# Adaptive denoising diffusion modelling via random time reversal

MFO Mini-Workshop on *Probabilistic Perspectives in Neural Network-Based Machine Learning*

Lukas Trottner
based on joint work with Sören Christensen, Jan Kallsen and Claudia Strauch
29 October 2025

University of Stuttgart
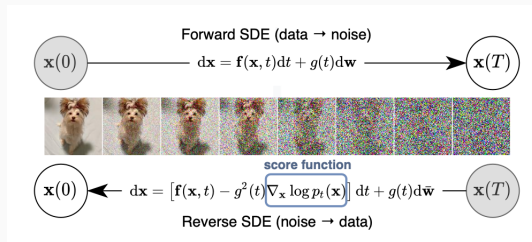Kiel University       Heidelberg University

**University of Stuttgart**
Germany

# Denoising diffusion models

- 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.*

**Denoising Diffusion Models**

- for some fixed time $T > 0$ consider the forward model

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

- under sufficient regularity conditions, the forward model has a solution $X = (X_t)_{t \in [0,T]}$ with marginal densities $p_t(x) = \int p_{0,t}(y, x)\,p_0(dy)$ such that the time reversal $\overleftarrow{X}_t = X_{T-t}$ solves

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

where

$$\overline{b}_i(t, x) = b_i(t, x) - \frac{1}{p_t(x)} \sum_{j,k=1}^{d} \frac{\partial}{\partial x_j}\left(p_t(x)\sigma_{ik}(t, x)\sigma_{jk}(t, x)\right)$$

$$= b_i(t, x) - (\nabla \cdot \Sigma(t, x))_i - (\nabla \log p_t(x))_i, \quad i = 1, \ldots, d, \Sigma = \sigma\sigma^\top$$

⤳ time-reversed process solves a time-inhomogeneous SDE, now with drift $-\overline{b}(T - \cdot, \cdot)$ involving the score $\nabla \log p_t$, which depends on the unknown data distribution $p_0$

⤳ score needs to be estimated from the data

## Denoising score matching

- denoising score matching:

$$\nabla \log p_t(x) = \frac{\int \nabla_x p_{0,t}(y,x) \, p_0(\mathrm{d}y)}{p_t(x)} = \int \nabla_x \log p_{0,t}(y,x) \underbrace{\frac{p_{0,t}(y,x) \, p_0(\mathrm{d}y)}{p_t(x)}}_{=\mathbb{P}(X_0 \in \mathrm{d}y \mid X_t = x)}$$

$$= \mathbb{E}[\nabla_2 \log p_{0,t}(X_0, X_t) \mid X_t = x]$$

and thus

$$\mathfrak{s} := \nabla \log p_t \in \underset{s \text{ meas.}}{\arg \min} \, \mathbb{E}\big[\|s(X_t) - \nabla_2 \log p_{0,t}(X_0, X_t)\|^2\big]$$

⇝ given data $(X_0^i)_{i \in [n]} \overset{\text{iid}}{\sim} p_0$ define the denoising score estimator

$$\hat{\mathfrak{s}} \in \underset{s \in \mathcal{S}}{\arg \min} \, \frac{1}{n} \sum_{i=1}^{n} \int_{\underline{T}}^{T} \|s(t, X_t^i) - \nabla_2 \log p_{0,t}(X_0^i, X_t^i)\|^2 \, \mathrm{d}t,$$

where $0 < \underline{T} \ll T$ and $\mathcal{S}$ is an approximating function class, e.g. space-time neural networks

## Generative process

On $[0, T - \underline{T}]$, simulate

$$\mathrm{d}Y_t = \left(-b(T-t, Y_t) + \nabla \cdot \Sigma(T-t, Y_t) + \Sigma(T-t, Y_t)\hat{\mathfrak{s}}(T-t, Y_t)\right)\mathrm{d}t + \sigma(T-t, Y_t)\,\mathrm{d}W_t, \quad \mathbb{P}^{Y_0}(\mathrm{d}y) \approx p_T(y)\,\mathrm{d}y$$

Output:

$$Y_{T-\underline{T}} \overset{d}{\approx} \overleftarrow{X}_{T-\underline{T}} = X_{\underline{T}} \overset{d}{\approx} X_0$$
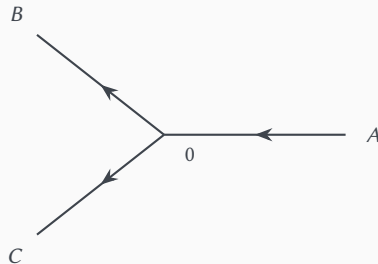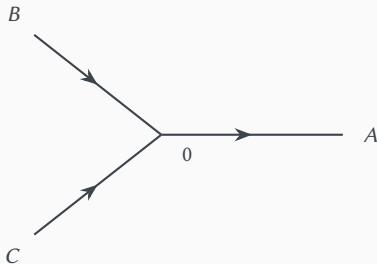
---

**Basic observations**

- time reversal at deterministic time $T$ forces the backward process to be time-inhomogeneous
- if $p_0$ has low-dimensional support $\mathcal{M}$, for small $t$ and $x$ close to $\mathcal{M}$, $\nabla \log p_t(x)$ is approximately orthogonal to $\mathcal{M}$ (Stanczuk et al., 2024)[1]
- initialising the generative process in a distribution that is not close to $\mathbb{P}^{X_T}$ and simulating for $T - \underline{T}$ time units will not give useful results ⤳ algorithm is not adaptive to the noise level in the data

---

[1] Stanczuk et al. (2024). Your diffusion model secretly knows the dimension of the data manifold. *ICML*.

## Homogeneous time reversal

- Markov property: "the past and future of a Markov process are conditionally independent given the present" ⇝ time-reversed Markov processes are Markov
- to ensure that a homogeneous Markov process remains homogeneous under time reversal, we need to reverse at a suitable random (life)time $\zeta$. This can be
  - a randomised stopping time such as an independent exponential time;
  - a last exit time;
  - a first hitting time;
  - any terminal time, that is, any stopping time $T$ such that $T = t + T \circ \theta_t$ on $\{T > t\}$
- retaining the strong Markov property under time reversal is a bit more tricky:

## *h*-transforms and time reversal

### *h*-transform

For a possibly killed, homogeneous strong Markov process $X$ with state space $S$, let $h$ be an excessive function, that is

$$\mathbb{E}_x[h(X_t)] \le h(x) \quad \text{and} \quad \lim_{t \to 0} \mathbb{E}_x[h(X_t)] = h(x).$$

Then,

$$P_t^h f(x) = \mathbb{E}_x\Big[\frac{h(X_t)}{h(x)} f(X_t) \mathbf{1}_{\{X_t \in S\}}\Big] \mathbf{1}_{(0,\infty)}(h(x)), \quad f \in \mathcal{B}_b(\mathbb{R}^d),$$

defines a sub-Markov semigroup. The corresponding Markov process $X^h$ is strong Markov and is called *h*-transform of $X$.

- suppose that $X$ is a continuous and self-dual Feller process (i.e., its generator satisfies $A = A^*$)
- if $X^h$ has a finite killing time $\zeta$, then the time-reversed process $\overleftarrow{X_t^h} = X_{\zeta-t}^h$ is homogeneous, strong Markov and is a $\overleftarrow{h}$-transform of $X$.

# *h*-transforming a killed diffusion

- consider a symmetric diffusion process

$$\mathrm{d}X_t = b(X_t)\,\mathrm{d}t + \sigma(X_t)\,\mathrm{d}W_t$$

with invariant measure $m$ and let $Z$ be its version killed at an independent exponential time with parameter $r > 0$

- as an excessive function for $Z$ use

$$h(x) = \int G_r(x, y)\,\kappa(\mathrm{d}y)$$

for the Green kernel $G_r(x, y) = \int_0^\infty e^{-rt} p_t(x, y)\,\mathrm{d}y$ and a representing measure $\kappa$

- $\kappa(\mathrm{d}y) = r\,\mathrm{d}y \rightsquigarrow h = 1$ and $Z^h = Z$
- $\kappa(\mathrm{d}y) = \frac{1}{G_r(x_0, y)} \beta(\mathrm{d}y) \rightsquigarrow Z$ conditioned to have distribution $\beta$ before killing if started in $x_0$
- $Z$ is a killed Brownian motion and $\kappa(\mathrm{d}y) = \sigma_R(\mathrm{d}y)$ for the surface measure $\sigma_R$ of an $R$-sphere $\mathbb{S}^{d-1}(R) \rightsquigarrow$ $Z^h$ is killed at last exit from $\mathbb{S}^{d-1}(R)$

## A time-homogeneous generative process

**Proposition**

1. $Z^h$ is an Itô-diffusion with dynamics

$$dZ_t^h = \left(b(Z_t^h) + \Sigma(X_t)\nabla \log h(X_t)\right) dt + \sigma(Z_t^h) dW_t$$

   outside supp $\kappa$ and its distribution at the lifetime is given by

$$\mathbb{P}_x(Z_{\zeta_-}^h \in dy) = \frac{G_r(x,y)}{h(x)}\kappa(dy)$$

2. Let $\alpha = \mathbb{P}^{Z_0^h}$. Then $\overleftarrow{Z_t^h}$ is an $\overleftarrow{h}$-transform of $Z$ with initial distribution $\mathbb{P}_\alpha(Z_{\zeta_-}^h \in dy)$ and

$$\overleftarrow{h}(x) := \int \frac{G_r(x,y)}{h(y)} \alpha(dy).$$

   In particular, $\overleftarrow{Z^h}$ has dynamics

$$d\overleftarrow{Z_t^h} = \left(b(\overleftarrow{Z_t^h}) + \Sigma(\overleftarrow{Z_t^h})\nabla \log \overleftarrow{h}(\overleftarrow{Z_t^h})\right) dt + \sigma(\overleftarrow{Z_t^h}) d\overline{W}_t,$$

   outside supp $\alpha =: \mathcal{M}$ and $\mathbb{P}_\alpha(\overleftarrow{Z_{\zeta_-}^h} \in dy \mid \overleftarrow{Z_0^h} = x) = \frac{G_r(x,y)}{\overleftarrow{h}(x)h(y)}\alpha(dy)$ for $\mathbb{P}_\alpha(Z_{\zeta_-}^h \in \cdot)$-a.e. $x$.

**A time-homogeneous generative process**

Idealised algorithm:

1. Initialise $Z_0^{\tilde{h}} \sim \tilde{\beta} \approx \mathbb{P}_\alpha(Z_{\zeta-}^h)$
   - for ergodic forward process with stationary distribution $\mu$ and small exponential killing rate $r > 0$, choose $\tilde{\beta} = \mu$   [$\leftrightarrow$ ergodic diffusion model]
   - for exponentially killed BM with small killing rate $r > 0$, choose $\tilde{\beta} = \text{Laplace}(0, (2r)^{-1/2}\mathbb{I}_d)$   [$\leftrightarrow$ variance exploding diffusion model]
   - for $\kappa(\mathrm{d}y) = \frac{1}{G_r(x_0, y)}\delta_z$, choose $\tilde{\beta} = \delta_z$

2. Simulate diffusion $Z^{\tilde{h}}$ until killing time and output $Z_{\zeta-}^{\tilde{h}}$

**Requirements for implementation**

1. learn $\nabla \log \overleftarrow{h}$ (only a function in space – no time component);
2. learn killing time $\zeta$ of $Z^{\overleftarrow{h}}$

## Learning to kill

**Polarity hypothesis**

Assume that $\mathcal{M} = \operatorname{supp} \alpha$ is polar for $X$, i.e., for any $x \in \mathbb{R}^d$, $\mathbb{P}_x(\inf\{t > 0 : X_t \in \mathcal{M}\} < \infty) = 0$.

**Theorem**

Under the polarity hypothesis, the backward process $\overleftarrow{Z^h}$ is killed at first entrance into $\mathcal{M}$.

Possible strategies to estimate a $\delta$-fattening $\mathcal{M}_\delta = \{x : \operatorname{dist}(x, \mathcal{M}) \leq \delta\}$ given data $X^1, \ldots, X^n \overset{\text{iid}}{\sim} \alpha$ and an estimator $\hat{\mathfrak{s}}$ of $\mathfrak{s} := \nabla \log \overleftarrow{h}$:

- plug-in approach: estimate $\mathcal{M}_\delta$ directly or indirectly by setting $\widehat{\mathcal{M}_\delta} = (\widehat{\mathcal{M}})_\delta$; then set
  $\hat{\zeta} := \inf\{t \geq 0 : Z_t^{\hat{\mathfrak{s}}} \in \widehat{\mathcal{M}_\delta}\}$
- use explosive behaviour of $\mathfrak{s}$ as $x \to \mathcal{M}$:

**Theorem**

Suppose that $\mathcal{M}$ is polar for $X$ and $Y$ solving $dY_t = \sigma(Y_t) dB_t$. Then, it a.s. holds that

$$\zeta = \inf\Big\{t \geq 0 : \sup_{s \leq t} |\mathfrak{s}(\overleftarrow{Z_s^h})| = \infty\Big\} = \inf\Big\{t \geq 0 : \|\mathfrak{s}(\overleftarrow{Z^h})\|_{L^2([0,t])} = \infty\Big\}.$$

**Denoising score matching**

- for $\mathbb{P}_\alpha(Z^h_{\zeta_-} \in \cdot)$-a.e. $x$

$$\mathfrak{s}(x) = \nabla \log \overleftarrow{h}(x) = \frac{1}{\overleftarrow{h}(x)} \int \nabla_x G_r(x, y) \frac{1}{h(y)} \, \alpha(\mathrm{d}y) = \int \nabla_x \log G_r(x, y) \frac{G_r(x, y)}{\overleftarrow{h}(x) h(y)} \, \alpha(\mathrm{d}y)$$

$$= \mathbb{E}\big[\nabla_x \log G_r(x, Z^{\overleftarrow{h}}_{\zeta_-}) \mid Z^{\overleftarrow{h}}_0 = x\big]$$

$$= \mathbb{E}_\alpha\big[\nabla_x \log G_r(x, Z^h_0) \mid Z^h_{\zeta_-} = x\big]$$

- this implies that on $\mathbb{R}^d \setminus \mathscr{M}_\delta$, $\mathfrak{s}$ agrees $\mathbb{P}_\alpha(Z^h_{\zeta_-} \in \cdot)$-a.e. with the minimiser of

$$\mathcal{B}(\mathbb{R}^d; \mathbb{R}^d) \ni s \mapsto \mathbb{E}_\alpha\Big[\big\| s(Z^h_{\zeta_-}) - \nabla \log G_r(Z^h_0, Z^h_{\zeta_-})\big\|^2 \mathbf{1}_{\{\|Z^h_{\zeta_-} - Z^h_0\| > \delta\}}\Big]$$

- note that if $Z^h = Z$, then $\zeta \sim \mathrm{Exp}(r)$ independent of $X$, $Z_{\zeta_-} = X_\zeta$ has full support and we have

$$\mathbb{E}_\alpha\Big[\big\| s(Z^h_{\zeta_-}) - \nabla \log G_r(Z^h_0, Z^h_{\zeta_-})\big\|^2 \mathbf{1}_{\{\|Z^h_{\zeta_-} - Z^h_0\| > \delta\}}\Big] = r \mathbb{E}_\alpha\Big[\int_0^\zeta \big\| s(Z^h_t) - \nabla \log G_r(Z^h_0, Z^h_t)\big\|^2 \mathbf{1}_{\{\|Z^h_t - Z^h_0\| > \delta\}} \, \mathrm{d}t\Big]$$

**Projection learning**

- we don't have to start the backward process approximately in $\mathbb{P}_\alpha(Z^h_{\zeta_-} \in dy)$: it will always be killed on the data support $\mathcal{M}$ and different initialisations will yield different output distributions supported on $\mathcal{M} \rightsquigarrow$ natural conditioning
- a natural question is therefore what happens if we don't start the generative process from pure noise but something more informative, say a masked or moderately noised picture



- it turns out that the natural conditioning aspect entails a blessing of dimensionality

## Projection learning

Let $Z$ be an exponentially killed Brownian motion. Then,

$$\overleftarrow{h}(x) = \int G_r(x, y)\,\alpha(\mathrm{d}y), \quad G_r(x, y) = 2(2\pi)^{-d/2} r\Big(\frac{\sqrt{2r}}{|x-y|}\Big)^{\frac{d-2}{2}} K_{\frac{d-2}{2}}\Big(\frac{\sqrt{2r}}{|x-y|}\Big).$$

For large $d$,

$$\nabla \log \overleftarrow{h}(x) \approx d\,\frac{\int \frac{x-y}{|x-y|^d}\,\alpha(\mathrm{d}y)}{\int |x-y|^{2-d}\,\alpha(\mathrm{d}y)}$$

and thus, if there is a unique projection $x^* \in \arg\min_{y \in \mathscr{M}} |x - y|$ of $x$ onto $\mathscr{M}$, then

$$\nabla \log \overleftarrow{h}(x) \approx d\,\frac{x^* - x}{|x^* - x|^2} = d\,\frac{\mathrm{sign}(x^* - x)}{|x^* - x|}$$

### Theorem

Let $\delta, \varepsilon > 0$ and fix an observation $x \in \mathbb{R}^d$. If $\alpha(B(x, r)) > \varepsilon$ for some ball $B(x, r)$ with radius $r > 0$ around $y$, then

$$\mathbb{P}\Big(Z_{\zeta-}^{\overleftarrow{h}} \in \mathscr{M} \cap B(x, (1+\delta)r) \mid Z_0^{\overleftarrow{h}} = x\Big) \geq 1 - \frac{1}{\varepsilon}(1+\delta)^{2-d}.$$

## Projection learning

Consider now estimators $\hat{\mathfrak{s}}_n$, an independent Brownian motion $W$ and let $\widehat{Z}^{\hat{\mathfrak{s}}_n}$ be the process solving

$$\mathrm{d}\widehat{Z}_t^{\hat{\mathfrak{s}}_n} = \hat{\mathfrak{s}}_n(\widehat{Z}_t^{\hat{\mathfrak{s}}_n})\mathbf{1}_{\{t \leq \hat{\zeta}\}}\,\mathrm{d}t + \mathbf{1}_{\{t \leq \hat{\zeta}\}}\,\mathrm{d}W_t, \quad \hat{\zeta} := \inf\Big\{t \geq 0 \, : \, \|\widehat{Z}^{\hat{\mathfrak{s}}_n}\|_{L^2[0,t]} > M\Big\}.$$

### Theorem

Fix an observation $x \in \mathbb{R}^d$. Suppose that

- for any $\tilde{\delta}, \delta, \varepsilon > 0$ it holds for sufficiently large $n$ that

$$\mathbb{P}\left(\Big\|\big(\hat{\mathfrak{s}}_n(Z^{\tilde{h}}) - \mathfrak{s}(Z^{\tilde{h}})\big)\mathbf{1}_{\{Z^{\tilde{h}} \notin \mathscr{M}_{\tilde{\delta}}\}}\Big\|_{L^2(\zeta)} > \delta \,\Big|\, Z_0^{\tilde{h}} = x\right) < \varepsilon$$

- for any $n \in \mathbb{N}$ and $\tilde{\delta} > 0$, the function $\hat{\mathfrak{s}}_n$ is $L_{\tilde{\delta}}$-Lipschitz on $\mathscr{M}_{\tilde{\delta}}^{\mathsf{c}}$

Let $\delta, \varepsilon, \tilde{\delta}, \tilde{\varepsilon} > 0$. If $\alpha(B(x, r)) > \varepsilon$, then, for sufficiently large $M > 0$ and $n \in \mathbb{N}$,

$$\mathbb{P}\big(\widehat{Z}_{\hat{\zeta}}^{\hat{\mathfrak{s}}_n} \in \mathscr{M}_{\tilde{\delta}} \cap B(x, (1+\delta)r) \,\big|\, \widehat{Z}_0^{\hat{\mathfrak{s}}_n} = x\big) \, > \, 1 - \frac{1}{\varepsilon}(1+\delta)^{2-d} - \tilde{\varepsilon}.$$

Thank you for your attention!