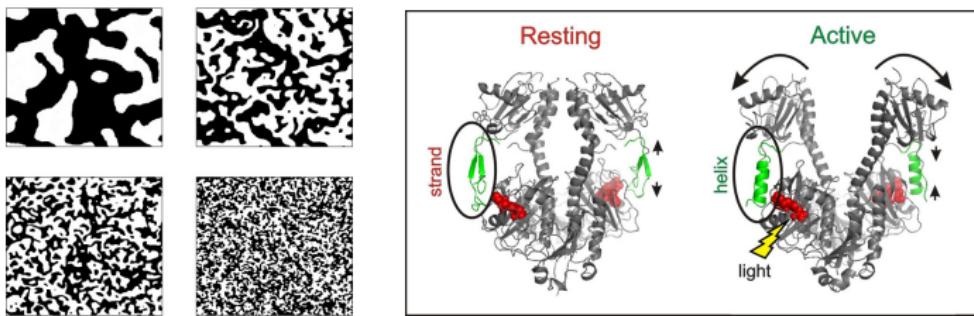


Boltzmann Generators

F. Noé¹

Group Seminar, Nov 6, 2018

Limitations of Monte Carlo Sampling

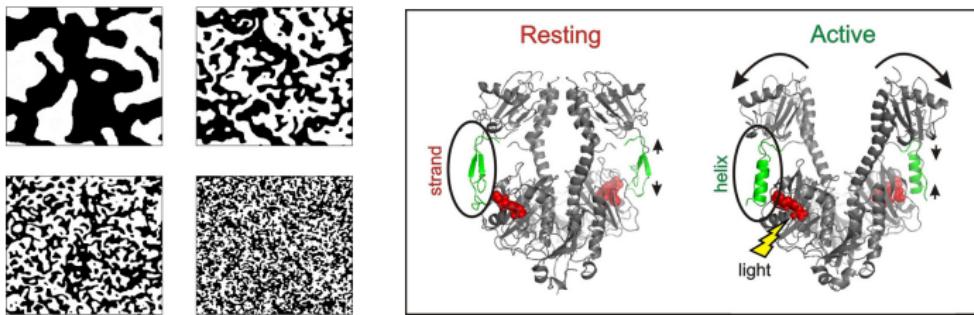


- **Input:** Reduced Potential Energy $u(\mathbf{x})$ in coordinates $\mathbf{x} \in \mathbb{R}^n$.
- **Aim:** Sample Equilibrium Distribution.

$$\mu(\mathbf{x}) \propto e^{-u(\mathbf{x})}$$

- **Problem 1:** For increasing n , the subvolume of low-energy configurations is vanishingly small compared to \mathbb{R}^n and has a complex shape.
 - Direct MD sampling (with rejection or reweighting) in configuration space is hopeless .
- **Problem 2:** Metastable states or phases
 - MCMC or MD methods with small steps converge very slowly
 - Guessing large MCMC proposal steps is hard and problem-specific.

Limitations of Monte Carlo Sampling

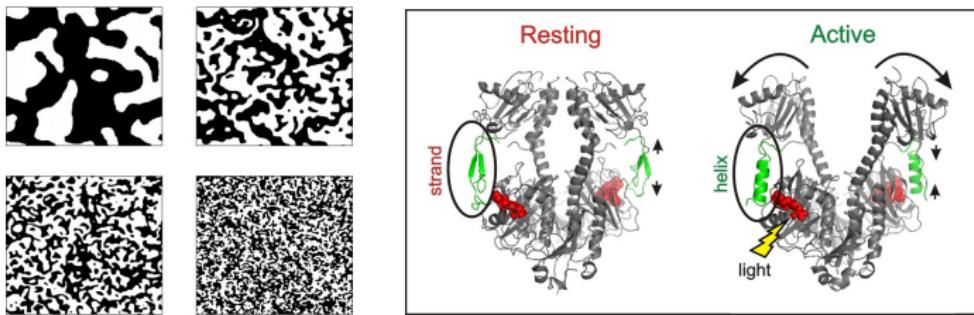


- **Input:** Reduced Potential Energy $u(\mathbf{x})$ in coordinates $\mathbf{x} \in \mathbb{R}^n$.
- **Aim:** Sample Equilibrium Distribution.

$$\mu(\mathbf{x}) \propto e^{-u(\mathbf{x})}$$

- **Problem 1:** For increasing n , the subvolume of low-energy configurations is vanishingly small compared to \mathbb{R}^n and has a complex shape.
 - Direct MD sampling (with rejection or reweighting) in configuration space is hopeless .
- **Problem 2:** Metastable states or phases
 - MCMC or MD methods with small steps converge very slowly
 - Guessing large MCMC proposal steps is hard and problem-specific.

Limitations of Monte Carlo Sampling

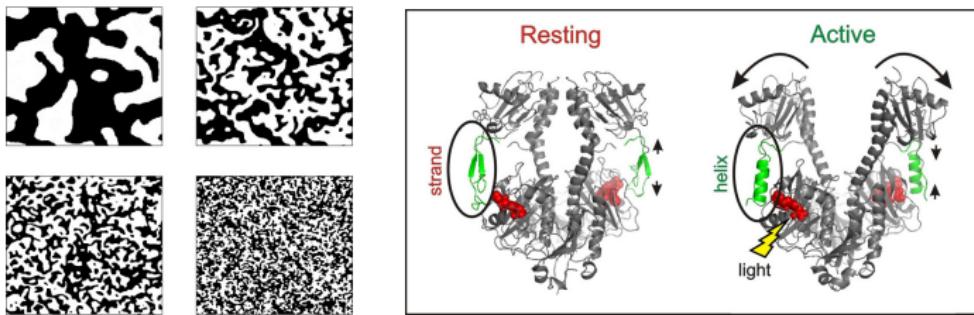


- **Input:** Reduced Potential Energy $u(\mathbf{x})$ in coordinates $\mathbf{x} \in \mathbb{R}^n$.
- **Aim:** Sample Equilibrium Distribution.

$$\mu(\mathbf{x}) \propto e^{-u(\mathbf{x})}$$

- **Problem 1:** For increasing n , the subvolume of low-energy configurations is vanishingly small compared to \mathbb{R}^n and has a complex shape.
 - Direct MD sampling (with rejection or reweighting) in configuration space is hopeless .
- **Problem 2:** Metastable states or phases
 - MCMC or MD methods with small steps converge very slowly
 - Guessing large MCMC proposal steps is hard and problem-specific.

Limitations of Monte Carlo Sampling

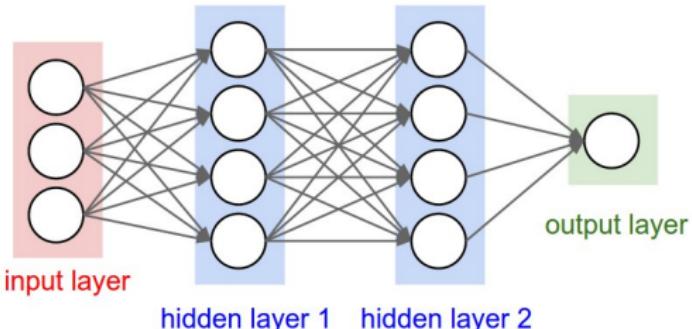


- **Input:** Reduced Potential Energy $u(\mathbf{x})$ in coordinates $\mathbf{x} \in \mathbb{R}^n$.
- **Aim:** Sample Equilibrium Distribution.

$$\mu(\mathbf{x}) \propto e^{-u(\mathbf{x})}$$

- **Problem 1:** For increasing n , the subvolume of low-energy configurations is vanishingly small compared to \mathbb{R}^n and has a complex shape.
 - Direct MD sampling (with rejection or reweighting) in configuration space is hopeless .
- **Problem 2:** Metastable states or phases
 - MCMC or MD methods with small steps converge very slowly
 - Guessing large MCMC proposal steps is hard and problem-specific.

Reminder: Multilayer Neural Networks

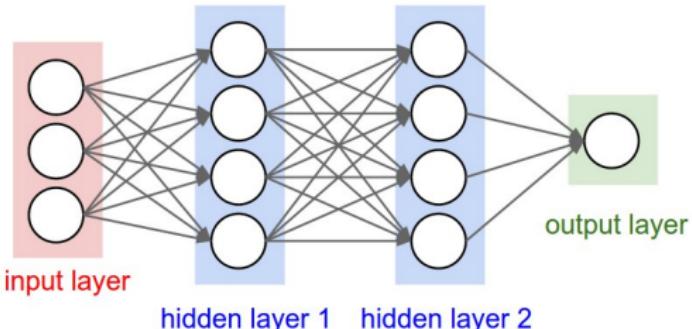


Sequence of linear and nonlinear transforms defined by the recursion:

$$\mathbf{x}^{l+1} = \sigma(\mathbf{W}^l \mathbf{x}^l + \mathbf{b}^l)$$

- L layers indexed by $l = 1, \dots, L$. Input vector $\mathbf{x}^{(0)}$ does not count as a layer.
- $\mathbf{x}^l \in \mathbb{R}^{n_l}$: Activations of n_l neurons at layer l .
- Trainable weights $\mathbf{W}^l \in \mathbb{R}^{n_{l-1} \times n_l}$ and biases $\mathbf{b}^l \in \mathbb{R}^{n_l}$ at each layer. W_{ij}^l is connecting neuron j of layer $l-1$ with neuron i of layer l .
- Nonlinear function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$.
- Output vector: $\hat{\mathbf{y}} = \mathbf{x}^L$ ($\hat{\mathbf{y}}$: network predictions. \mathbf{y} : training values)

Reminder: Multilayer Neural Networks

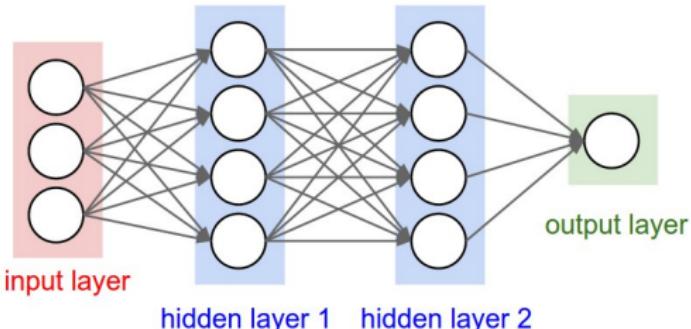


Sequence of linear and nonlinear transforms defined by the recursion:

$$\mathbf{x}^{l+1} = \sigma(\mathbf{W}^l \mathbf{x}^l + \mathbf{b}^l)$$

- L layers indexed by $l = 1, \dots, L$. Input vector $\mathbf{x}^{(0)}$ does not count as a layer.
- $\mathbf{x}^l \in \mathbb{R}^{n_l}$: Activations of n_l neurons at layer l
- Trainable weights $\mathbf{W}^l \in \mathbb{R}^{n_{l-1} \times n_l}$ and biases $\mathbf{b}^l \in \mathbb{R}^{n_l}$ at each layer. W_{ij}^l is connecting neuron j of layer $l-1$ with neuron i of layer l .
- Nonlinear function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$.
- Output vector: $\hat{\mathbf{y}} = \mathbf{x}^L$ ($\hat{\mathbf{y}}$: network predictions. \mathbf{y} : training values)

Reminder: Multilayer Neural Networks

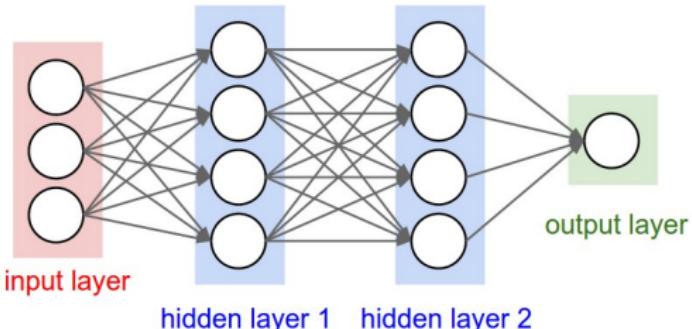


Sequence of linear and nonlinear transforms defined by the recursion:

$$\mathbf{x}^{l+1} = \sigma(\mathbf{W}^l \mathbf{x}^l + \mathbf{b}^l)$$

- L layers indexed by $l = 1, \dots, L$. Input vector $\mathbf{x}^{(0)}$ does not count as a layer.
- $\mathbf{x}^l \in \mathbb{R}^{n_l}$: Activations of n_l neurons at layer l
- Trainable weights $\mathbf{W}^l \in \mathbb{R}^{n_{l-1} \times n_l}$ and biases $\mathbf{b}^l \in \mathbb{R}^{n_l}$ at each layer. W_{ij}^l is connecting neuron j of layer $l-1$ with neuron i of layer l .
- Nonlinear function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$.
- Output vector: $\hat{\mathbf{y}} = \mathbf{x}^L$ ($\hat{\mathbf{y}}$: network predictions. \mathbf{y} : training values)

Reminder: Multilayer Neural Networks

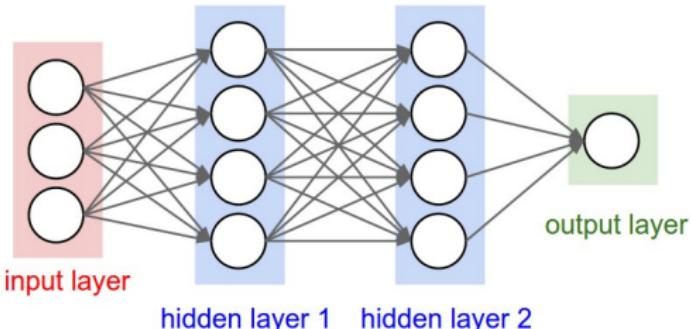


Sequence of linear and nonlinear transforms defined by the recursion:

$$\mathbf{x}^{l+1} = \sigma(\mathbf{W}^l \mathbf{x}^l + \mathbf{b}^l)$$

- L layers indexed by $l = 1, \dots, L$. Input vector $\mathbf{x}^{(0)}$ does not count as a layer.
- $\mathbf{x}^l \in \mathbb{R}^{n_l}$: Activations of n_l neurons at layer l
- Trainable weights $\mathbf{W}^l \in \mathbb{R}^{n_{l-1} \times n_l}$ and biases $\mathbf{b}^l \in \mathbb{R}^{n_l}$ at each layer. W_{ij}^l is connecting neuron j of layer $l-1$ with neuron i of layer l .
- Nonlinear function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$.
- Output vector: $\hat{\mathbf{y}} = \mathbf{x}^L$ ($\hat{\mathbf{y}}$: network predictions. \mathbf{y} : training values)

Reminder: Multilayer Neural Networks



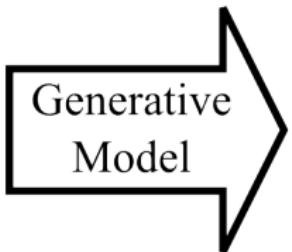
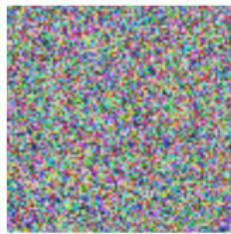
Sequence of linear and nonlinear transforms defined by the recursion:

$$\mathbf{x}^{l+1} = \sigma(\mathbf{W}^l \mathbf{x}^l + \mathbf{b}^l)$$

- L layers indexed by $l = 1, \dots, L$. Input vector $\mathbf{x}^{(0)}$ does not count as a layer.
- $\mathbf{x}^l \in \mathbb{R}^{n_l}$: Activations of n_l neurons at layer l
- Trainable weights $\mathbf{W}^l \in \mathbb{R}^{n_{l-1} \times n_l}$ and biases $\mathbf{b}^l \in \mathbb{R}^{n_l}$ at each layer. W_{ij}^l is connecting neuron j of layer $l-1$ with neuron i of layer l .
- Nonlinear function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$.
- Output vector: $\hat{\mathbf{y}} = \mathbf{x}^L$ ($\hat{\mathbf{y}}$: network predictions. \mathbf{y} : training values)

Reminder: Generative Neural Networks

Noise $\sim N(0,1)$



Reminder: Generative Neural Networks

- Idea: Learn to sample intractable $p(\mathbf{x})$ by sampling tractable latent distribution

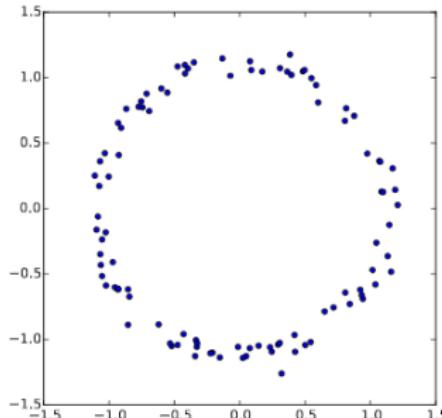
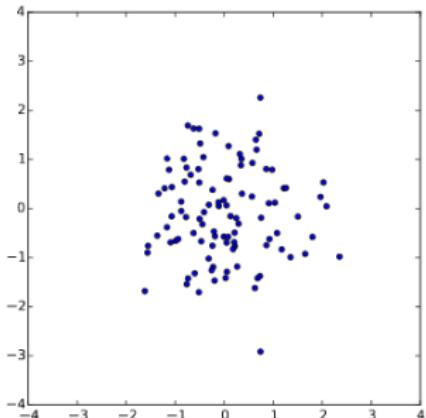
$$\mathbf{z} \sim p(\mathbf{z})$$

and perform a linear transformation to a desired distribution:

$$\mathbf{x} = G(\mathbf{z}, \theta) \sim p(\mathbf{x}).$$

- Example:

- Left: Samples from normal distribution, $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
- Right: Samples mapped through $G(\mathbf{z}) = \frac{\mathbf{z}}{10} + \frac{\mathbf{z}}{\|\mathbf{z}\|}$ to form a ring.



Reminder: Generative Neural Networks

- Idea: Learn to sample intractable $p(\mathbf{x})$ by sampling tractable latent distribution

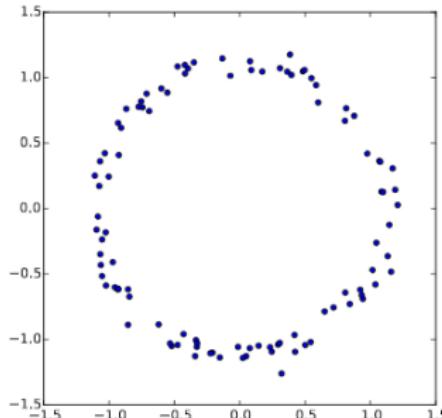
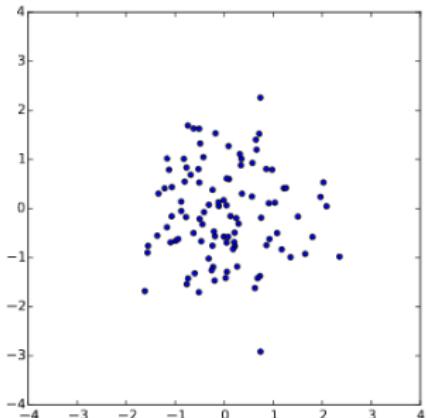
$$\mathbf{z} \sim p(\mathbf{z})$$

and perform a linear transformation to a desired distribution:

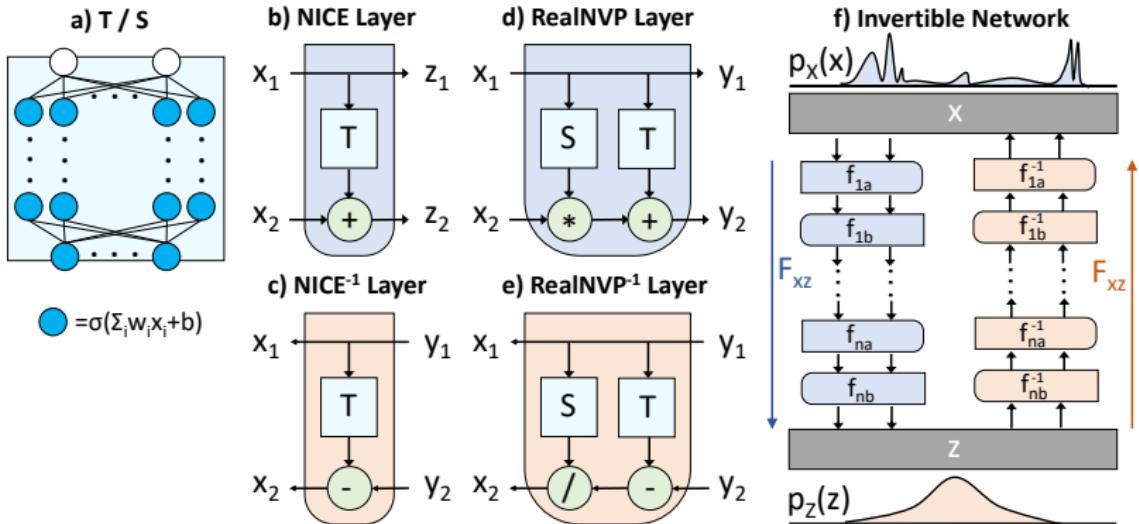
$$\mathbf{x} = G(\mathbf{z}, \theta) \sim p(\mathbf{x}).$$

- Example:

- Left: Samples from normal distribution, $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
- Right: Samples mapped through $G(\mathbf{z}) = \frac{\mathbf{z}}{10} + \frac{\mathbf{z}}{\|\mathbf{z}\|}$ to form a ring.



Boltzmann Generators



- **NICE**: Dinh, Krueger, Y. Bengio, ICLR 2015
- **RealNVP**: Dinh, Sohl-Dickstein, S. Bengio, ICLR 2017

Transformation of random variables

- Invertible transformation:

$$\mathbf{z} = F_{xz}(\mathbf{x}; \theta)$$

$$\mathbf{x} = F_{zx}(\mathbf{z}; \theta).$$

- Jacobians:

$$\mathbf{J}_{zx}(\mathbf{z}; \theta) = \left[\frac{\partial F_{zx}(\mathbf{z}; \theta)}{\partial z_1}, \dots, \frac{\partial F_{zx}(\mathbf{z}; \theta)}{\partial z_n} \right]$$

$$\mathbf{J}_{xz}(\mathbf{x}; \theta) = \left[\frac{dF_{xz}(\mathbf{x}; \theta)}{dx_1}, \dots, \frac{dF_{xz}(\mathbf{x}; \theta)}{dx_n} \right]$$

- Transformation of random variables:

$$\begin{aligned} p_X(\mathbf{x}) &= p_Z(\mathbf{z}) |\det \mathbf{J}_{zx}(\mathbf{z})|^{-1} \\ &= p_Z(T_{xz}(\mathbf{x})) |\det \mathbf{J}_{xz}(\mathbf{x})| \end{aligned} \tag{1}$$

$$\begin{aligned} p_Z(\mathbf{z}) &= p_X(\mathbf{x}) |\det \mathbf{J}_{xz}(\mathbf{x})|^{-1} \\ &= p_X(T_{zx}(\mathbf{z})) |\det \mathbf{J}_{zx}(\mathbf{z})| \end{aligned} \tag{2}$$

Transformation of random variables

- Invertible transformation:

$$\mathbf{z} = F_{xz}(\mathbf{x}; \theta)$$

$$\mathbf{x} = F_{zx}(\mathbf{z}; \theta).$$

- Jacobians:

$$\mathbf{J}_{zx}(\mathbf{z}; \theta) = \left[\frac{\partial F_{zx}(\mathbf{z}; \theta)}{\partial z_1}, \dots, \frac{\partial F_{zx}(\mathbf{z}; \theta)}{\partial z_n} \right]$$

$$\mathbf{J}_{xz}(\mathbf{x}; \theta) = \left[\frac{dF_{xz}(\mathbf{x}; \theta)}{dx_1}, \dots, \frac{dF_{xz}(\mathbf{x}; \theta)}{dx_n} \right]$$

- Transformation of random variables:

$$\begin{aligned} p_X(\mathbf{x}) &= p_Z(\mathbf{z}) |\det \mathbf{J}_{zx}(\mathbf{z})|^{-1} \\ &= p_Z(T_{xz}(\mathbf{x})) |\det \mathbf{J}_{xz}(\mathbf{x})| \end{aligned} \tag{1}$$

$$\begin{aligned} p_Z(\mathbf{z}) &= p_X(\mathbf{x}) |\det \mathbf{J}_{xz}(\mathbf{x})|^{-1} \\ &= p_X(T_{zx}(\mathbf{z})) |\det \mathbf{J}_{zx}(\mathbf{z})| \end{aligned} \tag{2}$$

Transformation of random variables

- Invertible transformation:

$$\mathbf{z} = F_{xz}(\mathbf{x}; \theta)$$

$$\mathbf{x} = F_{zx}(\mathbf{z}; \theta).$$

- Jacobians:

$$\mathbf{J}_{zx}(\mathbf{z}; \theta) = \left[\frac{\partial F_{zx}(\mathbf{z}; \theta)}{\partial z_1}, \dots, \frac{\partial F_{zx}(\mathbf{z}; \theta)}{\partial z_n} \right]$$

$$\mathbf{J}_{xz}(\mathbf{x}; \theta) = \left[\frac{dF_{xz}(\mathbf{x}; \theta)}{dx_1}, \dots, \frac{dF_{xz}(\mathbf{x}; \theta)}{dx_n} \right]$$

- Transformation of random variables:

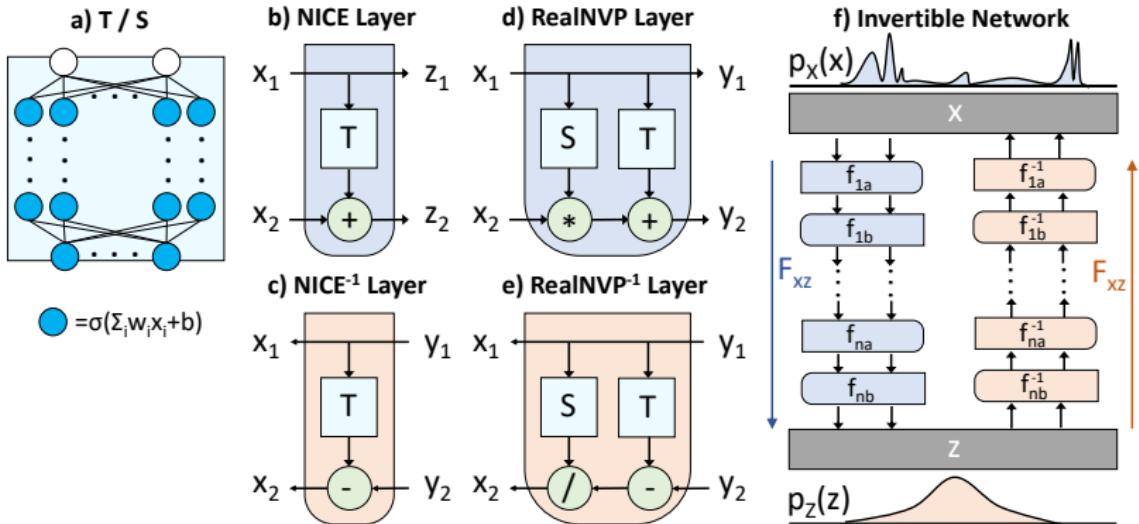
$$\begin{aligned} p_X(\mathbf{x}) &= p_Z(\mathbf{z}) |\det \mathbf{J}_{zx}(\mathbf{z})|^{-1} \\ &= p_Z(T_{xz}(\mathbf{x})) |\det \mathbf{J}_{xz}(\mathbf{x})| \end{aligned} \tag{1}$$

$$\begin{aligned} p_Z(\mathbf{z}) &= p_X(\mathbf{x}) |\det \mathbf{J}_{xz}(\mathbf{x})|^{-1} \\ &= p_X(T_{zx}(\mathbf{z})) |\det \mathbf{J}_{zx}(\mathbf{z})| \end{aligned} \tag{2}$$

Invertible network components

Layer	f_{xz}	$ \det \mathbf{J}_{xz} $	f_{zx}	$ \det \mathbf{J}_{zx} $
NICE	$\begin{aligned} \mathbf{z}_1 &= \mathbf{x}_1 \\ \mathbf{z}_2 &= \mathbf{x}_2 + T(\mathbf{x}_1; \theta) \end{aligned}$	1	$\begin{aligned} \mathbf{x}_1 &= \mathbf{z}_1 \\ \mathbf{x}_2 &= \mathbf{z}_2 - T(\mathbf{y}_1; \theta) \end{aligned}$	1
Scaling, Exp	$\mathbf{z} = e^{\mathbf{k}} \circ \mathbf{x}$	$e^{\sum_i k_i}$	$\mathbf{x} = e^{-\mathbf{k}} \circ \mathbf{z}$	$e^{-\sum_i k_i}$
RealNVP	$\begin{aligned} \mathbf{z}_1 &= \mathbf{x}_1 \\ \mathbf{z}_2 &= \mathbf{x}_2 \odot \exp(S(\mathbf{x}_1; \theta)) \quad e^{\sum_i S_i(\mathbf{x}_1; \theta)} \\ &\quad + T(\mathbf{x}_1; \theta) \end{aligned}$	$\begin{aligned} \mathbf{x}_1 &= \mathbf{z}_1 \\ \mathbf{x}_2 &= (\mathbf{z}_2 - T(\mathbf{x}_1; \theta)) \\ &\quad \odot \exp(-S(\mathbf{z}_1; \theta)) \end{aligned}$	$e^{-\sum_i S_i(\mathbf{z}_1; \theta)}$	

Boltzmann Generators



- **NICE**: Dinh, Krueger, Y. Bengio, ICLR 2015
- **RealNVP**: Dinh, Sohl-Dickstein, S. Bengio, ICLR 2017

Transformation of random variables

- Distributions:

$$\text{Prior } q_Z(\mathbf{z}) \xrightarrow{F_{zx}} p_X(\mathbf{x}) \text{ generated}$$

$$\text{Boltzmann } \mu_X(\mathbf{x}) \xrightarrow{F_{xz}} p_Z(\mathbf{z}) \text{ generated}$$

- Aim: sample configurations \mathbf{x} from **Boltzmann distribution**

$$\mu_X(\mathbf{x}) = Z_X^{-1} e^{-u(\mathbf{x})} \quad (3)$$

- Reduced energy at temperature T :

$$u(\mathbf{x}) = \frac{U(\mathbf{x})}{k_B T}$$

- Prior distribution: Sample input in \mathbf{z} from isotropic Gaussian:

$$q_Z^k(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \sigma_k^2 \mathbf{I}) = Z_Z^{-1} e^{-\frac{1}{2} \|\mathbf{z}\|^2 / \sigma_k^2}, \quad (4)$$

Prior energy:

$$u_Z^k(\mathbf{z}) = -\log q_Z^k(\mathbf{z}) = \frac{1}{2\sigma_k^2} \|\mathbf{z}\|^2 + \text{const}$$

Transformation of random variables

- Distributions:
Prior $q_Z(\mathbf{z}) \xrightarrow{F_{zx}} p_X(\mathbf{x})$ generated
Boltzmann $\mu_X(\mathbf{x}) \xrightarrow{F_{xz}} p_Z(\mathbf{z})$ generated
- Aim: sample configurations \mathbf{x} from **Boltzmann distribution**

$$\mu_X(\mathbf{x}) = Z_X^{-1} e^{-u(\mathbf{x})} \quad (3)$$

- Reduced energy at temperature T :

$$u(\mathbf{x}) = \frac{U(\mathbf{x})}{k_B T}$$

- Prior distribution: Sample input in \mathbf{z} from isotropic Gaussian:

$$q_Z^k(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \sigma_k^2 \mathbf{I}) = Z_Z^{-1} e^{-\frac{1}{2} \|\mathbf{z}\|^2 / \sigma_k^2}, \quad (4)$$

Prior energy:

$$u_Z^k(\mathbf{z}) = -\log q_Z^k(\mathbf{z}) = \frac{1}{2\sigma_k^2} \|\mathbf{z}\|^2 + \text{const}$$

Transformation of random variables

- Distributions:
Prior $q_Z(\mathbf{z}) \xrightarrow{F_{zx}} p_X(\mathbf{x})$ generated
Boltzmann $\mu_X(\mathbf{x}) \xrightarrow{F_{xz}} p_Z(\mathbf{z})$ generated
- Aim: sample configurations \mathbf{x} from **Boltzmann distribution**

$$\mu_X(\mathbf{x}) = Z_X^{-1} e^{-u(\mathbf{x})} \quad (3)$$

- Reduced energy at temperature T :

$$u(\mathbf{x}) = \frac{U(\mathbf{x})}{k_B T}$$

- Prior distribution: Sample input in \mathbf{z} from isotropic Gaussian:

$$q_Z^k(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \sigma_k^2 \mathbf{I}) = Z_Z^{-1} e^{-\frac{1}{2} \|\mathbf{z}\|^2 / \sigma_k^2}, \quad (4)$$

Prior energy:

$$u_Z^k(\mathbf{z}) = -\log q_Z^k(\mathbf{z}) = \frac{1}{2\sigma_k^2} \|\mathbf{z}\|^2 + \text{const}$$

Transformation of random variables

- Distributions:
Prior $q_Z(\mathbf{z}) \xrightarrow{F_{zx}} p_X(\mathbf{x})$ generated
Boltzmann $\mu_X(\mathbf{x}) \xrightarrow{F_{xz}} p_Z(\mathbf{z})$ generated
- Aim: sample configurations \mathbf{x} from **Boltzmann distribution**

$$\mu_X(\mathbf{x}) = Z_X^{-1} e^{-u(\mathbf{x})} \quad (3)$$

- Reduced energy at temperature T :

$$u(\mathbf{x}) = \frac{U(\mathbf{x})}{k_B T}$$

- Prior distribution:** Sample input in \mathbf{z} from isotropic Gaussian:

$$q_Z^k(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \sigma_k^2 \mathbf{I}) = Z_Z^{-1} e^{-\frac{1}{2} \|\mathbf{z}\|^2 / \sigma_k^2}, \quad (4)$$

Prior energy:

$$u_Z^k(\mathbf{z}) = -\log q_Z^k(\mathbf{z}) = \frac{1}{2\sigma_k^2} \|\mathbf{z}\|^2 + \text{const}$$

Latent KL divergence

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback-Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- KL divergence between q and p , entropy H_q :

$$\begin{aligned} \text{KL}(q \parallel p) &= \int q(\mathbf{x}) [\log q(\mathbf{x}) - \log p(\mathbf{x})] d\mathbf{x}, \\ &= H_q - \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}, \end{aligned}$$

- KL divergence between generated $p_X(\mathbf{x})$ and Boltzmann distribution:

$$\begin{aligned} \text{KL}_\theta [q_Z \parallel p_Z] &= H_Z - \int q_Z(\mathbf{z}) \log p_Z(\mathbf{z}; \theta) d\mathbf{z}, \\ &= H_Z - \int q_Z(\mathbf{z}) [\log \mu_X(F_{zx}(\mathbf{z}; \theta)) + \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|] d\mathbf{z}, \\ &= \underbrace{H_Z + \log Z_X}_{\text{const}} + \mathbb{E}_{\mathbf{z} \sim q_Z(\mathbf{z})} [u(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|] \end{aligned}$$

- KL loss:

$$J_{KL} = \mathbb{E}_{\mathbf{z} \sim q_Z(\mathbf{z})} [u(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|]. \quad (5)$$

Latent KL divergence

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback-Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- KL divergence between q and p , entropy H_q :

$$\begin{aligned} \text{KL}(q \parallel p) &= \int q(\mathbf{x}) [\log q(\mathbf{x}) - \log p(\mathbf{x})] d\mathbf{x}, \\ &= H_q - \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}, \end{aligned}$$

- KL divergence between generated $p_X(\mathbf{x})$ and Boltzmann distribution:

$$\begin{aligned} \text{KL}_\theta [q_Z \parallel p_Z] &= H_Z - \int q_Z(\mathbf{z}) \log p_Z(\mathbf{z}; \theta) d\mathbf{z}, \\ &= H_Z - \int q_Z(\mathbf{z}) [\log \mu_X(F_{zx}(\mathbf{z}; \theta)) + \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|] d\mathbf{z}, \\ &= \underbrace{H_Z + \log Z_X}_{\text{const}} + \mathbb{E}_{\mathbf{z} \sim q_Z(\mathbf{z})} [u(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|] \end{aligned}$$

- KL loss:

$$J_{KL} = \mathbb{E}_{\mathbf{z} \sim q_Z(\mathbf{z})} [u(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|]. \quad (5)$$

Latent KL divergence

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback-Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- KL divergence between q and p , entropy H_q :

$$\begin{aligned} \text{KL}(q \parallel p) &= \int q(\mathbf{x}) [\log q(\mathbf{x}) - \log p(\mathbf{x})] d\mathbf{x}, \\ &= H_q - \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}, \end{aligned}$$

- KL divergence between generated $p_X(\mathbf{x})$ and Boltzmann distribution:

$$\begin{aligned} \text{KL}_\theta [q_Z \parallel p_Z] &= H_Z - \int q_Z(\mathbf{z}) \log p_Z(\mathbf{z}; \theta) d\mathbf{z}, \\ &= H_Z - \int q_Z(\mathbf{z}) [\log \mu_X(F_{zx}(\mathbf{z}; \theta)) + \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|] d\mathbf{z}, \\ &= \underbrace{H_Z + \log Z_X}_{\text{const}} + \mathbb{E}_{\mathbf{z} \sim q_Z(\mathbf{z})} [u(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|] \end{aligned}$$

- KL loss:

$$J_{KL} = \mathbb{E}_{\mathbf{z} \sim q_Z(\mathbf{z})} [u(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|]. \quad (5)$$

Latent KL divergence

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback–Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- KL divergence between q and p , entropy H_q :

$$\begin{aligned} \text{KL}(q \parallel p) &= \int q(\mathbf{x}) [\log q(\mathbf{x}) - \log p(\mathbf{x})] d\mathbf{x}, \\ &= H_q - \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}, \end{aligned}$$

- KL divergence between generated $p_X(\mathbf{x})$ and Boltzmann distribution:

$$\begin{aligned} \text{KL}_\theta [q_Z \parallel p_Z] &= H_Z - \int q_Z(\mathbf{z}) \log p_Z(\mathbf{z}; \theta) d\mathbf{z}, \\ &= H_Z - \int q_Z(\mathbf{z}) [\log \mu_X(F_{zx}(\mathbf{z}; \theta)) + \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|] d\mathbf{z}, \\ &= \underbrace{H_Z + \log Z_X}_{\text{const}} + \mathbb{E}_{\mathbf{z} \sim q_Z(\mathbf{z})} [u(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|] \end{aligned}$$

- KL loss:

$$J_{KL} = \mathbb{E}_{\mathbf{z} \sim q_Z(\mathbf{z})} [u(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)|]. \quad (5)$$

Multi- T latent KL divergence

- **Boltzmann distribution** $e^{-u^k(\mathbf{x})}$ at **multiple temperatures** (T^1, \dots, T^K) . Reference temperature T^0 , relative temperature τ_k :

$$u^k(\mathbf{x}) = \frac{T^0}{T^k} u^0(\mathbf{x}) = \frac{u^0(\mathbf{x})}{\tau_k}$$

- **Prior distribution** $q_Z^k(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \sigma_k^2 \mathbf{I}) = Z_Z^{-1} e^{-\frac{1}{2}\|\mathbf{z}\|^2/\sigma_k^2}$, **Prior energy**:

$$u_Z^k(\mathbf{z}) = -\log q_Z^k(\mathbf{z}) = \frac{1}{2\sigma_k^2} \|\mathbf{z}\|^2 + \text{const}$$

→ variance takes same role as relative temperature. Setting variance=1 at standard temperature, we obtain:

$$\sigma_k^2 = \tau_k.$$

- **Multi-temperature KL loss**:

$$J_{KL}^{T^1, \dots, T^K} = \sum_{k=1}^K \mathbb{E}_{\mathbf{z} \sim q_Z^k(\mathbf{z})} \left[u^k(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)| \right].$$

Multi- T latent KL divergence

- **Boltzmann distribution** $e^{-u^k(\mathbf{x})}$ at **multiple temperatures** (T^1, \dots, T^K) . Reference temperature T^0 , relative temperature τ_k :

$$u^k(\mathbf{x}) = \frac{T^0}{T^k} u^0(\mathbf{x}) = \frac{u^0(\mathbf{x})}{\tau_k}$$

- **Prior distribution** $q_Z^k(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \sigma_k^2 \mathbf{I}) = Z_Z^{-1} e^{-\frac{1}{2}\|\mathbf{z}\|^2/\sigma_k^2}$, **Prior energy**:

$$u_Z^k(\mathbf{z}) = -\log q_Z^k(\mathbf{z}) = \frac{1}{2\sigma_k^2} \|\mathbf{z}\|^2 + \text{const}$$

→ variance takes same role as relative temperature. Setting variance=1 at standard temperature, we obtain:

$$\sigma_k^2 = \tau_k.$$

- **Multi-temperature KL loss**:

$$J_{KL}^{T^1, \dots, T^K} = \sum_{k=1}^K \mathbb{E}_{\mathbf{z} \sim q_Z^k(\mathbf{z})} \left[u^k(F_{zx}(\mathbf{z}; \theta)) - \log |J_{zx}(\mathbf{z}; \theta)| \right].$$

Multi- T latent KL divergence

- **Boltzmann distribution** $e^{-u^k(\mathbf{x})}$ at **multiple temperatures** (T^1, \dots, T^K) . Reference temperature T^0 , relative temperature τ_k :

$$u^k(\mathbf{x}) = \frac{T^0}{T^k} u^0(\mathbf{x}) = \frac{u^0(\mathbf{x})}{\tau_k}$$

- **Prior distribution** $q_Z^k(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \sigma_k^2 \mathbf{I}) = Z_Z^{-1} e^{-\frac{1}{2}\|\mathbf{z}\|^2/\sigma_k^2}$, **Prior energy**:

$$u_Z^k(\mathbf{z}) = -\log q_Z^k(\mathbf{z}) = \frac{1}{2\sigma_k^2} \|\mathbf{z}\|^2 + \text{const}$$

→ variance takes same role as relative temperature. Setting variance=1 at standard temperature, we obtain:

$$\sigma_k^2 = \tau_k.$$

- **Multi-temperature KL loss**:

$$J_{KL}^{T^1, \dots, T^K} = \sum_{k=1}^K \mathbb{E}_{\mathbf{z} \sim q_Z^k(\mathbf{z})} \left[u^k(F_{zx}(\mathbf{z}; \theta)) - \log |\mathbf{J}_{zx}(\mathbf{z}; \theta)| \right].$$

Configuration KL divergence

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback-Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- KL divergence between generated distribution $q_Z(z)$ and Gaussian:

$$\begin{aligned} \text{KL}_\theta(\mu_X \parallel p_X) &= H_X - \int \mu_X(x) \log p_X(x; \theta) dx \\ &= H_X - \int \mu_X(x) [\log q_Z(F_{xz}(x; \theta)) + \log |\mathbf{J}_{xz}(z; \theta)|] dx. \\ &= \underbrace{H_X + \log Z_Z}_{\text{const}} + \mathbb{E}_{x \sim \mu(x)} \left[\frac{1}{\sigma^2} \|F_{xz}(x; \theta)\|^2 - \log |\mathbf{J}_{xz}(x; \theta)| \right]. \end{aligned}$$

Problem: sampling $x \sim \mu(x)$ is difficult and our goal!

- Maximum Likelihood loss:

$$\begin{aligned} J_{ML} &= -\mathbb{E}_{x \sim \rho(x)} [\log p_X(x; \theta)] \\ &= \mathbb{E}_{x \sim \rho(x)} \left[\frac{1}{\sigma^2} \|F_{xz}(x; \theta)\|^2 - \log |\mathbf{J}_{xz}(x; \theta)| \right] \end{aligned}$$

Configuration KL divergence

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback-Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- KL divergence between generated distribution $q_Z(\mathbf{z})$ and Gaussian:

$$\begin{aligned} \text{KL}_\theta(\mu_X \parallel p_X) &= H_X - \int \mu_X(\mathbf{x}) \log p_X(\mathbf{x}; \theta) d\mathbf{x} \\ &= H_X - \int \mu_X(\mathbf{x}) [\log q_Z(F_{xz}(\mathbf{x}; \theta)) + \log |\mathbf{J}_{xz}(\mathbf{z}; \theta)|] d\mathbf{x}. \\ &= \underbrace{H_X + \log Z_Z}_{\text{const}} + \mathbb{E}_{\mathbf{x} \sim \mu(\mathbf{x})} \left[\frac{1}{\sigma^2} \|F_{xz}(\mathbf{x}; \theta)\|^2 - \log |\mathbf{J}_{xz}(\mathbf{x}; \theta)| \right]. \end{aligned}$$

Problem: sampling $\mathbf{x} \sim \mu(\mathbf{x})$ is difficult and our goal!

- Maximum Likelihood loss:

$$\begin{aligned} J_{ML} &= -\mathbb{E}_{\mathbf{x} \sim \rho(\mathbf{x})} [\log p_X(\mathbf{x}; \theta)] \\ &= \mathbb{E}_{\mathbf{x} \sim \rho(\mathbf{x})} \left[\frac{1}{\sigma^2} \|F_{xz}(\mathbf{x}; \theta)\|^2 - \log |\mathbf{J}_{xz}(\mathbf{x}; \theta)| \right] \end{aligned}$$

Configuration KL divergence

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback-Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- KL divergence between generated distribution $q_Z(\mathbf{z})$ and Gaussian:

$$\begin{aligned}\text{KL}_\theta(\mu_X \parallel p_X) &= H_X - \int \mu_X(\mathbf{x}) \log p_X(\mathbf{x}; \theta) d\mathbf{x} \\ &= H_X - \int \mu_X(\mathbf{x}) [\log q_Z(F_{xz}(\mathbf{x}; \theta)) + \log |\mathbf{J}_{xz}(\mathbf{z}; \theta)|] d\mathbf{x}. \\ &= \underbrace{H_X + \log Z_Z}_{\text{const}} + \mathbb{E}_{\mathbf{x} \sim \mu(\mathbf{x})} \left[\frac{1}{\sigma^2} \|F_{xz}(\mathbf{x}; \theta)\|^2 - \log |\mathbf{J}_{xz}(\mathbf{x}; \theta)| \right].\end{aligned}$$

Problem: sampling $\mathbf{x} \sim \mu(\mathbf{x})$ is difficult and our goal!

- **Maximum Likelihood** loss:

$$\begin{aligned}J_{ML} &= -\mathbb{E}_{\mathbf{x} \sim \rho(\mathbf{x})} [\log p_X(\mathbf{x}; \theta)] \\ &= \mathbb{E}_{\mathbf{x} \sim \rho(\mathbf{x})} \left[\frac{1}{\sigma^2} \|F_{xz}(\mathbf{x}; \theta)\|^2 - \log |\mathbf{J}_{xz}(\mathbf{x}; \theta)| \right]\end{aligned}$$

Reaction coordinate loss

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback–Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- Reaction coordinate loss:

$$\begin{aligned} J_{RC} &= \int p(R(x)) \log p(R(x)) dR(x) \\ &= \mathbb{E}_{x \sim p_X(x)} \log p(R(x)). \end{aligned}$$

- Implementation:

- Reaction coordinate function R is user input
- Min and max bounds are given
- $p(R(x))$ is computed as a batchwise kernel density estimate.

Reaction coordinate loss

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback–Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

- Reaction coordinate loss:

$$\begin{aligned} J_{RC} &= \int p(R(\mathbf{x})) \log p(R(\mathbf{x})) dR(\mathbf{x}) \\ &= \mathbb{E}_{\mathbf{x} \sim p_X(\mathbf{x})} \log p(R(\mathbf{x})). \end{aligned}$$

- Implementation:

- Reaction coordinate function R is user input
- Min and max bounds are given
- $p(R(\mathbf{x}))$ is computed as a batchwise kernel density estimate.

Reaction coordinate loss

- Loss function:

$$J = \underbrace{w_{ML} J_{ML}}_{\text{max likelihood}} + \underbrace{w_{KL} J_{KL}}_{\text{Kullback–Leibler}} + \underbrace{w_{RC} J_{RC}}_{\text{reaction coordinate}}$$

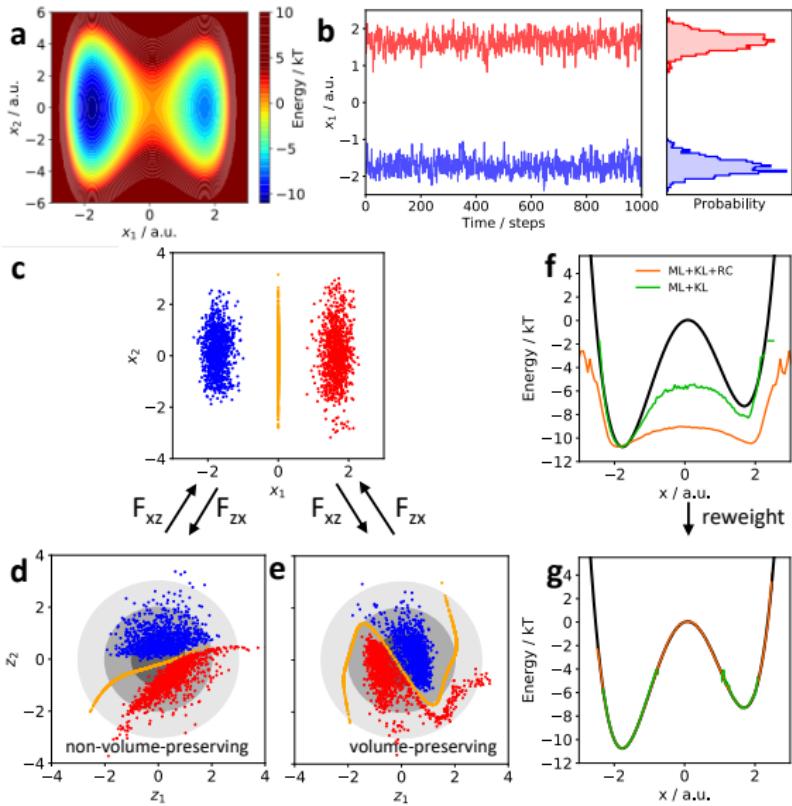
- Reaction coordinate loss:

$$\begin{aligned} J_{RC} &= \int p(R(\mathbf{x})) \log p(R(\mathbf{x})) dR(\mathbf{x}) \\ &= \mathbb{E}_{\mathbf{x} \sim p_X(\mathbf{x})} \log p(R(\mathbf{x})). \end{aligned}$$

- Implementation:

- Reaction coordinate function R is user input
- Min and max bounds are given
- $p(R(\mathbf{x}))$ is computed as a batchwise kernel density estimate.

Boltzmann Generator - 1D Example



Reweighting

- Probability reweighting: assign to each generated configuration \mathbf{x} the statistical weight:

$$w_X(\mathbf{x} \mid \mathbf{z}) = \frac{\mu_X(\mathbf{x})}{p_X(\mathbf{x})} = \frac{p_Z(\mathbf{z})}{q_Z(\mathbf{z})}. \quad (6)$$
$$\propto e^{-u_X(T_{ZX}(\mathbf{z})) + u_Z(\mathbf{z}) + \log|\det(\mathbf{J}_{ZX}(\mathbf{z}))|}$$

- Weighted expectation values

$$\mathbb{E}[O] \approx \frac{\sum_{i=1}^N w_X(\mathbf{x}_i) O(\mathbf{x}_i)}{\sum_{i=1}^N w_X(\mathbf{x}_i)}.$$

Reweighting

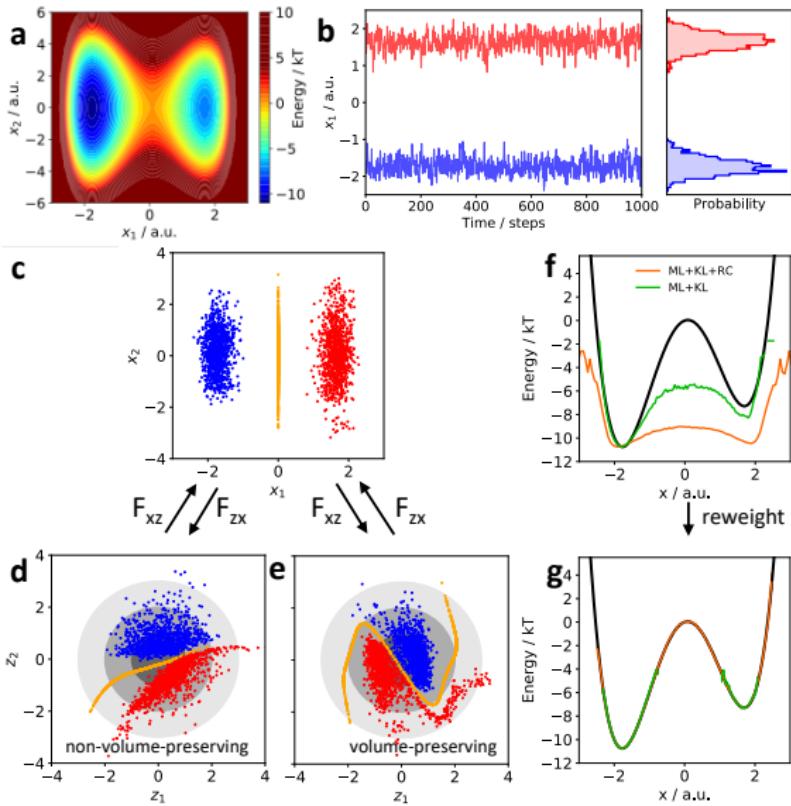
- Probability reweighting: assign to each generated configuration \mathbf{x} the statistical weight:

$$w_X(\mathbf{x} \mid \mathbf{z}) = \frac{\mu_X(\mathbf{x})}{p_X(\mathbf{x})} = \frac{p_Z(\mathbf{z})}{q_Z(\mathbf{z})}. \quad (6)$$
$$\propto e^{-u_X(T_{ZX}(\mathbf{z})) + u_Z(\mathbf{z}) + \log|\det(\mathbf{J}_{ZX}(\mathbf{z}))|}$$

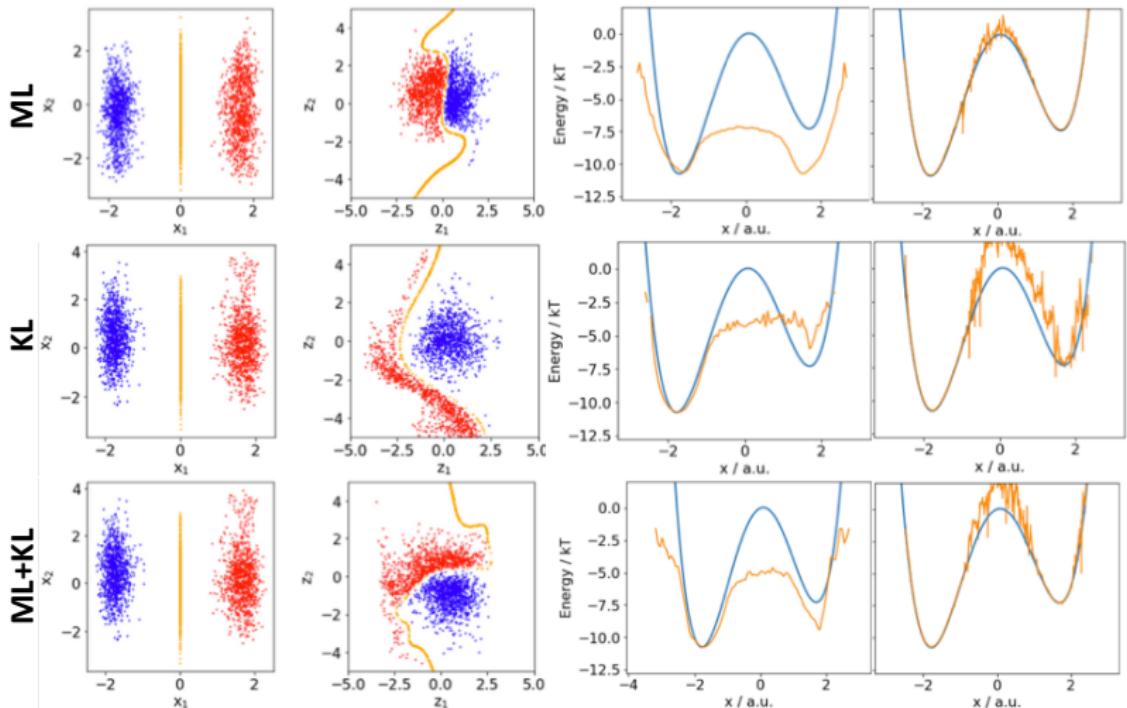
- Weighted expectation values

$$\mathbb{E}[O] \approx \frac{\sum_{i=1}^N w_X(\mathbf{x}_i) O(\mathbf{x}_i)}{\sum_{i=1}^N w_X(\mathbf{x}_i)}.$$

Boltzmann Generator - 1D Example

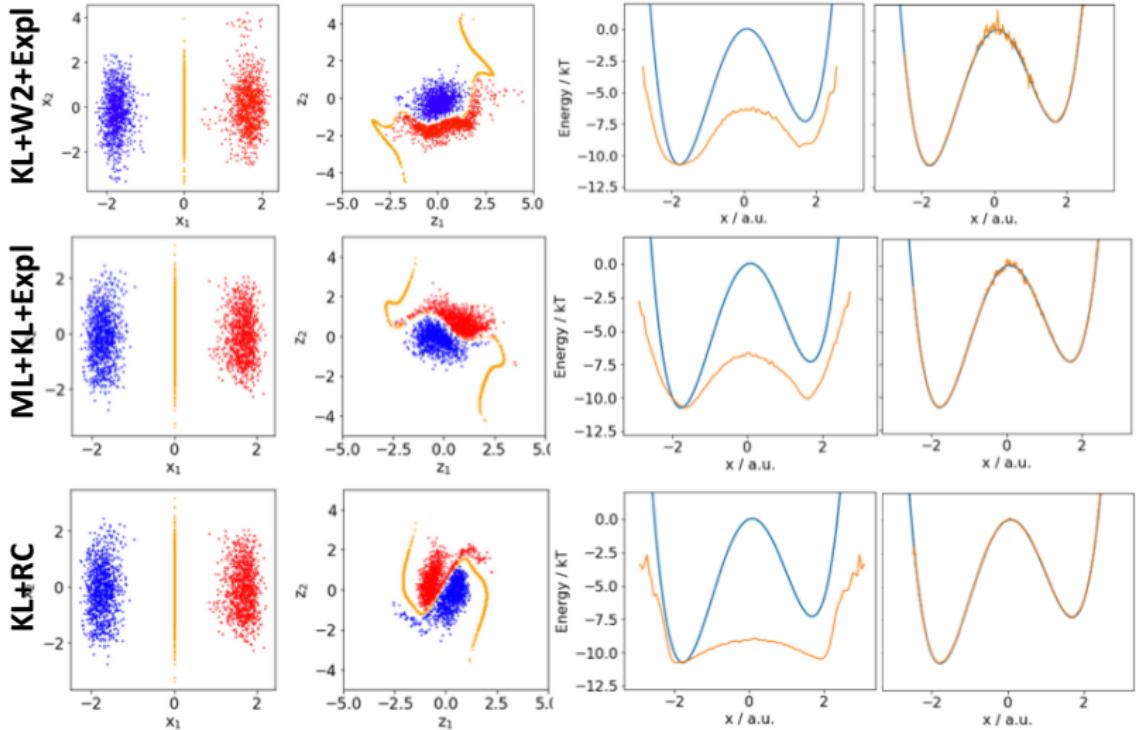


Different training methods



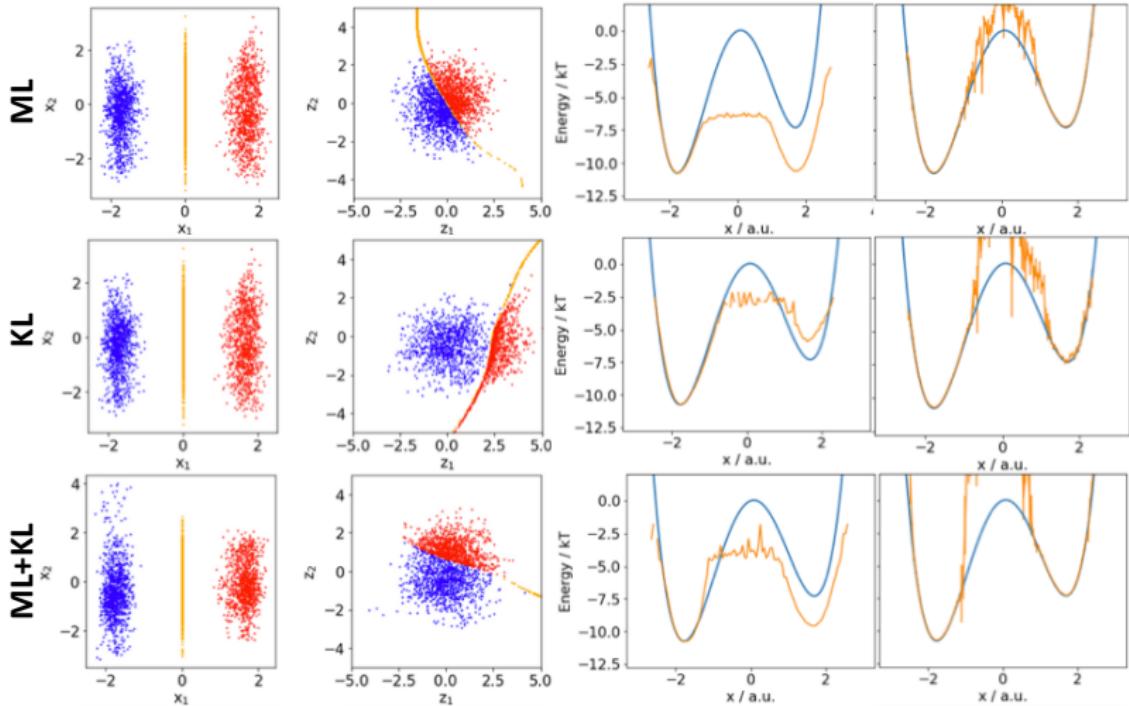
Architecture N₁₀S, $nl_{layers} = 2$, $nl_{hidden} = 100$, $\sigma = \tanh$

Different training methods



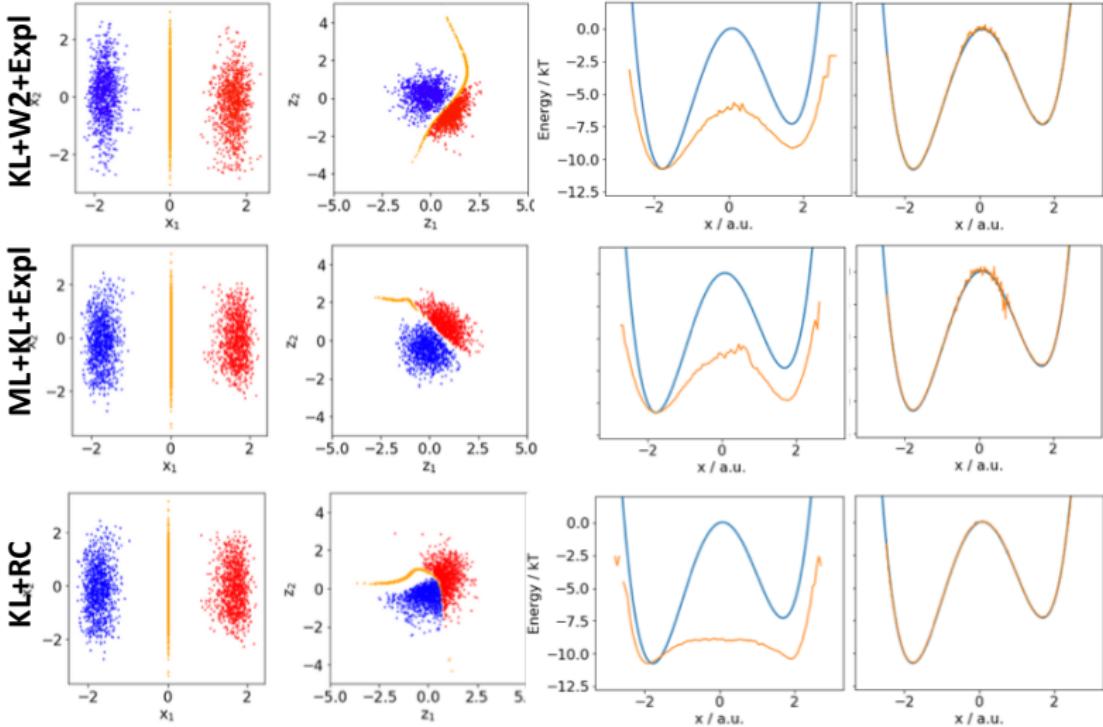
Architecture N₁₀S, $nl_{layers} = 2$, $nl_{hidden} = 100$, $\sigma = \tanh$

Different training methods



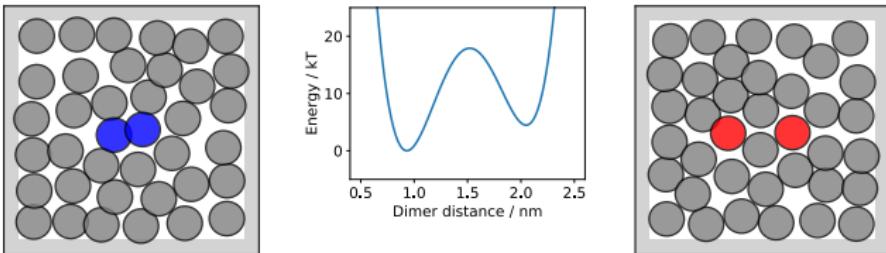
Architecture R_{10} , $nI_{layers} = 2$, $nI_{hidden} = 100$, $\sigma = \tanh$

Different training methods



Architecture R_{10} , $nl_{layers} = 2$, $nl_{hidden} = 100$, $\sigma = \tanh$

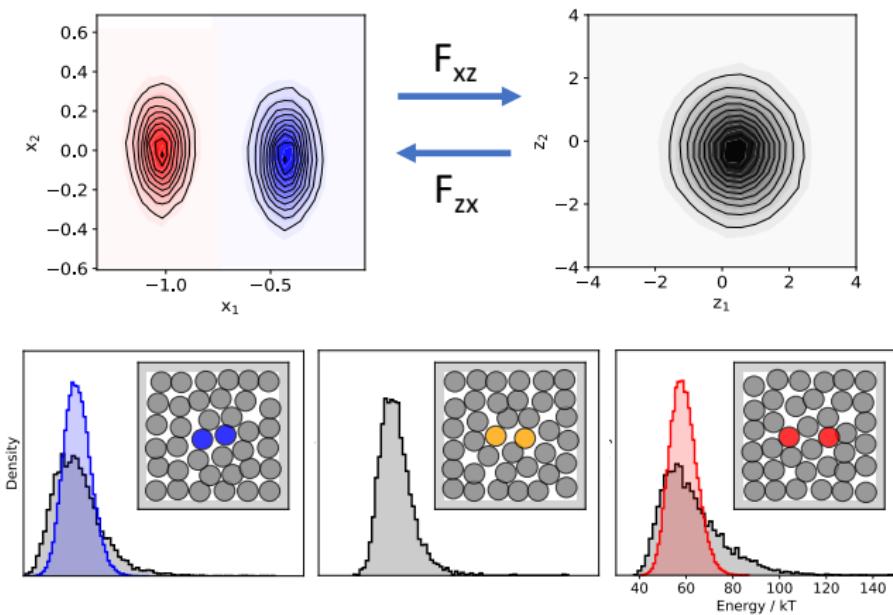
Boltzmann Generator - Particle dimer



- Configuration: $\mathbf{x} = [\mathbf{x}_{1x}, \mathbf{x}_{1y}, \mathbf{x}_{2x}, \mathbf{x}_{2y}, \dots, \mathbf{x}_{(n_s+2)x}, \mathbf{x}_{(n_s+2)y}]$
- Dimer distance $d = \|\mathbf{x}_1 - \mathbf{x}_2\|$, Heavyside function h , **potential energy**:

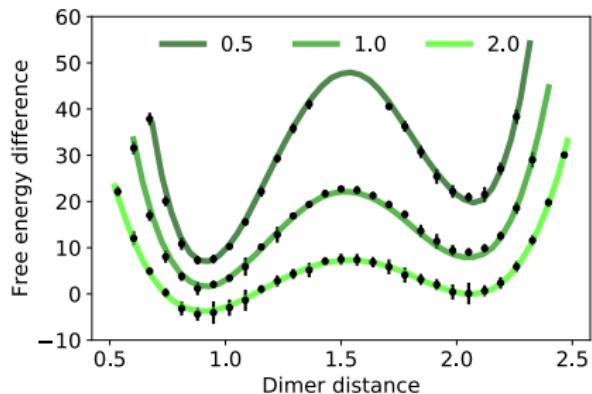
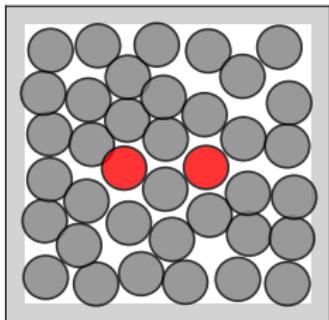
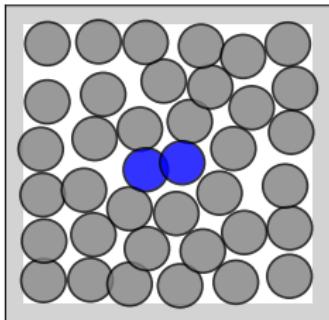
$$\begin{aligned} U(\mathbf{x}) = & k_d (\mathbf{x}_{1x} + \mathbf{x}_{2x})^2 + k_d \mathbf{x}_{1y}^2 + k_d \mathbf{x}_{2y}^2 \\ & + \frac{1}{4} a(d - d_0)^4 - \frac{1}{2} b(d - d_0)^2 + c(d - d_0)^4 \\ & + \sum_{i=1}^{n+2} h(-\mathbf{x}_{1x} - l_{\text{box}}) k_{\text{box}} (-\mathbf{x}_{1x} - l_{\text{box}})^2 + \sum_{i=1}^{n+2} h(\mathbf{x}_{1x} - l_{\text{box}}) k_{\text{box}} (\mathbf{x}_{1x} - l_{\text{box}})^2 \\ & + \sum_{i=1}^{n+2} h(-\mathbf{x}_{1y} - l_{\text{box}}) k_{\text{box}} (-\mathbf{x}_{1y} - l_{\text{box}})^2 + \sum_{i=1}^{n+2} h(\mathbf{x}_{1y} - l_{\text{box}}) k_{\text{box}} (\mathbf{x}_{1y} - l_{\text{box}})^2 \\ & + \varepsilon \sum_{i=1}^{n+1} \sum_{j=i+1, j \neq 2}^{n+2} \left(\frac{\sigma}{\|\mathbf{x}_i - \mathbf{x}_j\|} \right)^{12} \end{aligned}$$

Boltzmann Generator - Particle dimer

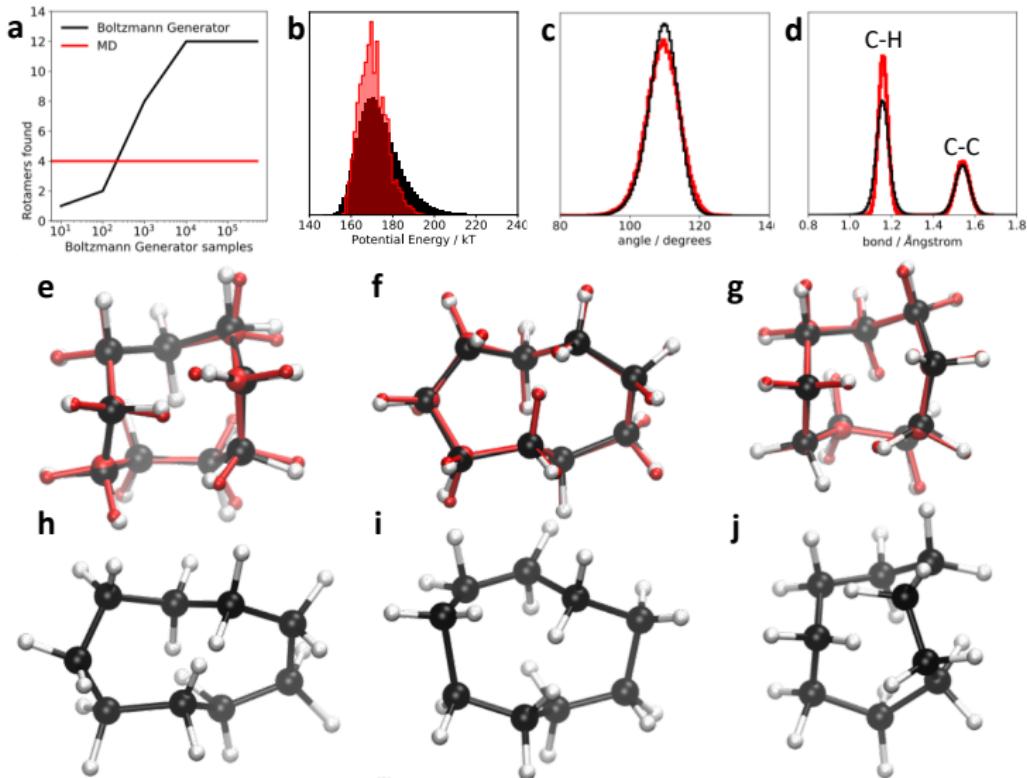


Architecture R₈, $n_{layers} = 3$, $n_{hidden} = 200$, $\sigma = \tanh$, $w_{KL} = 1$,
 $w_{ML} = 0.01$, $w_{RC} = 10$

Boltzmann Generator - Particle dimer

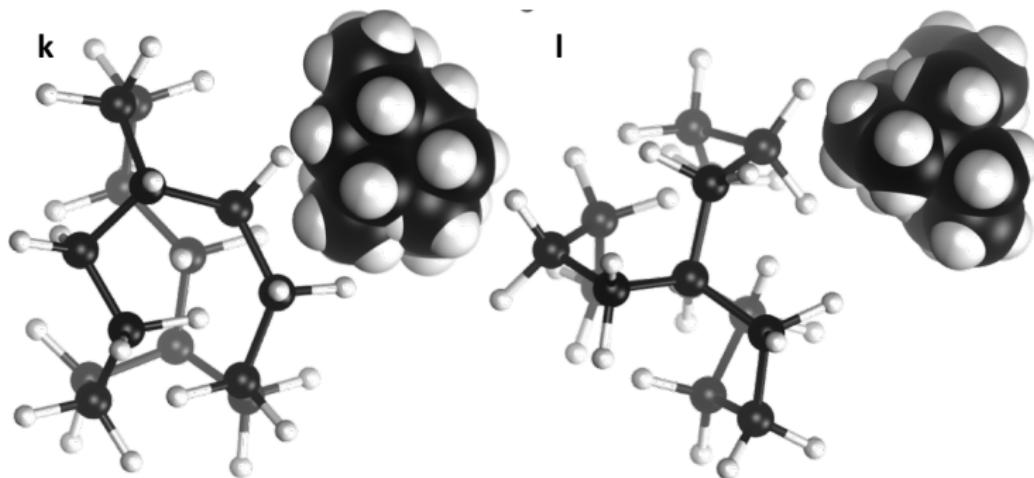


Boltzmann Generator - Hydrocarbons



Architecture R₄, $n_{\text{layers}} = 2$, $n_{\text{hidden}} = 100$, $\sigma = \tanh$, $w_{KL} = 0.01$, $w_{ML} = 1$

Boltzmann Generator - Hydrocarbons



Architecture R_4 , $nI_{layers} = 2$, $nI_{hidden} = 100$, $\sigma = \tanh$, $w_{KL} = 0.01$,
 $w_{ML} = 1$