

# Statistical guarantees for denoising reflected diffusion models

MFO Mini-Workshop on Statistical Challenges for Deep Generative Models

---

Claudia Strauch and [Lukas Trottner](#)

joint work with [Asbjørn Holk Thomsen](#)

17 February 2025

Heidelberg University

[University of Birmingham](#)

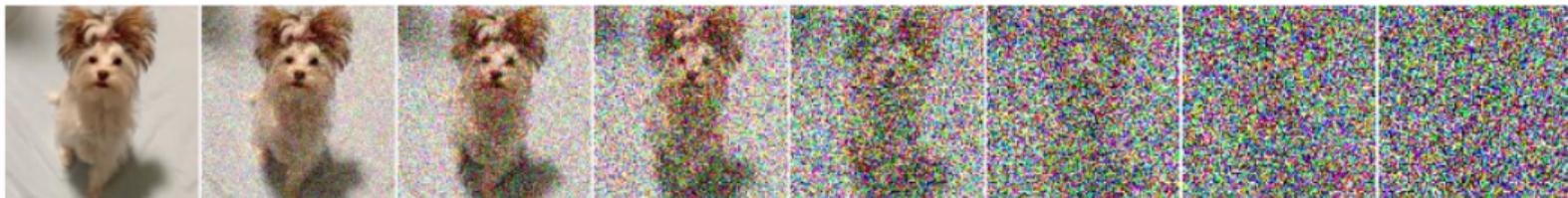
[Aarhus University](#)



UNIVERSITY OF  
BIRMINGHAM

## Motivation:

*“Creating noise from data is easy; creating data from noise is **generative modeling**.”*



Source: Song et al. (2021). Score based generative modeling through stochastic differential equations. *ICLR*.

# Generative modelling

- ↪ involves **learning the underlying distribution** of a dataset to **generate new, similar data points**
- ↪ aims to model the data distribution  $p(x)$  for observed data  $x$ , allowing the **generation of new samples**  $x'$  that resemble the original data
- ↪ essential in applications like **image synthesis**, text generation, and data augmentation

## Core tasks:

1. **Density estimation:** Learning the probability distribution  $p(x)$  or its properties.
2. **Sampling:** Drawing new samples  $x'$  from the learned  $p(x)$ .

**Examples:** Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs) and normalizing flows

# Generative modelling

- ↪ involves **learning the underlying distribution** of a dataset to **generate new, similar data points**
- ↪ aims to model the data distribution  $p(x)$  for observed data  $x$ , allowing the **generation of new samples**  $x'$  that resemble the original data
- ↪ essential in applications like **image synthesis**, text generation, and data augmentation

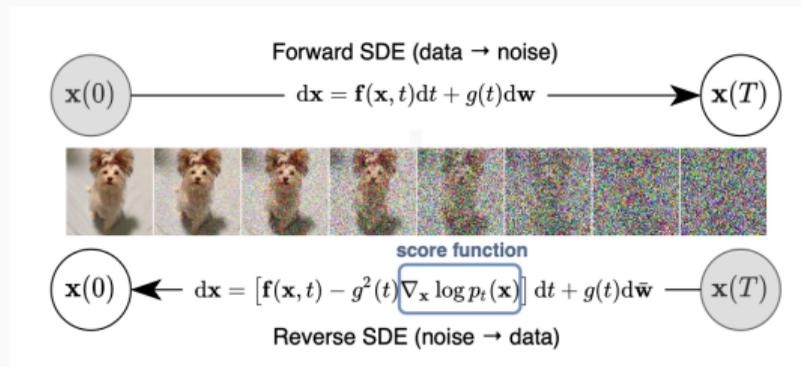
## Core tasks:

1. **Density estimation:** Learning the probability distribution  $p(x)$  or its properties.
2. **Sampling:** Drawing new samples  $x'$  from the learned  $p(x)$ .

**Examples:** Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs) and normalizing flows

## Denoising diffusion models (DDMs)

- provide an **iterative generative algorithm** to create new samples that approximately match the target distribution  $p_0$ , given a finite number of samples corresponding to an unknown  $p_0$
- **general idea**: find a **stochastic process** that perturbs  $p_0$  to a new distribution  $p_T$  such that
  - 1)  $p_T$  or a good approximation thereof is **easy to sample from**, and
  - 2) the perturbation is **reversible** in the sense that we know how to **simulate the time-reversed process**



Source: Song et al. (2021). Score based generative modeling through stochastic differential equations. *ICLR*.

## Statistical challenges in generative modelling

### Unifying principle of generative modelling:

transform noise to **create new data** that matches a given training data set

↪ transformations must **adapt to the information contained in the training data**, which is **high-dimensional** in typical machine learning applications

Generative models have demonstrated remarkable **empirical success** across diverse domains, including images, videos, and text, despite their differences in methodology:

- models like **Generative Adversarial Networks (GANs)** aim to **directly approximate** the transformation from noise to data using adversarial training
- **Denosing Diffusion Models (DDMs)** **dynamically evolve noise into data** by approximating the characteristics of a **stochastic process**

*Under what conditions do these models ensure that the generated distribution converges to the target distribution at a (minimax) optimal rate?*

# Statistical challenges in generative modelling

## Unifying principle of generative modelling:

transform noise to **create new data** that matches a given training data set

↪ transformations must **adapt to the information contained in the training data**, which is **high-dimensional** in typical machine learning applications

Generative models have demonstrated remarkable **empirical success** across diverse domains, including images, videos, and text, despite their differences in methodology:

- models like **Generative Adversarial Networks (GANs)** aim to **directly approximate** the transformation from noise to data using adversarial training
- **Denosing Diffusion Models (DDMs)** **dynamically evolve noise into data** by approximating the characteristics of a **stochastic process**

*Under what conditions do these models ensure that the **generated distribution** converges to the **target distribution** at a (minimax) optimal rate?*

## Classical Denoising Diffusion Models (DDMs)

- for some fixed time  $T > 0$  and suitable drift  $b : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  and dispersion  $\sigma : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ , consider the **forward model**

$$dX_t = b(t, X_t) dt + \sigma(t, X_t) dW_t, \quad t \in [0, T], X_0 \sim p_0,$$

$W = (W_t)_{t \in [0, T]}$  some standard  $d$ -dimensional Brownian motion

- under sufficient regularity conditions, the forward model has a solution  $X = (X_t)_{t \in [0, T]}$  with marginal densities  $(p_t)_{t \in [0, T]}$  such that the **time-reversed process**  $\tilde{X}_t = X_{T-t}$ ,  $t \in [0, T]$ , solves

$$d\tilde{X}_t = -\bar{b}(T-t, \tilde{X}_t) dt + \sigma(T-t, \tilde{X}_t) d\bar{W}_t, \quad t \in [0, T], \tilde{X}_0 \sim p_T,$$

for some Brownian motion  $(\bar{W}_t)_{t \in [0, T]}$  and drift  $\bar{b} : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  given by

$$\bar{b}_i(t, x) = b_i(t, x) - \frac{1}{p_t(x)} \sum_{j,k=1}^d \frac{\partial}{\partial x_j} [p_t(x) \sigma_{ik}(t, x) \sigma_{jk}(t, x)], \quad i = 1, \dots, d$$

- ↪ time-reversed process solves a **time-inhomogeneous SDE**, now with drift  $-\bar{b}(T - \cdot, \cdot)$  and dispersion  $\sigma(T - \cdot, \cdot)$

## Classical Denoising Diffusion Models (DDMs)

- for some fixed time  $T > 0$  and suitable drift  $b : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  and dispersion  $\sigma : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ , consider the **forward model**

$$dX_t = b(t, X_t) dt + \sigma(t, X_t) dW_t, \quad t \in [0, T], X_0 \sim p_0,$$

$W = (W_t)_{t \in [0, T]}$  some standard  $d$ -dimensional Brownian motion

- under sufficient regularity conditions, the forward model has a solution  $X = (X_t)_{t \in [0, T]}$  with marginal densities  $(p_t)_{t \in [0, T]}$  such that the **time-reversed process**  $\tilde{X}_t = X_{T-t}$ ,  $t \in [0, T]$ , solves

$$d\tilde{X}_t = -\bar{b}(T-t, \tilde{X}_t) dt + \sigma(T-t, \tilde{X}_t) d\bar{W}_t, \quad t \in [0, T], \tilde{X}_0 \sim p_T,$$

for some Brownian motion  $(\bar{W}_t)_{t \in [0, T]}$  and drift  $\bar{b} : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  given by

$$\bar{b}_i(t, x) = b_i(t, x) - \frac{1}{p_t(x)} \sum_{j, k=1}^d \frac{\partial}{\partial x_j} [p_t(x) \sigma_{ik}(t, x) \sigma_{jk}(t, x)], \quad i = 1, \dots, d$$

- ↪ time-reversed process solves a **time-inhomogeneous SDE**, now with drift  $-\bar{b}(T - \cdot, \cdot)$  and dispersion  $\sigma(T - \cdot, \cdot)$

## Classical DDMs

- **standard convention:** set  $\sigma(t, x) = \gamma(t)\mathbb{I}_d$  for some scalar function  $\gamma$
- ↪ **forward model** is given by a (possibly time-inhomogeneous) **Ornstein–Uhlenbeck process** with explicit transition densities, and the backward drift becomes

$$\bar{b}(t, x) = b(t, x) - \gamma^2(t) \underbrace{\nabla \log p_t(x)}$$

“**score**” of the forward model

- ↪ **backwards process** has the dynamics

$$d\tilde{X}_t = (-b(T-t, \tilde{X}_t) + \gamma^2(T-t) \nabla \log p_{T-t}(\tilde{X}_t)) dt + \gamma(T-t) d\bar{W}_t \quad t \in [0, T], \tilde{X}_0 \sim p_T$$

- as  $t \rightarrow T$ , the density of  $\tilde{X}_t$  approaches  $p_0 \implies$  **simulating the reverse process** generates new data **samples corresponding to the target  $p_0$**
- **note:** we are free to choose the coefficients of our forward process (i.e.,  $b$  and  $\sigma$ ), but the **score function**  $\nabla \log p_t$  depends on  $p_0$  ↪ needs to be estimated from the data (“**score matching**”)

## Classical DDMs

- **standard convention:** set  $\sigma(t, x) = \gamma(t)\mathbb{I}_d$  for some scalar function  $\gamma$
- ↪ **forward model** is given by a (possibly time-inhomogeneous) **Ornstein–Uhlenbeck process** with explicit transition densities, and the backward drift becomes

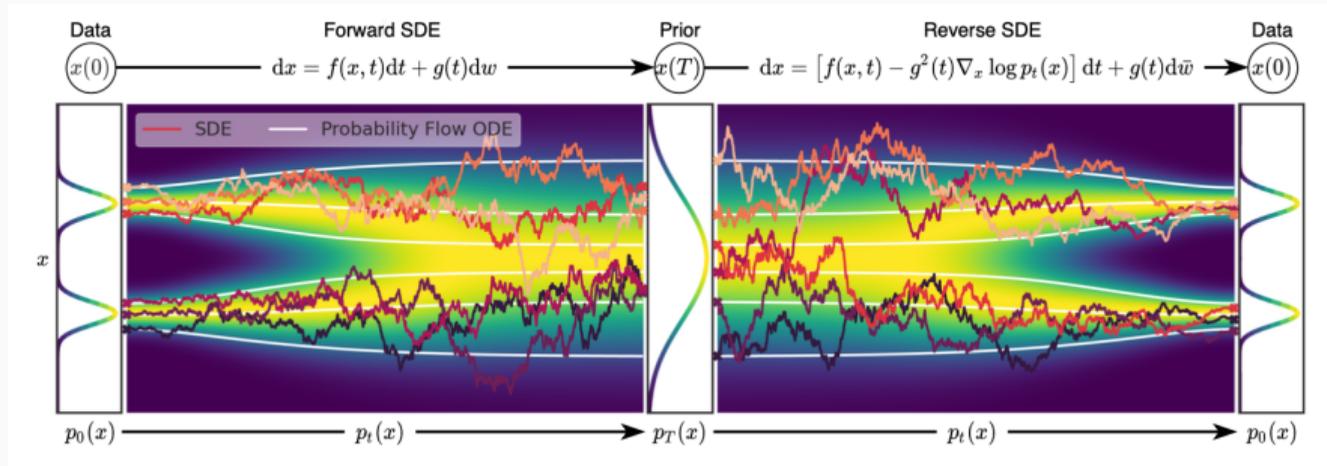
$$\bar{b}(t, x) = b(t, x) - \gamma^2(t) \underbrace{\nabla \log p_t(x)}_{\text{“score” of the forward model}}$$

- ↪ **backwards process** has the dynamics

$$d\tilde{X}_t = \left( -b(T-t, \tilde{X}_t) + \gamma^2(T-t) \nabla \log p_{T-t}(\tilde{X}_t) \right) dt + \gamma(T-t) d\bar{W}_t \quad t \in [0, T], \tilde{X}_0 \sim p_T$$

- as  $t \rightarrow T$ , the density of  $\tilde{X}_t$  approaches  $p_0 \implies$  **simulating the reverse process** generates new data **samples corresponding to the target  $p_0$**
- **note:** we are free to choose the coefficients of our forward process (i.e.,  $b$  and  $\sigma$ ), but the **score function**  $\nabla \log p_t$  depends on  $p_0$  ↪ needs to be estimated from the data (“**score matching**”)

# Statistical aspects of denoising diffusion models



Source: Song et al. (2021). Score based generative modeling through stochastic differential equations. *ICLR*.

## Statistical questions:

1. are diffusion models **minimax learners** (in terms of smoothness assumptions on  $p_0$ )?
2. how can empirical lack of curse of dimensionality be explained?  $\rightsquigarrow$  **submanifold hypothesis**
3. alternative model designs with better theoretical/experimental justification?

# Diffusion models are minimax optimal distribution estimators<sup>1</sup>

*“Is diffusion modeling a **good distribution estimator**? In other words, how can the **estimation error of the generated data distribution** be explicitly bounded by the number of the training data and in a data structure dependent way?”*

Assumptions on the **initial distribution with density**  $p_0$  can be summarised by three key components:

- (i)  $p_0$  is **compactly supported** on a  $d$ -dimensional hypercube;
- (ii)  $p_0$  is **bounded away from zero** on its support;
- (iii)  $p_0$  has **Besov smoothness of order  $s$  away from the support boundary** (where  $s$  is allowed to be sufficiently small to not necessarily imply continuity of  $p_0$ ) and is **infinitely differentiable close to the boundary**.

Under these conditions, [Oko et al. \(2024\)](#) show that generated data distribution achieves the **nearly minimax optimal estimation rate**  $n^{-\frac{s}{2s+d}}(\log n)^8$  in total variation distance.

---

<sup>1</sup>K. Oko, S. Akiyama, and T. Suzuki (2023). **Diffusion Models are Minimax Optimal Distribution Estimators**. *ICML*.

## Convergence of diffusion models under the manifold hypothesis

Convergence rates (even optimal ones) expressed **in terms of the ambient dimension  $d$**  fall short of capturing the empirical success of DDMs

↪ gap is related to the **manifold hypothesis**: real-world high-dimensional data often reside on **lower-dimensional manifolds**, to which well-trained **generative models are believed to adapt**

Tang and Yang (2024)<sup>2</sup> establish (up to log factors) the **minimax convergence rate**  $C(d)n^{-\frac{s+1}{2s+d}}$  in Wasserstein-1 distance for distributions  $p_0$  such that

- (i)  $p_0$  is supported on a **compact and  $\beta$ -smooth  $\tilde{d}$ -dimensional** submanifold  $\mathcal{M}$ , where  $\beta \geq 2$ ;
- (ii)  $p_0$  is **bounded away from zero** on  $\mathcal{M}$ ;
- (iii)  $p_0$  has **smoothness of order**  $s \in [0, \beta - 1]$  w.r.t. the volume measure on  $\mathcal{M}$ .

---

<sup>2</sup>R. Tang and Y. Yang (2024). **Adaptivity of Diffusion Models to Manifold Structures**. *AISTATS*.

# Convergence of diffusion models under the manifold hypothesis

Convergence rates (even optimal ones) expressed **in terms of the ambient dimension  $d$**  fall short of capturing the empirical success of DDMs

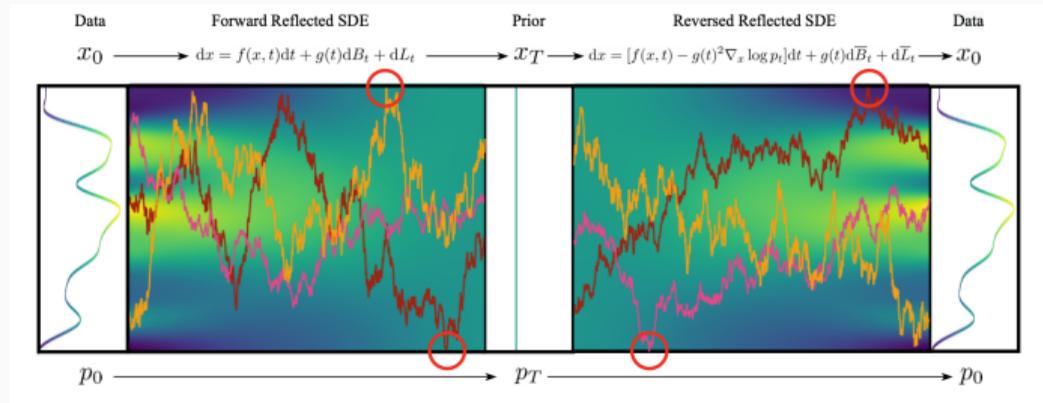
↪ gap is related to the **manifold hypothesis**: real-world high-dimensional data often reside on **lower-dimensional manifolds**, to which well-trained **generative models are believed to adapt**

- **multiplicative factor  $C(d)$**  in Tang and Yang's convergence rate is **of order  $d^{s+\tilde{d}/2}$**  and thus potentially very large for high ambient dimension  $d$
- most recently, [Azangulov et al. \(2024\)](#)<sup>2</sup> show that this **multiplicative factor can be significantly reduced to the order  $\sqrt{d}$**

---

<sup>2</sup>I. Azangulov, G. Deligiannidis and J. Rousseau (2024). **Convergence of Diffusion Models Under the Manifold Hypothesis in High-Dimensions**. arXiv:2409.18804

# Denoising reflected diffusion models



Source: Lou and Ermon (2023). Reflected Diffusion Models. *ICML*.

## Questions:

1. are diffusion models minimax learners (in terms of smoothness assumptions on  $p_0$ )?
2. how can empirical lack of curse of dimensionality be explained?  $\rightsquigarrow$  submanifold hypothesis
3. **alternative model designs** with better theoretical/experimental justification?

## Denoising Reflected Diffusion Models (DRDMs) in a nutshell

- extend DDMs by constraining both forward and backward processes to a **bounded domain**  $D \subset \mathbb{R}^d$
- **forward process** includes a **reflection term** to enforce boundary constraints,

$$dX_t = b(X_t) dt + \sigma(X_t) dW_t + \nu(X_t) d\ell_t, \quad X_0 \in \bar{D},$$

where  $\ell_t$  is the **local time** at the boundary  $\partial D$  and  $\nu$  determines the direction of reflection

- under technical conditions (Cattiaux, 1988)<sup>3</sup>, **time reversed process** is reflected at the boundary as well and solves

$$d\tilde{X}_t = -\bar{b}(t, \tilde{X}_t) dt + \sigma(\tilde{X}_t) d\bar{W}_t + \nu(\tilde{X}_t) d\bar{\ell}_t, \quad \tilde{X}_0 \sim p_T$$

- retains Markov properties with **constrained state space** and specific Neumann boundary conditions

---

<sup>3</sup>Cattiaux (1988). Time reversal of diffusion processes with a boundary condition. SPA

## Denoising Reflected Diffusion Models (DRDMs) in a nutshell

- extend DDMs by constraining both forward and backward processes to a **bounded domain**  $D \subset \mathbb{R}^d$
- **forward process** includes a **reflection term** to enforce boundary constraints,

$$dX_t = b(X_t) dt + \sigma(X_t) dW_t + \nu(X_t) d\ell_t, \quad X_0 \in \bar{D},$$

where  $\ell_t$  is the **local time** at the boundary  $\partial D$  and  $\nu$  determines the direction of reflection

- under technical conditions (Cattiaux, 1988)<sup>3</sup>, **time reversed process** is reflected at the boundary as well and solves

$$d\tilde{X}_t = -\bar{b}(t, \tilde{X}_t) dt + \sigma(\tilde{X}_t) d\bar{W}_t + \nu(\tilde{X}_t) d\bar{\ell}_t, \quad \tilde{X}_0 \sim p_T$$

- retains Markov properties with **constrained state space** and specific Neumann boundary conditions

---

<sup>3</sup>Cattiaux (1988). **Time reversal of diffusion processes with a boundary condition.** *SPA*

## Comparison: DDMs versus their reflected counterparts

- **domain:** DDMs operate on  $\mathbb{R}^d$ , while DRDMs are constrained to a bounded domain  $D \subset \mathbb{R}^d$   
↪ DRDMs include **reflection terms** to ensure dynamics remain in  $\overline{D}$ , while DDMs do not account for spatial constraints
- **implementation complexity:** DRDMs require managing **boundary local times** and **Neumann conditions**, introducing additional complexity
- **applications:** DRDMs are better suited for generating data **confined to specific domains** or **bounded physical spaces**



Source: Lou and Ermon (2023). Reflected Diffusion Models. *ICML*.

## Generative modelling with **reflected** diffusions: Forward process

- assume that  $D \subseteq \mathbb{R}^d$  is an **open, connected and bounded set with  $\mathcal{C}^\infty$  boundary  $\partial D$**
- consider the **reflected time-homogeneous forward model**

$$dX_t = b(X_t) dt + \sigma(X_t) dW_t + \nu(X_t) d\ell_t, \quad X_0 \in \bar{D},$$

with smooth and bounded coefficients  $b: \bar{D} \rightarrow \mathbb{R}^d$ ,  $\sigma: \bar{D} \rightarrow \mathbb{R}^{d \times d}$  and **conormal reflection** determined by

$$\nu(x) := \frac{1}{2} \sigma \sigma^\top(x) n(x), \quad x \in \partial D$$

$\rightsquigarrow$   $n$  is the inward unit normal vector at the boundary  $\partial D$ ,  $(\ell_t)_{t \geq 0}$  is the local time at  $\partial D$  satisfying

$$\ell_t = \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \int_0^t \mathbf{1}_{(\partial D)_\varepsilon}(X_s) ds$$

$\rightsquigarrow$  boundary reflection process reflects  $X$  in a conormal direction whenever it hits the boundary  $\partial D$ , thus constraining the state space of the diffusion to the compact set  $\bar{D}$

## Forward process and SDE

- **key requirement:** precise understanding of the forward process's limiting behaviour to determine the runtime needed for the backward initialisation to approximate the **true terminal forward distribution**  $p_T$
- **subsequent simplification:** choose  $b = \nabla f$  and  $\sigma = \sqrt{2f} \mathbb{I}_{d \times d}$  for some **diffusivity**  
 $\mathcal{C}^\infty(\bar{D}) \ni f : \mathbb{R}^d \rightarrow [f_{\min}, \infty) \subset (0, \infty)$

↪ time-homogeneous **forward dynamics** are described by the **divergence form**  $L^2$ -generator

$$\mathcal{A} = \nabla \cdot f \nabla = \langle \nabla f, \nabla \cdot \rangle + f \Delta,$$

corresponding to the constrained SDE

$$dX_t = \nabla f(X_t) dt + \sqrt{2f(X_t)} dW_t + \nu(X_t) d\ell_t$$

- ↪ both the **reflected forward and backward SDEs** exhibit **normal reflection** at the boundary
- divergence theorem  $\implies$  **invariant distribution** of the forward Markov process  $X$  is the (easy-to-sample-from) **uniform distribution** on  $\bar{D}$ , i.e.,  $\mu = \text{Leb}|_{\bar{D}} / \text{Leb}(\bar{D})$

## Forward process and SDE

- **key requirement:** precise understanding of the forward process's limiting behaviour to determine the runtime needed for the backward initialisation to approximate the **true terminal forward distribution**  $p_T$
  - **subsequent simplification:** choose  $b = \nabla f$  and  $\sigma = \sqrt{2f} \mathbb{I}_{d \times d}$  for some **diffusivity**  $\mathcal{C}^\infty(\bar{D}) \ni f : \mathbb{R}^d \rightarrow [f_{\min}, \infty) \subset (0, \infty)$
- ↪ time-homogeneous **forward dynamics** are described by the **divergence form**  $L^2$ -generator

$$\mathcal{A} = \nabla \cdot f \nabla = \langle \nabla f, \nabla \cdot \rangle + f \Delta,$$

corresponding to the constrained SDE

$$dX_t = \nabla f(X_t) dt + \sqrt{2f(X_t)} dW_t + \nu(X_t) d\ell_t$$

- ↪ both the **reflected forward and backward SDEs** exhibit **normal reflection at the boundary**
- divergence theorem  $\implies$  **invariant distribution** of the forward Markov process  $X$  is the (easy-to-sample-from) **uniform distribution** on  $\bar{D}$ , i.e.,  $\mu = \text{Leb}|_{\bar{D}} / \text{Leb}(\bar{D})$

## Spectral properties

- under the given assumptions, there exist orthonormal eigenpairs  $(\lambda_j, e_j)_{j \geq 0}$  of the operator  $-\nabla \cdot f \nabla$  satisfying

$$0 = \lambda_0 < \lambda_1 \leq \lambda_2 \leq \dots, \quad \lambda_j \asymp j^{2/d}$$

- transition densities can be expressed as

$$q_t(x, y) = \sum_{j \geq 0} e^{-t\lambda_j} e_j(x) e_j(y), \quad x, y \in D$$

- bounds on eigenfunctions:

$$\|e_j\|_{H^k} \lesssim \lambda_j^{k/2} \asymp j^{k/d}, \quad j \geq 1, \quad \|e_j\|_{\infty} \lesssim j^{\tau}, \text{ for } \tau > 1/2$$

$\rightsquigarrow$  smoothing property of densities:

$$\|p_t\|_{H^k} \lesssim \|p_0\|_{\infty} e^{-tj^{2/d}} j^{\tau+k/d}, \quad t > 0,$$

for arbitrary  $\tau > 1/2 \rightsquigarrow p_t \in C^{\infty}(D)$  for any bounded initial density  $p_0$

## Backward process and score approximation

- backward dynamics becomes

$$d\tilde{X}_t = (\nabla f(\tilde{X}_t) + 2f(\tilde{X}_t)\nabla \log p_{T-t}(\tilde{X}_t)) dt + \sqrt{2f(\tilde{X}_t)} d\bar{W}_t + v(\tilde{X}_t) d\bar{e}_t,$$

with initialisation  $\tilde{X}_0 \sim p_T$

- spectral decomposition of the transition densities  $\implies$  score is explicitly given by

$$\nabla \log p_t(x) = \frac{\sum_{j \geq 0} e^{-t\lambda_j} \langle p_0, e_j \rangle_{L^2} \nabla e_j(x)}{\sum_{j \geq 0} e^{-t\lambda_j} \langle p_0, e_j \rangle_{L^2} e_j(x)}, \quad x \in D, t > 0$$

$\rightsquigarrow$  will be instrumental in analysing the score approximation properties of neural networks underlying the algorithm

## Neural network classes

- use ReLU activation function  $\sigma(y) = y \vee 0$ , and, for any  $b, x \in \mathbb{R}^m$ , let  $\sigma_b(x) = \begin{bmatrix} \sigma(x_1 - b_1) \\ \sigma(x_2 - b_2) \\ \vdots \\ \sigma(x_m - b_m) \end{bmatrix}$
- consider functions of the form

$$\varphi(x) = A_L \sigma_{b_L} A_{L-1} \sigma_{b_{L-1}} \cdots A_1 \sigma_{b_1} A_0 x,$$

where  $A_i \in \mathbb{R}^{W_{i+1} \times W_i}$ ,  $b_i \in \mathbb{R}^{W_{i+1}}$  for  $i = 0, \dots, L$ , and where there are at most a total of  $S$  non-zero entries of the  $A_i$ 's and  $b_i$ 's and all entries are numerically at most  $B$

↪ class of networks

$$\Phi(L, W, S, B) := \left\{ A_L \sigma_{b_L} A_{L-1} \sigma_{b_{L-1}} \cdots A_1 \sigma_{b_1} A_0 \mid A_i \in \mathbb{R}^{W_{i+1} \times W_i}, b_i \in \mathbb{R}^{W_{i+1}}, \right. \\ \left. \sum_{i=0}^L (\|A_i\|_0 + \|b_i\|_0) \leq S, \max_{i \in \{0, \dots, L\}} (\|A_i\|_\infty \vee \|b_i\|_\infty) \leq B \right\}$$

## Generative modeling with reflected diffusions

- denote the **true score** by  $\mathfrak{s}^\circ(x, t) := \nabla \log p_t(x)$ , and assume we are given **data samples**  $(X_{0,i})_{i \in [n]} \stackrel{\text{i.i.d.}}{\sim} p_0$
- for a **hypothesis class**  $\mathcal{S}$  of neural networks and  $\mathfrak{s} \in \mathcal{S} \cup \{\mathfrak{s}^\circ\}$ , define

$$L_{\mathfrak{s}}(x) := \mathbb{E} \left[ \int_{\underline{T}}^{\overline{T}} |\mathfrak{s}(X_t, t) - \nabla_y \log q_t(x, X_t)|^2 \mid X_0 = x \right]$$

$\rightsquigarrow$   $\overline{T}$  is the terminal runtime of the reflected forward process

$\rightsquigarrow$   $\underline{T} \in (0, \overline{T})$  is such that we run the reflected generative process, which is initialised with distribution  $\mathcal{U}(D)$ , until  $\overline{T} - \underline{T}$

- denote the **empirical denoising score matching loss** associated to  $\mathfrak{s}$  by

$$\hat{L}_{\mathfrak{s},n} := \frac{1}{n} \sum_{i=1}^n L_{\mathfrak{s}}(X_{0,i}),$$

and define the **empirical score minimiser** by

$$\hat{\mathfrak{s}}_n := \arg \min_{\mathfrak{s} \in \mathcal{S}} \hat{L}_{\mathfrak{s},n}$$

- let  $\bar{X}^{\hat{s}}$  be a solution of the reflected SDE

$$d\bar{X}_t^{\hat{s}} = (\nabla f(\bar{X}_t^{\hat{s}}) + 2f(\bar{X}_t^{\hat{s}})s(\bar{X}_t^{\hat{s}}, t)) dt + \sqrt{2f(\bar{X}_t^{\hat{s}})} d\bar{W}_t + v(\bar{X}_t^{\hat{s}}) d\bar{\ell}_t, \quad t \in [0, \bar{T}, -\underline{T}],$$

$$\bar{X}_0^{\hat{s}} \sim \mathcal{U}(\bar{D}),$$

for some Brownian motion  $(\bar{W}_t)_{t \in [0, \bar{T} - \underline{T}]}$  and local time  $(\bar{\ell}_t)_{t \in [0, \bar{T} - \underline{T}]}$  at the boundary  $\partial D$ , and denote its density at time  $t$  by  $\bar{p}_t^{\hat{s}}$

- initialisation  $X_0^{\hat{s}_n} \sim \mathcal{U}(\bar{D})$
- $(\bar{p}_t^{\hat{s}_n})_{t \in [0, \bar{T}]}$  are the densities of the backward process driven by the score estimate  $\hat{s}_n$
- assessing the quality of the generated samples boils down to analysing the distance between the distribution induced by  $p_0$  and the (random) distribution induced by  $\bar{p}_{\bar{T} - \underline{T}}^{\hat{s}_n}$

- let  $\bar{X}^{\hat{s}}$  be a solution of the reflected SDE

$$d\bar{X}_t^{\hat{s}} = (\nabla f(\bar{X}_t^{\hat{s}}) + 2f(\bar{X}_t^{\hat{s}})s(\bar{X}_t^{\hat{s}}, t)) dt + \sqrt{2f(\bar{X}_t^{\hat{s}})} d\bar{W}_t + v(\bar{X}_t^{\hat{s}}) d\bar{\ell}_t, \quad t \in [0, \bar{T}, -\underline{T}],$$

$$\bar{X}_0^{\hat{s}} \sim \mathcal{U}(\bar{D}),$$

for some Brownian motion  $(\bar{W}_t)_{t \in [0, \bar{T} - \underline{T}]}$  and local time  $(\bar{\ell}_t)_{t \in [0, \bar{T} - \underline{T}]}$  at the boundary  $\partial D$ , and denote its density at time  $t$  by  $\bar{p}_t^{\hat{s}}$

- initialisation  $X_0^{\hat{s}_n} \sim \mathcal{U}(\bar{D})$
- $\rightsquigarrow (\bar{p}_t^{\hat{s}_n})_{t \in [0, \bar{T}]}$  are the densities of the backward process driven by the score estimate  $\hat{s}_n$
- $\rightsquigarrow$  assessing the **quality of the generated samples** boils down to analysing the distance between the **distribution induced by  $p_0$**  and the (random) **distribution induced by  $\bar{p}_{\bar{T} - \underline{T}}^{\hat{s}_n}$**

## Main result

**Theorem** (Holk, CS and LT (2024))

Assume  $p_0 = \tilde{p}_0 + \alpha$  with  $\tilde{p}_0 \in H_c^s(D)$ ,  $\alpha > 0$ ,  $s \in \mathbb{N} \cap (d/2, \infty)$ , and let

$$\underline{T} \asymp n^{-\frac{2s}{((2-d/s)\wedge 1)(2s+d)}}, \quad \bar{T} = \frac{s}{\lambda_1(2s+d)} \log n.$$

Then, there exists a class of feed forward ReLU neural networks  $\mathcal{S}$ , with explicit size constraints in terms of  $n$ ,  $d$  and  $s$ , such that

$$\mathbb{E}[\text{TV}(p_0, \hat{p}_{\bar{T}-\underline{T}}^{\hat{\mathcal{S}}_n})] \lesssim n^{-\frac{s}{2s+d}} (\log n)^3 (\log \log n)^{1/2}.$$

Letting  $\bar{p}_t^{\mathcal{S}}$  be the density at time  $t$  of the time-reversed forward process

$$\mathbb{E}[\text{TV}(p_0, \hat{p}_{\bar{T}-\underline{T}}^{\hat{\mathcal{S}}_n})] \leq \underbrace{\text{TV}(p_0, p_{\underline{T}})}_{=: \text{(I)}} + \underbrace{\text{TV}(\mathbb{P}(X_{\bar{T}} \in \cdot \mid X_0 \sim p_0), \mathcal{U}(\bar{D}))}_{=: \text{(II)}} + \underbrace{\mathbb{E}[\text{TV}(\bar{p}_{\bar{T}-\underline{T}}^{\mathcal{S}^\circ}, \hat{p}_{\bar{T}-\underline{T}}^{\hat{\mathcal{S}}_n})]}_{=: \text{(III)}}.$$

## Main result

**Theorem** (Holk, CS and LT (2024))

Assume  $p_0 = \tilde{p}_0 + \alpha$  with  $\tilde{p}_0 \in H_c^s(D)$ ,  $\alpha > 0$ ,  $s \in \mathbb{N} \cap (d/2, \infty)$ , and let

$$\underline{T} \asymp n^{-\frac{2s}{((2-d/s)\wedge 1)(2s+d)}}, \quad \bar{T} = \frac{s}{\lambda_1(2s+d)} \log n.$$

Then, there exists a class of feed forward ReLU neural networks  $\mathcal{S}$ , with explicit size constraints in terms of  $n$ ,  $d$  and  $s$ , such that

$$\mathbb{E}[\text{TV}(p_0, \hat{p}_{\bar{T}-\underline{T}}^{\hat{\mathfrak{S}}_n})] \lesssim n^{-\frac{s}{2s+d}} (\log n)^3 (\log \log n)^{1/2}.$$

Letting  $\bar{p}_t^{\mathfrak{S}}$  be the density at time  $t$  of the time-reversed forward process

$$\mathbb{E}[\text{TV}(p_0, \hat{p}_{\bar{T}-\underline{T}}^{\hat{\mathfrak{S}}_n})] \leq \underbrace{\text{TV}(p_0, p_{\underline{T}})}_{=:(\text{I})} + \underbrace{\text{TV}(\mathbb{P}(X_{\bar{T}} \in \cdot \mid X_0 \sim p_0), \mathcal{U}(\bar{D}))}_{=:(\text{II})} + \underbrace{\mathbb{E}[\text{TV}(\bar{p}_{\bar{T}-\underline{T}}^{\mathfrak{S}^\circ}, \hat{p}_{\bar{T}-\underline{T}}^{\hat{\mathfrak{S}}_n})]}_{=:(\text{III})}.$$

## Error decomposition

$$\mathbb{E}[\text{TV}(p_0, \hat{p}_{\bar{T}-\underline{T}}^{\hat{\xi}_n})] \leq \underbrace{\text{TV}(p_0, p_{\underline{T}})}_{=: \text{(I)}} + \underbrace{\text{TV}(\mathbb{P}(X_{\bar{T}} \in \cdot \mid X_0 \sim p_0), \mathcal{U}(\bar{D}))}_{=: \text{(II)}} + \underbrace{\mathbb{E}[\text{TV}(\bar{p}_{\bar{T}-\underline{T}}^{\hat{\xi}_n}, \hat{p}_{\bar{T}-\underline{T}}^{\hat{\xi}_n})]}_{=: \text{(III)}}$$

(I) represents the error induced by **stopping early the backward process** initialised by the true forward terminal density  $p_{\bar{T}}$  at time  $\bar{T} - \underline{T}$

↪ controlled via **small time heat kernel bounds** for the transition densities

↪ relies on Hölder continuity of  $p_0$ : for  $\beta \in [(2 - d/s) \wedge 1, 1]$ ,

$$|p_0(x) - p_0(y)| \leq c_\beta |x - y|^\beta, \quad x, y \in D$$

**Lemma** (Holk, CS and LT (2024))

There exists a constant  $C$  depending only on  $f, d, D, \beta$  and  $c_\beta$  such that

$$\text{TV}(p_0, p_{\underline{T}}) \leq C \underline{T}^{\beta/2}, \quad \underline{T} \leq 1.$$

## Error decomposition

$$\mathbb{E}[\text{TV}(p_0, \hat{p}_{\bar{T}-\underline{T}}^{\hat{g}_n})] \leq \underbrace{\text{TV}(p_0, p_{\bar{T}})}_{=: \text{(I)}} + \underbrace{\text{TV}(\mathbb{P}(X_{\bar{T}} \in \cdot \mid X_0 \sim p_0), \mathcal{U}(\bar{D}))}_{=: \text{(II)}} + \underbrace{\mathbb{E}[\text{TV}(\bar{p}_{\bar{T}-\underline{T}}^{\hat{g}_n}, \hat{p}_{\bar{T}-\underline{T}}^{\hat{g}_n})]}_{=: \text{(III)}}$$

(II) is the error associated to **starting the backward process in its stationary distribution** instead of  $p_{\bar{T}}$

↪ controlled in terms of the **spectral gap  $\lambda_1$**  of  $\mathcal{A}$ , which can be lower bounded by

$\lambda_1 \geq f_{\min}/C_P(D)$ , where  $C_P(D)$  is the Poincaré constant of the domain  $D$

**Lemma** (Holk, CS and LT (2024))

It holds that

$$\text{TV}(\mathbb{P}(X_{\bar{T}} \in \cdot \mid X_0 \sim p_0), \mathcal{U}(\bar{D})) \leq \frac{\sqrt{\text{Leb}(D)}}{2} \|p_0\|_{L^2} e^{-\lambda_1 \bar{T}}, \quad \bar{T} > 0.$$

## Error decomposition

$$\mathbb{E}[\text{TV}(p_0, \bar{p}_{\bar{T}-\underline{T}}^{\hat{\mathfrak{s}}_n})] \leq \underbrace{\text{TV}(p_0, p_{\underline{T}})}_{=:(\text{I})} + \underbrace{\text{TV}(\mathbb{P}(X_{\bar{T}} \in \cdot \mid X_0 \sim p_0), \mathcal{U}(\bar{D}))}_{=:(\text{II})} + \underbrace{\mathbb{E}[\text{TV}(\bar{p}_{\bar{T}-\underline{T}}^{\mathfrak{s}^\circ}, \bar{p}_{\bar{T}-\underline{T}}^{\hat{\mathfrak{s}}_n})]}_{=:(\text{III})}$$

(III) quantifies the error coming from running the backward process with the drift determined by the estimated score  $\hat{\mathfrak{s}}_n$  instead of the true score  $\mathfrak{s}^\circ$

↪ by Girsanov's theorem and Pinsker's inequality, controlled by

$$\mathbb{E}\left[\int_{\underline{T}}^{\bar{T}} \int_D |\hat{\mathfrak{s}}(x, t) - \nabla \log p_t(x)|^2 p_t(x) dx dt\right]$$

↪ key to bounding it: equivalence between **explicit** and **denoising score matching**, i.e.,

$$\int_{\underline{T}}^{\bar{T}} \int_D |\mathfrak{s}(y, t) - \nabla \log p_t(y)|^2 p_t(y) dy dt = \mathbb{E}[L_{\mathfrak{s}}(X_0)] + C,$$

where  $C \leq 0$  is a constant that is independent of  $\mathfrak{s}$

## Bounding the score matching error (III)

Generalisation loss can be bounded in terms of the **minimal score approximation error** over the class  $\mathcal{S}$  and the **complexity of the induced function class**  $\mathcal{L} := \{L_{\mathfrak{s}} : \mathfrak{s} \in \mathcal{S}\}$  for a desired precision level  $\delta$ :

**Theorem** (Oko et. al (2023))

Suppose that  $\sup_{\mathfrak{s} \in \mathcal{S}} \|L_{\mathfrak{s}} - L_{\mathfrak{s}^*}\|_{\infty} \leq C(\mathcal{L}) < \infty$ . Then, for any  $\delta > 0$  such that  $\mathcal{N}(\mathcal{L}, \|\cdot\|_{\infty}, \delta) \geq 3$ , it holds that

$$\begin{aligned} & \mathbb{E} \left[ \int_{\underline{I}}^{\overline{T}} \int_D |\hat{\mathfrak{s}}(x, t) - \nabla \log p_t(x)|^2 p_t(x) dx dt \right] \\ & \leq 2 \inf_{\mathfrak{s} \in \mathcal{S}} \int_{\underline{I}}^{\overline{T}} \int_D |\mathfrak{s}(x, t) - \nabla \log p_t(x)|^2 p_t(x) dx dt + 2 \frac{C(\mathcal{L})}{n} \left( \frac{37}{9} \log \mathcal{N}(\mathcal{L}, \|\cdot\|_{\infty}, \delta) + 32 \right) + 3\delta. \end{aligned}$$

- ↪ **next step:** control both the uniform loss upper bound  $C(\mathcal{L})$  and the covering number  $\mathcal{N}(\mathcal{L}, \|\cdot\|_{\infty}, \delta)$  for  $\delta = n^{-2s/(2s+d)}$  [✓ since  $\log \mathcal{N}(\mathcal{L}, \|\cdot\|_{\infty}, \delta) \leq \log \mathcal{N}(\Phi(L, W, S, B), \|\cdot\|_{\infty}, \frac{\delta}{C\overline{T}}) \lesssim LS \log \left( \frac{LWB\overline{T}}{\delta} \right)$ ]
- ↪ **final and most fundamental question:** treatment of the explicit score approximation error

# Strategy for bounding the approximation error

## 1. Truncate

$$p_t(x) = \sum_{j=0}^{\infty} e^{-\lambda_j t} \langle p_0, e_j \rangle e_j(x) \approx \sum_{j=0}^N e^{-\lambda_j t} \langle p_0, e_j \rangle e_j(x) =: h_N(x, t)$$

and use  $\nabla h_N(x, t) \approx \nabla p_t(x)$

$$\rightsquigarrow \int_{\underline{T}}^{\overline{T}} \int_{\mathbb{R}^d} \left| s^\circ(x, t) - \frac{\nabla h_N(x, t)}{h_N(x, t)} \right|^2 p_t(x) dx dt \lesssim N^{-2s/d} \log \underline{T} \implies N \asymp n^{d/(2s+d)}$$

- for an appropriately chosen discrete set of time points  $\{t_j\}$ , use the spatial smoothness of  $h_N(x, t_j)$  induced by the Sobolev smoothness of  $p_0$  to obtain an efficient neural network approximation of  $h_N(\cdot, t_j)$ , based on general approximation results from Suzuki (2019)<sup>4</sup>;
- approximate the space-time functions  $h_N(x, t)$ ,  $\nabla h_N(x, t)$  by constructing a neural network approximation of a **polynomial time interpolation** of the neural networks from Step 2., where the interpolation degree is adapted to the parameters  $N$ ,  $s$  and  $d$ .

---

<sup>4</sup>Suzuki (2019). **Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality.** *ICLR*

# Strategy for bounding the approximation error

## 1. Truncate

$$p_t(x) = \sum_{j=0}^{\infty} e^{-\lambda_j t} \langle p_0, e_j \rangle e_j(x) \approx \sum_{j=0}^N e^{-\lambda_j t} \langle p_0, e_j \rangle e_j(x) =: h_N(x, t)$$

and use  $\nabla h_N(x, t) \approx \nabla p_t(x)$

$$\rightsquigarrow \int_{\underline{T}}^{\overline{T}} \int_{\mathbb{R}^d} \left| s^\circ(x, t) - \frac{\nabla h_N(x, t)}{h_N(x, t)} \right|^2 p_t(x) dx dt \lesssim N^{-2s/d} \log \underline{T} \implies N \asymp n^{d/(2s+d)}$$

- for an appropriately chosen discrete set of time points  $\{t_j\}$ , use the spatial smoothness of  $h_N(x, t_j)$  induced by the Sobolev smoothness of  $p_0$  to obtain an efficient neural network approximation of  $h_N(\cdot, t_j)$ , based on general approximation results from [Suzuki \(2019\)](#)<sup>4</sup>;
- approximate the space-time functions  $h_N(x, t)$ ,  $\nabla h_N(x, t)$  by constructing a neural network approximation of a [polynomial time interpolation](#) of the neural networks from Step 2., where the interpolation degree is adapted to the parameters  $N$ ,  $s$  and  $d$ .

---

<sup>4</sup>Suzuki (2019). **Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality.** *ICLR*

# Strategy for bounding the approximation error

## 1. Truncate

$$p_t(x) = \sum_{j=0}^{\infty} e^{-\lambda_j t} \langle p_0, e_j \rangle e_j(x) \approx \sum_{j=0}^N e^{-\lambda_j t} \langle p_0, e_j \rangle e_j(x) =: h_N(x, t)$$

and use  $\nabla h_N(x, t) \approx \nabla p_t(x)$

$$\rightsquigarrow \int_{\underline{T}}^{\overline{T}} \int_{\mathbb{R}^d} \left| s^\circ(x, t) - \frac{\nabla h_N(x, t)}{h_N(x, t)} \right|^2 p_t(x) dx dt \lesssim N^{-2s/d} \log \underline{T} \implies N \asymp n^{d/(2s+d)}$$

- for an appropriately chosen discrete set of time points  $\{t_j\}$ , use the spatial smoothness of  $h_N(x, t_j)$  induced by the Sobolev smoothness of  $p_0$  to obtain an efficient neural network approximation of  $h_N(\cdot, t_j)$ , based on general approximation results from [Suzuki \(2019\)](#)<sup>4</sup>;
- approximate the space-time functions  $h_N(x, t)$ ,  $\nabla h_N(x, t)$  by constructing a neural network approximation of a [polynomial time interpolation](#) of the neural networks from Step 2., where the interpolation degree is adapted to the parameters  $N$ ,  $s$  and  $d$ .

---

<sup>4</sup>Suzuki (2019). **Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality.** *ICLR*

## Theorem (Holk, CS and LT (2024))

Let  $0 < \underline{T} < \bar{T}$  and  $n \in \mathbb{N}$  sufficiently large be given with  $\underline{T} \in \text{Poly}(n^{-1})$ . Then, there exists a neural network  $\mathfrak{g} \in \Phi(L(n), W(n), S(n), B(n))$  satisfying

$$\int_{\underline{T}}^{\bar{T}} \int_D |\mathfrak{g}(x, t) - \nabla_x \log p_t(x)|^2 p_t(x) dx dt \lesssim n^{-\frac{2s}{2s+d}} (\log n)^2 (\bar{T} + \log(\underline{T}^{-1})).$$

The size of the network is evaluated as

$$\begin{aligned} L(n) &\lesssim \log n \log \log n, \\ \|W(n)\|_{\infty} &\lesssim Mn^{\frac{d}{2s+d}} \log n, \\ S(n) &\lesssim Mn^{\frac{d}{2s+d}} (\log n)^2, \quad \text{and} \\ B(n) &\lesssim n^{\frac{1}{2s+d}} \vee \frac{1}{\underline{T}}, \end{aligned}$$

where  $M \in O(|\log \frac{\bar{T}}{\underline{T}}|)$ . Furthermore, the network can be chosen such that there exists a constant  $C < \infty$  depending only on  $p_0$  and  $D$  such that  $|\mathfrak{g}(x, t)| \leq \frac{C}{\sqrt{t}}$  for all  $t \in [\underline{T}, \bar{T}]$  and  $x \in D$ .

## Future Research

- extending the DRDM framework to **data supported on lower-dimensional submanifolds**  $\rightsquigarrow$  challenging because  $L^2$ -techniques don't translate naturally, but explicit form of Skorokhod map for reflected BM in  $D = [0, 1]^d$  gives a possible starting point
- Use score approximation techniques to unify the statistical analysis in the framework of **denoising Markov models** (Benton et al., 2024)<sup>5</sup> for appropriate self-adjoint forward Markov processes
- $\rightsquigarrow$  efficient sampling methods for the **generative reflected process with estimated score**?
  - natural sampling schemes for reflected diffusions combine an Euler–Maryuama discretisation with a projection (Śłomiński, 1994)<sup>6</sup> or rejection (Fishman et al., 2024)<sup>7</sup> step
- how to alter the training objective to **avoid training data replication** while maintaining statistical optimality?
  - Vardanyan et. al (2024)<sup>8</sup> suggest penalised Wasserstein-1-GAN generator obtained from

$$\hat{g}_n \in \arg \min_{g \in \mathcal{G}} \left\{ W_1 \left( g \# \mathcal{U}_d, \frac{1}{n} \sum_{i=1}^n \delta_{X_i} \right) + \lambda \min_{h \in \mathcal{H}} \int_{[0,1]^d} |h \circ g(u) - u|^2 du \right\},$$

---

<sup>5</sup>Benton, Shi, De Bortoli, Delegiannidis and Doucet (2024). **From denoising diffusions to denoising Markov models**. *JRSS B*

<sup>6</sup>Śłomiński (1994). **On approximation of solutions of multidimensional SDE's with reflecting boundary conditions**. *SPA*

<sup>7</sup>Fishman, Klarner, Mathieu, Hutchinson and De Bortoli (2024). **Metropolis Sampling for Constrained Diffusion Models**.

*NeurIPS*

<sup>8</sup>Vardanyan, Hunanyan, Galstyan, Minasyan and Dalalyan (2024). **Statistically Optimal Generative Modeling with Maximum Deviation from the Empirical Distribution**. *ICML*

## Future Research

- extending the DRDM framework to **data supported on lower-dimensional submanifolds**  $\rightsquigarrow$  challenging because  $L^2$ -techniques don't translate naturally, but explicit form of Skorokhod map for reflected BM in  $D = [0, 1]^d$  gives a possible starting point
- Use score approximation techniques to unify the statistical analysis in the framework of **denoising Markov models** (Benton et al., 2024)<sup>5</sup> for appropriate self-adjoint forward Markov processes
- $\rightsquigarrow$  efficient sampling methods for the **generative reflected process with estimated score**
  - natural sampling schemes for reflected diffusions combine an Euler–Maryuama discretisation with a projection (Śłomiński, 1994)<sup>6</sup> or rejection (Fishman et al., 2024)<sup>7</sup> step
- how to alter the training objective to **avoid training data replication** while maintaining statistical optimality?
  - Vardanyan et. al (2024)<sup>8</sup> suggest penalised Wasserstein-1-GAN generator obtained from

$$\hat{g}_n \in \arg \min_{g \in \mathcal{G}} \left\{ W_1 \left( g \# \mathcal{U}_d, \frac{1}{n} \sum_{i=1}^n \delta_{X_i} \right) + \lambda \min_{h \in \mathcal{H}} \int_{[0,1]^d} |h \circ g(u) - u|^2 du \right\},$$

---

<sup>5</sup>Benton, Shi, De Bortoli, Delegiannidis and Doucet (2024). **From denoising diffusions to denoising Markov models.** *JRSS B*

<sup>6</sup>Śłomiński (1994). **On approximation of solutions of multidimensional SDE's with reflecting boundary conditions.** *SPA*

<sup>7</sup>Fishman, Klarner, Mathieu, Hutchinson and De Bortoli (2024). **Metropolis Sampling for Constrained Diffusion Models.**

*NeurIPS*

<sup>8</sup>Vardanyan, Hunanyan, Galstyan, Minasyan and Dalalyan (2024). **Statistically Optimal Generative Modeling with Maximum Deviation from the Empirical Distribution.** *ICML*

## Future Research

- extending the DRDM framework to **data supported on lower-dimensional submanifolds**  $\rightsquigarrow$  challenging because  $L^2$ -techniques don't translate naturally, but explicit form of Skorokhod map for reflected BM in  $D = [0, 1]^d$  gives a possible starting point
- Use score approximation techniques to unify the statistical analysis in the framework of **denoising Markov models** (Benton et al., 2024)<sup>5</sup> for appropriate self-adjoint forward Markov processes
- $\rightsquigarrow$  efficient sampling methods for the **generative reflected process with estimated score**?
  - natural sampling schemes for reflected diffusions combine an Euler–Maryuama discretisation with a projection (Śłomiński, 1994)<sup>6</sup> or rejection (Fishman et al., 2024)<sup>7</sup> step
- how to alter the training objective to **avoid training data replication** while maintaining statistical optimality?
  - Vardanyan et al (2024)<sup>8</sup> suggest penalised Wasserstein-1-GAN generator obtained from

$$\hat{g}_n \in \arg \min_{g \in \mathcal{G}} \left\{ W_1 \left( g \# \mathcal{U}_d, \frac{1}{n} \sum_{i=1}^n \delta_{X_i} \right) + \lambda \min_{h \in \mathcal{H}} \int_{[0,1]^d} |h \circ g(u) - u|^2 du \right\},$$

---

<sup>5</sup>Benton, Shi, De Bortoli, Delegiannidis and Doucet (2024). **From denoising diffusions to denoising Markov models**. *JRSS B*

<sup>6</sup>Śłomiński (1994). **On approximation of solutions of multidimensional SDE's with reflecting boundary conditions**. *SPA*

<sup>7</sup>Fishman, Klarner, Mathieu, Hutchinson and De Bortoli (2024). **Metropolis Sampling for Constrained Diffusion Models**. *NeurIPS*

<sup>8</sup>Vardanyan, Hunanyan, Galstyan, Minasyan and Dalalyan (2024). **Statistically Optimal Generative Modeling with Maximum Deviation from the Empirical Distribution**. *ICML*