

# Critically-Damped Langevin Score-based Generative Models:

*introduction, motivation and convergence.*

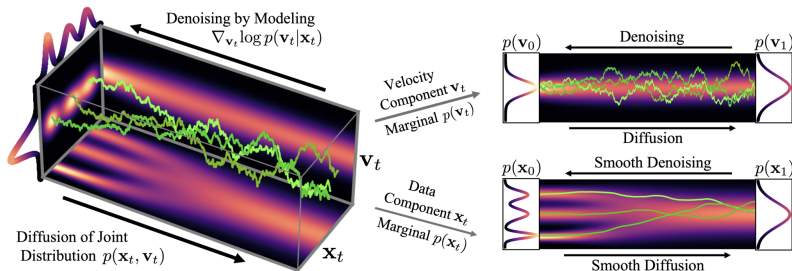
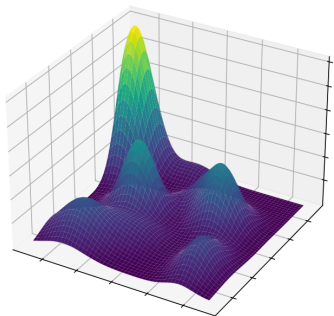


Image taken from [Dockhorn et al. \(2022\)](#).

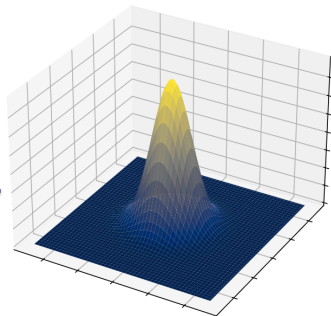
# Generative modeling framework.

- ▶  $\mathcal{D} = \{x_i\}_{i=1}^n \in (\mathbb{R}^d)^n$  a collection of i.i.d. samples from an **unknown** distribution  $\pi_{\text{data}}$
- ▶ Goal: **generate new samples from**  $\pi_{\text{data}}$  (i.e. find a proba  $\pi_{\infty}$  and a simulable kernel  $Q$  such that  $\pi_{\text{data}} \simeq \pi_{\infty} Q$ ).

Complex data distribution  $\pi_{\text{data}}$



Easy-to-sample distribution  $\pi_{\infty}$

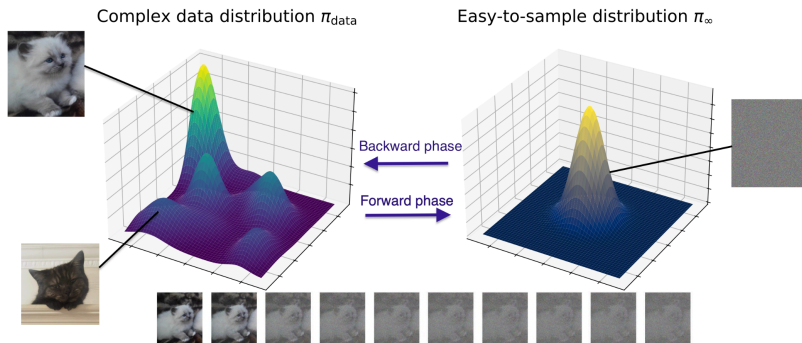


$\pi_{\infty} Q$



# SGMs Philosophy.

- ▶ “Creating noise from data is easy; creating data from noise is generative modeling.” (Song et al., 2021)



# Table of Contents

## **1. Introduction to SGMs**

- 1.1 Forward-backward SDEs
- 1.2 SGMs in practice.

## **2. Diffusion on an extended phase space.**

- 2.1 Can we make other choices of forward ?
- 2.2 Turn CLD into a generative model.

## **3. Theoretical result**

- 3.1 Convergence results for VP and VE processes.
- 3.2 Convergence of CLD.

# Table of Contents

## 1. Introduction to SGMs

- 1.1 Forward-backward SDEs
- 1.2 SGMs in practice.

## 2. Diffusion on an extended phase space.

- 2.1 Can we make other choices of forward ?
- 2.2 Turn CLD into a generative model.

## 3. Theoretical result

- 3.1 Convergence results for VP and VE processes.
- 3.2 Convergence of CLD.

# From data to noise: the forward process.

- ▶  $(\vec{X}_t)_{t \in [0, T]}$  is solution to an **Ornstein–Uhlenbeck** process:

$$d\vec{X}_t = -\vec{X}_t dt + \sqrt{2} dB_t, \quad \vec{X}_0 \sim \pi_{\text{data}}.$$

- ▶ If unfamiliar with SDEs: limit of a discrete-time process given by

$$X_{k+h} = \sqrt{1-2h} X_k + \sqrt{2h} Z_k, \quad Z_k \sim \mathcal{N}(0, I_d), \quad h \rightarrow 0.$$

- ▶ **Intuition:** destroys signal via Gaussian noise and rescaling.

## SGMs through SDE: more on the forward process.

- ▶ The noising procedure implies a scaling down of the data points  $d\vec{X}_t = -\vec{X}_t dt$ ,

## SGMs through SDE: more on the forward process.

- ... and a Gaussian noising process  $d\vec{X}_t = \sqrt{2}dB_t$ ,



SGMs through SDE: more on the forward process.

# From noise to data: the backward process.

- ▶ This forward process admits a **time-reversed process** (Anderson, 1982; Cattiaux et al., 2021), i.e.

$$\left(\overleftarrow{X}_t\right)_{t \in [0, T]} \stackrel{\mathcal{L}}{=} \left(\overrightarrow{X}_{T-t}\right)_{t \in [0, T]}$$

with,

$$d\overleftarrow{X}_t = \left( \overleftarrow{X}_t + 2 \underbrace{\nabla \log p_{T-t}(\overleftarrow{X}_t)}_{\text{score function}} \right) dt + \sqrt{2} dB_t, \quad \overleftarrow{X}_0 \sim p_T.$$

with  $p_t$  the p.d.f. of  $\overrightarrow{X}_t$ .

- ▶ The score term drives the backward process towards **regions of high probability**.
- ▶ This is (**almost**) a **generative model**:  $\overleftarrow{X}_T \sim \pi_{\text{data}}$ .

# Table of Contents

## 1. Introduction to SGMs

- 1.1 Forward-backward SDEs
- 1.2 SGMs in practice.

## 2. Diffusion on an extended phase space.

- 2.1 Can we make other choices of forward ?
- 2.2 Turn CLD into a generative model.

## 3. Theoretical result

- 3.1 Convergence results for VP and VE processes.
- 3.2 Convergence of CLD.

# SGMs in Practice I: mixing time.

- ▶ Let  $Q_t$  be the semigroup of  $\overleftarrow{X}_t$ :

$$Q_t(x, dy) = \mathbb{P} \left( \overleftarrow{X}_t \in dy \mid \overleftarrow{X}_0 = x \right) .$$

- ▶ **Time-reversal holds when  $\overleftarrow{X}_0 \sim p_T$ , i.e.**

$$\pi_{\text{data}} = p_T Q_T .$$

- ▶ But  $p_t$  depends on  $\pi_{\text{data}}$ :

$$p_t(x_t) = \int_{\mathbb{R}^d} \underbrace{p_t(x_t | x_0)}_{\text{p.d.f. of } \overrightarrow{X}_t | \overrightarrow{X}_0} \pi_{\text{data}}(dx_0) .$$

- ▶ In practice, one wants an **independent** and **easy-to-sample** probability  $\pi_\infty$  to initialize the generative model.

# SGMs in Practice I: mixing time.

- ▶ 💡 leverage the ergodicity of the O-U kernel.
- ▶ **Forward process** admits time marginal with  $Z \sim \mathcal{N}(0, I_d)$  and  $Z \perp X_0$ :

$$\vec{X}_t = e^{-t} \vec{X}_0 + \sqrt{1 - e^{-2t}} Z$$

- ▶ For  $T$  large, the initial conditions are forgotten:

$$p_T \approx \pi_\infty \sim \mathcal{N}(0, I_d) .$$



**Mixing Time Error:**  $\pi_{\text{data}} \simeq \pi_\infty Q_T$

## SGMs in practice II: learn the score function.

- ▶ The backward process depends on the score function  $\nabla \log p_t(x)$ .
- ▶ The forward process marginals can be sampled exactly.
- ▶ Train a **deep neural network**  $s_\theta : [0, T] \times \mathbb{R}^d \mapsto \mathbb{R}^d$  to minimize:

$$\mathcal{L}_{\text{naive}}(\theta) = \mathbb{E} \left[ \left\| s_\theta \left( \tau, \vec{X}_\tau \right) - \nabla \log p_\tau \left( \vec{X}_\tau \right) \right\|^2 \right],$$

with  $\tau \sim \mathcal{U}(0, T)$  independent of the forward process  $(\vec{X}_t)_{t \geq 0}$ .

- ▶ But  $p_\tau(x)$  is **unknown** !

## SGMs in practice II: learn the score function.

- ▶ 💡 its **conditional version** shares the same optimum (Hyvärinen and Dayan, 2005; Vincent, 2011):

$$\mathcal{L}_{\text{score}}(\theta) = \mathbb{E} \left[ \left\| s_{\theta} \left( \tau, \vec{X}_{\tau} \right) - \nabla \log p_{\tau} \left( \vec{X}_{\tau} | \vec{X}_0 \right) \right\|^2 \right] .$$

- ▶ The conditional score is explicit:

$$\nabla \log p_{\tau}(\vec{X}_{\tau} | \vec{X}_0) = \frac{m_{\tau} \vec{X}_0 - \vec{X}_{\tau}}{\sigma_{\tau}^2} = -\frac{Z}{\sigma_{\tau}}$$

with  $m_{\tau} = e^{-\tau}$  and  $\sigma_{\tau} = \sqrt{1 - m_{\tau}^2}$ .

- ▶ Score matching Neural Networks writes as,

$$\mathcal{L}_{\text{score}}(\theta) = \mathbb{E} \left[ \left\| s_{\theta} \left( \tau, \vec{X}_{\tau} \right) + \frac{Z}{\sigma_{\tau}} \right\|^2 \right] .$$



**Approximation error:**  $\pi_{\text{data}} \approx \pi_{\infty} Q_T^{\theta}$

## SGMs in practice III: simulate from the backward kernel.

- ▶ The backward drift is **non-linear**: non-Gaussian.
- ▶ 💡 discretize  $[0, T]$  in  $N$  steps with  $t_k = kh$ ,  $h = T/N$ .
- ▶ Euler–Maruyama discretization:

$$\bar{X}_{t_{k+1}} = \bar{X}_{t_k} + h \left( \bar{X}_{t_k} + 2 s_\theta(T - t_k, \bar{X}_{t_k}) \right) + \sqrt{2h} Z_k$$

- ▶ Other approaches preserving the time marginals exist (e.g. ODE sampling).

⚠ **Discretization error:**  $\pi_{\text{data}} \approx \pi_\infty Q_{T,N}^\theta := \hat{\pi}_{\infty,N}^\theta$



# Table of Contents

## 1. Introduction to SGMs

- 1.1 Forward-backward SDEs
- 1.2 SGMs in practice.

## 2. Diffusion on an extended phase space.

- 2.1 Can we make other choices of forward ?
- 2.2 Turn CLD into a generative model.

## 3. Theoretical result

- 3.1 Convergence results for VP and VE processes.
- 3.2 Convergence of CLD.

# Draw inspiration from MCMC.

- ▶ In **sampling**, one wants to sample from  $\pi \propto e^{-U}$ .
  - ▶ When  $U : \mathbb{R}^d \rightarrow \mathbb{R}$  is smooth and typically strongly convex,

$$dX_t = -\nabla U(X_t)dt + \sqrt{2}dB_t$$

admits  $\pi$  as invariant measure.

- ▶ Sampling can be done by discretization (ULA) or accept-reject corrections (MALA).
- ▶ This can be extended to a **kinetic setting**:

$$d \begin{pmatrix} X_t \\ V_t \end{pmatrix} = \begin{pmatrix} V_t \\ -(V_t + \nabla U(X_t)) \end{pmatrix} dt + \sqrt{2} \begin{pmatrix} 0 \\ 1 \end{pmatrix} dB_t$$

- ▶ Stationary distribution  $\pi(dx, dv) \propto e^{-U(x) - \frac{\|v\|^2}{2}} dx dv$ .

# Extending the Phase Space of SGMs.

- 💡 Augment data space with a **velocity component**  $(V_t)_{t \in [0, T]}$ .
- 💡  $X_t$  and  $V_t$  are coupled through *Hamiltonian-like* interactions.
- 💡 **Noise injection** only on the **velocity** component.

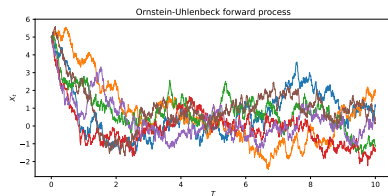
**Forward process:** for  $\vec{\mathbf{U}}_t = (\vec{X}_t, \vec{V}_t)^\top \in \mathbb{R}^2$  and  $B_t \in \mathbb{R}^2$ ,

$$d \begin{pmatrix} X_t \\ V_t \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ -1 & -2 \end{pmatrix} \begin{pmatrix} X_t \\ V_t \end{pmatrix} dt + \begin{pmatrix} 0 & 0 \\ 0 & \sigma \end{pmatrix} dB_t, \quad (X_0, V_0) \sim \pi_{\text{data}} \otimes \pi_v,$$

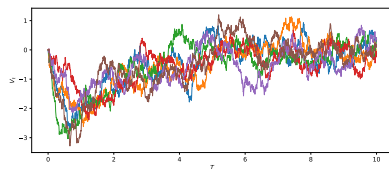
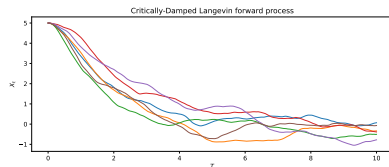
where  $\pi_v \sim \mathcal{N}(0, v^2)$ . We use compact matrix notation,

$$d\vec{\mathbf{U}}_t = A\vec{\mathbf{U}}_t dt + \Sigma dB_t, \quad \vec{\mathbf{U}}_0 \sim \pi_{\text{data}} \otimes \pi_v.$$

# Noising process comparison.



**OU process**



**CLD process (position & velocity)**

# Table of Contents

## 1. Introduction to SGMs

- 1.1 Forward-backward SDEs
- 1.2 SGMs in practice.

## 2. Diffusion on an extended phase space.

- 2.1 Can we make other choices of forward ?
- 2.2 Turn CLD into a generative model.

## 3. Theoretical result

- 3.1 Convergence results for VP and VE processes.
- 3.2 Convergence of CLD.

# What makes a forward SDE a generative model?

- ▶ A well-chosen **noising process** must satisfy three properties:
  1. **Interpolation:** transforms the data distribution  $\pi_{\text{data}}$  into an easy-to-sample prior  $\pi_{\infty}$ .
  2. **Learnability:** time-dependent score functions  $\nabla \log p_t(\cdot)$  can be learned.
  3. **Efficient sampling:** marginals can be efficiently simulated.
- ▶ Examples: Variance-Preserving (VP), Variance-Exploding (VE), flow matching and **Critically-Damped Langevin (CLD) diffusions**.

# 1. Interpolate between the data distribution and a prior.

- ▶ The forward process evolves in the **extended phase space**  $\vec{\mathbf{U}}_t = (\vec{X}_t, \vec{V}_t)^\top \in \mathbb{R}^2$  as

$$\vec{\mathbf{U}}_t = e^{tA} \vec{\mathbf{U}}_0 + \int_0^t e^{(t-s)A} \Sigma dB_s. \quad (1)$$

and converges to  $\pi_\infty \sim \mathcal{N}(0_2, \Sigma_\infty)$ .

- ▶ Time reversal property applies on the extended space, *i.e.*

$$(\vec{X}_t, \vec{V}_t)_{t \in [0, T]} = (\overleftarrow{X}_{T-t}, \overleftarrow{V}_{T-t})_{t \in [0, T]}$$

- ▶ leading to the backward SDE:

$$d\overleftarrow{\mathbf{U}}_t = -A\overleftarrow{\mathbf{U}}_t dt + \Sigma^2 \nabla \log p_{T-t}(\overleftarrow{\mathbf{U}}_t) dt + \Sigma dB_t,$$

with  $p_t(x, v)$  de p.d.f of (1).

## 2. Score function can be learned

- ▶ As before  $s_\theta$  trained to learn conditional score but on the whole phase-space state  $\vec{\mathbf{U}}_t = (\vec{X}_t, \vec{V}_t)$ :

$$\mathcal{L}_{\text{DSM}}(\theta) = \mathbb{E} \left[ \left\| s_\theta(t, \vec{\mathbf{U}}_t) - \nabla \log p_t(\vec{\mathbf{U}}_t \mid \vec{\mathbf{U}}_0) \right\|^2 \right] .$$

- ▶ However, we know that  $\vec{V}_0 \sim \mathcal{N}(0, v^2)$ , so we can marginalize  $\vec{\mathbf{U}}_0 = (\vec{X}_0, \vec{V}_0)^\top$  over  $V_0$ , leading to a **closed-form expression** of  $\nabla \log p_t(\vec{\mathbf{U}}_t \mid \vec{X}_0)$ :

$$\mathcal{L}_{\text{HSM}}(\theta) = \mathbb{E} \left[ \left\| s_\theta(t, \vec{\mathbf{U}}_t) - \nabla \log p_t(\vec{\mathbf{U}}_t \mid \vec{X}_0) \right\|^2 \right] ,$$

yielding **more stable training** objective.



### 3. Marginals can be sampled.

Different numerical schemes:

- ▶ Euler–Maruyama (standard baseline);
- ▶ Symplectic integrators design for position-velocity state-spaces.

Combining all this leads to **better numerical performance** (Dockhorn et al., 2022):

Table 1: Unconditional CIFAR-10 generative performance.

Class	Model	NLL↓	FID↓
Score	CLD-SGM (Prob. Flow) ( <i>ours</i> )	≤3.31	2.25
	CLD-SGM (SDE) ( <i>ours</i> )	-	2.23
Score	DDPM++, VPSDE (Prob. Flow) (Song et al., 2021c)	3.13	3.08
	DDPM++, VPSDE (SDE) (Song et al., 2021c)	-	2.41
	DDPM++, sub-VP (Prob. Flow) (Song et al., 2021c)	2.99	2.92
	DDPM++, sub-VP (SDE) (Song et al., 2021c)	-	2.41
	NCSN++, VESDE (SDE) (Song et al., 2021c)	-	2.20
	LSGM (Vahdat et al., 2021)	≤3.43	2.10
	LSGM-100M (Vahdat et al., 2021)	≤2.96	4.60
	DDPM (Ho et al., 2020)	≤3.75	3.17
	NCSN (Song & Ermon, 2019)	-	25.3
	Adversarial DSM (Jolicœur-Martineau et al., 2021b)	-	6.10
	Likelihood SDE (Song et al., 2021b)	2.84	2.87
	DDIM (100 steps) (Song et al., 2021a)	-	4.16
	FastDDPM (100 steps) (Kong & Ping, 2021)	-	2.86
	Improved DDPM (Nichol & Dhariwal, 2021)	3.37	2.90
	VDM (Kingma et al., 2021)	≤2.49	7.41 (4.00)
	UDM (Kim et al., 2021)	3.04	2.33
	D3PM (Austin et al., 2021)	≤3.44	7.34
	Gotta Go Fast (Jolicœur-Martineau et al., 2021a)	-	2.44
	DDPM Distillation (Luhman & Luhman, 2021)	-	9.36

# Table of Contents

## 1. Introduction to SGMs

- 1.1 Forward-backward SDEs
- 1.2 SGMs in practice.

## 2. Diffusion on an extended phase space.

- 2.1 Can we make other choices of forward ?
- 2.2 Turn CLD into a generative model.

## 3. Theoretical result

- 3.1 Convergence results for VP and VE processes.
- 3.2 Convergence of CLD.

# A variety of convergence results

- ▶ Assume score function is appropriately learned e.g.

$$\|s_{\theta}(t, \mathbf{U}_t) - \nabla \log p_t(\mathbf{U}_t)\|_{L_2} \leq M$$

where the expectation is taken under some appropriately chosen stochastic process.

- ▶ Using this framework a **variety of upper bounds** to the **distance between the data distribution and the generated distribution**  $d(\pi_{\text{data}}, \hat{\pi})$  have been established for various metrics:
  - ▶ For the total variation distance and Kullback-Leibler divergence: De Bortoli et al. (2021); Conforti et al. (2023); Bortoli et al. (2023); Chen et al. (2023); Chen (2023).
  - ▶ For the Wasserstein distance: Lee et al. (2022, 2023); Bruno et al. (2023); Gao et al. (2023).

## Wasserstein-2: upper bounds.

- ▶ The  $\mathcal{W}_2$  distance is defined as

$$\mathcal{W}_2^2(\pi_{\text{data}}, \hat{\pi}_{\infty, N}^{\theta}) = \inf \left\{ \mathbb{E} \left[ \left\| \vec{X}_0 - \bar{X}_{\infty, N}^{\theta} \right\|^2 \right], \vec{X}_0 \sim \pi_{\text{data}}, \bar{X}_{\infty, N}^{\theta} \sim \hat{\pi}_{\infty, N}^{\theta} \right\}$$

- ▶ Control the errors already presented:

$$\begin{aligned} \mathcal{W}_2(\pi_{\text{data}}, \hat{\pi}_{\infty, N}^{\theta}) &\leq \underbrace{\mathcal{W}_2(\mathcal{L}(\bar{X}_T), \mathcal{L}(\bar{X}_N))}_{\text{Discretization}} + \underbrace{\mathcal{W}_2(\mathcal{L}(\bar{X}_N), \mathcal{L}(\bar{X}_{\infty, N}))}_{\text{Mixing time}} \\ &\quad + \underbrace{\mathcal{W}_2(\mathcal{L}(\bar{X}_{\infty, N}), \mathcal{L}(\bar{X}_{\infty, N}^{\theta}))}_{\text{Score approx.}} \\ &\leq e^{-T} c_1 + M c_2 + \sqrt{h} c_3, \end{aligned}$$

with  $T > 0$  the diffusion time,  $M$  the score approximation quality and  $h = T/N$  the discretization step size.

## Backward contraction for O.U forward I.

- Proof relies **mostly** on **contraction for Euclidean norm**:

$$\mathcal{W}_2^2(\pi_{\text{data}}, \hat{\pi}_{\infty, N}^{\theta}) \leq \left\| \vec{X}_0 - \vec{X}_{\infty, N}^{\theta} \right\|_{L_2}^2$$

- Fix  $x, y \in (\mathbb{R})^2$ :

$$d\overleftarrow{X}_t^x = \left( \overleftarrow{X}_t^x + 2\nabla \log p_{T-t}(\overleftarrow{X}_t^x) \right) dt + \sqrt{2}dB_t, \quad \overleftarrow{X}_0 = x \text{ p.s.}$$

$$d\overleftarrow{X}_t^y = \left( \overleftarrow{X}_t^y + 2\nabla \log p_{T-t}(\overleftarrow{X}_t^y) \right) dt + \sqrt{2}dB_t, \quad \overleftarrow{X}_0 = y \text{ p.s.}$$

- Consider a synchronous coupling and introduce the difference ODE  $Z_t = \overleftarrow{X}_t^x - \overleftarrow{X}_t^y$ , which satisfies

$$dZ_t = Z_t + 2 \underbrace{\left( \nabla \log p_{T-t}(\overleftarrow{X}_t^x) - \nabla \log p_{T-t}(\overleftarrow{X}_t^y) \right)}_{:=\Delta_t} dt$$

- and study

$$\frac{d}{dt} \|Z_t\|^2 = 2Z_t^\top (dZ_t) = 2 \left( \|Z_t\|^2 + 2 \langle Z_t, \Delta_t \rangle \right).$$

## Backward contraction for O.U forward II.

- ▶ If  $p_{T-t}$  is  $\lambda$ -**log-concave**, then, there exists  $\lambda > 0$ , such that
  - ▶  $\langle Z_t, \nabla \log p_{T-t}(\overleftarrow{X}_t^x) - \nabla \log p_{T-t}(\overleftarrow{X}_t^y) \rangle \leq -\lambda \|Z_t\|^2$  ;
  - ▶  $\nabla^2 \log p_{T-t} \preccurlyeq -\lambda I_d$ .
- ▶ Therefore, using Grönwall inequality,

$$\begin{aligned} \frac{d}{dt} \|Z_t\|^2 &\leq 2(1 - 2\lambda) \|Z_t\|^2 \\ &\leq e^{2(1-2\lambda)t} \|Z_0\|^2. \end{aligned}$$

- ▶ **Takeaway:** strong log-concavity ( $\lambda > 1/2$ ) **gives contraction** for  $\|\cdot\|$  which implies  $\mathcal{W}_2$  contraction.

# Table of Contents

## 1. Introduction to SGMs

- 1.1 Forward-backward SDEs
- 1.2 SGMs in practice.

## 2. Diffusion on an extended phase space.

- 2.1 Can we make other choices of forward ?
- 2.2 Turn CLD into a generative model.

## 3. Theoretical result

- 3.1 Convergence results for VP and VE processes.
- 3.2 Convergence of CLD.

# Log-concavity is not enough for CLD

$$d\overleftarrow{\mathbf{U}}_t = -A\overleftarrow{\mathbf{U}}_t dt + \Sigma^2 \nabla \log p_{T-t}(\overleftarrow{\mathbf{U}}_t) dt + \Sigma dB_t,$$

- Consider a synchronous coupling and study the stability of the difference ODE  $\mathbf{Z}_t = \overleftarrow{\mathbf{U}}_t^x - \overleftarrow{\mathbf{U}}_t^y$ , which satisfies

$$\frac{d}{dt}(\|\mathbf{Z}_t\|^2) = -2\mathbf{Z}_t^\top A \mathbf{Z}_t + 2\mathbf{Z}_t^\top \Sigma^2 H_t \mathbf{Z}_t.$$

with

$$\Sigma = \begin{pmatrix} 0 & 0 \\ 0 & \sigma \end{pmatrix}$$

where we used the mean value theorem with

$$H_t = \nabla^2 \log p_{T-t} = \begin{pmatrix} H_t^{xx} & H_t^{xv} \\ H_t^{vx} & H_t^{vv} \end{pmatrix}.$$



## Log-concavity is not enough for CLD (counterexample).

Take  $p_t(x, v)$  that is  $(1 - c)$ -**log-concave** with,

$$H_t = \begin{pmatrix} -1 & c \\ c & -1 \end{pmatrix}, \quad 0 < c < 1; \quad \text{Spec}(H_t) = -1 \pm c,$$

Then with  $\Sigma = \text{diag}(0, \sigma)$  and  $\mathbf{Z} = (Z_x, Z_v)$ ,

$$2\mathbf{Z}^\top \Sigma^2 H_t \mathbf{Z} = 2\sigma^2(-c Z_x Z_v - Z_v^2).$$

Choose  $Z_x > 0$ ,  $Z_v < 0$  with ratio  $|Z_x|/|Z_v| > 1/c$ . Then the RHS becomes **positive**.

**Takeaway:** even though  $p_t$  is log-concave, the projected curvature  $\Sigma^2 H_t$  is *not* negative semidefinite. Uniform contraction is hopeless.

# Solution 1: Long-term regularity of the renormalized score

**Idea.** Introduce a *renormalized* formulation of the backward process:

$$d\overleftarrow{\mathbf{U}}_t = \tilde{A}\overleftarrow{\mathbf{U}}_t dt + \Sigma^2 \nabla \log \tilde{p}_{T-t}(\overleftarrow{\mathbf{U}}_t) dt + \Sigma dB_t, \quad \tilde{p}_t := \frac{p_t}{p_\infty}.$$

**Key properties.**

1.  $\tilde{A}$  is **negative definite**.
2.  $\tilde{p}_t$  "quantifies" **deviation from equilibrium**  $p_\infty$ .
3. Its curvature  $\nabla^2 \log \tilde{p}_t$  characterizes the **regularity of the score**.

# Structure & regularity assumptions on $p_{\text{data}}$

## Finite relative Fisher information

$$\mathcal{I}(p_{\text{data}} \mid p_{\infty}) = \int_{\mathbb{R}^d} \left\| \nabla \log \frac{p_{\text{data}}}{p_{\infty}}(x) \right\|^2 p_{\text{data}}(x) dx < \infty .$$

## Log-Lipschitz perturbation of a strongly log-concave base

$$p_{\text{data}}(x) \propto \exp \left( - [V(x) + H(x)] \right) ,$$

with for all  $x, y \in (\mathbb{R}^d)^2$ :

- ▶  $\exists \alpha > 0$  such that  $\alpha I_d \preceq \nabla^2 V(x)$  ;
- ▶  $|H(x) - H(y)| \leq L \|x - y\|$  .

## One-sided Lipschitz score

$$-(\nabla \log p_{\text{data}}(x) - \nabla \log p_{\text{data}}(y))^{\top} (x - y) \leq L_0 \|x - y\|^2, \quad \forall x, y \in \mathbb{R}^d .$$

# Regularity of the renormalized score

- ▶ Under the previous hypotheses, there exists a constant  $C > 0$  such that, for all  $t \in (0, T]$ ,

$$\|\nabla^2 \log \tilde{p}_t(\cdot)\| \leq C \left(1 + \frac{1}{\sqrt{t}}\right) e^{-2t} = \tilde{L}_t.$$

- ▶ **Interpretation.**
  - ▶ **Short times** ( $t \rightarrow 0$ ): the singularity in  $1/\sqrt{t}$  remains integrable.
  - ▶ **Long times** ( $t \rightarrow \infty$ ): exponential decay.
- ▶ **Takeaway:** The renormalized score function **regularizes over time**.

## Solution 1: Long-term regularity of the renormalized score

► Informally, there exists  $\mathfrak{M}$  PSD matrix and  $\eta > 0$  such that

$$\begin{aligned}\frac{d}{dt} \|\mathbf{Z}_t\|_{\mathfrak{M}}^2 &\leq -2\mathbf{Z}_t^\top \mathfrak{M} \mathbf{A} \mathbf{Z}_t + 2\mathbf{Z}_t^\top \mathfrak{M} \Sigma^2 \left( \nabla \log p_{T-t} \left( \overleftarrow{\mathbf{U}}_t^x \right) - \nabla \log p_{T-t} \left( \overleftarrow{\mathbf{U}}_t^y \right) \right) \\ &\leq 2(-\eta + \sigma^2 \tilde{L}_t) \|\mathbf{Z}_t\|_{\mathfrak{M}}^2 .\end{aligned}$$

Using Grönwall's lemma, there exists  $C > 0$ , such that,

$$\begin{aligned}\|\mathbf{Z}_t\|_{\mathfrak{M}}^2 &\leq e^{-2\eta t + \sigma^2 \int_0^t \tilde{L}_s ds} \|\mathbf{Z}_0\|_{\mathfrak{M}}^2 \\ &\leq C e^{-2\eta t} \|\mathbf{Z}_0\|_{\mathfrak{M}}^2 ,\end{aligned}$$

 **Contraction!**

## Final $\mathcal{W}_2$ upper bound

- From contraction in the extended phase space, the three sources of errors (mixing, approximation, and discretization) can be controlled jointly:

$$\mathcal{W}_2\left(\pi_{\text{data}} \otimes \pi_v, \mathcal{L}\left(\bar{\mathbf{U}}_T^\theta\right)\right) \leq c_1 e^{-c_2 T} \mathcal{W}_2\left(\pi_{\text{data}} \otimes \pi_v, \pi_\infty\right) + c_1 \sigma^2 M + c_1 \sqrt{h}.$$

- Projecting onto the position component  $X$  ( $P_X(x, v) = x$ ) preserves the  $\mathcal{W}_2$  distance, since  $P_X$  is 1-Lipschitz:

$$\mathcal{W}_2\left(\pi_{\text{data}}, \mathcal{L}\left(\bar{X}_T^\theta\right)\right) \leq \mathcal{W}_2\left(\pi_{\text{data}} \otimes \pi_v, \mathcal{L}\left(\bar{\mathbf{U}}_T^\theta\right)\right).$$

## Solution 2: restore ellipticity

**Idea.** Inject a small amount of noise on *all* coordinates:

$$\Sigma = \begin{pmatrix} \varepsilon & 0 \\ 0 & \sigma \end{pmatrix}, \quad \varepsilon > 0.$$

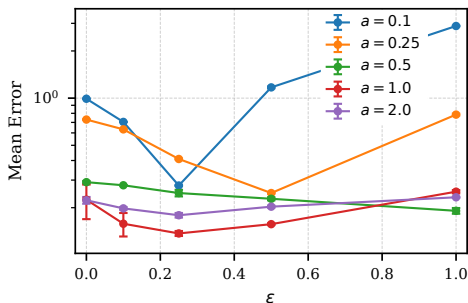
### Consequences.

- ▶ **Uniform ellipticity:** (multi-dimensional O.U. structure).
- ▶ **More quantitative bounds :** standard log-concave tools apply.
- ▶ **Practice:**  $\varepsilon$  provides a new parameters to control the regularity of the sample paths.

## Solution 2: Numerical aspects

**Empirics (Funnel dataset,  $d = 100$ ).**

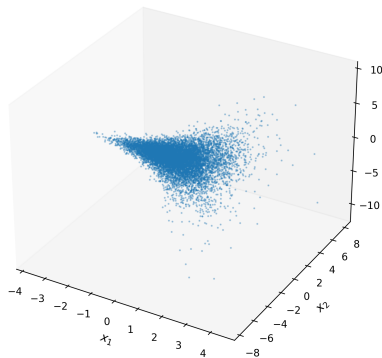
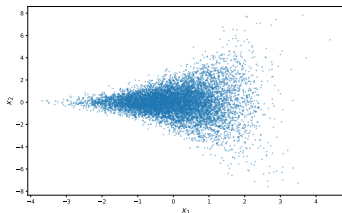
- ▶ Small  $\varepsilon$  often **improves** sliced- $\mathcal{W}_2$  vs.  $\varepsilon = 0$  (CLD baseline).
- ▶ **Trade-off:** slight sensitivity to other hyperparameters.



Mean  $\mathcal{W}_2$  over 5 runs; error bars represent  $\pm$  one standard deviation.



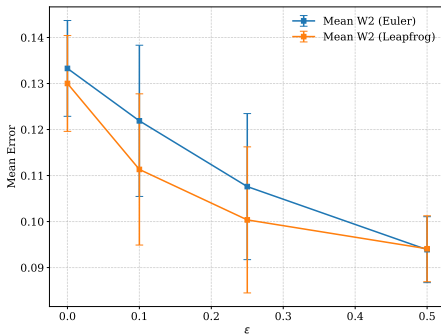
# Funnel distribution scatter plot



**Figure:** 10 000 samples from a funnel distribution in dimension 50. Plot of the 1st and 2nd dimension (left) and plot of the 1st, 2nd and 3rd dimension (right).

## Solution 2: Numerical aspects

- ▶ Even with a small  $\varepsilon > 0$ , **structure-preserving integrators** can further improve performance.
- ▶ **But:** higher computational cost — the network must learn full gradients  $\nabla \log p_t(x, v)$  instead of velocity-only terms  $\nabla_v \log p_t(v)$ , **doubling the effective dimension**.



ooo

- B. D. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- V. D. Bortoli, J. Thornton, J. Heng, and A. Doucet. Diffusion schrödinger bridge with applications to score-based generative modeling, 2023.
- S. Bruno, Y. Zhang, D.-Y. Lim, Ö. D. Akyildiz, and S. Sabanis. On diffusion-based generative models and their error bounds: The log-concave case with full convergence estimates. *arXiv preprint arXiv:2311.13584*, 2023.
- P. Cattiaux, G. Conforti, I. Gentil, and C. Léonard. Time reversal of diffusion processes under a finite entropy condition. *arXiv preprint arXiv:2104.07708*, 2021.
- S. Chen, S. Chewi, J. Li, Y. Li, A. Salim, and A. R. Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions, 2023.
- T. Chen. On the importance of noise scheduling for diffusion models. *arXiv preprint arXiv:2301.10972*, 2023.

- G. Conforti, A. Durmus, and M. G. Silveri. Score diffusion models without early stopping: finite fisher information is all you need, 2023.
- V. De Bortoli, J. Thornton, J. Heng, and A. Doucet. Diffusion schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34:17695–17709, 2021.
- T. Dockhorn, A. Vahdat, and K. Kreis. Score-based generative modeling with critically-damped langevin diffusion. In *International Conference on Learning Representations (ICLR)*, 2022.
- X. Gao, H. M. Nguyen, and L. Zhu. Wasserstein convergence guarantees for a general class of score-based generative models, 2023.
- A. Hyvärinen and P. Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- H. Lee, J. Lu, and Y. Tan. Convergence for score-based generative modeling with polynomial complexity. *Advances in Neural Information Processing Systems*, 35:22870–22882, 2022.

- H. Lee, J. Lu, and Y. Tan. Convergence of score-based generative modeling for general data distributions. In *International Conference on Algorithmic Learning Theory*, pages 946–985. PMLR, 2023.
- Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations (ICLR)*, 2021.
- P. Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011. doi: 10.1162/NECO\_a\_00142.