

Pursuits and Challenges Towards Simulation-free Training of Neural Sampler

Jiajun He & Yuanqi Du

@MILA Sampling Reading Group

Collaborators and Supervisors:

Francisco Vargas, Dinghuai Zhang, Shreyas Padhy, RuiKang OuYang, Carla Gomes, José Miguel Hernández-Lobato

Sampling

Unnormalized density function:

$$p_{\text{target}}(x) = \frac{\tilde{p}(x)}{Z}, \quad Z = \int \tilde{p}(x) dx$$

Obtain sample $x \sim p_{\text{target}}$.

Sampling

Unnormalized density function:

$$p_{\text{target}}(x) = \frac{\tilde{p}(x)}{Z}, \quad Z = \int \tilde{p}(x) dx$$

Obtain sample $x \sim p_{\text{target}}$.

👉 Bayesian inference: $p_{\text{target}} \propto \text{likelihood} \times \text{prior}$

Sampling

Unnormalized density function:

$$p_{\text{target}}(x) = \frac{\tilde{p}(x)}{Z}, \quad Z = \int \tilde{p}(x) dx$$

Obtain sample $x \sim p_{\text{target}}$.

- 👉 Bayesian inference: $p_{\text{target}} \propto \text{likelihood} \times \text{prior}$
- 👉 Boltzmann distribution (molecules, etc): $p_{\text{target}} \propto \exp(-\beta U)$

Sampling

Unnormalized density function:

$$p_{\text{target}}(x) = \frac{\tilde{p}(x)}{Z}, \quad Z = \int \tilde{p}(x) dx$$

Obtain sample $x \sim p_{\text{target}}$.

- 👉 Bayesian inference: $p_{\text{target}} \propto \text{likelihood} \times \text{prior}$
- 👉 Boltzmann distribution (molecules, etc): $p_{\text{target}} \propto \exp(-\beta U)$
- 👉 Rare event: $p_{\text{target}}(x) \propto \mathbf{1}_B(x)p(x)$

Sampling – classical approach

Markov chain Monte Carlo (MCMC)

Sampling – classical approach

Markov chain Monte Carlo (MCMC)

For example, unadjusted Langevin dynamics:

$$dX_t = \nabla \log \tilde{p}(X_t) dt + \sqrt{2} dW_t$$

Sampling – classical approach

Markov chain Monte Carlo (MCMC)

For example, unadjusted Langevin dynamics:

$$dX_t = \nabla \log \tilde{p}(X_t) dt + \sqrt{2} dW_t$$

 dependent samples; auto-correlation reduces efficiency sample size

Sampling – classical approach

Markov chain Monte Carlo (MCMC)

For example, unadjusted Langevin dynamics:

$$dX_t = \nabla \log \tilde{p}(X_t) dt + \sqrt{2} dW_t$$

- 😞 dependent samples; auto-correlation reduces efficiency sample size
- 😞 ergodicity; only guarantee convergence with infinite steps

Neural samplers

Train a neural network to amortize the sampling process

Neural samplers

Train a neural network to amortize the sampling process

😊 independent samples!

😊 can mix in finite time

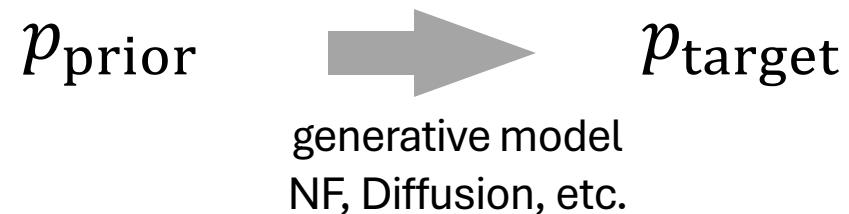
Neural samplers

Train a neural network to amortize the sampling process

😊 independent samples!

😊 can mix in finite time

Neural samplers are in fact generative models:



Diffusion Neural samplers

Train a diffusion (like) model

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t,$$

Diffusion Neural samplers

Train a diffusion (like) model

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t,$$

transporting samples from p_{prior} to p_{target} :

$$X_0 \sim p_{\text{prior}}, \text{ and want } X_T \sim p_{\text{target}}$$

Diffusion Neural samplers

1. Time-reversal sampler

Diffusion Neural samplers

1. Time-reversal sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

Want: $X_T \sim p_{\text{target}}$

Diffusion Neural samplers

1. Time-reversal sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

Set a target process: $dY_t = g(Y_t, t)dt + \sigma\sqrt{2}dW_t, Y_0 \sim p_{\text{target}},$

Diffusion Neural samplers

1. Time-reversal sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

Set a target process: $dY_t = \underbrace{g(Y_t, t)}_{\text{a simple function, e.g., 0}}dt + \sigma\sqrt{2}dW_t, Y_0 \sim p_{\text{target}},$

a simple function, e.g., 0

Diffusion Neural samplers

1. Time-reversal sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

align

Set a target process: $dY_t = \underbrace{g(Y_t, t)}_{\text{a simple function, e.g., 0}}dt + \sigma\sqrt{2}dW_t, Y_0 \sim p_{\text{target}},$

a simple function, e.g., 0

Diffusion Neural samplers

1. Time-reversal sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

align
 $X_t \sim Y_{T-t}$

Set a target process: $dY_t = \underbrace{g(Y_t, t)}_{\text{a simple function, e.g., 0}}dt + \sigma\sqrt{2}dW_t, Y_0 \sim p_{\text{target}},$

a simple function, e.g., 0

Diffusion Neural samplers

1. Time-reversal sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

align

$$X_t \sim Y_{T-t}$$

$$X_T \sim p_{\text{target}}$$

Set a target process:

$$dY_t = \underbrace{g(Y_t, t)}_{\text{a simple function, e.g., } 0}dt + \sigma\sqrt{2}dW_t, Y_0 \sim p_{\text{target}},$$

a simple function, e.g., 0

Diffusion Neural samplers

1. Time-reversal sampler

Want a sample process (prior to target),

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

align

$$X_t \sim Y_{T-t}$$

$$X_T \sim p_{\text{target}}$$

Set a target process:

$$dY_t = \underbrace{g(Y_t, t)}_{\text{a simple function, e.g., 0}}dt + \sigma\sqrt{2}dW_t, Y_t \sim p_{\text{target}},$$

a simple function, e.g., 0

Diffusion Neural samplers

1. Time-reversal sampler

Want a sample process (prior to target),

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

To be the **time-reversal**,

align

$$X_t \sim Y_{T-t}$$

$$X_T \sim p_{\text{target}}$$

Set a target process:

$$dY_t = \underbrace{g(Y_t, t)}_{\text{a simple function, e.g., } 0}dt + \sigma\sqrt{2}dW_t, Y_t \sim p_{\text{target}},$$

a simple function, e.g., 0

Diffusion Neural samplers

1. Time-reversal sampler

Want a sample process (prior to target),

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

To be the **time-reversal**,

of a simple process (target to prior)

align
 $X_t \sim Y_{T-t}$
 p_{target}

Set a target process:

$$dY_t = \underbrace{g(Y_t, t)}_{\text{a simple function, e.g., } 0}dt + \sigma\sqrt{2}dW_t, Y_t \sim p_{\text{target}},$$

a simple function, e.g., 0

Diffusion Neural samplers

1. Time-reversal sampler

This includes

- (1) DDS (denoising diffusion sampler)
 - (2) PIS (path integral sampler)
 - (3) DIS (diffusion time-reversal sampler)
 - (4) GFlowNet (generative flow network)
 - (5) iDEM (iterated denoising energy matching)
 - (6) RDMC (reversal diffusion monte carlo)
 - (7) PINN (physics-informed neural networks) sampler
- ...

Diffusion Neural samplers

Any other ways?

Diffusion Neural samplers

2. Escorted transport sampler

Diffusion Neural samplers

2. Escorted transport sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

Want: $X_T \sim p_{\text{target}}$

Diffusion Neural samplers

2. Escorted transport sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

define a sequence of interpolants $\pi_t: \pi_0 = p_{\text{prior}}, \pi_T = p_{\text{target}}$

Diffusion Neural samplers

2. Escorted transport sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

define a sequence of interpolants $\pi_t: \pi_0 = p_{\text{prior}}, \pi_T = p_{\text{target}}$

If marginal of $X_t \sim \pi_t \rightarrow X_T \sim p_{\text{target}}$

Diffusion Neural samplers

2. Escorted transport sampler

Want a sample process (prior to target),

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

define a sequence of interpolants $\pi_t: \pi_0 = p_{\text{prior}}, \pi_T = p_{\text{target}}$

If marginal of $X_t \sim \pi_t \implies X_T \sim p_{\text{target}}$

Diffusion Neural samplers

2. Escorted transport sampler

Want a sample process (prior to target),

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

whose **marginal densities aligns with**

define a sequence of interpolants $\pi_t: \pi_0 = p_{\text{prior}}, \pi_T = p_{\text{target}}$

If marginal of $X_t \sim \pi_t \implies X_T \sim p_{\text{target}}$

Diffusion Neural samplers

2. Escorted transport sampler

Want a sample process (prior to target),

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

whose **marginal densities aligns with**

define a sequence of interpolants $\pi_t: \pi_0 = p_{\text{prior}}, \pi_T = p_{\text{target}}$

pre-defined interpolants between prior and target

If marginal of $X_t \sim \pi_t \implies X_T \sim p_{\text{target}}$

Diffusion Neural samplers

2. Escorted transport sampler

This includes

- (1) CMCD (Controlled Monte Carlo Diffusions)
- (2) NETS (non-equilibrium transport sampler)
- (3) PINN (physics-informed neural networks) sampler
- (4) LFIS (Liouville Flow Importance Sampler)

...

Diffusion Neural samplers

Any other ways?

Diffusion Neural samplers

3. Annealed variance reduction sampler

Diffusion Neural samplers

3. Annealed variance reduction sampler

$$dX_t = f_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

Want: $X_T \sim p_{\text{target}}$

Diffusion Neural samplers

3. Annealed variance reduction sampler

What if we do not train it?

$$dX_t = \overbrace{f_\theta(X_t, t)} dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

Want: $X_T \sim p_{\text{target}}$

Diffusion Neural samplers

3. Annealed variance reduction sampler

What if we do not train it?

$$dX_t = \overbrace{f(X_t, t)} dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

$$X_T \neq p_{\text{target}}$$

Diffusion Neural samplers

3. Annealed variance reduction sampler

$$dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}},$$

$$dX_t = g(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}},$$

Diffusion Neural samplers

3. Annealed variance reduction sampler

$$dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}}, \rightarrow \vec{\mathbf{Q}}(X)$$

$$dX_t = g(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}}, \rightarrow \vec{\mathbf{P}}(X)$$

Diffusion Neural samplers

3. Annealed variance reduction sampler

$$dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}}, \quad \rightarrow \quad \vec{\mathbf{Q}}(X)$$

$$dX_t = g(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}}, \quad \rightarrow \quad \overleftarrow{\mathbf{P}}(X)$$

$$\text{Importance weight: } \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}(X)}$$

Diffusion Neural samplers

3. Annealed variance reduction sampler

$$dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}}, \rightarrow \vec{\mathbf{Q}}(X)$$

$$dX_t = g_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}}, \rightarrow \overleftarrow{\mathbf{P}}_\theta(X)$$

$$\text{Importance weight: } \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}_\theta(X)}$$

Diffusion Neural samplers

3. Annealed variance reduction sampler

Align $\left(\begin{array}{l} dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}}, \rightarrow \vec{\mathbf{Q}}(X) \\ dX_t = g_{\theta}(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}}, \rightarrow \overleftarrow{\mathbf{P}}_{\theta}(X) \end{array} \right.$

Importance weight: $\frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}_{\theta}(X)}$

Diffusion Neural samplers

3. Annealed variance reduction sampler

Align $\left(\begin{array}{l} dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}}, \rightarrow \vec{\mathbf{Q}}(X) \\ dX_t = g_{\theta}(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}}, \rightarrow \overleftarrow{\mathbf{P}}_{\theta}(X) \end{array} \right.$

Small variance

$$\text{Importance weight: } \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}_{\theta}(X)}$$

Diffusion Neural samplers

3. Annealed variance reduction sampler

Predefine a sample process (prior to target),

Align $\left(\begin{array}{l} dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}}, \Rightarrow \vec{Q}(X) \\ dX_t = g_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}}, \Rightarrow \overleftarrow{P}_\theta(X) \end{array} \right.$

Small variance

Importance weight: $\frac{d\vec{Q}(X)}{d\overleftarrow{P}_\theta(X)}$

Diffusion Neural samplers

3. Annealed variance reduction sampler

Predefine a sample process (prior to target),

Align $\left(\begin{array}{l} dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}}, \Rightarrow \vec{Q}(X) \\ \text{define or train a backward process (target to prior),} \\ dX_t = g_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}}, \Rightarrow \overleftarrow{P}_\theta(X) \end{array} \right.$

Small variance

Importance weight: $\frac{d\vec{Q}(X)}{d\overleftarrow{P}_\theta(X)}$

Diffusion Neural samplers

3. Annealed variance reduction sampler

Predefine a sample process (prior to target),

Align $\left(\begin{array}{l} dX_t = f(X_t, t)dt + \sigma\sqrt{2}dW_t, X_0 \sim p_{\text{prior}}, \Rightarrow \vec{Q}(X) \\ \text{define or train a backward process (target to prior),} \\ dX_t = g_\theta(X_t, t)dt + \sigma\sqrt{2}dW_t^t, X_T \sim p_{\text{target}}, \Rightarrow \overleftarrow{P}_\theta(X) \\ \text{perform importance sampling} \end{array} \right.$

Small variance

$$\text{Importance weight: } \frac{d\vec{Q}(X)}{d\overleftarrow{P}_\theta(X)}$$

Diffusion Neural samplers

3. Annealed variance reduction sampler

This includes

- (1) AIS (Annealed Importance Sampling)
- (2) MCD (Monte Carlo Diffusion)
- (3) LDVI (Langevin Diffusion Variational Inference)

...

Diffusion Neural samplers

1. Time-reversal sampler
2. Escorted transport sampler
3. Annealed variance reduction sampler

Diffusion Neural samplers

1. Time-reversal sampler
2. Escorted transport sampler
3. Annealed variance reduction sampler

These are design choices for the sampling processes

Diffusion Neural samplers

1. Time-reversal sampler
2. Escorted transport sampler
3. Annealed variance reduction sampler

These are design choices for the sampling processes
But **how to train them?**

Diffusion Neural samplers

1. Time-reversal sampler
2. Escorted transport sampler
3. Annealed variance reduction sampler

These are design choices for the sampling processes
But **how to train them?**

Diffusion Neural samplers

a. path-measure alignment

Diffusion Neural samplers

a. path-measure alignment

For any desired process $\vec{Q}(X)$ or $\vec{Q}_\theta(X)$

Diffusion Neural samplers

a. path-measure alignment

For any desired process $\vec{Q}(X)$ or $\vec{Q}_\theta(X)$

we can write down its desired time-reversal $\overleftarrow{P}(X)$ or $\overleftarrow{P}_\theta(X)$,

Diffusion Neural samplers

a. path-measure alignment

For any desired process $\vec{\mathbf{Q}}(X)$ or $\vec{\mathbf{Q}}_{\theta}(X)$

we can write down its desired time-reversal $\overleftarrow{\mathbf{P}}(X)$ or $\overleftarrow{\mathbf{P}}_{\theta}(X)$,

) align

Diffusion Neural samplers

a. path-measure alignment

$$D_{\text{KL}}[\vec{\mathbf{Q}}|\overleftarrow{\mathbf{P}}] = E_{\vec{\mathbf{Q}}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}(X)} \right]$$

$$D_{\text{LV}}[\vec{\mathbf{Q}}|\overleftarrow{\mathbf{P}}] = \text{Var}_{\vec{\pi}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}(X)} \right]$$

$$D_{\text{TB}}[\vec{\mathbf{Q}}|\overleftarrow{\mathbf{P}}] = E_{\vec{\pi}} \left[\left(\log \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}(X)} - k \right)^2 \right]$$

Other choices exist, including sub-TB, DB, etc...

Diffusion Neural samplers

b. marginal alignment

Diffusion Neural samplers

b. marginal alignment

For any desired process $\vec{Q}(X)$ or $\vec{Q}_\theta(X)$

Diffusion Neural samplers

b. marginal alignment


For any desired process $\vec{Q}(X)$ or $\vec{Q}_\theta(X)$

we can write down its desired marginal $q_t(X_t)$,

Diffusion Neural samplers

b. marginal alignment

For any desired process $\vec{Q}(X)$ or $\vec{Q}_\theta(X)$
we can write down its desired marginal $q_t(X_t)$,



align

Diffusion Neural samplers

b. marginal alignment

Score matching with Score estimator

PINN

Action matching

...

Diffusion Neural samplers

	Time-reversal sampler	Escorted transport sampler	Annealed Variance Reduction Sampler
Path measure alignment	DDS, DIS, PIS, GFN	CMCD, SLCD	MCD
Marginal alignment	iDEM, RDMC, PINN-sampler	NETS, PINN-sampler, LFIS	

Diffusion Neural samplers - Desiderata

Let's look at the loss again, for example:

Diffusion Neural samplers - Desiderata

Let's look at the loss again, for example:

$$D_{\text{KL}}[\vec{\mathbf{Q}}|\overleftarrow{\mathbf{P}}] = E_{\vec{\mathbf{Q}}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}(X)} \right]$$

$$D_{\text{LV}}[\vec{\mathbf{Q}}|\overleftarrow{\mathbf{P}}] = \text{Var}_{\vec{\pi}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}(X)} \right]$$

$$D_{\text{TB}}[\vec{\mathbf{Q}}|\overleftarrow{\mathbf{P}}] = E_{\vec{\pi}} \left[\left(\log \frac{d\vec{\mathbf{Q}}(X)}{d\overleftarrow{\mathbf{P}}(X)} - k \right)^2 \right]$$

Diffusion Neural samplers - Desiderata

Let's look at the loss again, for example:

$$D_{\text{KL}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \mathbb{E}_{\vec{\mathbf{Q}}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} \right]$$

$$D_{\text{LV}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \text{Var}_{\vec{\pi}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} \right]$$

$$D_{\text{TB}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \mathbb{E}_{\vec{\pi}} \left[\left(\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} - k \right)^2 \right]$$

Diffusion Neural samplers - Desiderata

Let's look at the loss again, for example:

$$D_{\text{KL}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \mathbb{E}_{\vec{\mathbf{Q}}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} \right]$$

$$D_{\text{LV}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \text{Var}_{\vec{\pi}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} \right]$$

$$D_{\text{TB}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \mathbb{E}_{\vec{\pi}} \left[\left(\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} - k \right)^2 \right]$$

🙄 need to simulate the trajectory – expensive!

Diffusion Neural samplers - Desiderata

Let's look at the loss again, for example:

$$D_{\text{KL}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \mathbb{E}_{\vec{\mathbf{Q}}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} \right]$$

$$D_{\text{LV}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \text{Var}_{\vec{\pi}} \left[\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} \right]$$

$$D_{\text{TB}}[\vec{\mathbf{Q}}|\vec{\mathbf{P}}] = \mathbb{E}_{\vec{\pi}} \left[\left(\log \frac{d\vec{\mathbf{Q}}(X)}{d\vec{\mathbf{P}}(X)} - k \right)^2 \right]$$

🙄 need to simulate the trajectory – expensive!

Any ways for “simulation-free” training?

Simulation-free training of Diffusion Neural samplers

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$

The second way of sampling

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$

The second way of sampling $X_0 = F_\theta(Z, 0), \quad Z \sim p_{\text{base}}$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), Z \sim p_{\text{base}}$

The second way of sampling $X_0 = F_\theta(Z, 0), Z \sim p_{\text{base}}$
 $dX_t = \partial_t F_\theta(Z, t) dt$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$

The second way of sampling $X_0 = F_\theta(Z, 0), Z \sim p_{\text{base}}$
 $dX_t = \partial_t F_\theta(\underbrace{Z, t}) dt$
 $Z = F_\theta^{-1}(X_t, t)$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), Z \sim p_{\text{base}}$

The second way of sampling $X_0 = F_\theta(Z, 0), Z \sim p_{\text{base}}$
 $dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$

The second way of sampling $X_0 = F_\theta(Z, 0), Z \sim p_{\text{base}}$
$$dX_t = \underbrace{\partial_t F_\theta(F_\theta^{-1}(X_t, t), t)}_{\text{Standard form of ODE}} dt$$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), Z \sim p_{\text{base}}$

The second way of sampling $X_0 = F_\theta(Z, 0), Z \sim p_{\text{base}}$
 $dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt$
 $dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t \sqrt{2}dW_t$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$

The second way of sampling $X_0 = F_\theta(Z, 0), Z \sim p_{\text{base}}$
 $dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt$
 $dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling $X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$

The second way of sampling $X_0 = F_\theta(Z, 0), Z \sim p_{\text{base}}$
 $dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt$ Easily obtained by NF

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t$$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling

$$X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$$

directly sample from time t

The second way of sampling

$$X_0 = F_\theta(Z, 0), \quad Z \sim p_{\text{base}}$$

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t) dt$$

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t) dt + \sigma_t^2 \nabla \log q_\theta(X_t, t) dt + \sigma_t \sqrt{2} dW_t$$

Simulation-free training of Diffusion Neural samplers

avoid simulating the trajectory (entirely) during training.



using a time-dependent normalizing flow

Define $F_\theta(\cdot, t)$ as an invertible function

The first way of sampling

$$X_t = F_\theta(Z, t), \quad Z \sim p_{\text{base}}$$

directly sample from time t

The second way of sampling

$$X_0 = F_\theta(Z, 0), \quad Z \sim p_{\text{base}}$$
$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t) dt$$

**Calculate the same loss
as other diffusion samplers**

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t) dt + \sigma_t^2 \nabla \log q_\theta(X_t, t) dt + \sigma_t \sqrt{2} dW_t$$

Simulation-free training of Diffusion Neural samplers

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t, X_0 \sim p_{\text{prior}}$$

Simulation-free training of Diffusion Neural samplers

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t, X_0 \sim p_{\text{prior}}$$

Align

$$dX_t = \underbrace{g(X_t)}dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim p_{\text{target}}$$

a simple function, e.g., 0

Simulation-free training of Diffusion Neural samplers

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t, X_0 \sim p_{\text{prior}}$$



time-reversal

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt - \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim q_\theta(\cdot, T)$$

$$dX_t = \underbrace{g(X_t)}dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim p_{\text{target}}$$

a simple function, e.g., 0

Align

Simulation-free training of Diffusion Neural samplers

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t, X_0 \sim p_{\text{prior}}$$



time-reversal

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt - \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim q_\theta(\cdot, T)$$

Align

$$dX_t = \underbrace{g(X_t)}dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim p_{\text{target}}$$

a simple function, e.g., 0

Simulation-free training of Diffusion Neural samplers

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t, X_0 \sim p_{\text{prior}}$$



time-reversal

Align

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt - \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim q_\theta(\cdot, T)$$

$$dX_t = \underbrace{g(X_t)}dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim p_{\text{target}}$$

a simple function, e.g., 0

✓ same direction – Girsanov Theorem applicable

Simulation-free training of Diffusion Neural samplers

$$dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt + \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t, X_0 \sim p_{\text{prior}}$$



time-reversal

Align

$$\left(\begin{array}{l} dX_t = \partial_t F_\theta(F_\theta^{-1}(X_t, t), t)dt - \sigma_t^2 \nabla \log q_\theta(X_t, t)dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim q_\theta(\cdot, T) \\ \\ dX_t = \underbrace{g(X_t)}dt + \sigma_t \sqrt{2}dW_t^-, X_T \sim p_{\text{target}} \end{array} \right.$$

a simple function, e.g., 0

✓ same direction – Girsanov Theorem applicable

✓ simulation-free evaluation – can always obtain sample by 1-step $X_t = F_\theta(Z, t), Z \sim p_{\text{base}}$

Simulation-free training of Diffusion Neural samplers

 Great! How does it perform?

Simulation-free training of Diffusion Neural samplers

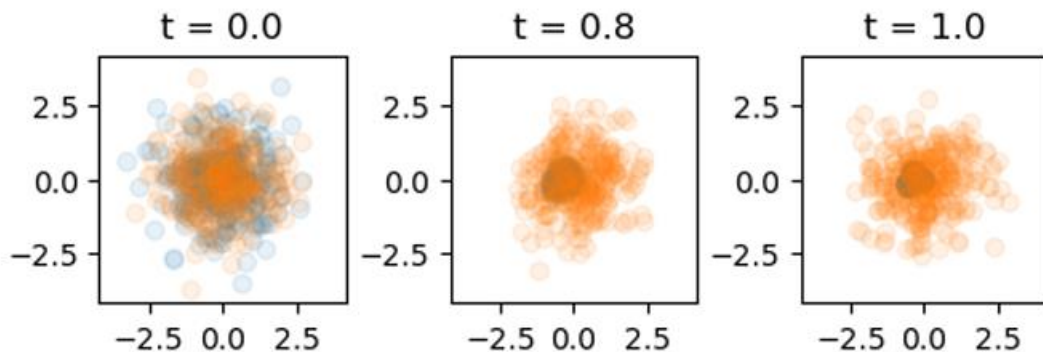
 Great! How does it perform?

 unfortunately...

Simulation-free training of Diffusion Neural samplers

😊 Great! How does it perform?

😭 unfortunately...

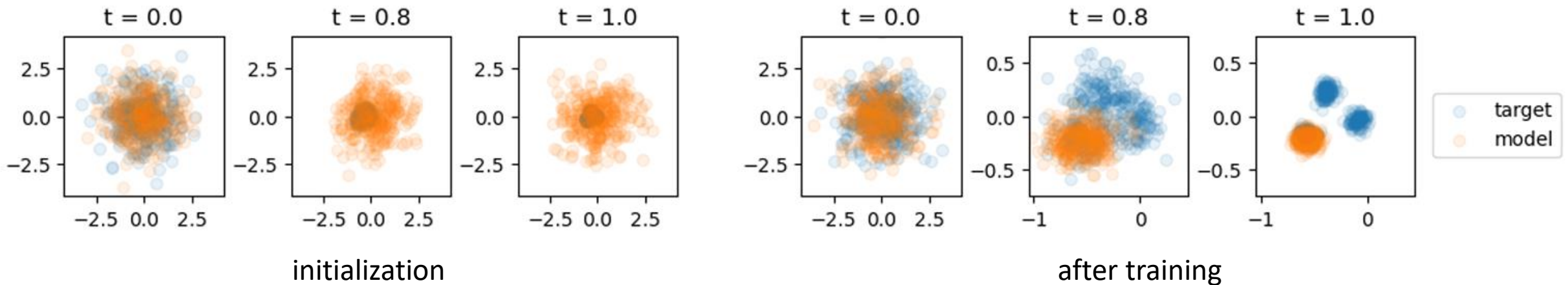


initialization

Simulation-free training of Diffusion Neural samplers

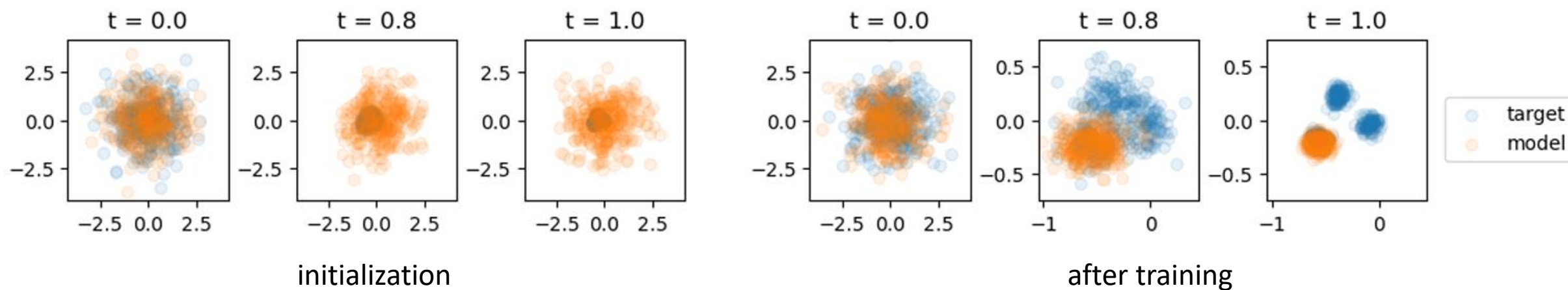
😊 Great! How does it perform?

😞 unfortunately...



Simulation-free training of Diffusion Neural samplers

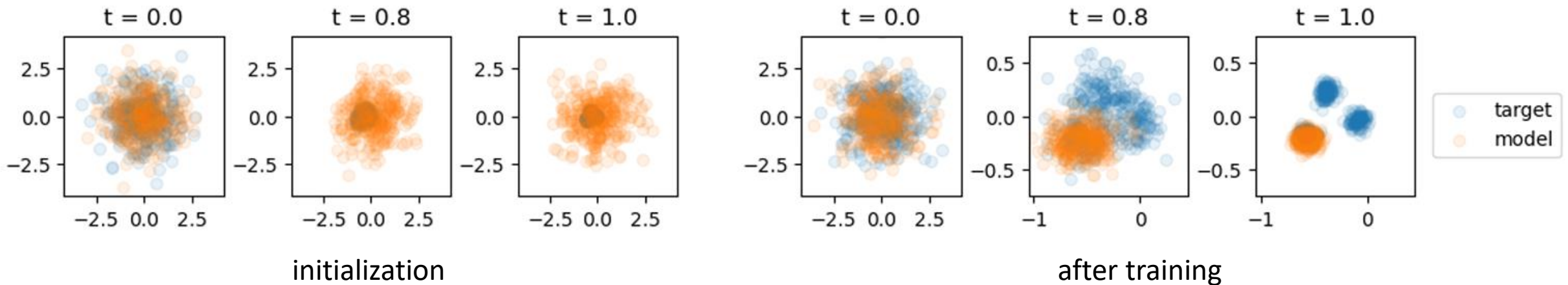
Why?



Simulation-free training of Diffusion Neural samplers

Why?

Objective? 🤔 same as DDS

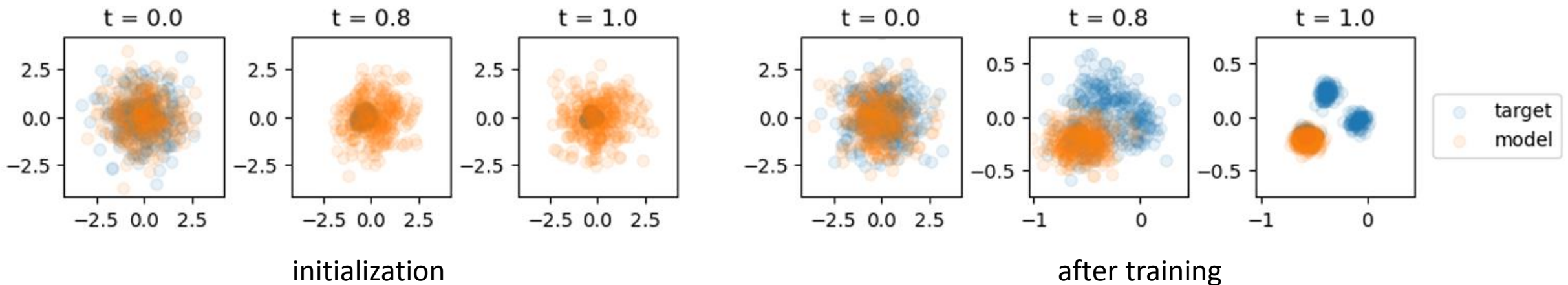


Simulation-free training of Diffusion Neural samplers

Why?

Objective? 🤔 same as DDS

Capacity? 🤔 target is so simple



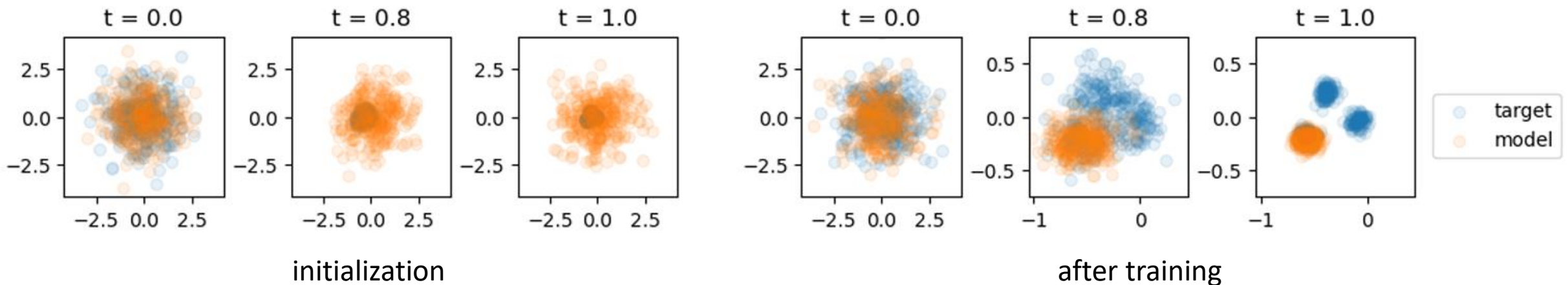
Simulation-free training of Diffusion Neural samplers

Why?

Objective? 🤔 same as DDS

Capacity? 🤔 target is so simple

Network parameterization? 🤔 might be the reason



Langevin Preconditioning

Langevin Preconditioning

a. DDS/PIS

warm-up initialization

$$f_{\theta}(\cdot, t) = \underbrace{\text{NN}_{1,\theta}(\cdot, t)}_{\approx 0} + \text{NN}_{2,\theta}(t) \circ \nabla \log p_{\text{target}}(\cdot)$$

Langevin Preconditioning

a. DDS/PIS

warm-up initialization

$$f_{\theta}(\cdot, t) = \text{NN}_{1,\theta}(\cdot, t) + \text{NN}_{2,\theta}(t) \circ \nabla \log p_{\text{target}}(\cdot)$$

≈ 0

Langevin gradient

b. CMCD/NETS

$$\pi_t(\cdot) = p_{\text{prior}}^{1-\beta}(\cdot) p_{\text{target}}^{\beta}(\cdot) \quad \text{Langevin gradient}$$

$$dX_t = (f_{\theta}(X_t, t) + \sigma_t^2 \nabla \log \pi_t(X_t)) dt + \sqrt{2} \sigma_t dW_t \quad \overrightarrow{\mathbf{Q}}_{\theta}(X)$$

$$dX_t = (f_{\theta}(X_t, t) - \sigma_t^2 \nabla \log \pi_t(X_t)) dt + \sqrt{2} \sigma_t dW_t^- \quad \overleftarrow{\mathbf{P}}_{\theta}(X)$$

Langevin Preconditioning

a. DDS/PIS

warm-up initialization

$$f_{\theta}(\cdot, t) = \text{NN}_{1,\theta}(\cdot, t) + \text{NN}_{2,\theta}(t) \circ \nabla \log p_{\text{target}}(\cdot)$$

≈ 0

Langevin gradient

What if we remove this Langevin?

b. CMCD/NETS

$$\pi_t(\cdot) = p_{\text{prior}}^{1-\beta}(\cdot) p_{\text{target}}^{\beta}(\cdot) \quad \text{Langevin gradient}$$

$$dX_t = (f_{\theta}(X_t, t) + \sigma_t^2 \nabla \log \pi_t(X_t)) dt + \sqrt{2} \sigma_t dW_t \quad \overrightarrow{\mathbf{Q}}_{\theta}(X)$$

$$dX_t = (f_{\theta}(X_t, t) - \sigma_t^2 \nabla \log \pi_t(X_t)) dt + \sqrt{2} \sigma_t dW_t^- \quad \overleftarrow{\mathbf{P}}_{\theta}(X)$$

How to remove Langevin?

a. DDS

$$f_{\theta}(\cdot, t) = \text{NN}_{1,\theta}(\cdot, t) + \text{NN}_{2,\theta}(t) \circ \nabla \log p_{\text{target}}(\cdot)$$

b. CMCD's Optimality condition (Nelson's relation)

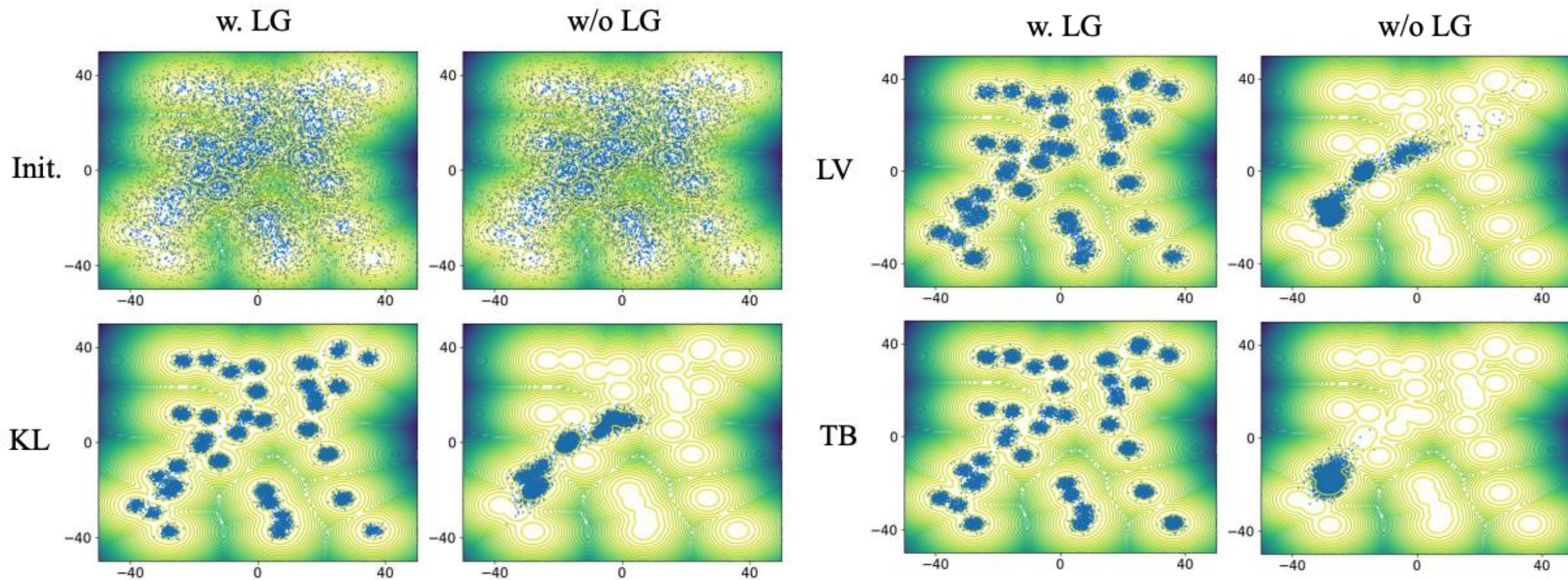
$$D(\overrightarrow{Q_{\theta}^{p_{\text{prior}}, f_{\theta} + \sigma_t^2 \nabla \log \pi_t}}, \overleftarrow{P_{\theta}^{p_{\text{target}}, f_{\theta} - \sigma_t^2 \nabla \log \pi_t}})$$

But it is not necessary!

$$D(\overrightarrow{Q_{\theta}^{p_{\text{prior}}, f_{\theta}}}, \overleftarrow{P_{\theta}^{p_{\text{target}}, f_{\theta} - 2\sigma_t^2 \nabla \log \pi_t}})$$

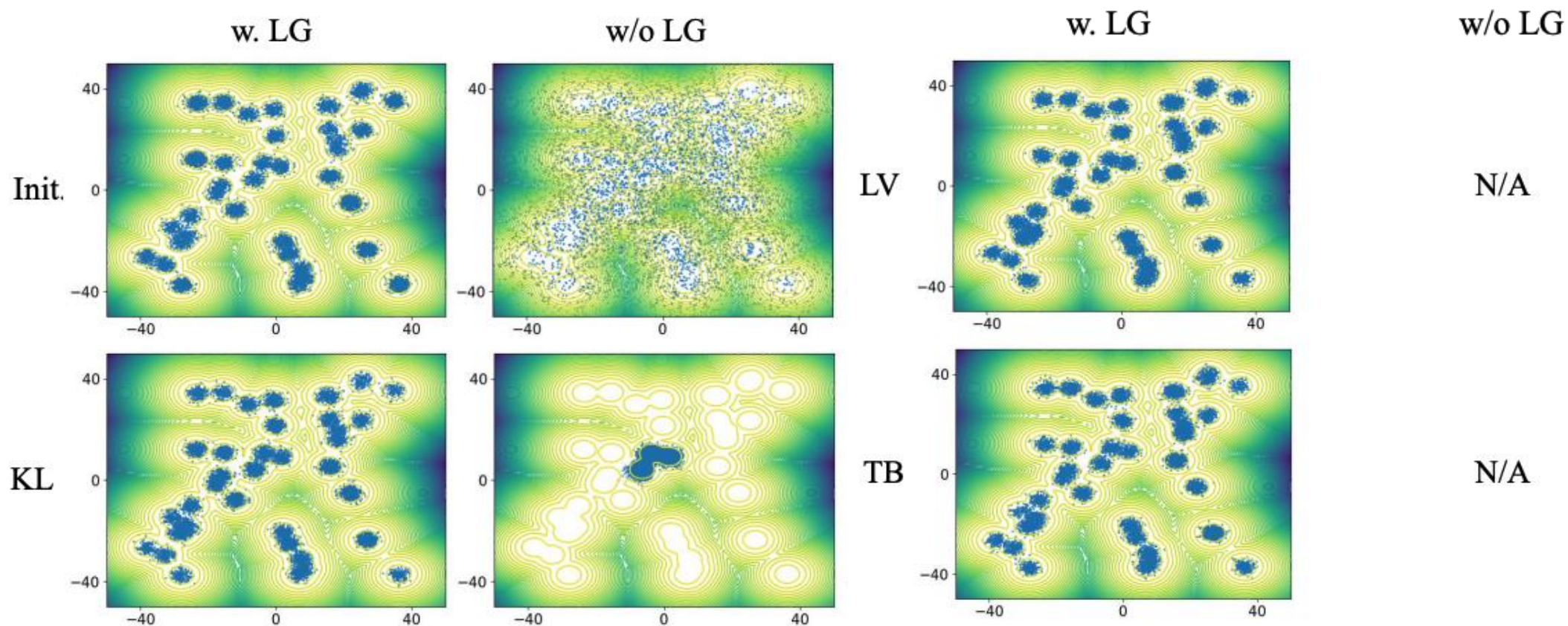
Empirical Results

a. Langevin precondition is necessary to prevent mode collapse



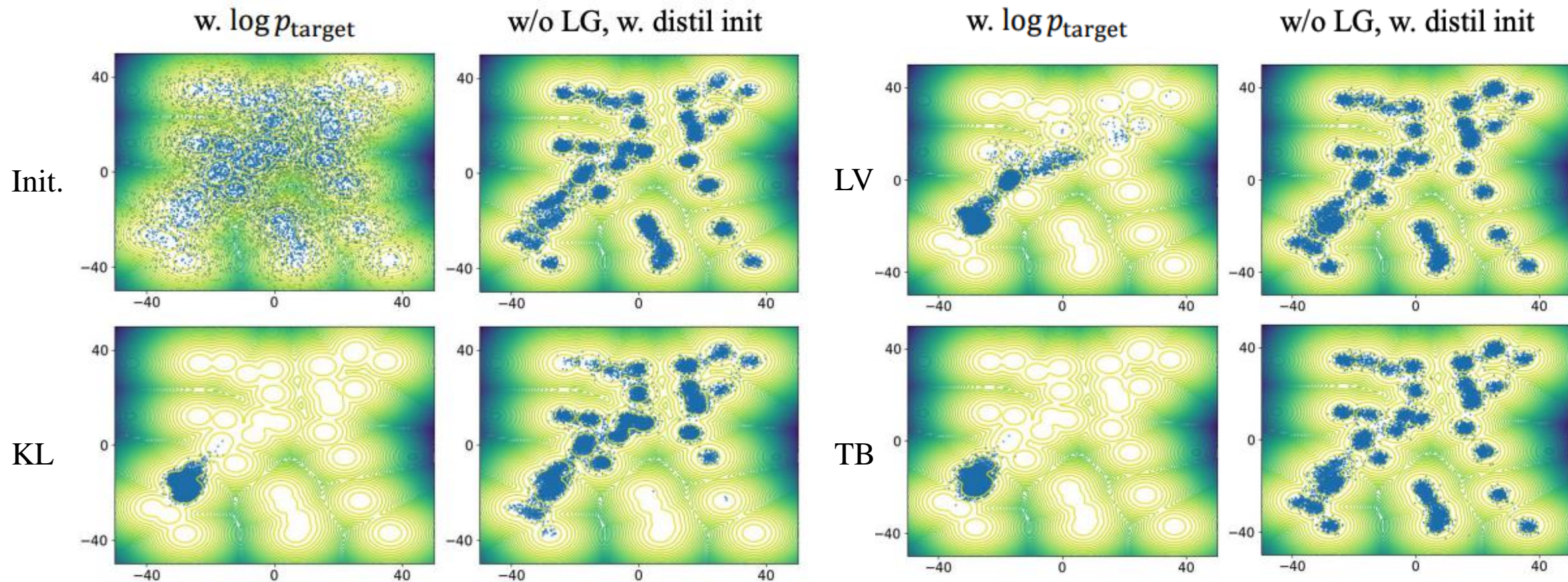
Empirical Results

a. Langevin precondition is necessary to prevent mode collapse



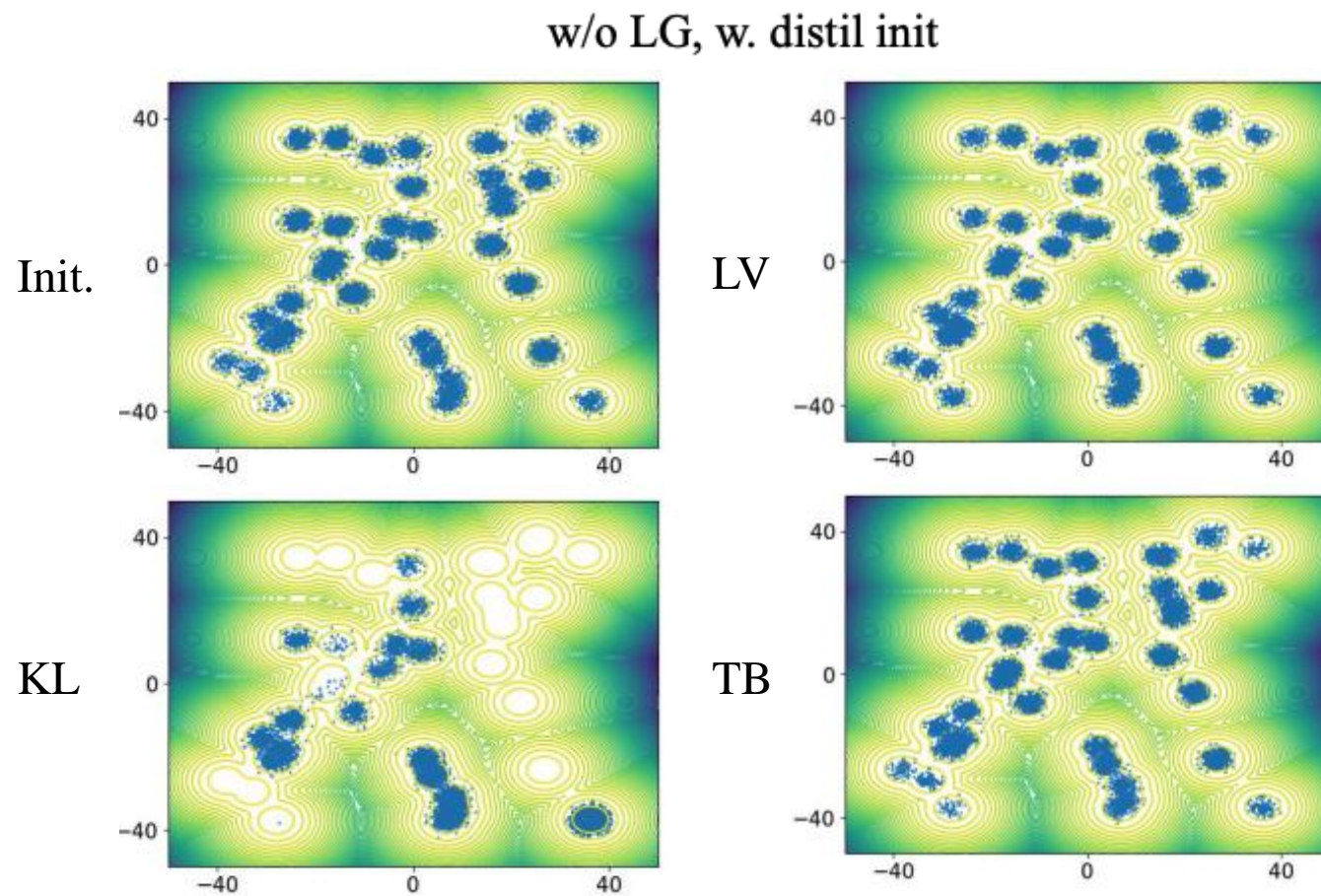
Empirical Results

b. Does other ways of incorporating the target information help?



Empirical Results

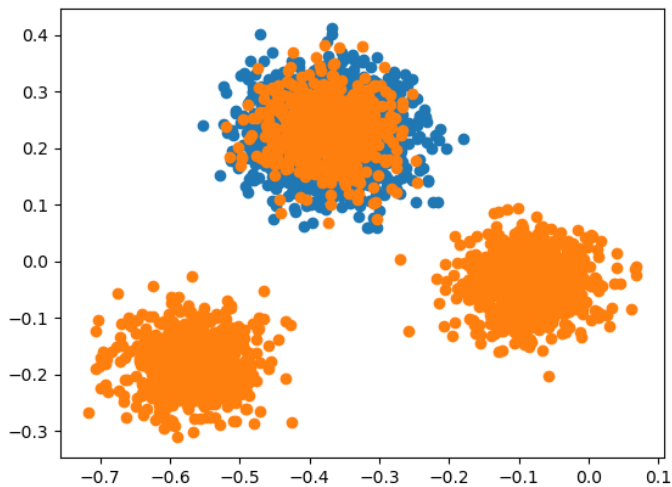
b. Do other ways of incorporating the target information help?



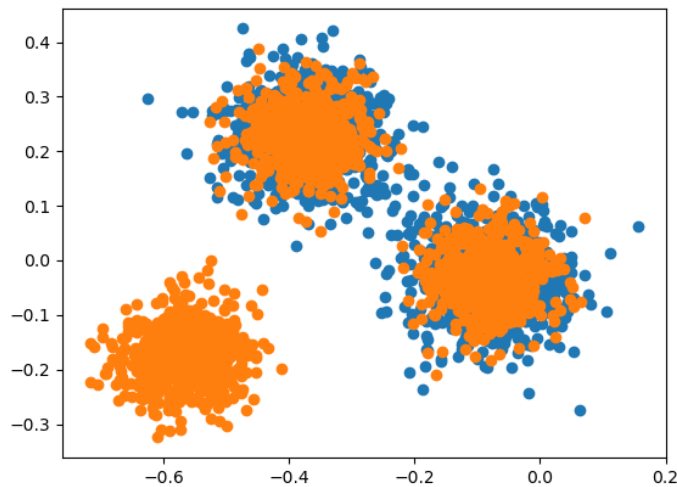
Empirical Results

DDS w/o Langevin for GMM-3:

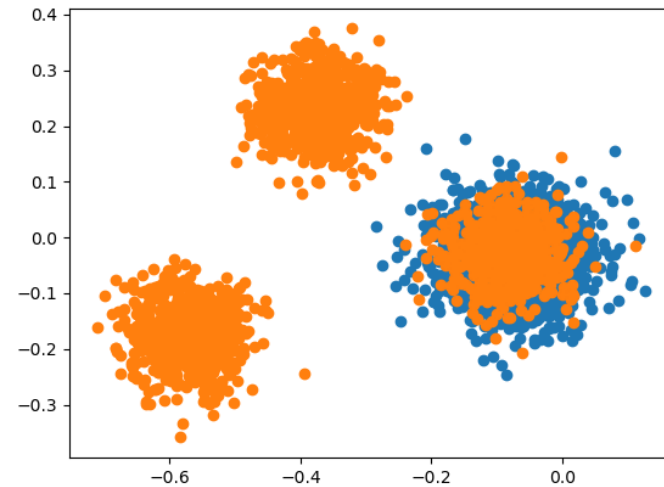
KL



LV



TB



Empirical Results

c. How about sample efficiency?

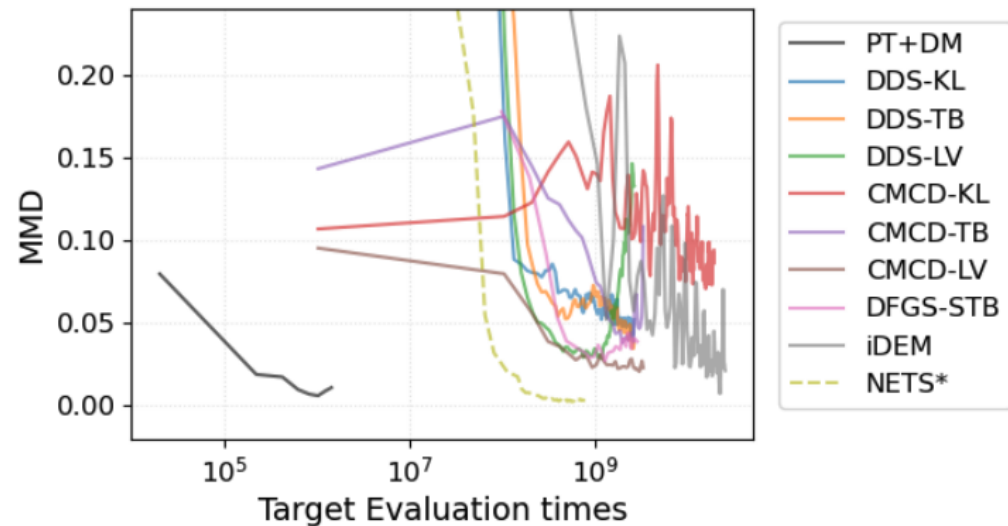
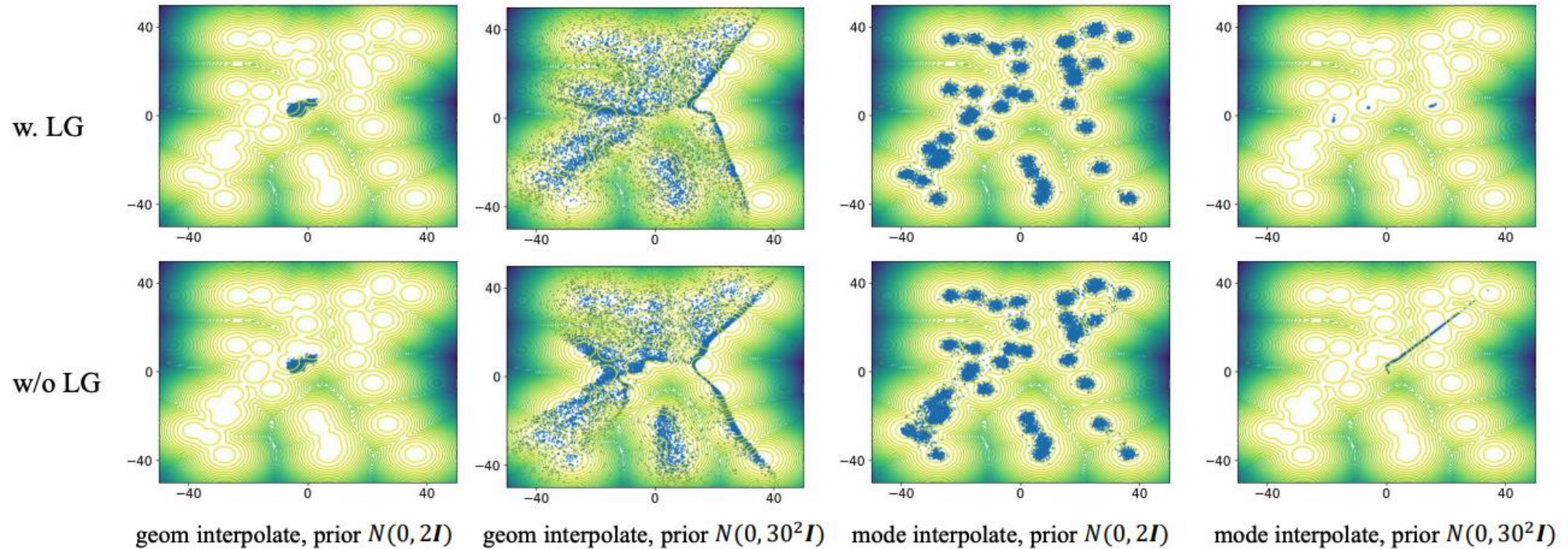


Figure 2: Sample quality vs target evaluation times for different approaches with different objectives on GMM-40 target. *NETS uses mode interpolation, which is distinct from that employed in others.

Empirical Results

d. PINN objective is different

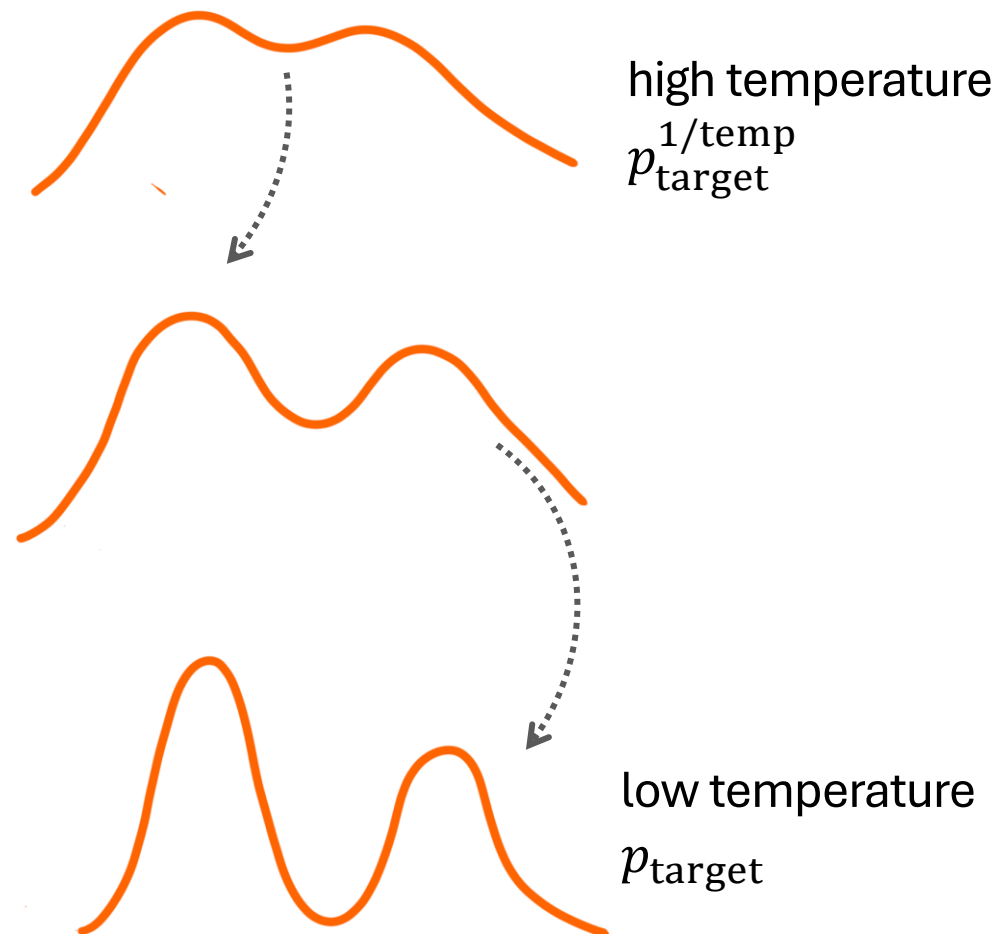
1. different interpolant
2. different prior size
3. “consistent” behavior



Parallel Tempering/Replica Exchange

😊 SOTA MCMC in MD simulation

😊 Highly parallel

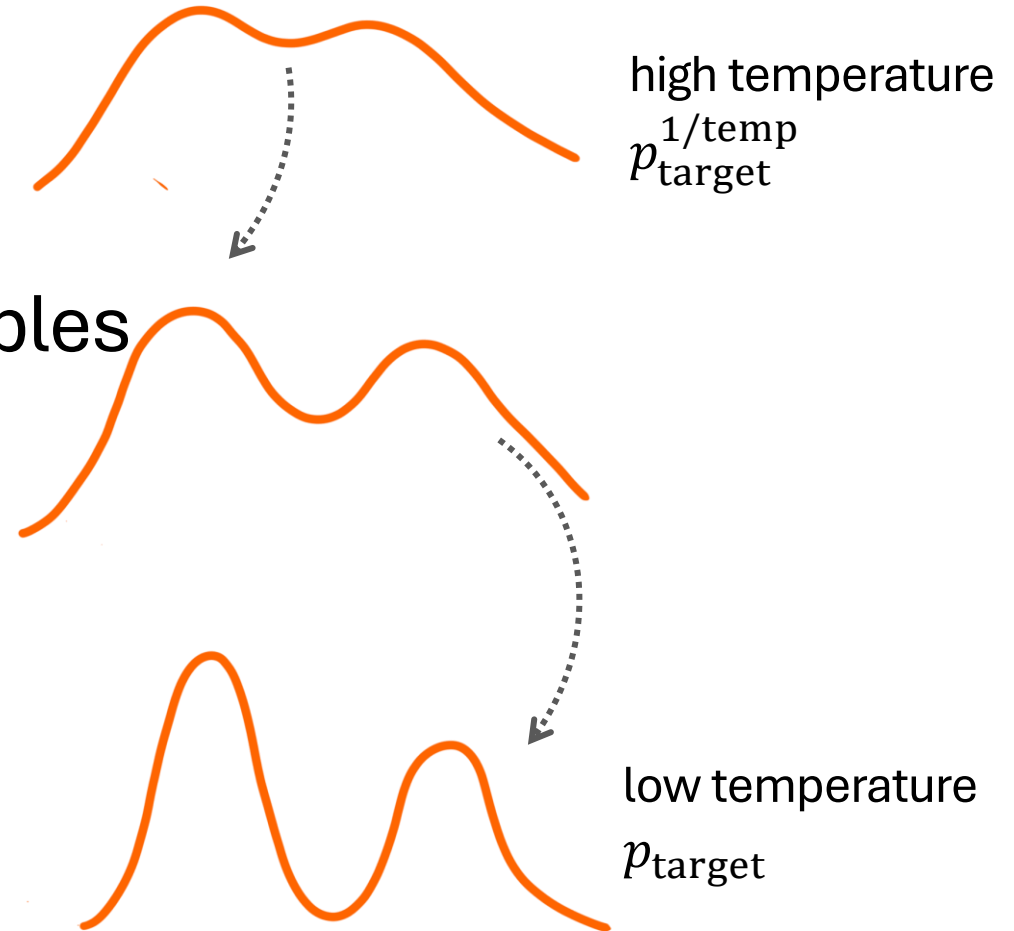


Parallel Tempering/Replica Exchange

Different use from neural samplers...

😞 Correlated samples

😞 Need more simulation for new samples



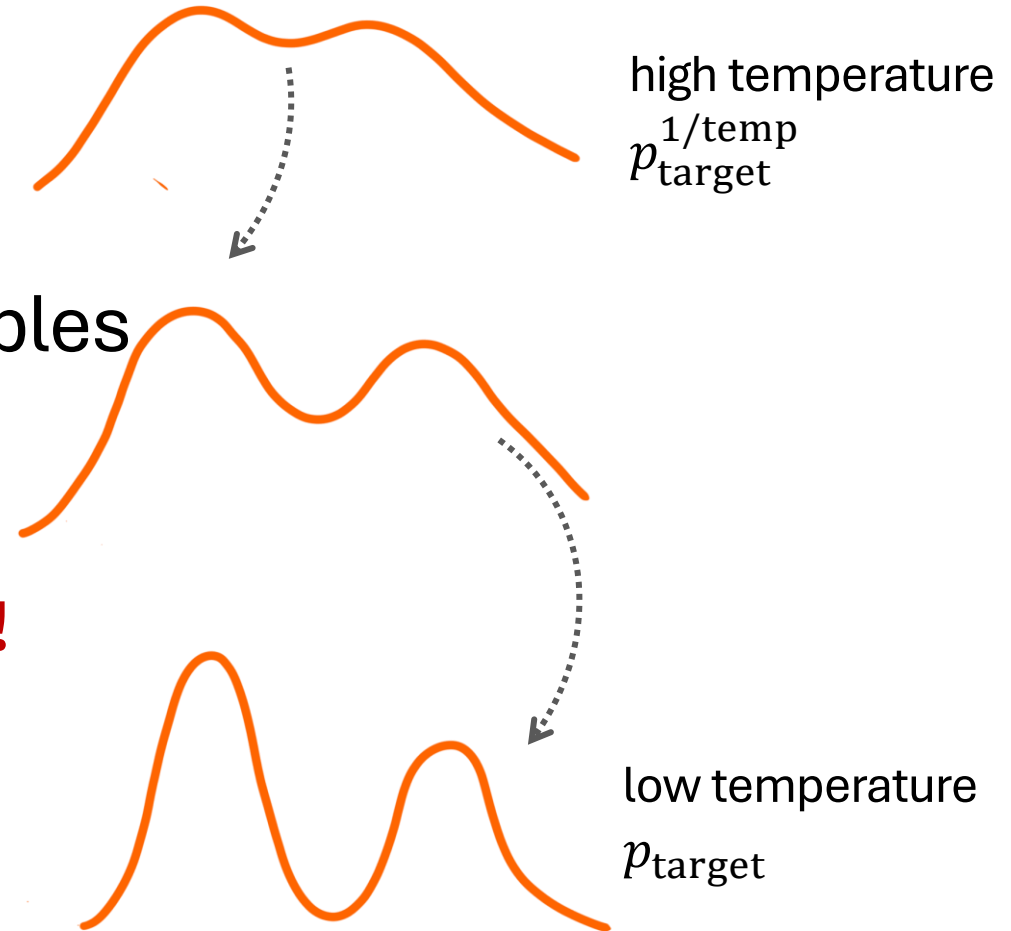
Parallel Tempering/Replica Exchange

Different use from neural samplers...

😞 Correlated samples

😞 Need more simulation for new samples

Generative models can easily address them!
But is it worth it?



Parallel Tempering/Replica Exchange

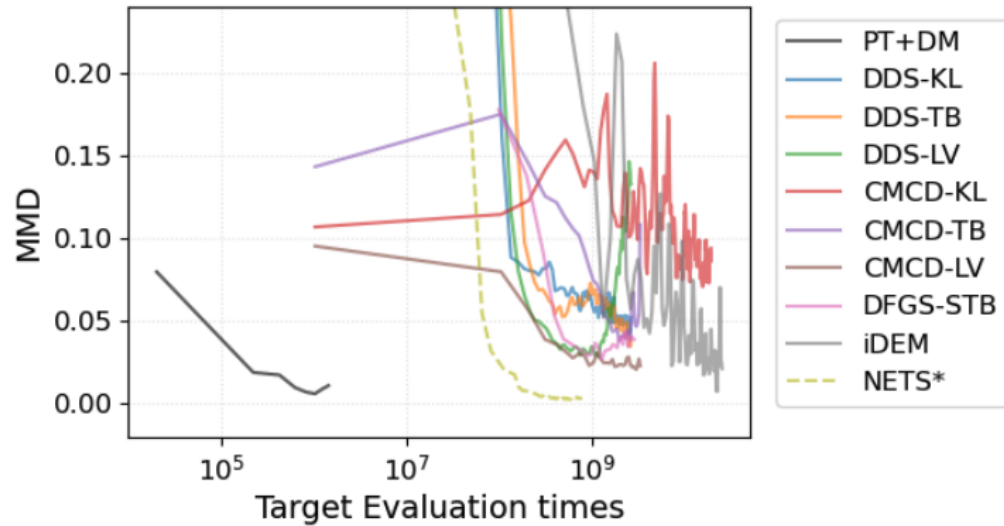
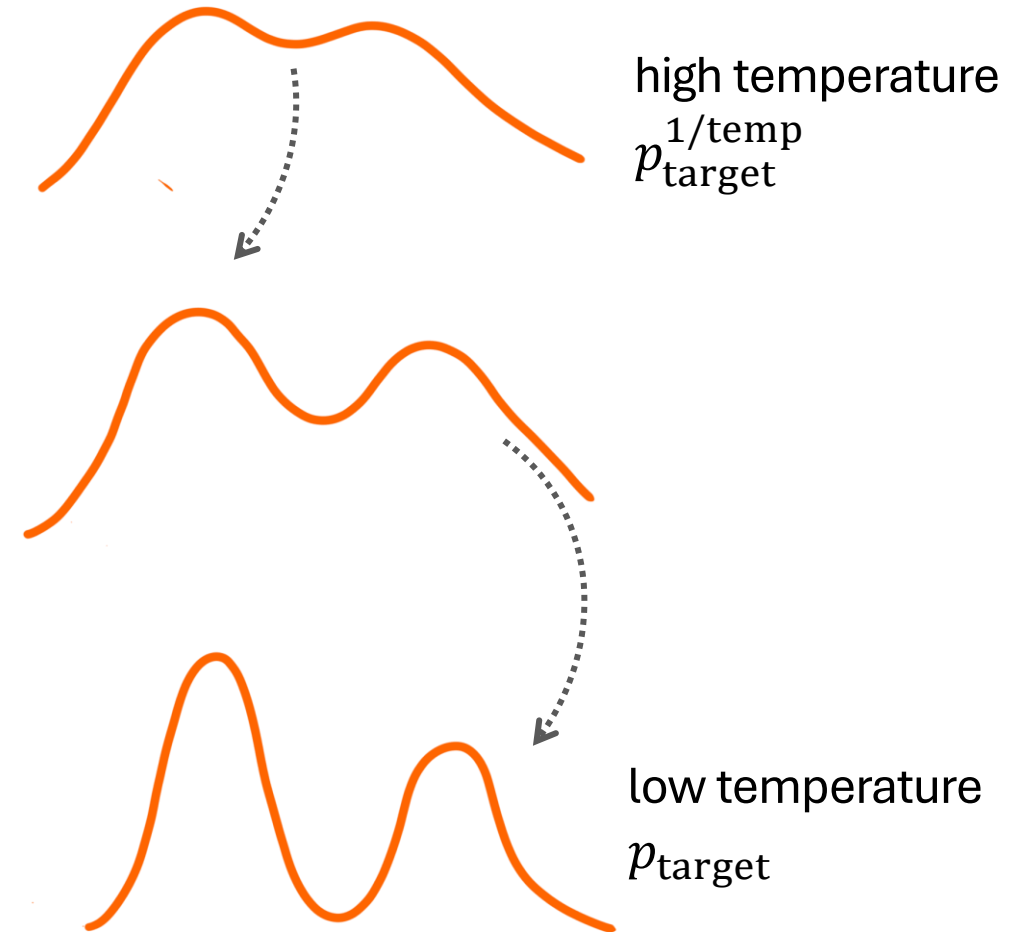


Figure 2: Sample quality vs target evaluation times for different approaches with different objectives on GMM-40 target. *NETS uses mode interpolation, which is distinct from that employed in others.



Discussion & Takeaway

1. We should not hide **Langevin** gradient used

Discussion & Takeaway

1. We **should not hide Langevin** gradient used
2. If we need Langevin gradient anyway, we need to **talk about sample efficiency** (we should also be open to initialize using data)

Discussion & Takeaway

1. We **should not hide Langevin** gradient used
2. If we need Langevin gradient anyway, we need to **talk about sample efficiency** (we should also be open to initialize using data)
3. Improving **PT** is a promising direction (solve challenges with neural network ansatz)

Discussion & Takeaway

1. We **should not hide Langevin** gradient used
2. If we need Langevin gradient anyway, we need to **talk about sample efficiency** (we should also be open to initialize using data)
3. Improving **PT** is a promising direction (solve challenges with neural network ansatz)
4. **Better prior, interpolant, explorative objectives** still needed