

# A dynamical systems perspective on measure transport and generative modeling

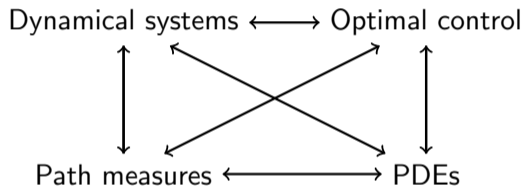
Lorenz Richter

Machine-Learning-Based Sampling in Lattice Field Theory and Quantum Chemistry  
Bethe Center for Theoretical Physics, Bonn

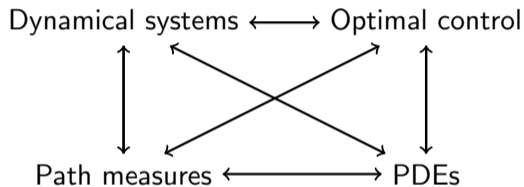
October, 2024



- Sampling via measure transport can be seen from different perspectives:

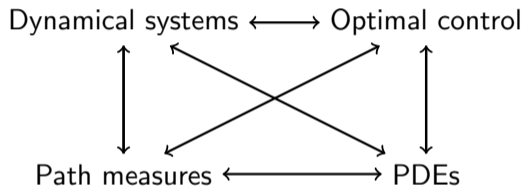


- Sampling via measure transport can be seen from different perspectives:



- The different perspectives will eventually allow us to propose new numerical algorithms.

- Sampling via measure transport can be seen from different perspectives:



- The different perspectives will eventually allow us to propose new numerical algorithms.
- This is joint work with Julius Berner (Caltech), Jingtong Sun (Caltech), Denis Blessing (KIT) and Nikolas Nüsken (King's College).

## Task

Sample from a complex (high-dimensional, multimodal) distribution  $\mathcal{D}$ .

## Task

Sample from a complex (high-dimensional, multimodal) distribution  $\mathcal{D}$ .

$\mathcal{D}$  can be given in the form of:

1. **samples**  $X^{(i)} \sim \mathcal{D}$  (images, text, audio, ...).

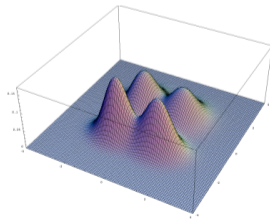


## Task

Sample from a complex (high-dimensional, multimodal) distribution  $\mathcal{D}$ .

$\mathcal{D}$  can be given in the form of:

1. **samples**  $X^{(i)} \sim \mathcal{D}$  (images, text, audio, ...).
2. an (unnormalized) **density** (e.g., in Bayesian statistics, computational physics and chemistry).



- Impressive results for the first case:





- Impressive results for the first case:

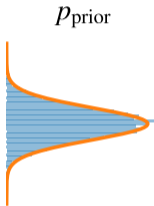


- The second case is a focus of (our) current research.

**Goal:** We want to sample from distribution  $p_{\text{target}} = \rho/\mathcal{Z}$ .

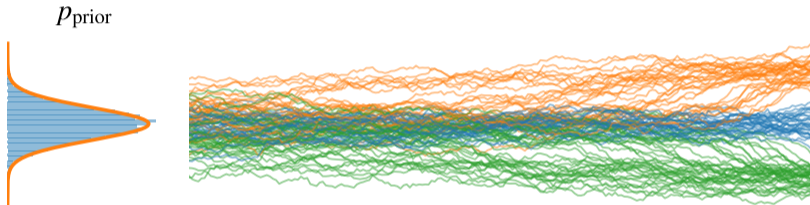
# Sampling via dynamical systems

**Goal:** We want to sample from distribution  $p_{\text{target}} = \rho / \mathcal{Z}$ .



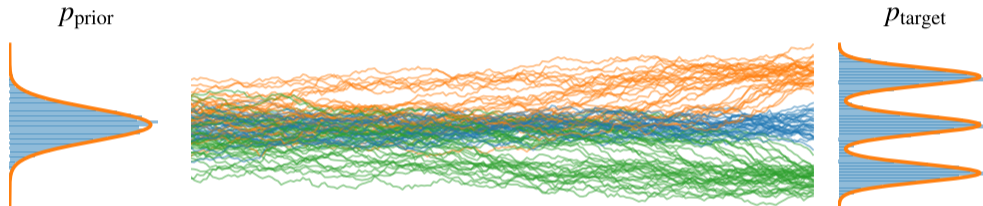
# Sampling via dynamical systems

**Goal:** We want to sample from distribution  $p_{\text{target}} = \rho / \mathcal{Z}$ .



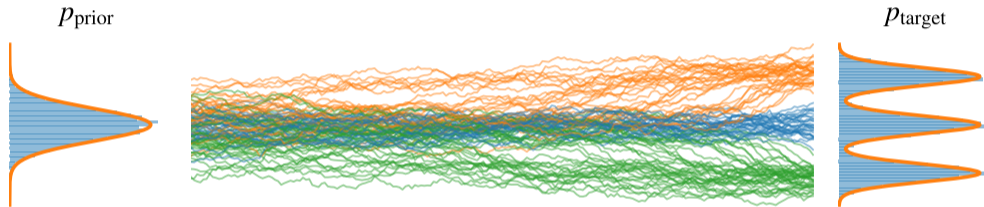
# Sampling via dynamical systems

**Goal:** We want to sample from distribution  $p_{\text{target}} = \rho / \mathcal{Z}$ .



# Sampling via dynamical systems

**Goal:** We want to sample from distribution  $p_{\text{target}} = \rho / \mathcal{Z}$ .



**SDE**

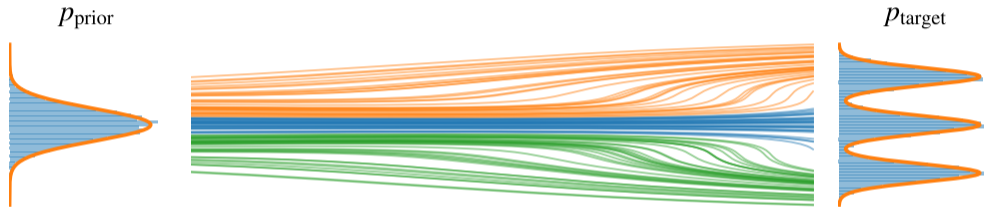
**Setting:**

$$X_0 \sim p_{\text{prior}}$$

$$dX_s = \mu(X_s, s) ds + \sigma(s) dW_s$$

# Sampling via dynamical systems

**Goal:** We want to sample from distribution  $p_{\text{target}} = \rho / \mathcal{Z}$ .



**Setting:**

$$X_0 \sim p_{\text{prior}}$$

**SDE**

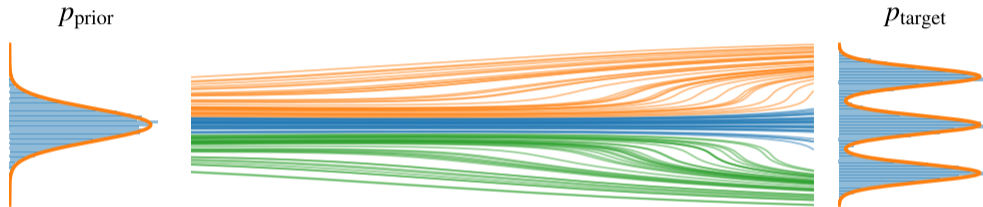
$$dX_s = \mu(X_s, s) ds + \sigma(s) dW_s$$

**ODE**

$$dX_s = \mu(X_s, s) ds$$

# Sampling via dynamical systems

**Goal:** We want to sample from distribution  $p_{\text{target}} = \rho/\mathcal{Z}$ .



**SDE**

**ODE**

**Setting:**

$$X_0 \sim p_{\text{prior}}$$

$$dX_s = \mu(X_s, s) ds + \sigma(s) dW_s$$

$$dX_s = \mu(X_s, s) ds$$

**Idea:** Learn  $\mu$  s.t.  $X_T \sim p_{\text{target}}$ .



## Attempt I: PDE perspective

- Considering the density of  $X_t$ , denoted by  $p_X(\cdot, t)$ , leads to the following PDEs:
  - ▶ **SDE:** Fokker-Planck equation

$$\partial_t p_X + \operatorname{div}(p_X \mu) - \frac{1}{2} \operatorname{Tr}(\sigma \sigma^\top \nabla^2 p_X) = 0,$$

- Considering the density of  $X_t$ , denoted by  $p_X(\cdot, t)$ , leads to the following PDEs:

- ▶ **SDE:** Fokker-Planck equation

$$\partial_t p_X + \operatorname{div}(p_X \mu) - \frac{1}{2} \operatorname{Tr}(\sigma \sigma^\top \nabla^2 p_X) = 0,$$

- ▶ **ODE:** Continuity equation

$$\partial_t p_X + \operatorname{div}(p_X \mu) = 0,$$

- Considering the density of  $X_t$ , denoted by  $p_X(\cdot, t)$ , leads to the following PDEs:

- ▶ **SDE:** Fokker-Planck equation

$$\partial_t p_X + \operatorname{div}(p_X \mu) - \frac{1}{2} \operatorname{Tr}(\sigma \sigma^\top \nabla^2 p_X) = 0,$$

- ▶ **ODE:** Continuity equation

$$\partial_t p_X + \operatorname{div}(p_X \mu) = 0,$$

with boundary conditions  $p_X(\cdot, 0) = p_{\text{prior}}$  and  $p_X(\cdot, T) = p_{\text{target}}$ .

- Considering the density of  $X_t$ , denoted by  $p_X(\cdot, t)$ , leads to the following PDEs:

- ▶ **SDE:** Fokker-Planck equation

$$\partial_t p_X + \operatorname{div}(p_X \mu) - \frac{1}{2} \operatorname{Tr}(\sigma \sigma^\top \nabla^2 p_X) = 0,$$

- ▶ **ODE:** Continuity equation

$$\partial_t p_X + \operatorname{div}(p_X \mu) = 0,$$

with boundary conditions  $p_X(\cdot, 0) = p_{\text{prior}}$  and  $p_X(\cdot, T) = p_{\text{target}}$ .

- **Idea:** Identify pairs  $(\mu, p_X)$  that fulfill the above PDEs.

- **Variational formulation of PDEs:** Consider loss functionals

$$\mathcal{L} : C(\mathbb{R}^d \times [0, T], \mathbb{R}^d) \times C(\mathbb{R}^d \times [0, T], \mathbb{R}) \rightarrow \mathbb{R}_{\geq 0}$$

that are zero if and only if a pair  $(\mu, p_X)$  fulfills the corresponding PDE.

- **Variational formulation of PDEs:** Consider loss functionals

$$\mathcal{L} : C(\mathbb{R}^d \times [0, T], \mathbb{R}^d) \times C(\mathbb{R}^d \times [0, T], \mathbb{R}) \rightarrow \mathbb{R}_{\geq 0}$$

that are zero if and only if a pair  $(\mu, p_X)$  fulfills the corresponding PDE.

- For numerical stability, we consider the PDEs in log-space,  $V := \log p_X$ , yielding

- ▶ **SDE:**

$$\mathcal{R}_{\log\text{FP}}(\mu, V) := \partial_t V + \text{div}(\mu) + \nabla V \cdot \mu - \frac{1}{2} \|\sigma^\top \nabla V\|^2 - \frac{1}{2} \text{Tr}(\sigma \sigma^\top \nabla^2 V) = 0,$$

- ▶ **ODE:**

$$\mathcal{R}_{\log\text{CE}}(\mu, V) := \partial_t V + \text{div}(\mu) + \nabla V \cdot \mu = 0.$$

## Attempt I: PDE perspective – general vs. constrained dynamics

- Learning  $\mu$  and  $p_X$  simultaneously typically leads to non-unique solutions.



## Attempt I: PDE perspective – general vs. constrained dynamics

- Learning  $\mu$  and  $p_X$  simultaneously typically leads to non-unique solutions.
- Adding constraints allows for unique solutions:
  - ▶ **Annealing:** Fix  $p_X$  (i.e. a path of densities) and consider

$$\mathcal{R}_{\log\text{FP}}^{\text{anneal}}(\tilde{\mu}) := \mathcal{R}_{\log\text{FP}}(\tilde{\mu}, V), \quad \mathcal{R}_{\log\text{CE}}^{\text{anneal}}(\tilde{\mu}) := \mathcal{R}_{\log\text{CE}}(\tilde{\mu}, V).$$

# Attempt I: PDE perspective – general vs. constrained dynamics

- Learning  $\mu$  and  $p_X$  simultaneously typically leads to non-unique solutions.
- Adding constraints allows for unique solutions:

- ▶ **Annealing:** Fix  $p_X$  (i.e. a path of densities) and consider

$$\mathcal{R}_{\log\text{FP}}^{\text{anneal}}(\tilde{\mu}) := \mathcal{R}_{\log\text{FP}}(\tilde{\mu}, V), \quad \mathcal{R}_{\log\text{CE}}^{\text{anneal}}(\tilde{\mu}) := \mathcal{R}_{\log\text{CE}}(\tilde{\mu}, V).$$

- ▶ **Score-based generative modeling:** Fix  $\mu = \sigma\sigma^\top \nabla V - f$  and consider

$$\mathcal{R}_{\text{score}}(\tilde{V}) := \mathcal{R}_{\log\text{FP}}(\sigma\sigma^\top \nabla \tilde{V} - f, \tilde{V}).$$

# Attempt I: PDE perspective – general vs. constrained dynamics

- Learning  $\mu$  and  $p_X$  simultaneously typically leads to non-unique solutions.
- Adding constraints allows for unique solutions:

- ▶ **Annealing:** Fix  $p_X$  (i.e. a path of densities) and consider

$$\mathcal{R}_{\log\text{FP}}^{\text{anneal}}(\tilde{\mu}) := \mathcal{R}_{\log\text{FP}}(\tilde{\mu}, V), \quad \mathcal{R}_{\log\text{CE}}^{\text{anneal}}(\tilde{\mu}) := \mathcal{R}_{\log\text{CE}}(\tilde{\mu}, V).$$

- ▶ **Score-based generative modeling:** Fix  $\mu = \sigma\sigma^\top \nabla V - f$  and consider

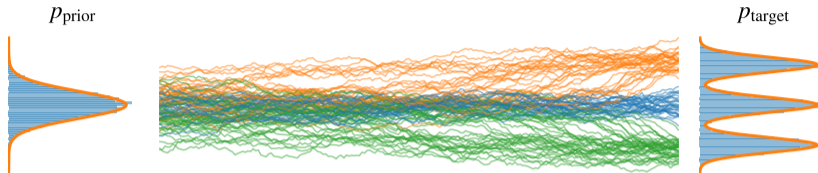
$$\mathcal{R}_{\text{score}}(\tilde{V}) := \mathcal{R}_{\log\text{FP}}(\sigma\sigma^\top \nabla \tilde{V} - f, \tilde{V}).$$

- ▶ **Optimal transport & Schrödinger bridges:** Additionally minimize  $\mathbb{E} \left[ \frac{1}{2} \int_0^T \|\mu(X_s, s)\|^2 ds \right]$ .  
Find  $\mu = \nabla \Phi$ , where  $\Phi$  solves

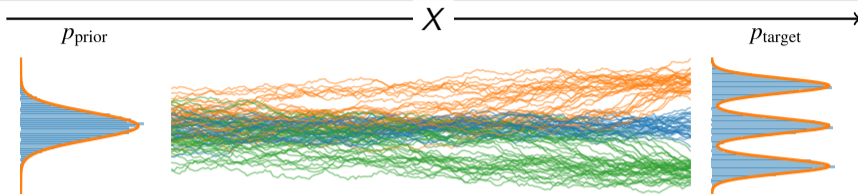
$$\mathcal{R}_{\text{HJB}}^{\text{SB}}(\Phi) := \partial_t \Phi + \frac{1}{2} \|\nabla \Phi\|^2 + \frac{1}{2} \text{Tr}(\sigma\sigma^\top \nabla^2 \Phi) = 0, \quad \mathcal{R}_{\text{HJB}}^{\text{OT}}(\Phi) := \partial_t \Phi + \frac{1}{2} \|\nabla \Phi\|^2 = 0.$$

## Attempt II: Time-reversals and path space measures

# Attempt II: Time-reversals and path space measures



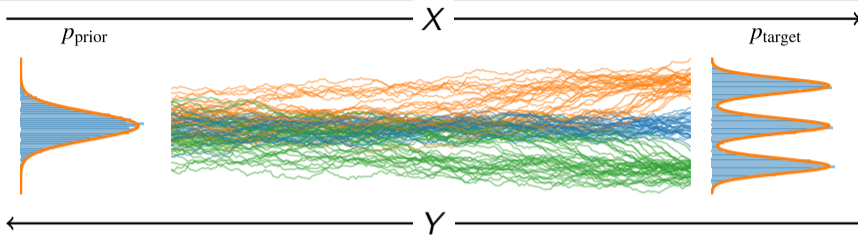
## Attempt II: Time-reversals and path space measures



- **Setting:** Consider forward and reverse SDE:

$$dX_s = \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, \quad X_0 \sim p_{\text{prior}},$$

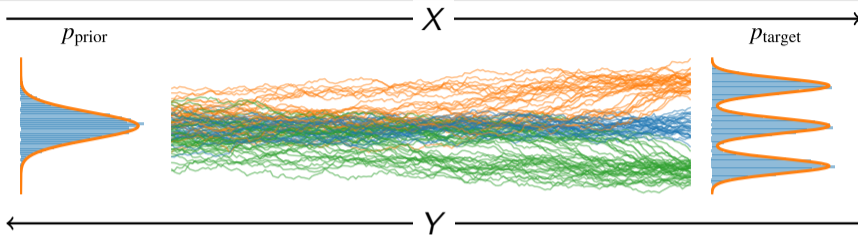
## Attempt II: Time-reversals and path space measures



- **Setting:** Consider forward and reverse SDE:

$$\begin{aligned}dX_s &= \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, & X_0 &\sim p_{\text{prior}}, \\dY_s &= \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, & Y_T &\sim p_{\text{target}}.\end{aligned}$$

## Attempt II: Time-reversals and path space measures



- **Setting:** Consider forward and reverse SDE:

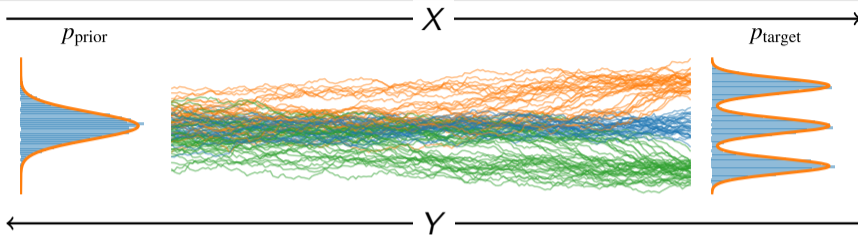
$$dX_s = \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, \quad X_0 \sim p_{\text{prior}},$$

$$dY_s = \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, \quad Y_T \sim p_{\text{target}}.$$

- **Idea:** Learn  $\tilde{\mu}_F, \tilde{\mu}_B$  s.t.  $X$  is time-reversal of  $Y$ , implying  $X_T \sim p_{\text{target}}, Y_0 \sim p_{\text{prior}}$ .



## Attempt II: Time-reversals and path space measures



- **Setting:** Consider forward and reverse SDE:

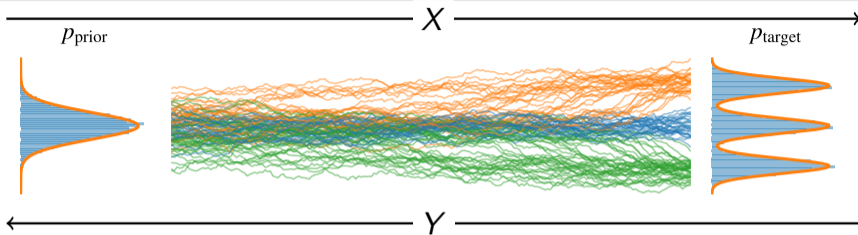
$$dX_s = \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, \quad X_0 \sim p_{\text{prior}},$$

$$dY_s = \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, \quad Y_T \sim p_{\text{target}}.$$

- **Idea:** Learn  $\tilde{\mu}_F, \tilde{\mu}_B$  s.t.  $X$  is time-reversal of  $Y$ , implying  $X_T \sim p_{\text{target}}, Y_0 \sim p_{\text{prior}}$ .

- ▶ Learn  $\tilde{\mu}_F$  and  $\tilde{\mu}_B$  simultaneously:  $\mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B)$

# Attempt II: Time-reversals and path space measures



- **Setting:** Consider forward and reverse SDE:

$$dX_s = \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, \quad X_0 \sim p_{\text{prior}},$$

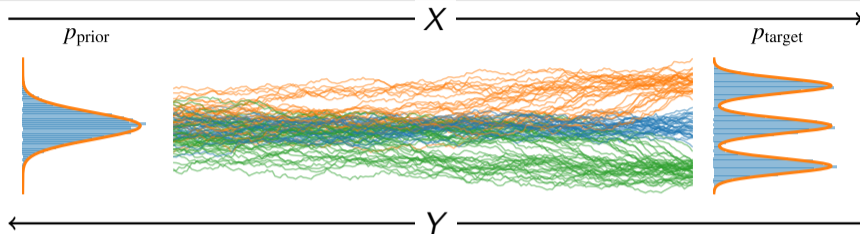
$$dY_s = \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, \quad Y_T \sim p_{\text{target}}.$$

- **Idea:** Learn  $\tilde{\mu}_F, \tilde{\mu}_B$  s.t.  $X$  is time-reversal of  $Y$ , implying  $X_T \sim p_{\text{target}}, Y_0 \sim p_{\text{prior}}$ .

- ▶ Learn  $\tilde{\mu}_F$  and  $\tilde{\mu}_B$  simultaneously:  $\mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B)$

- ▶ Fixe annealing  $p_X$ :  $\mathcal{L}_{\text{CMCD}}^{\text{BSDE}}(\tilde{\mu}_F)$

# Attempt II: Time-reversals and path space measures



- **Setting:** Consider forward and reverse SDE:

$$dX_s = \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, \quad X_0 \sim p_{\text{prior}},$$

$$dY_s = \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, \quad Y_T \sim p_{\text{target}}.$$

- **Idea:** Learn  $\tilde{\mu}_F, \tilde{\mu}_B$  s.t.  $X$  is time-reversal of  $Y$ , implying  $X_T \sim p_{\text{target}}, Y_0 \sim p_{\text{prior}}$ .

- ▶ Learn  $\tilde{\mu}_F$  and  $\tilde{\mu}_B$  simultaneously:  $\mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B)$
- ▶ Fix annealing  $p_X$ :  $\mathcal{L}_{\text{CMCD}}^{\text{BSDE}}(\tilde{\mu}_F)$
- ▶ Fix  $\tilde{\mu}_B$  suitably:  $\mathcal{L}_{\text{DIS}}^{\text{BSDE}}(\tilde{\mu}_F)$

## Attempt II: Time-reversals and path space measures

- We consider the SDEs

$$\begin{aligned}dX_s &= \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, & X_0 &\sim p_{\text{prior}}, \\dY_s &= \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, & Y_T &\sim p_{\text{target}}.\end{aligned}$$

## Attempt II: Time-reversals and path space measures

- We consider the SDEs

$$\begin{aligned}dX_s &= \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, & X_0 &\sim p_{\text{prior}}, \\dY_s &= \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, & Y_T &\sim p_{\text{target}}.\end{aligned}$$

- **Path space perspective:** Consider path measures  $\mathbb{P}_{X^{\mu_F}}$  and  $\mathbb{P}_{\bar{Y}^{\mu_B}}$ .

## Attempt II: Time-reversals and path space measures

- We consider the SDEs

$$\begin{aligned}dX_s &= \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, & X_0 &\sim p_{\text{prior}}, \\dY_s &= \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, & Y_T &\sim p_{\text{target}}.\end{aligned}$$

- **Path space perspective:** Consider path measures  $\mathbb{P}_{X^{\mu_F}}$  and  $\mathbb{P}_{\tilde{Y}^{\mu_B}}$ .
- Identify drifts  $\mu_F, \mu_B$  via divergence of those measures

$$\mu_F, \mu_B \in \arg \min_{\tilde{\mu}_F, \tilde{\mu}_B} D(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}).$$

## Attempt II: Time-reversals and path space measures

- We consider the SDEs

$$\begin{aligned}dX_s &= \tilde{\mu}_F(X_s, s) ds + \sigma(s) dW_s, & X_0 &\sim p_{\text{prior}}, \\dY_s &= \tilde{\mu}_B(Y_s, s) ds + \sigma(s) \bar{d}W_s, & Y_T &\sim p_{\text{target}}.\end{aligned}$$

- **Path space perspective:** Consider path measures  $\mathbb{P}_{X^{\mu_F}}$  and  $\mathbb{P}_{\tilde{Y}^{\mu_B}}$ .
- Identify drifts  $\mu_F, \mu_B$  via divergence of those measures

$$\mu_F, \mu_B \in \arg \min_{\tilde{\mu}_F, \tilde{\mu}_B} D(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}).$$

### Proposition (Log-likelihood for path measures)

$$\begin{aligned}\log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) &= \int_0^T \left( \sigma^{-2}(\tilde{\mu}_F + \tilde{\mu}_B) \cdot \left( \tilde{\mu}_R + \frac{\tilde{\mu}_B - \tilde{\mu}_F}{2} \right) + \nabla \cdot \tilde{\mu}_B \right)(X_s^{\tilde{\mu}_R}, s) ds \\ &\quad + \int_0^T \sigma^{-1}(\tilde{\mu}_F + \tilde{\mu}_B)(X_s^{\tilde{\mu}_R}, s) \cdot dW_s + \log \frac{p_{\text{prior}}(X_0^{\tilde{\mu}_R})}{p_{\text{target}}(X_T^{\tilde{\mu}_R})}\end{aligned}$$

Connections and equivalences: divergences and loss functions



# Connections and equivalences: Time-reversals and BSDEs

- We need to choose a loss  $\mathcal{L}$  or a divergence  $D$ .

# Connections and equivalences: Time-reversals and BSDEs

- We need to choose a loss  $\mathcal{L}$  or a divergence  $D$ .
- First choice: BSDE-based losses for stochastic evolutions.

# Connections and equivalences: Time-reversals and BSDEs

- We need to choose a loss  $\mathcal{L}$  or a divergence  $D$ .
- First choice: BSDE-based losses for stochastic evolutions.
- **BSDE loss:** stochastic representation of PDE via Itô's formula. For the process

$$dX_s = \mu(X_s, s)ds + \sigma(s) dW_s.$$

and a PDE

$$\partial_t V + \frac{1}{2} \text{Tr} \left( \sigma \sigma^\top \nabla^2 V \right) + \mu \cdot \nabla V + h(\cdot, \cdot, V, \nabla V) = 0$$

it holds

$$\mathcal{R}_{\text{BSDE}}(V) = V(X_0, 0) - V(X_T, T) + \int_0^T \left( \partial_s V + \frac{1}{2} \text{Tr}(\sigma \sigma^\top \nabla^2 V) + \mu \cdot \nabla V \right)(X_s, s) ds + \int_0^T \sigma^\top \nabla V(X_s, s) \cdot dW_s = 0.$$

# Connections and equivalences: Time-reversals and BSDEs

- We need to choose a loss  $\mathcal{L}$  or a divergence  $D$ .
- First choice: BSDE-based losses for stochastic evolutions.
- **BSDE loss**: stochastic representation of PDE via Itô's formula. For the process

$$dX_s = \mu(X_s, s)ds + \sigma(s) dW_s.$$

and a PDE

$$\partial_t V + \frac{1}{2} \text{Tr} \left( \sigma \sigma^\top \nabla^2 V \right) + \mu \cdot \nabla V + h(\cdot, \cdot, V, \nabla V) = 0$$

it holds

$$\mathcal{R}_{\text{BSDE}}(V) = V(X_0, 0) - V(X_T, T) - \int_0^T h(\cdot, \cdot, V, \nabla V)(X_s, s) ds + \int_0^T \sigma^\top \nabla V(X_s, s) \cdot dW_s = 0.$$

# Connections and equivalences: Time-reversals and BSDEs

- We need to choose a loss  $\mathcal{L}$  or a divergence  $D$ .
- First choice: BSDE-based losses for stochastic evolutions.
- **BSDE loss**: stochastic representation of PDE via Itô's formula. For the process

$$dX_s = \mu(X_s, s)ds + \sigma(s) dW_s.$$

and a PDE

$$\partial_t V + \frac{1}{2} \text{Tr} \left( \sigma \sigma^\top \nabla^2 V \right) + \mu \cdot \nabla V + h(\cdot, \cdot, V, \nabla V) = 0$$

it holds

$$\mathcal{R}_{\text{BSDE}}(V) = V(X_0, 0) - V(X_T, T) - \int_0^T h(\cdot, \cdot, V, \nabla V)(X_s, s) ds + \int_0^T \sigma^\top \nabla V(X_s, s) \cdot dW_s = 0.$$

- We can now consider the loss

$$\mathcal{L}_{\text{BSDE}}(\tilde{V}) = \mathbb{E} \left[ \left( \mathcal{R}_{\text{BSDE}}(\tilde{V})(X) \right)^2 \right],$$

where the expectation is over different realizations of the process  $X$ .

# Connections and equivalences: Time-reversals and BSDEs

- BSDE-based losses are equivalent to a particular divergence between path space measures:

$$D_{\text{BSDE}}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) = \mathbb{E} \left[ \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)^2 \right] = \mathcal{L}_{\log\text{FP}}^{\text{BSDE}}(\tilde{\mu}, \tilde{V}) = \mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B)$$

# Connections and equivalences: Time-reversals and BSDEs

- BSDE-based losses are equivalent to a particular divergence between path space measures:

$$D_{\text{BSDE}}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) = \mathbb{E} \left[ \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)^2 \right] = \mathcal{L}_{\log\text{FP}}^{\text{BSDE}}(\tilde{\mu}, \tilde{V}) = \mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B)$$

## Proposition (Equivalence to trajectory-based methods)

*The BSDE versions of our losses are equivalent to previously existing losses:*

# Connections and equivalences: Time-reversals and BSDEs

- BSDE-based losses are equivalent to a particular divergence between path space measures:

$$D_{\text{BSDE}}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) = \mathbb{E} \left[ \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)^2 \right] = \mathcal{L}_{\log\text{FP}}^{\text{BSDE}}(\tilde{\mu}, \tilde{V}) = \mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B)$$

## Proposition (Equivalence to trajectory-based methods)

*The BSDE versions of our losses are equivalent to previously existing losses:*

- 1 Assuming the reparametrization  $\sigma\sigma^\top \nabla \tilde{V} = \tilde{\mu}_F - \tilde{\mu}_B$ , it holds

$$\mathcal{L}_{\log\text{FP}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{V}) = \mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B).$$



# Connections and equivalences: Time-reversals and BSDEs

- BSDE-based losses are equivalent to a particular divergence between path space measures:

$$D_{\text{BSDE}}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) = \mathbb{E} \left[ \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)^2 \right] = \mathcal{L}_{\log\text{FP}}^{\text{BSDE}}(\tilde{\mu}, \tilde{V}) = \mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B)$$

## Proposition (Equivalence to trajectory-based methods)

*The BSDE versions of our losses are equivalent to previously existing losses:*

- 1 Assuming the reparametrization  $\sigma\sigma^\top \nabla \tilde{V} = \tilde{\mu}_F - \tilde{\mu}_B$ , it holds

$$\mathcal{L}_{\log\text{FP}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{V}) = \mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B).$$

- 2 It holds

$$\mathcal{L}_{\log\text{FP}}^{\text{anneal,BSDE}}(\tilde{\mu}_F) = \mathcal{L}_{\text{CMCD}}^{\text{BSDE}}(\tilde{\mu}_F).$$

# Connections and equivalences: Time-reversals and BSDEs

- BSDE-based losses are equivalent to a particular divergence between path space measures:

$$D_{\text{BSDE}}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) = \mathbb{E} \left[ \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)^2 \right] = \mathcal{L}_{\log\text{FP}}^{\text{BSDE}}(\tilde{\mu}, \tilde{V}) = \mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B)$$

## Proposition (Equivalence to trajectory-based methods)

*The BSDE versions of our losses are equivalent to previously existing losses:*

- 1 Assuming the reparametrization  $\sigma\sigma^\top \nabla \tilde{V} = \tilde{\mu}_F - \tilde{\mu}_B$ , it holds

$$\mathcal{L}_{\log\text{FP}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{V}) = \mathcal{L}_{\text{Bridge}}^{\text{BSDE}}(\tilde{\mu}_F, \tilde{\mu}_B).$$

- 2 It holds

$$\mathcal{L}_{\log\text{FP}}^{\text{anneal,BSDE}}(\tilde{\mu}_F) = \mathcal{L}_{\text{CMCD}}^{\text{BSDE}}(\tilde{\mu}_F).$$

- 3 Assuming the reparametrization  $\sigma\sigma^\top \nabla \tilde{V} = \tilde{\mu}_F - f$ , it holds

$$\mathcal{L}_{\text{score}}^{\text{BSDE}}(\tilde{V}) = \mathcal{L}_{\text{DIS}}^{\text{BSDE}}(\tilde{\mu}_F).$$

## Connections and equivalences: Path measures and optimal control

- $D = D_{\text{KL}}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{\gamma}^{\tilde{\mu}_B}}) = \mathbb{E} \left[ \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{\gamma}^{\tilde{\mu}_B}}} (X^{\tilde{\mu}_F}) \right]$  leads to stochastic optimal control:

# Connections and equivalences: Path measures and optimal control

- $D = D_{\text{KL}}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{\gamma}^{\tilde{\mu}_B}}) = \mathbb{E} \left[ \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{\gamma}^{\tilde{\mu}_B}}} (X^{\tilde{\mu}_F}) \right]$  leads to stochastic optimal control:

Proposition (Verification theorem, time-reversed diffusion sampling (DIS))

Set  $\tilde{\mu}_F := f + \sigma u$ , i.e. let  $X^u$  be defined by

$$dX_s^u = (f + \sigma u)(X_s^u, s) ds + \sigma(s) dW_s,$$

and fix  $\tilde{\mu}_B = f$ . Consider the loss

$$\mathcal{L}(u) = D_{\text{KL}}(\mathbb{P}_{X^u} | \mathbb{P}_{X^{u^*}}) = D_{\text{KL}}(\mathbb{P}_{X^u} | \mathbb{P}_{\tilde{\gamma}}) - D_{\text{KL}}(\mathbb{P}_{X_0^u} | \mathbb{P}_{Y_T}),$$

where  $\mathbb{P}_{X^u}$  denotes the path space measure of  $X^u$  etc. Then it holds that

$$-\log \mathcal{Z} = \min_{u \in \mathcal{U}} \mathcal{L}(u) := \min_{u \in \mathcal{U}} \mathbb{E} \left[ \int_0^T \left( \frac{1}{2} \|u\|^2 - \text{div}(f) \right) (X_s^u, s) ds + \log \frac{\rho_{Y_T}(X_0^u)}{\rho(X_T^u)} \right],$$

where the unique minimum is attained by  $u^* := \sigma^\top \nabla \log \bar{p}_Y$ .

- **Detour:** let us consider the case with available data samples, but no density  $\rho$

- **Detour:** let us consider the case with available data samples, but no density  $\rho$
- The previous loss is not feasible  $\rightarrow$  we cannot minimize a divergence directly

- **Detour:** let us consider the case with available data samples, but no density  $\rho$
- The previous loss is not feasible  $\rightarrow$  we cannot minimize a divergence directly
- **Trick:** instead of  $D_{\text{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{\tilde{Y}})$  let us consider  $D_{\text{KL}}(\mathbb{P}_Y|\mathbb{P}_{\tilde{X}^u}) = \mathbb{E} \left[ \log \frac{d\mathbb{P}_Y}{d\mathbb{P}_{\tilde{X}^u}}(Y) \right]$   
(which is possible since we have data samples)

- **Detour:** let us consider the case with available data samples, but no density  $\rho$
- The previous loss is not feasible  $\rightarrow$  we cannot minimize a divergence directly
- **Trick:** instead of  $D_{\text{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{\tilde{Y}})$  let us consider  $D_{\text{KL}}(\mathbb{P}_Y|\mathbb{P}_{\tilde{X}^u}) = \mathbb{E} \left[ \log \frac{d\mathbb{P}_Y}{d\mathbb{P}_{\tilde{X}^u}}(Y) \right]$  (which is possible since we have data samples), yielding

$$\mathcal{L}(u) = \min_{u \in \mathcal{U}} \mathbb{E} \left[ \int_0^T \left( \frac{1}{2} \|u\|^2 + \text{div}(\sigma u - f) \right) (Y_s, s) ds + \log \frac{\rho_{X_T^u}(Y_0)}{\rho_{X_0^u}(Y_T)} \right].$$



- **Detour:** let us consider the case with available data samples, but no density  $\rho$
- The previous loss is not feasible  $\rightarrow$  we cannot minimize a divergence directly
- **Trick:** instead of  $D_{\text{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{\tilde{Y}})$  let us consider  $D_{\text{KL}}(\mathbb{P}_Y|\mathbb{P}_{\tilde{X}^u}) = \mathbb{E} \left[ \log \frac{d\mathbb{P}_Y}{d\mathbb{P}_{\tilde{X}^u}}(Y) \right]$  (which is possible since we have data samples), yielding

$$\mathcal{L}(u) = \min_{u \in \mathcal{U}} \mathbb{E} \left[ \int_0^T \left( \frac{1}{2} \|u\|^2 + \text{div}(\sigma u - f) \right) (Y_s, s) ds + \log \frac{p_{X_T^u}(Y_0)}{p_{X_0^u}(Y_T)} \right].$$

- Not tractable since  $p_{X_T^u}$  is not known

## Path space perspective: Recovering score-based generative modeling

- **Detour:** let us consider the case with available data samples, but no density  $\rho$
- The previous loss is not feasible  $\rightarrow$  we cannot minimize a divergence directly
- **Trick:** instead of  $D_{\text{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{\tilde{Y}})$  let us consider  $D_{\text{KL}}(\mathbb{P}_Y|\mathbb{P}_{\tilde{X}^u}) = \mathbb{E} \left[ \log \frac{d\mathbb{P}_Y}{d\mathbb{P}_{\tilde{X}^u}}(Y) \right]$  (which is possible since we have data samples), yielding

$$\mathcal{L}(u) = \min_{u \in \mathcal{U}} \mathbb{E} \left[ \int_0^T \left( \frac{1}{2} \|u\|^2 + \text{div}(\sigma u - f) \right) (Y_s, s) ds + \log \frac{p_{X_T^u}(Y_0)}{p_{X_0^u}(Y_T)} \right].$$

- Not tractable since  $p_{X_T^u}$  is not known  $\rightarrow$  as a remedy, maximize ELBO

$$\mathbb{E} \left[ \log p_{X_T^u}(Y_0) \right] \geq \mathbb{E} \left[ \log p_{X_0^u}(Y_T) - \int_0^T \left( \frac{1}{2} \|u\|^2 + \text{div}(\sigma u - f) \right) (Y_s, s) ds \right]$$

- **Denoising score matching:** Rewrite the divergence and apply Monte Carlo approximation of the time-integral, using  $\tau \sim \mathcal{U}([0, T])$ , so that no trajectories are needed anymore:

- **Denoising score matching:** Rewrite the divergence and apply Monte Carlo approximation of the time-integral, using  $\tau \sim \mathcal{U}([0, T])$ , so that no trajectories are needed anymore:

$$\mathbb{E}[\log p_{X_\tau^u}(Y_0)] \geq \mathbb{E}\left[\log p_{X_0^u}(Y_T) - \int_0^T \left(\frac{1}{2}\|u\|^2 + \operatorname{div}(\sigma u - f)\right)(Y_s, s) ds\right]$$

# Path space perspective: Recovering score-based generative modeling

- **Denoising score matching:** Rewrite the divergence and apply Monte Carlo approximation of the time-integral, using  $\tau \sim \mathcal{U}([0, T])$ , so that no trajectories are needed anymore:

Divergence theorem:  $u \cdot \sigma^\top \nabla \log p_{Y_s|Y_0}$  (in expectation)

$$\mathbb{E}[\log p_{X_\tau^u}(Y_0)] \geq \mathbb{E} \left[ \log p_{X_0^u}(Y_T) - \int_0^T \left( \frac{1}{2} \|u\|^2 + \text{div}(\sigma u) - \text{div}(f) \right)(Y_s, s) ds \right]$$

(Annotations: 'const.' with arrows pointing to  $\log p_{X_0^u}(Y_T)$ ,  $\text{div}(\sigma u)$ , and  $\text{div}(f)$ ; a vertical arrow points from the divergence theorem text to  $\text{div}(\sigma u)$ )

# Path space perspective: Recovering score-based generative modeling

- **Denosing score matching:** Rewrite the divergence and apply Monte Carlo approximation of the time-integral, using  $\tau \sim \mathcal{U}([0, T])$ , so that no trajectories are needed anymore:

$$\begin{aligned} & \text{Divergence theorem: } u \cdot \sigma^\top \nabla \log p_{Y_s|Y_0} \text{ (in expectation)} \\ & \begin{array}{ccc} \text{const.} & & \text{const.} \\ \downarrow & & \downarrow \end{array} \\ \mathbb{E}[\log p_{X_\tau^u}(Y_0)] & \geq \mathbb{E} \left[ \log p_{X_0^u}(Y_\tau) - \int_0^\tau \left( \frac{1}{2} \|u\|^2 + \text{div}(\sigma u) - \text{div}(f) \right)(Y_s, s) ds \right] \\ & = \frac{T}{2} \underbrace{\mathbb{E} \left[ \|u(Y_\tau, \tau) - \sigma^\top(\tau) \nabla \log p_{Y_\tau|Y_0}(Y_\tau|Y_0)\|^2 \right]}_{\text{denosing score matching}} + \text{const.} \end{aligned}$$

- We propose a novel divergence:

Definition (Log-variance divergence)

$$D_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) := \text{Var} \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)$$

- We propose a novel divergence:

## Definition (Log-variance divergence)

$$D_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) := \text{Var} \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)$$

- In principle arbitrary choice for  $\tilde{\mu}_R$  allows to balance exploration and exploitation.



- We propose a novel divergence:

## Definition (Log-variance divergence)

$$D_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) := \text{Var} \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)$$

- In principle arbitrary choice for  $\tilde{\mu}_R$  allows to balance exploration and exploitation.
- No differentiation through the SDE solver.

# Path space perspective: Novel divergences

- We propose a novel divergence:

## Definition (Log-variance divergence)

$$D_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) := \text{Var} \left( \log \frac{d\mathbb{P}_{X^{\tilde{\mu}_F}}}{d\mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}}(X^{\tilde{\mu}_R}) \right)$$

- In principle arbitrary choice for  $\tilde{\mu}_R$  allows to balance exploration and exploitation.
- No differentiation through the SDE solver.

## Proposition (Equivalence with KL divergence)

$$\frac{1}{2} \left( \frac{\delta}{\delta \tilde{\mu}_F} D_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) \Big|_{\tilde{\mu}_R = \tilde{\mu}_F} \right) = \frac{\delta}{\delta \tilde{\mu}_F} D_{KL}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}})$$

## Proposition (Control variate)

$$\frac{1}{2} \left( \frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{Y^{\tilde{\mu}_B}}) \Big|_{\tilde{\mu}_R = \tilde{\mu}_F} \right) \text{ is a control variate version of } \frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{KL}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{Y^{\tilde{\mu}_B}}).$$

## Proposition (Control variate)

$$\frac{1}{2} \left( \frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) \Big|_{\tilde{\mu}_R = \tilde{\mu}_F} \right) \text{ is a control variate version of } \frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{KL}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}).$$

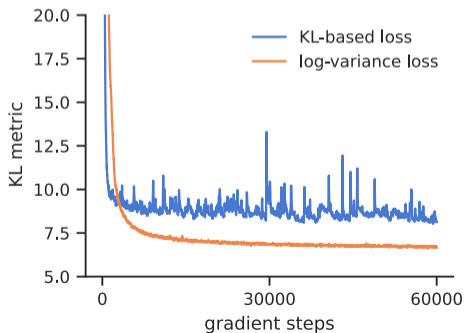
- $X^{\text{CV}} = X + C$ , where  $\mathbb{E}[C] = 0$

# The log-variance divergence: Variance reduction

## Proposition (Control variate)

$\frac{1}{2} \left( \frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) \Big|_{\tilde{\mu}_R = \tilde{\mu}_F} \right)$  is a control variate version of  $\frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{KL}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}})$ .

- $X^{CV} = X + C$ , where  $\mathbb{E}[C] = 0$

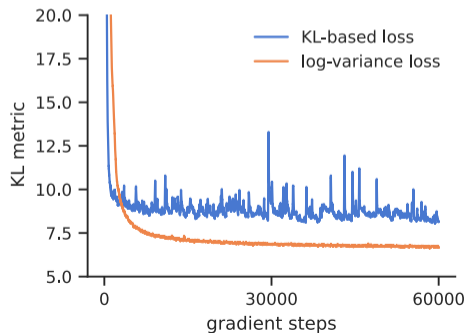


# The log-variance divergence: Variance reduction

## Proposition (Control variate)

$\frac{1}{2} \left( \frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) \Big|_{\tilde{\mu}_R = \tilde{\mu}_F} \right)$  is a control variate version of  $\frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{KL}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}})$ .

- $X^{CV} = X + C$ , where  $\mathbb{E}[C] = 0$
- This leads to variance reduction in the estimated gradient.

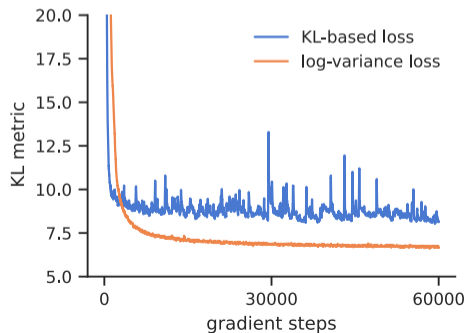


# The log-variance divergence: Variance reduction

## Proposition (Control variate)

$\frac{1}{2} \left( \frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}}) \Big|_{\tilde{\mu}_R = \tilde{\mu}_F} \right)$  is a control variate version of  $\frac{\delta}{\delta \tilde{\mu}_F} \widehat{D}_{KL}(\mathbb{P}_{X^{\tilde{\mu}_F}} | \mathbb{P}_{\tilde{Y}^{\tilde{\mu}_B}})$ .

- $X^{CV} = X + C$ , where  $\mathbb{E}[C] = 0$
- This leads to variance reduction in the estimated gradient.
- Usually implying faster and better convergence of gradient based optimization.



# The log-variance divergence: Robustness properties

## Proposition (Robustness at solution)

$$\text{Var} \left( \frac{\delta}{\delta \tilde{\mu}_F} \Big|_{\tilde{\mu}_F = \mu_F} \widehat{D}_{\text{LV}}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{Y^{\tilde{\mu}_B}}) \right) = 0, \quad \text{Var} \left( \frac{\delta}{\delta \tilde{\mu}_B} \Big|_{\tilde{\mu}_B = \mu_B} \widehat{D}_{\text{LV}}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{Y^{\tilde{\mu}_F}}) \right) = 0.$$



# The log-variance divergence: Robustness properties

## Proposition (Robustness at solution)

$$\text{Var} \left( \frac{\delta}{\delta \tilde{\mu}_F} \Big|_{\tilde{\mu}_F = \mu_F} \widehat{D}_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{Y^{\tilde{\mu}_B}}) \right) = 0, \quad \text{Var} \left( \frac{\delta}{\delta \tilde{\mu}_B} \Big|_{\tilde{\mu}_B = \mu_B} \widehat{D}_{LV}^{\tilde{\mu}_R}(\mathbb{P}_{X^{\tilde{\mu}_F}}, \mathbb{P}_{Y^{\tilde{\mu}_F}}) \right) = 0.$$

## Proposition (Robustness in high dimensions)

$$\frac{\sqrt{\text{Var} \left( \widehat{D}_{LV}^w \left( \bigotimes_{i=1}^d \mathbb{P}_i, \bigotimes_{i=1}^d \mathbb{Q}_i \right) \right)}}{D_{LV}^w \left( \bigotimes_{i=1}^d \mathbb{P}_i, \bigotimes_{i=1}^d \mathbb{Q}_i \right)}$$

*can be bounded uniformly in  $d$ .*

- **Alternative:** PINN-based losses for stochastic and deterministic evolutions.

# A simulation-free attempt based on PDEs: PINN-based losses

- **Alternative:** PINN-based losses for stochastic and deterministic evolutions.
- Consider parametrization (encoding the boundary conditions)

$$\tilde{V}_{\varphi,z}(\cdot, t) = \frac{t}{T} \log \frac{\rho_{\text{target}}}{z(t)} + \left(1 - \frac{t}{T}\right) \log p_{\text{prior}} + \frac{t}{T} \left(1 - \frac{t}{T}\right) \varphi(\cdot, t),$$

where we learn  $z$  and  $\varphi$ , cf. Máté & Fleuret, 2023.

- **Alternative:** PINN-based losses for stochastic and deterministic evolutions.
- Consider parametrization (encoding the boundary conditions)

$$\tilde{V}_{\varphi, z}(\cdot, t) = \frac{t}{T} \log \frac{\rho_{\text{target}}}{z(t)} + \left(1 - \frac{t}{T}\right) \log p_{\text{prior}} + \frac{t}{T} \left(1 - \frac{t}{T}\right) \varphi(\cdot, t),$$

where we learn  $z$  and  $\varphi$ , cf. Máté & Fleuret, 2023.

- We can then minimize

$$\mathcal{L}(\tilde{\mu}, \tilde{V}) = \mathbb{E} \left[ \left( \mathcal{R}(\tilde{\mu}, \tilde{V})(\xi, \tau) \right)^2 \right],$$

where  $(\xi, \tau) \sim \nu$  are sampled from a measure  $\nu$ , e.g.  $\nu = \text{Unif}(\Omega \times [0, T])$ ,  $\Omega \subset \mathbb{R}^d$ .

- BSDE-based losses:

# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.

# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.
  - ▶ Gradients of the solutions (usually corresponding to the learned drift) can be learned directly.

# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.
  - ▶ Gradients of the solutions (usually corresponding to the learned drift) can be learned directly.
  - ▶ Only stochastic dynamics can be approached.



# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.
  - ▶ Gradients of the solutions (usually corresponding to the learned drift) can be learned directly.
  - ▶ Only stochastic dynamics can be approached.
- PINN-based losses:

# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.
  - ▶ Gradients of the solutions (usually corresponding to the learned drift) can be learned directly.
  - ▶ Only stochastic dynamics can be approached.
- PINN-based losses:
  - ▶ Can be readily applied to deterministic evolutions.

# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.
  - ▶ Gradients of the solutions (usually corresponding to the learned drift) can be learned directly.
  - ▶ Only stochastic dynamics can be approached.
- PINN-based losses:
  - ▶ Can be readily applied to deterministic evolutions.
  - ▶ Simulation-free, no time-discretization.

# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.
  - ▶ Gradients of the solutions (usually corresponding to the learned drift) can be learned directly.
  - ▶ Only stochastic dynamics can be approached.
- PINN-based losses:
  - ▶ Can be readily applied to deterministic evolutions.
  - ▶ Simulation-free, no time-discretization.
  - ▶ Off-policy training comes by design.

# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.
  - ▶ Gradients of the solutions (usually corresponding to the learned drift) can be learned directly.
  - ▶ Only stochastic dynamics can be approached.
- PINN-based losses:
  - ▶ Can be readily applied to deterministic evolutions.
  - ▶ Simulation-free, no time-discretization.
  - ▶ Off-policy training comes by design.
  - ▶ Need to know “essential support” of target density.

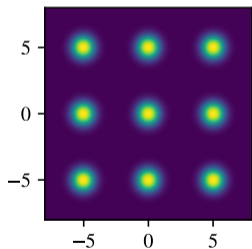
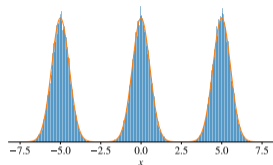
# Numerical implications of the losses

- BSDE-based losses:
  - ▶ Neither second-order nor time derivatives have to be computed.
  - ▶ Gradients of the solutions (usually corresponding to the learned drift) can be learned directly.
  - ▶ Only stochastic dynamics can be approached.
- PINN-based losses:
  - ▶ Can be readily applied to deterministic evolutions.
  - ▶ Simulation-free, no time-discretization.
  - ▶ Off-policy training comes by design.
  - ▶ Need to know “essential support” of target density.
  - ▶ Training is sensitive to hyperparameter tuning.

- We consider the following losses:

Method	Stochastic	Deterministic	BSDE version	Unique
General bridge	$\mathcal{L}_{\log\text{FP}}(\tilde{\mu}, \tilde{V})$	$\mathcal{L}_{\log\text{CE}}(\tilde{\mu}, \tilde{V})$	Bridge	✗
Prescribed bridge	$\mathcal{L}_{\log\text{FP}}^{\text{anneal}}(\tilde{\mu})$	$\mathcal{L}_{\log\text{CE}}^{\text{anneal}}(\tilde{\mu})$	CMCD	✓
Score-based	$\mathcal{L}_{\text{score}}(\tilde{V})$		DIS	✓
SB & OT	$\mathcal{L}_{\text{SB}}(\tilde{\mu}, \tilde{V})$	$\mathcal{L}_{\text{OT}}(\tilde{\mu}, \tilde{V})$		✓

# Numerical examples: Gaussian mixture ( $d = 2, 9$ modes)

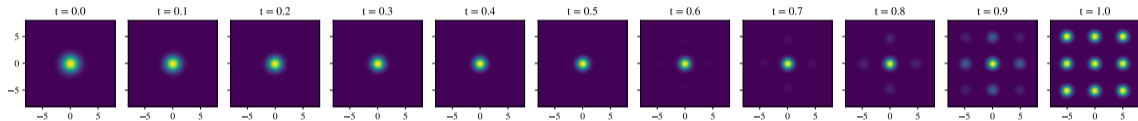


Problem	Method	Loss	$\Delta \log Z \downarrow$	$\mathcal{W}_\gamma^2 \downarrow$	ESS $\uparrow$	$\Delta \text{std} \downarrow$	sec./it. $\downarrow$
GMM ( $d = 2$ )	PIS-KL		1.094	0.467	0.0051	1.937	0.503
	PIS-LV		0.046	<b>0.020</b>	0.9093	0.023	0.500
	DIS-KL		1.551	0.064	0.0226	2.522	0.565
	DIS-LV		0.056	<b>0.020</b>	0.8660	0.004	0.536
	SDE	$\mathcal{L}_{\log\text{FP}}$	<b>0.000</b>	<b>0.020</b>	<b>1.0000</b>	0.004	0.011
	SDE-anneal	$\mathcal{L}_{\log\text{FP}}^{\text{anneal}}$	5.364	0.172	0.1031	0.209	0.062
	SDE-score	$\mathcal{L}_{\text{score}}$	0.009	<b>0.020</b>	0.9818	0.096	0.013
	SB	$\mathcal{L}_{\text{SB}}$	0.002	<b>0.020</b>	0.9959	0.050	0.017
	ODE	$\mathcal{L}_{\log\text{CE}}$	<b>0.000</b>	<b>0.020</b>	<b>1.0000</b>	<b>0.003</b>	<b>0.008</b>
ODE-anneal	$\mathcal{L}_{\log\text{CE}}^{\text{anneal}}$	4.227	0.044	0.1427	0.753	0.020	
OT	$\mathcal{L}_{\text{OT}}$	0.005	0.057	0.9932	0.065	0.080	



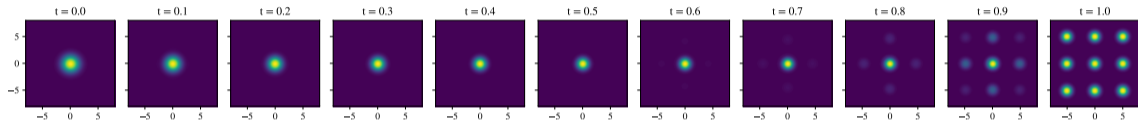
# Numerical examples: Gaussian mixture ( $d = 2, 9$ modes)

- Geometric annealing path can be suboptimal ( $\mathcal{L}_{\log\text{CE}}^{\text{anneal}}$ ):

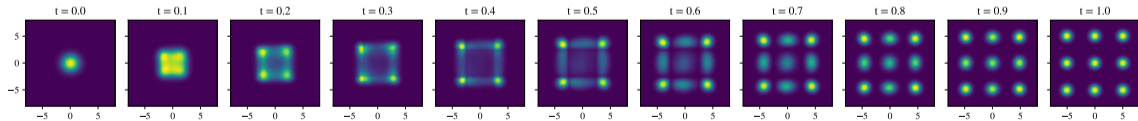


# Numerical examples: Gaussian mixture ( $d = 2, 9$ modes)

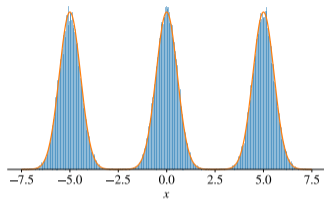
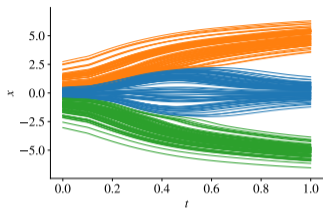
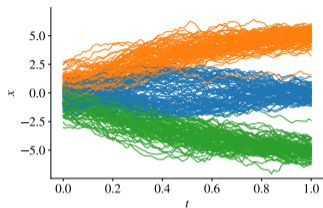
- Geometric annealing path can be suboptimal ( $\mathcal{L}_{\log\text{CE}}^{\text{anneal}}$ ):



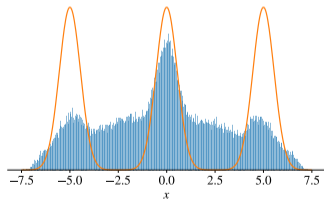
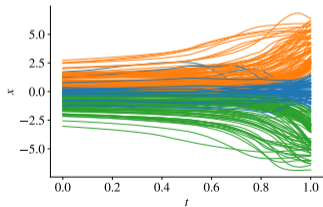
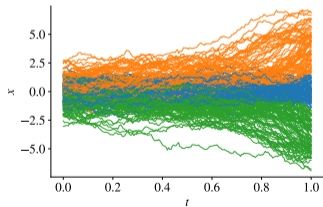
- The learned path seems to be more appropriate ( $\mathcal{L}_{\log\text{CE}}$ ):



# Numerical examples: Gaussian mixture ( $d = 2, 9$ modes)

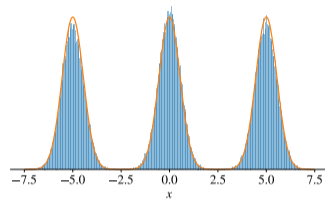
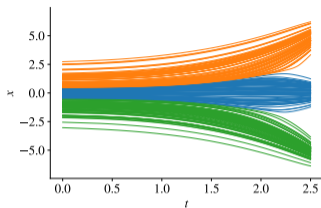
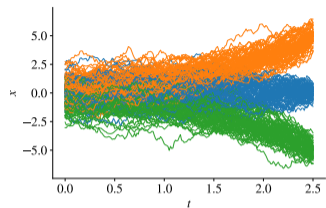


$\mathcal{L}_{\log\text{FP}}$

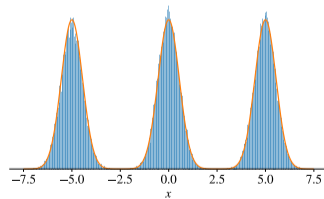
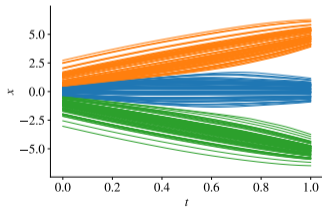
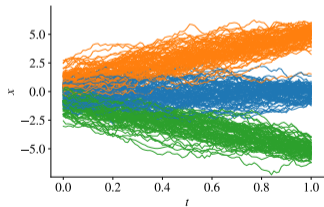


$\mathcal{L}_{\text{anneal}\log\text{FP}}$

# Numerical examples: Gaussian mixture ( $d = 2, 9$ modes)



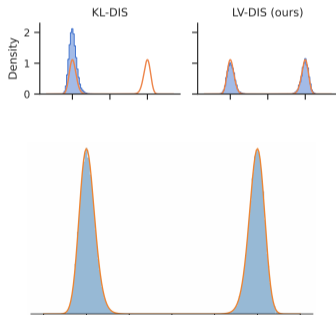
$\mathcal{L}_{\text{score}}$



$\mathcal{L}_{\text{SB}}$

# Numerical examples: Double well ( $d = 5, 32$ modes)

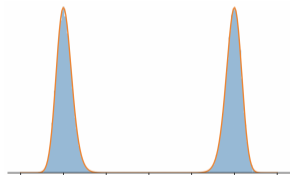
$$\rho(x) := \exp\left(-\sum_{i=1}^5 (x_i^2 - 4)^2\right)$$



Problem	Method	Loss	$\Delta \log Z \downarrow$	$\mathcal{W}_\gamma^2 \downarrow$	ESS $\uparrow$	$\Delta \text{std} \downarrow$	sec./it. $\downarrow$
MW ( $d = 5, m = 5, \delta = 4$ )	PIS-KL		3.567	1.699	0.0004	1.409	0.441
	PIS-LV		0.214	0.121	0.6744	0.001	0.402
	DIS-KL		1.462	1.175	0.0012	0.431	0.490
	DIS-LV		0.375	0.120	0.4519	0.001	0.437
	SDE	$\mathcal{L}_{\log \text{FP}}$	0.161	0.123	0.8167	0.016	0.017
	SDE-anneal	$\mathcal{L}_{\log \text{FP}}^{\text{anneal}}$	0.842	0.257	0.3464	0.004	0.014
	SDE-score	$\mathcal{L}_{\text{score}}$	3.969	0.427	0.0124	0.004	0.026
	SB	$\mathcal{L}_{\text{SB}}$	7.855	0.328	0.0314	0.045	0.029
	ODE	$\mathcal{L}_{\log \text{CE}}$	<b>0.000</b>	<b>0.118</b>	<b>0.9993</b>	<b>0.000</b>	<b>0.008</b>
	ODE-anneal	$\mathcal{L}_{\log \text{CE}}^{\text{anneal}}$	0.025	0.121	0.9506	0.005	0.010
OT	$\mathcal{L}_{\text{OT}}$	0.010	0.120	0.9862	0.002	0.020	

# Numerical examples: Double well ( $d = 50, 32$ modes)

$$\rho(x) := \exp \left( - \sum_{i=1}^5 (x_i^2 - 2)^2 - \frac{1}{2} \sum_{i=6}^{50} x_i^2 \right)$$



Problem	Method	Loss	$\Delta \log Z \downarrow$	$\mathcal{W}_\gamma^2 \downarrow$	ESS $\uparrow$	$\Delta \text{std} \downarrow$	sec./it. $\downarrow$
MW ( $d = 50, m = 5, \delta = 2$ )	PIS-KL		0.101	6.821	0.8172	0.001	0.479
	PIS-LV		0.087	6.823	0.8453	<b>0.000</b>	0.416
	DIS-KL		1.785	6.854	0.0225	0.009	0.522
	DIS-LV		1.783	6.855	0.0227	0.009	0.450
	SDE	$\mathcal{L}_{\log \text{FP}}$	0.038	6.820	0.9511	0.001	0.050
	SDE-anneal	$\mathcal{L}_{\log \text{FP}}^{\text{anneal}}$	0.270	6.899	0.9171	0.021	0.067
	SDE-score	$\mathcal{L}_{\text{score}}$	1.989	<b>6.803</b>	0.1065	0.016	0.053
	SB	$\mathcal{L}_{\text{SB}}$	189.71	7.552	0.0106	0.051	0.053
	ODE	$\mathcal{L}_{\log \text{CE}}$	<b>0.003</b>	6.815	<b>0.9937</b>	0.002	<b>0.023</b>
	ODE-anneal	$\mathcal{L}_{\log \text{CE}}^{\text{anneal}}$	1.759	6.821	0.2100	0.017	0.043
OT	$\mathcal{L}_{\text{OT}}$	0.104	6.824	0.9027	0.001	0.043	

$$\rho(\phi) = \exp \left( - \sum_{x \in \Lambda} \left( -2\kappa \sum_{\hat{\mu}=1}^2 \phi(x)\phi(x + \hat{\mu}) + (1 - 2\lambda)\phi(x)^2 + \lambda\phi(x)^4 \right) \right)$$

$$\rho(x) = \exp \left( - \sum_{i=1}^d (-2\kappa(x_i x_{i-L} + x_i x_{i+1}) + (1 - 2\lambda)x_i^2 + \lambda x_i^4) \right)$$



$$\rho(x) = \exp \left( - \sum_{i=1}^d (-2\kappa(x_i x_{i-L} + x_i x_{i+1}) + (1 - 2\lambda)x_i^2 + \lambda x_i^4) \right)$$

- Difficulty depends non-trivially on the choices of  $\kappa$  and  $\lambda$ .

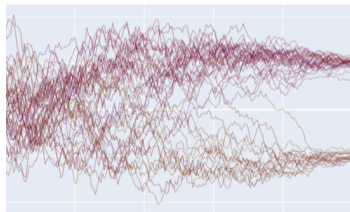
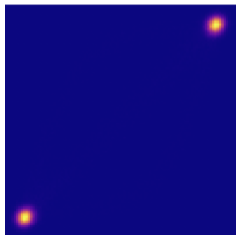
$$\rho(x) = \exp \left( - \sum_{i=1}^d (-2\kappa(x_i x_{i-L} + x_i x_{i+1}) + (1 - 2\lambda)x_i^2 + \lambda x_i^4) \right)$$

- Difficulty depends non-trivially on the choices of  $\kappa$  and  $\lambda$ .
- Preliminary results, not incorporating any symmetries (with  $\lambda = 0.022$ ), using Hamiltonian dynamics and learned priors.

# Numerical examples: $\phi^4$ field theory

$$\rho(x) = \exp \left( - \sum_{i=1}^d (-2\kappa(x_i x_{i-L} + x_i x_{i+1}) + (1 - 2\lambda)x_i^2 + \lambda x_i^4) \right)$$

- Difficulty depends non-trivially on the choices of  $\kappa$  and  $\lambda$ .
- Preliminary results, not incorporating any symmetries (with  $\lambda = 0.022$ ), using Hamiltonian dynamics and learned priors.



$\kappa$	dimension	$\Delta \log Z \downarrow$	ESS $\uparrow$
0.2	$16 \times 8$	0.0002	0.91
	$64 \times 8$	0.0057	0.08
0.5	$16 \times 8$	0.0176	0.79
	$64 \times 8$	0.0169	0.05

- We have established optimal control, PDE and path space perspectives on generative modeling.

- We have established optimal control, PDE and path space perspectives on generative modeling.
- This allows to carry over respective methods and theory to generative modeling.

- We have established optimal control, PDE and path space perspectives on generative modeling.
- This allows to carry over respective methods and theory to generative modeling.
- We introduced algorithms to sample from an (unnormalized) density, which are already competitive to MCMC/SMC.

- We have established optimal control, PDE and path space perspectives on generative modeling.
- This allows to carry over respective methods and theory to generative modeling.
- We introduced algorithms to sample from an (unnormalized) density, which are already competitive to MCMC/SMC.
- The log-variance divergence outperforms the KL divergence.

- We have established optimal control, PDE and path space perspectives on generative modeling.
- This allows to carry over respective methods and theory to generative modeling.
- We introduced algorithms to sample from an (unnormalized) density, which are already competitive to MCMC/SMC.
- The log-variance divergence outperforms the KL divergence.
- PINNs seem to be suitable for learning dynamical systems for sampling.



- We have established optimal control, PDE and path space perspectives on generative modeling.
- This allows to carry over respective methods and theory to generative modeling.
- We introduced algorithms to sample from an (unnormalized) density, which are already competitive to MCMC/SMC.
- The log-variance divergence outperforms the KL divergence.
- PINNs seem to be suitable for learning dynamical systems for sampling.
- Often, non-uniqueness helps to find a “better” solution.

- **General framework:** (stochastic) normalizing flows and GFlowNets can be incorporated, however, continuous-time perspective allows for more flexibility.

- **General framework:** (stochastic) normalizing flows and GFlowNets can be incorporated, however, continuous-time perspective allows for more flexibility.
- **SMC:** Annealed importance sampling and resampling can be naturally integrated. (Diffusion model version of CRAFT.)

- **General framework:** (stochastic) normalizing flows and GFlowNets can be incorporated, however, continuous-time perspective allows for more flexibility.
- **SMC:** Annealed importance sampling and resampling can be naturally integrated. (Diffusion model version of CRAFT.)
- **Hamiltonian dynamics:** underdamped versions can be considered and lead to improved performance.

# Thank you for your attention!

richter@zib.de

## References:

- N. Nüsken, L.R. *Solving high-dimensional Hamilton–Jacobi–Bellman PDEs using neural networks: perspectives from the theory of controlled diffusions and measures on path space*. Partial Differential Equations and Applications, 2021.
- J. Berner, L.R., K. Ullrich. *An optimal control perspective on diffusion-based generative modeling*. TMLR, 2024.
- L.R., J. Berner. *Improved sampling via learned diffusions*. ICLR, 2024.
- J. Sun, J. Berner, L.R. et al. *Dynamical measure transport and neural PDE solvers for sampling*. arXiv:2407.07873, 2024.