{\theta})] \\ &= -\beta \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} \left[ \underbrace{U(\mathbf{x}) - \frac{1}{\beta} \log r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi})}_{\tilde{U}(\mathbf{x}, \mathbf{z})} \right] + \mathbb{H}(q(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}))\end{aligned}$$

How to connect the objective with the physics, e.g. the atomistic potential  $U(\mathbf{x})$ ?

With  $p_{\text{target}} \propto -\beta U(\mathbf{x})$ :

$$\begin{aligned}\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi}) &= -\beta \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [U(\mathbf{x})] + \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [\log r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi}) - \log q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}; \boldsymbol{\theta})] \\ &= -\beta \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} \left[ \underbrace{U(\mathbf{x}) - \frac{1}{\beta} \log r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi})}_{\tilde{U}(\mathbf{x}, \mathbf{z})} \right] + \mathbb{H}(q(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}))\end{aligned}$$

Maximization of  $\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi})$  as balance between:

- Minimization of the averaged (joint) potential energy  $\tilde{U}(\mathbf{x}, \mathbf{z})$ .
- Maximization of the entropy of the generative model  $\mathbb{H}(q(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}))$ .

## Gradient estimation of the term involving $U(\mathbf{x})$

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}, \phi) = \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} \left[ -\beta U(\mathbf{x}) + \log r(\mathbf{z}|\mathbf{x}; \phi) - \log q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}; \boldsymbol{\theta}) \right].$$

Consider  $q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}), \mathbf{S}_{\boldsymbol{\theta}}(\mathbf{z}))$  and  $q(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\frac{\partial}{\partial \boldsymbol{\theta}} \left\langle \left\langle -\beta U(\mathbf{x}) \right\rangle_{q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})} \right\rangle_{q(\mathbf{z})}$$

## Gradient estimation of the term involving $U(\mathbf{x})$

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}, \phi) = \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} \left[ -\beta U(\mathbf{x}) + \log r(\mathbf{z}|\mathbf{x}; \phi) - \log q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}; \boldsymbol{\theta}) \right].$$

Consider  $q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}), \mathbf{S}_{\boldsymbol{\theta}}(\mathbf{z}))$  and  $q(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\frac{\partial}{\partial \boldsymbol{\theta}} \left\langle \left\langle -\beta U(\mathbf{x}) \right\rangle_{q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})} \right\rangle_{q(\mathbf{z})}$$

- Estimator highly affected by *noise*.
- Reparametrize  $\mathbf{x}$  by auxiliary random variable  $\boldsymbol{\epsilon}$  and differentiable transformation  $\mathbf{g}(\boldsymbol{\epsilon}; \mathbf{z})$ :

$$\begin{aligned} \mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z}) &= \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\epsilon}; \mathbf{z}) \\ &= \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}) + \boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{z}) \odot \boldsymbol{\epsilon} \quad \text{with } p(\boldsymbol{\epsilon}) = \mathcal{N}(\mathbf{0}, \mathbf{I}). \end{aligned}$$



## Gradient estimation of the term involving $U(\mathbf{x})$

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} \left[ -\beta U(\mathbf{x}) + \log r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi}) - \log q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}; \boldsymbol{\theta}) \right].$$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}} \left\langle \left\langle -\beta U(\mathbf{x}) \right\rangle_{q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})} \right\rangle_{q(\mathbf{z})} &= \left\langle \left\langle -\beta \frac{\partial U(\mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z}))}{\partial \boldsymbol{\theta}} \right\rangle_{p(\boldsymbol{\epsilon})} \right\rangle_{q(\mathbf{z})} \\ &= \left\langle \left\langle -\beta \frac{\partial U(\mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z}))}{\partial \mathbf{x}} \frac{\partial \mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z})}{\partial \boldsymbol{\theta}} \right\rangle_{p(\boldsymbol{\epsilon})} \right\rangle_{q(\mathbf{z})} \\ &= \left\langle \left\langle \beta \underbrace{F(\mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z}))}_{\text{Atomistic force-field}} \frac{\partial \mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z})}{\partial \boldsymbol{\theta}} \right\rangle_{p(\boldsymbol{\epsilon})} \right\rangle_{q(\mathbf{z})} \end{aligned}$$

## Gradient estimation of the term involving $U(\mathbf{x})$

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}, \phi) = \mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} \left[ -\beta U(\mathbf{x}) + \log r(\mathbf{z}|\mathbf{x}; \phi) - \log q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}; \boldsymbol{\theta}) \right].$$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}} \left\langle \left\langle -\beta U(\mathbf{x}) \right\rangle_{q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})} \right\rangle_{q(\mathbf{z})} &= \left\langle \left\langle -\beta \frac{\partial U(\mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z}))}{\partial \boldsymbol{\theta}} \right\rangle_{p(\boldsymbol{\epsilon})} \right\rangle_{q(\mathbf{z})} \\ &= \left\langle \left\langle -\beta \frac{\partial U(\mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z}))}{\partial \mathbf{x}} \frac{\partial \mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z})}{\partial \boldsymbol{\theta}} \right\rangle_{p(\boldsymbol{\epsilon})} \right\rangle_{q(\mathbf{z})} \\ &= \left\langle \left\langle \beta \underbrace{F(\mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z}))}_{\text{Atomistic force-field}} \frac{\partial \mathbf{x}(\boldsymbol{\epsilon}; \mathbf{z})}{\partial \boldsymbol{\theta}} \right\rangle_{p(\boldsymbol{\epsilon})} \right\rangle_{q(\mathbf{z})} \end{aligned}$$

### Learning by evaluating $U(\mathbf{x})$ and $F(\mathbf{x})$

Gradient estimation involves the *evaluation* of the force field  $F(\mathbf{x})$  at configurations  $\mathbf{x}^{(i)} \sim q(\mathbf{x})$ , *not* the *simulation* of  $p_{\text{target}}(\mathbf{x})$ .

## Target distribution

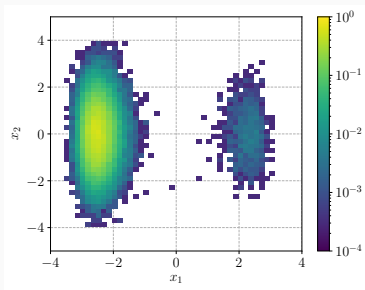


Figure 7: Reference histogram obtained by importance sampling.

$$p_{\text{target}}(\mathbf{x}) \propto e^{-\beta U(\mathbf{x})}$$

with  $\beta = 1$  and,

$$U(x_1, x_2) = \frac{1}{4}ax_1^4 - \frac{1}{2}bx_1^2 + cx_1 + \frac{1}{2}dx_2^2.$$

- Bistable in  $x_1$  and harmonic in  $x_2$ .
- Two distinct modes, one less pronounced.
- Random-walk MCMC is not able to capture both modes properly.

## Model details

Three components to specify:

1. Latent representation:

$$q(\mathbf{z}) = \mathcal{N}(\mathbf{0}, I)$$

2. Probabilistic **d**ecoder:

$$q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}), \mathbf{S}_{\boldsymbol{\theta}} = \text{diag}(\boldsymbol{\sigma}_{\boldsymbol{\theta}}^2(\mathbf{z})))$$

3. Probabilistic **e**ncoder:

$$r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi}) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\phi}}(\mathbf{x}), \mathbf{S}_{\boldsymbol{\phi}} = \text{diag}(\boldsymbol{\sigma}_{\boldsymbol{\phi}}^2(\mathbf{x})))$$

## Model details

Three components to specify:

1. Latent representation:

$$q(\mathbf{z}) = \mathcal{N}(\mathbf{0}, I)$$

2. Probabilistic **d**ecoder:

$$q(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}), \mathbf{S}_{\boldsymbol{\theta}} = \text{diag}(\boldsymbol{\sigma}_{\boldsymbol{\theta}}^2(\mathbf{z})))$$

3. Probabilistic **e**ncoder:

$$r(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi}) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\phi}}(\mathbf{x}), \mathbf{S}_{\boldsymbol{\phi}} = \text{diag}(\boldsymbol{\sigma}_{\boldsymbol{\phi}}^2(\mathbf{x})))$$

## Training details

- Tempering during training is employed, initially start with  $\beta_s = 1e - 10$  and  $\bar{p}_{\text{target}}(\mathbf{x}) \propto e^{-\beta_s U(\mathbf{x})}$ .
- Convergence with  $\beta_s$ .
- Increase  $\beta_s$  such that  $D_{\text{KL}}(q(\mathbf{x}, \mathbf{z}) || \bar{p}(\mathbf{x}, \mathbf{z}))$  does not exceed a threshold.
- Repeat until  $\beta_s = \beta$ .

## Bounds of the KL-divergence and predictions during training

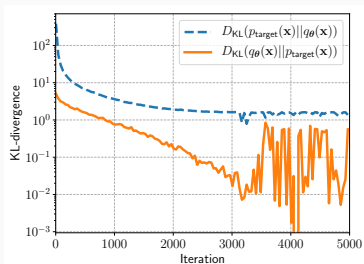


Figure 8: Upper- and lower bound of the training objective.

## Observable estimation

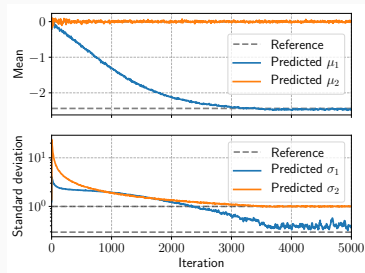


Figure 9: Mean and standard deviation compared to reference.

Predictive distribution  $q(\mathbf{x}) = \int q(\mathbf{x}|\mathbf{z})q(\mathbf{z})d\mathbf{z}$

Potential energy (e.g. at  $x_2 = 0$ )

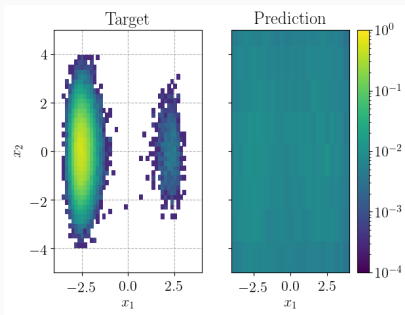


Figure 10: Target (left) and prediction (right) during learning.

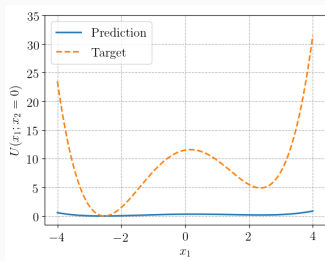
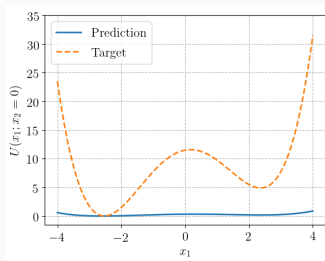


Figure 11: Reference potential energy  $U(\mathbf{x})$  and potential energy estimated by predictive distribution  $U_p(\mathbf{x}) \propto -\frac{1}{\beta} \log q(\mathbf{x})$  (noisy) at  $\{\mathbf{x}|x_2 = 0\}$

$$\text{Predictive distribution } q(\mathbf{x}) = \int q(\mathbf{x}|\mathbf{z})q(\mathbf{z})d\mathbf{z}$$

**Figure 10:** Target (left) and prediction (right) during learning.

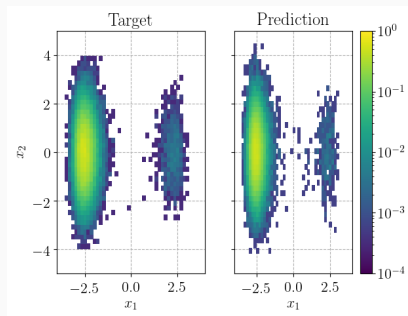
Potential energy (e.g. at  $x_2 = 0$ )



**Figure 11:** Reference potential energy  $U(\mathbf{x})$  and potential energy estimated by predictive distribution  $U_p(\mathbf{x}) \propto -\frac{1}{\beta} \log q(\mathbf{x})$  (noisy) at  $\{\mathbf{x}|x_2 = 0\}$

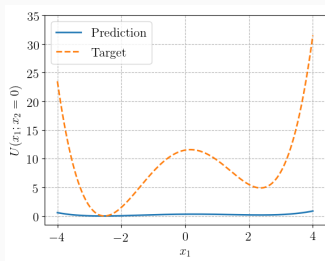


Predictive distribution  $q(x) = \int q(x|z)q(z)dz$



**Figure 10:** Target (left) and prediction (right) during learning.

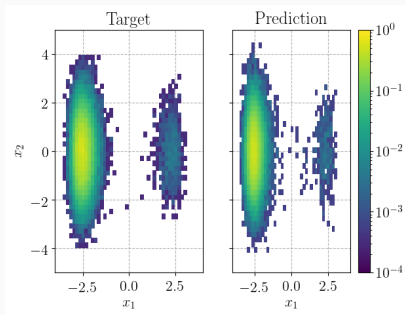
Potential energy (e.g. at  $x_2 = 0$ )



**Figure 11:** Reference potential energy  $U(x)$  and potential energy estimated by predictive distribution  $U_p(x) \propto -\frac{1}{\beta} \log q(x)$  (noisy) at  $\{x|x_2 = 0\}$

Predictive distribution  $q(x) = \int q(x|z)q(z)dz$

Potential energy (e.g. at  $x_2 = 0$ )

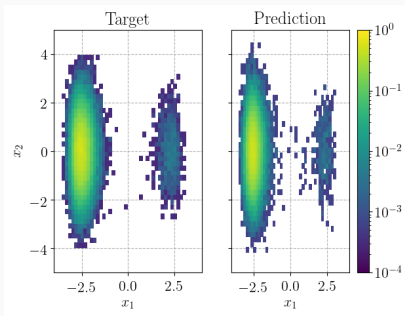


**Figure 10:** Target (left) and prediction (right) during learning.

**Figure 11:** Reference potential energy  $U(x)$  and potential energy estimated by predictive distribution

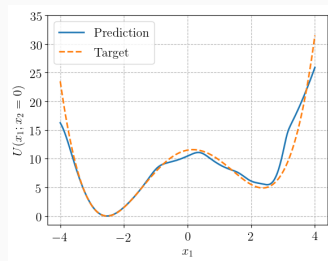
$U_p(x) \propto -\frac{1}{\beta} \log q(x)$  (noisy) at  $\{x|x_2 = 0\}$

Predictive distribution  $q(x) = \int q(x|z)q(z)dz$



**Figure 10:** Target (left) and prediction (right) during learning.

Potential energy (e.g. at  $x_2 = 0$ )



**Figure 11:** Reference potential energy  $U(x)$  and potential energy estimated by predictive distribution  $U_p(x) \propto -\frac{1}{\beta} \log q(x)$  (noisy) at  $\{x|x_2 = 0\}$

## What CVs are learned (without MD data)?

Latent embedding

CV-values  $z$  for *predicted*  $x$

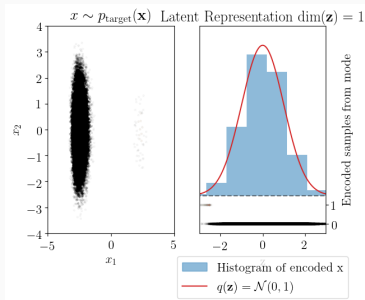


Figure 12: Reference samples encoded in the latent space.

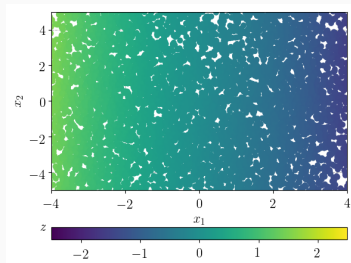
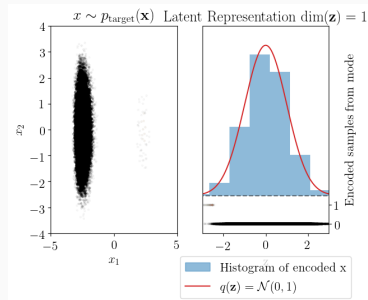


Figure 13: Samples  $x^{(i)} \sim q(\mathbf{x}; \theta)$  encoded with  $r(\mathbf{z}|\mathbf{x}; \phi)$ . Color of scatter indicate CV-values.

## What CVs are learned (without MD data)?

Latent embedding

CV-values  $z$  for *predicted*  $x$



**Figure 12:** Reference samples encoded in the latent space.

**Figure 13:** Samples  $x^{(i)} \sim q(x; \theta)$  encoded with  $r(z|x; \phi)$ . Color of scatter indicate CV-values.

## What CVs are learned (without MD data)?

Latent embedding

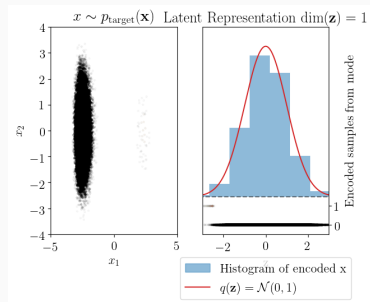


Figure 12: Reference samples encoded in the latent space.

CV-values  $z$  for *predicted*  $x$

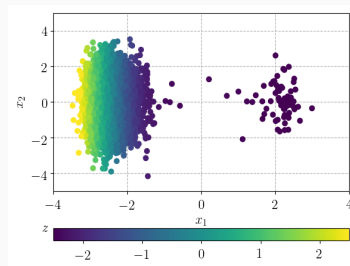


Figure 13: Samples  $x^{(i)} \sim q(x; \theta)$  encoded with  $r(\mathbf{z}|\mathbf{x}; \phi)$ . Color of scatter indicate CV-values.

Without using data, one is able to learn CVs capturing multi-modality.

## What CVs are learned (without MD data)?

Latent embedding

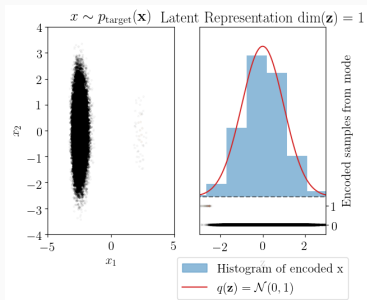


Figure 12: Reference samples encoded in the latent space.

CV-values  $z$  for predicted  $x$

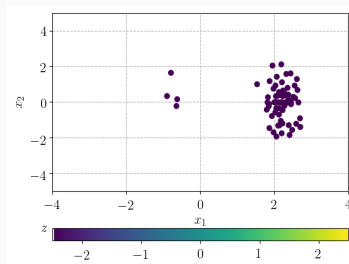


Figure 13: Samples  $x^{(i)} \sim q(x; \theta)$  encoded with  $r(z|x; \phi)$ . Color of scatter indicate CV-values.

## What CVs are learned (without MD data)?

Latent embedding

CV-values  $z$  for *predicted*  $x$

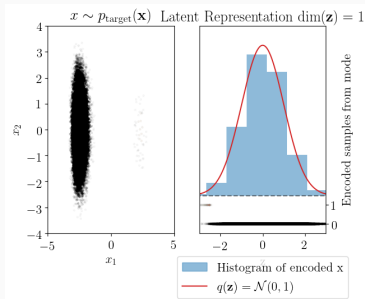


Figure 12: Reference samples encoded in the latent space.

Figure 13: Samples  $x^{(i)} \sim q(x; \theta)$  encoded with  $r(z|x; \phi)$ . Color of scatter indicate CV-values.



## What CVs are learned (without MD data)?

Latent embedding

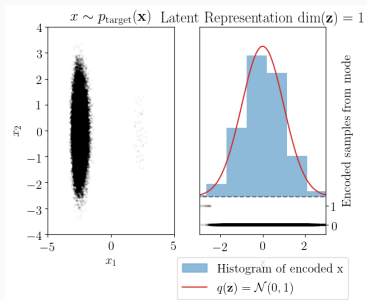


Figure 12: Reference samples encoded in the latent space.

CV-values  $z$  for *predicted*  $x$

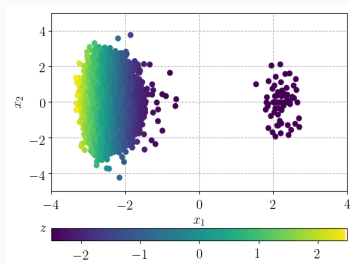


Figure 13: Samples  $x^{(i)} \sim q(x; \theta)$  encoded with  $r(z|x; \phi)$ . Color of scatter indicate CV-values.

## Significant modes for ALA-2

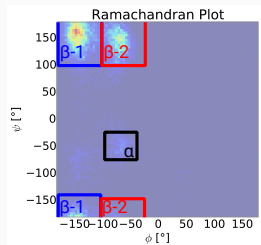
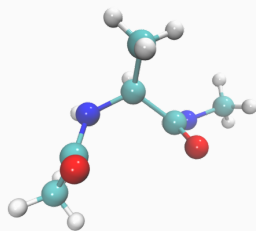
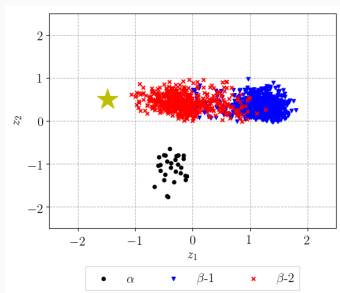


Figure 14: Reference modes

- What are relevant CVs of the system?
- The identification without data would facilitate the construction of enhanced sampling methods or biasing potentials.
- *Data-driven* approach: How to produce data accompanying whole configurational space while not being able to sample properly?

## Latent representation and prediction

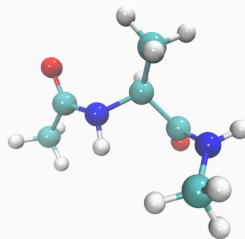
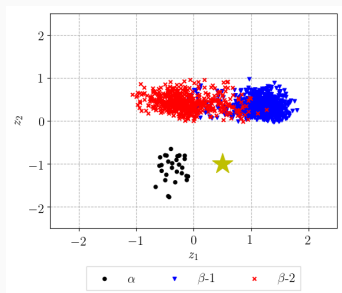


**Figure 15:** Latent representation of reference samples from multiple configurations (left). *Predicted* atomistic configurations (right), given the CV as indicated left.

## Latent representation and prediction

**Figure 15:** Latent representation of reference samples from multiple configurations (left). *Predicted* atomistic configurations (right), given the CV as indicated left.

## Latent representation and prediction



**Figure 15:** Latent representation of reference samples from multiple configurations (left). *Predicted* atomistic configurations (right), given the CV as indicated left.

**Can we learn characteristics just by *evaluating* the force-field?**

Characteristics (correlated to  $\phi - \psi$  angles) are learned and yield atomistic configurations.

## Data-Driven

$$D_{\text{KL}}(p_{\text{target}}(\mathbf{x})||q(\mathbf{x}))$$

- Robust framework of CG in small data regime.
- Requires data from different modes.
- Generative approach.
- Quantification of epistemic uncertainty.

## Outline

- Uncertainty quantification for variational approach.
- Mixed formulation:

$$F(\phi, \theta) = \alpha D_{\text{KL}}(p_{\text{target}}(\mathbf{x})||q(\mathbf{x})) + (1 - \alpha) D_{\text{KL}}(q(\mathbf{x})||p_{\text{target}}(\mathbf{x}))$$

## Embedded-Physics

$$D_{\text{KL}}(q(\mathbf{x})||p_{\text{target}}(\mathbf{x}))$$

- Generative approach.
- No physical insight presumed. Instead: reveal insight while learning.
- Explorative capabilities - **no data required**.
- Problem independent as long as force-field is accessible.

Thank you!

## Encoding network

Linear layer	Input dimension	Output dimension	Activation layer	Activation function
$l_{\phi}^{(1)}$	$\dim(\mathbf{x})$	$d_1$	$a^{(1)}$	SeLu <sup>1</sup>
$l_{\phi}^{(2)}$	$d_1$	$d_2$	$a^{(2)}$	SeLu
$l_{\phi}^{(3)}$	$d_2$	$d_3$	$a^{(3)}$	Log Sigmoid <sup>2</sup>
$l_{\phi}^{(4)}$	$d_3$	$\dim(\mathbf{z})$	None	-
$l_{\phi}^{(5)}$	$d_3$	$\dim(\mathbf{z})$	None	-

**Table 1:** Network specification of the encoding neural network with  $d_1 = 50$ ,  $d_2 = 100$ , and  $d_3 = 100$ .

## Decoding network

Linear layer	Input dimension	Output dimension	Activation layer	Activation function
$l_{\theta}^{(1)}$	$\dim(\mathbf{z})$	$d_3$	$\tilde{a}^{(1)}$	Tanh
$l_{\theta}^{(2)}$	$d_3$	$d_2$	$\tilde{a}^{(2)}$	Tanh
$l_{\theta}^{(3)}$	$d_2$	$d_1$	$\tilde{a}^{(3)}$	Tanh
$l_{\theta}^{(4)}$	$d_1$	$\dim(\mathbf{x})$	None	-

**Table 2:** Network specification of the decoding neural network with  $d_{\{1,2,3\}}$  as defined in Table 1.





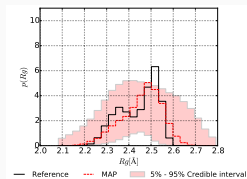
- Predictions accounting for uncertainty in  $\theta$ , with  $\theta \sim p(\theta|\mathbf{X})$ .
- Metropolis-within-Gibbs sampler [Mattei and Frelsen (2018)] corrects usage of approx. posterior  $q_\phi(\mathbf{z}|\mathbf{x})$ .

We illustrate this by computing the radius of gyration ( $R_g$ ) [Fluitt and de Pablo (2015); Carmichael and Shell (2012)] given as,

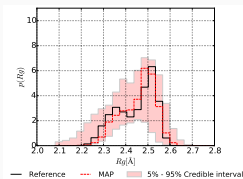
$$a_{R_g}(\mathbf{x}) = \sqrt{\frac{\sum_p m_p \|\mathbf{x}_p - \mathbf{x}_{\text{COM}}\|^2}{\sum_p m_p}}.$$

- The sum considers atoms  $p = 1, \dots, P$ .
- $m_p$  and  $\mathbf{x}_p$  denote the mass and the coordinates of each atom, respectively.
- $\mathbf{x}_{\text{COM}}$  denotes the center of mass of the peptide.

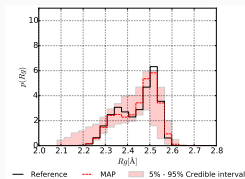
- Predictions accounting for uncertainty in  $\theta$ , with  $\theta \sim p(\theta|X)$ .
- Metropolis-within-Gibbs sampler [Mattei and Frelsen (2018)] corrects usage of approx. posterior  $q_\phi(\mathbf{z}|X)$ .



(a)  $N = 50$ .



(b)  $N = 200$ .



(c)  $N = 500$ .

**Figure 16:** Predicted radius of gyration with  $\dim(\mathbf{z}) = 2$  for various sizes  $N$  of the training dataset. The reference solution (black) is estimated by  $N = 10000$ . The shaded area represents the credible interval, reflecting the induced epistemic uncertainty from the limited amount of training data.



## Motivation

- How to avoid overfitting?
- How many  $\theta$ s are actually required?
- Can one search *across models*?

## Motivation

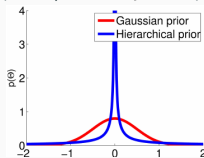
- How to avoid overfitting?
- How many  $\theta$ s are actually required?
- Can one search *across models*?

Sparsity-enforcing hierarchical prior (ARD, [MacKay 1994])

$$p(\boldsymbol{\theta}|\boldsymbol{\tau}) = \prod_j p(\theta_j|\tau_j)$$

$$\theta_j \sim \mathcal{N}(0, \tau_j^{-1})$$

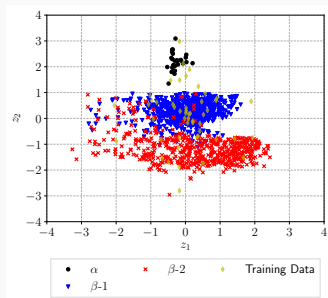
$$\tau_j \sim \text{Gamma}(a_0, b_0)$$



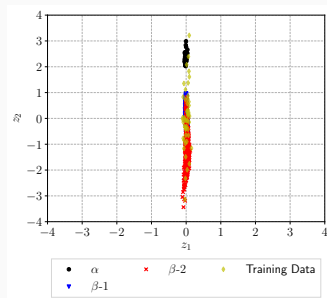
- Inner EM framework:
  - **E-step:** Estimate  $\langle \tau_j \rangle_{p(\tau_j|\theta_j)} = \frac{a_0+1/2}{b_0+\theta_j^2/2}$
  - **M-step:** Additive component to the derivative of the log-likelihood:
 
$$\frac{\partial}{\partial \tau_{c,j}} = - \langle \tau_j \rangle \theta_j$$



Sparsity prior alleviates learning physically meaningful CVs given low data (e.g.  $N=50$ )



(a) Active ARD prior.



(b) Without ARD prior.

**Figure 17:** Representation of the  $z$ -coordinates of the training data  $X$  with  $N = 50$  in the CV space (yellow diamonds). Using the trained model and the mean of  $q_\phi(\mathbf{z}|\mathbf{z})$  we computed the  $z$ -coordinates of 1527 test samples corresponding to different conformations of the alanine dipeptide to  $\alpha$  (black),  $\beta-1$  (blue), and  $\beta-2$  (red).

Lower bound the entropy  $\mathbb{H}(q(\mathbf{x}; \boldsymbol{\theta}))$

$$-\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] = -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} \left[ \log q(\mathbf{x}; \boldsymbol{\theta}) + \underbrace{D_{KL}(q_P(\mathbf{z}|\mathbf{x}) || q_P(\mathbf{z}|\mathbf{x}))}_{=0} \right]$$

Lower bound the entropy  $\mathbb{H}(q(\mathbf{x}; \boldsymbol{\theta}))$

$$\begin{aligned} -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] &= -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} \left[ \log q(\mathbf{x}; \boldsymbol{\theta}) + \underbrace{D_{\text{KL}}(q_{\text{P}}(\mathbf{z}|\mathbf{x}) || q_{\text{P}}(\mathbf{z}|\mathbf{x}))}_{=0} \right] \\ &\geq -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta}) + D_{\text{KL}}(q_{\text{P}}(\mathbf{z}|\mathbf{x}) || r(\mathbf{z}|\mathbf{x}))] \\ &= -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\mathbb{E}_{q_{\text{P}}(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{x}; \boldsymbol{\theta}) + \log q_{\text{P}}(\mathbf{z}|\mathbf{x}) - \log r(\mathbf{z}|\mathbf{x})]] . \end{aligned}$$

Lower bound the entropy  $\mathbb{H}(q(\mathbf{x}; \boldsymbol{\theta}))$ 

$$\begin{aligned}
 -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] &= -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} \left[ \log q(\mathbf{x}; \boldsymbol{\theta}) + \underbrace{D_{\text{KL}}(q_{\text{P}}(\mathbf{z}|\mathbf{x}) || q_{\text{P}}(\mathbf{z}|\mathbf{x}))}_{=0} \right] \\
 &\geq -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta}) + D_{\text{KL}}(q_{\text{P}}(\mathbf{z}|\mathbf{x}) || r(\mathbf{z}|\mathbf{x}))] \\
 &= -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\mathbb{E}_{q_{\text{P}}(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{x}; \boldsymbol{\theta}) + \log q_{\text{P}}(\mathbf{z}|\mathbf{x}) - \log r(\mathbf{z}|\mathbf{x})]] .
 \end{aligned}$$

Employ  $\log q_{\text{P}}(\mathbf{z}|\mathbf{x}) = \log q(\mathbf{z}) + \log q(\mathbf{x}|\mathbf{z}) - \log q(\mathbf{x})$ :

$$\begin{aligned}
 -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] &\geq -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\mathbb{E}_{q_{\text{P}}(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{x}; \boldsymbol{\theta}) + \log q(\mathbf{z}) + \log q(\mathbf{x}|\mathbf{z}) - \log q(\mathbf{x}) - \\
 &= -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\mathbb{E}_{q_{\text{P}}(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z}) + \log q(\mathbf{x}|\mathbf{z}) - \log r(\mathbf{z}|\mathbf{x})]]
 \end{aligned}$$



Lower bound the entropy  $\mathbb{H}(q(\mathbf{x}; \boldsymbol{\theta}))$ 

$$\begin{aligned}
 -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] &= -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} \left[ \log q(\mathbf{x}; \boldsymbol{\theta}) + \underbrace{D_{\text{KL}}(q_{\text{P}}(\mathbf{z}|\mathbf{x}) || q_{\text{P}}(\mathbf{z}|\mathbf{x}))}_{=0} \right] \\
 &\geq -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta}) + D_{\text{KL}}(q_{\text{P}}(\mathbf{z}|\mathbf{x}) || r(\mathbf{z}|\mathbf{x}))] \\
 &= -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\mathbb{E}_{q_{\text{P}}(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{x}; \boldsymbol{\theta}) + \log q_{\text{P}}(\mathbf{z}|\mathbf{x}) - \log r(\mathbf{z}|\mathbf{x})]] .
 \end{aligned}$$

Employ  $\log q_{\text{P}}(\mathbf{z}|\mathbf{x}) = \log q(\mathbf{z}) + \log q(\mathbf{x}|\mathbf{z}) - \log q(\mathbf{x})$ :

$$\begin{aligned}
 -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] &\geq -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\mathbb{E}_{q_{\text{P}}(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{x}; \boldsymbol{\theta}) + \log q(\mathbf{z}) + \log q(\mathbf{x}|\mathbf{z}) - \log q(\mathbf{x}) - \\
 &= -\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\mathbb{E}_{q_{\text{P}}(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z}) + \log q(\mathbf{x}|\mathbf{z}) - \log r(\mathbf{z}|\mathbf{x})]]
 \end{aligned}$$

$$-\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\theta})} [\log q(\mathbf{x}; \boldsymbol{\theta})] \geq -\mathbb{E}_{q(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})} [\log q(\mathbf{z}) + \log q(\mathbf{x}|\mathbf{z}) - \log r(\mathbf{z}|\mathbf{x})]$$

## References

---

*Reaction Rate Theory: Faraday Discussion 195 (Faraday Discussions).*

Royal Society of Chemistry, 2017. ISBN 178262483X. URL

<https://www.amazon.com/>

[Reaction-Rate-Theory-Discussion-Discussions/dp/178262483X?SubscriptionId=AKIAIOBINVZYXZQZ2U3A&tag=chimbiori05-20&linkCode=xm2&camp=2025&creative=165953&creativeASIN=178262483X](https://www.amazon.com/Reaction-Rate-Theory-Discussion-Discussions/dp/178262483X?SubscriptionId=AKIAIOBINVZYXZQZ2U3A&tag=chimbiori05-20&linkCode=xm2&camp=2025&creative=165953&creativeASIN=178262483X).

Matthew J. Beal and Zoubin Ghahramani. Variational bayesian learning of directed graphical models with hidden variables.

*Bayesian Anal.*, 1(4):793–831, 12 2006. doi: 10.1214/06-BA126. URL

<https://doi.org/10.1214/06-BA126>.

- I. Bilonis and P.S. Koutsourelakis. Free energy computations by minimization of kullback–leibler divergence: An efficient adaptive biasing potential method for sparse representations. *Journal of Computational Physics*, 231(9):3849 – 3870, 2012. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2012.01.033>. URL <http://www.sciencedirect.com/science/article/pii/S0021999112000630>.
- Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–3937, 2016. ISSN 0027-8424. doi: [10.1073/pnas.1517384113](https://doi.org/10.1073/pnas.1517384113). URL <https://www.pnas.org/content/113/15/3932>.

- Scott P. Carmichael and M. Scott Shell. A new multiscale algorithm and its application to coarse-grained peptide models for self-assembly. *The Journal of Physical Chemistry B*, 116(29): 8383–8393, 2012. doi: 10.1021/jp2114994. URL <https://doi.org/10.1021/jp2114994>. PMID: 22300263.
- Maghesree Chakraborty, Chenliang Xu, and Andrew D. White. Encoding and selecting coarse-grain mapping operators with hierarchical graphs, 2018.
- Pei-Yang Chen and Mark E. Tuckerman. Molecular dynamics based enhanced sampling of collective variables with very large time steps. *The Journal of Chemical Physics*, 148(2):024106, 2018. doi: 10.1063/1.4999447. URL <https://doi.org/10.1063/1.4999447>.

- Wei Chen and Andrew L Ferguson. Molecular enhanced sampling with autoencoders: On-the-fly collective variable discovery and accelerated free energy landscape exploration, 2017.
- Eric Darve, David Rodríguez-Gómez, and Andrew Pohorille. Adaptive biasing force method for scalar and vector free energy calculations. *The Journal of Chemical Physics*, 128(14):144120, 2008. doi: 10.1063/1.2829861. URL <https://doi.org/10.1063/1.2829861>.
- Andrew L. Ferguson. Bayeswham: A bayesian approach for free energy estimation, reweighting, and uncertainty quantification in the weighted histogram analysis method. *Journal of Computational Chemistry*, 38(18):1583–1605, 7 2017. ISSN 0192-8651. doi: 10.1002/jcc.24800.

- Andrew L. Ferguson, Athanassios Z. Panagiotopoulos, Pablo G. Debenedetti, and Ioannis G. Kevrekidis. Integrating diffusion maps with umbrella sampling: Application to alanine dipeptide. *The Journal of Chemical Physics*, 134(13):135103, 2011. doi: 10.1063/1.3574394. URL <https://doi.org/10.1063/1.3574394>.
- Aaron M. Fluitt and Juan J. de Pablo. An analysis of biomolecular force fields for simulations of polyglutamine in solution. *Biophysical Journal*, 109(5):1009 – 1018, 2015. ISSN 0006-3495. doi: <https://doi.org/10.1016/j.bpj.2015.07.018>. URL <http://www.sciencedirect.com/science/article/pii/S0006349515007249>.

Markos A. Katsoulakis and José Trashorras. Information loss in coarse-graining of stochastic particle dynamics. *Journal of Statistical Physics*, 122(1):115–135, Jan 2006. ISSN 1572-9613. doi: 10.1007/s10955-005-8063-1. URL <https://doi.org/10.1007/s10955-005-8063-1>.

Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2013.

Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 971–980. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/6698-self-normalizing-neural-networks.pdf>.

- Alessandro Laio and Michele Parrinello. Escaping free-energy minima. *Proceedings of the National Academy of Sciences*, 99(20): 12562–12566, 2002. ISSN 0027-8424. doi: 10.1073/pnas.202427399. URL <http://www.pnas.org/content/99/20/12562>.
- Pierre-Alexandre Mattei and Jes Frelsen. Leveraging the exact likelihood of deep latent variable models, 2018.
- W. G. Noid. Perspective: Coarse-grained models for biomolecular systems. *The Journal of Chemical Physics*, 139(9):090901, 2013. doi: 10.1063/1.4818908. URL <https://doi.org/10.1063/1.4818908>.
- Frank Noé and Hao Wu. Boltzmann generators - sampling equilibrium states of many-body systems with deep learning, 2018.



- Rajesh Ranganath, Dustin Tran, and David M. Blei. Hierarchical variational models. In *ICML*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 324–333. JMLR.org, 2016.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, pages 1278–1286, 2014. URL <http://jmlr.org/proceedings/papers/v32/rezende14.html>.
- Romelia Salomon-Ferrer, David A. Case, and Ross C. Walker. An overview of the amber biomolecular simulation package. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 3(2): 198–210. doi: 10.1002/wcms.1121. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/wcms.1121>.

- M. Scott Shell. The relative entropy is fundamental to multiscale and inverse thermodynamic problems. *The Journal of Chemical Physics*, 129(14):144108, 2008. doi: 10.1063/1.2992060. URL <https://doi.org/10.1063/1.2992060>.
- W. Clark Still, Anna Tempczyk, Ronald C. Hawley, and Thomas Hendrickson. Semianalytical treatment of solvation for molecular mechanics and dynamics. *Journal of the American Chemical Society*, 112(16):6127–6129, Aug 1990. ISSN 0002-7863. doi: 10.1021/ja00172a038. URL <https://doi.org/10.1021/ja00172a038>.

José Trashorras and Dimitrios Tsagkarogiannis. From mesoscale back to microscale: Reconstruction schemes for coarse-grained stochastic lattice systems. *SIAM Journal on Numerical Analysis*, 48(5):1647–1677, 2010. doi: 10.1137/080722382. URL <https://hal.archives-ouvertes.fr/hal-00275802>.

Christoph Wehmeyer and Frank Noé. Time-lagged autoencoders: Deep learning of slow collective variables for molecular kinetics. *The Journal of Chemical Physics*, 148(24):241703, 2018. doi: 10.1063/1.5011399. URL <https://doi.org/10.1063/1.5011399>.