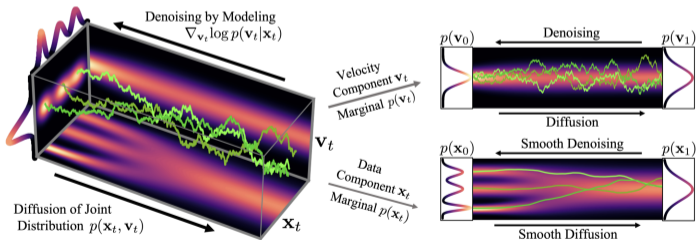


Wasserstein-2 Convergence of Critically Damped Diffusion Models

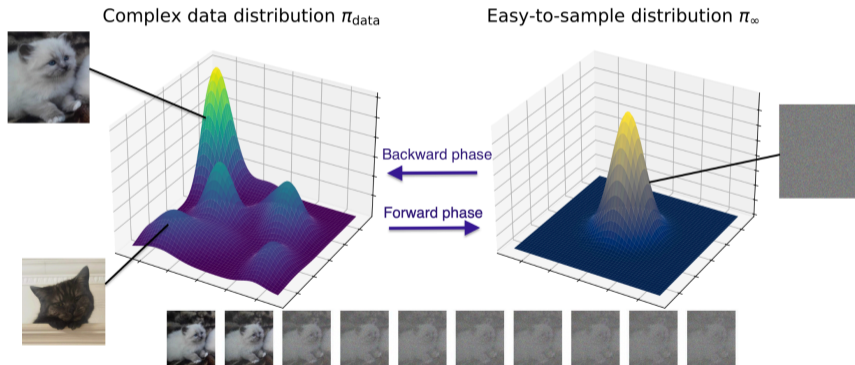
SGM errors, kinetic extension and stability.

Stanislas Strasman, S. Surendran, C. Boyer, S. Le Corff, V. Lemaire, A. Ocello.



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 π_{data} .
- Goal: **generate new samples from** π_{data} .



The Time-Reversal Perspective

Let p_t be the density of the forward process

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

Time reversal of diffusions

Fix $T > 0$, the backward process

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

is again a diffusion (Anderson, 1982) and solves

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

Ideal generative identity

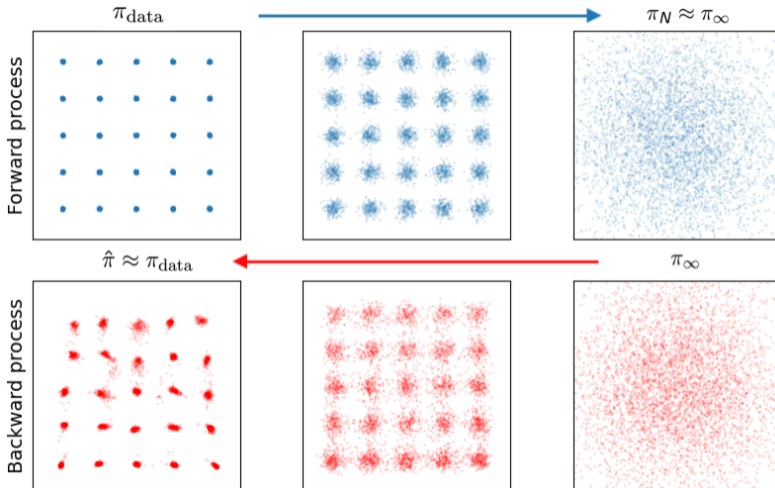
Let Q_T be the Markov transition kernel of the reverse process:

$$Q_T(x, dy) = \mathbb{P}(\overleftarrow{X}_T \in dy \mid \overleftarrow{X}_0 = x).$$

Then time reversal gives

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

Synthetic example: 2-dimensional mixture of 25 Gaussian.



SGMs in practice: three approximations



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

is almost generative, but ...

1. Initialization

p_T depends on the data. The exact terminal law is

$$\vec{X}_T \stackrel{\mathcal{L}}{=} e^{-T} \vec{X}_0 + \sqrt{1 - e^{-2T}} Z$$

For large T ,

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

$$\pi_{\text{data}} \approx \pi_\infty Q_T$$

2. Score approximation

The reverse drift uses the unknown score

$$\nabla \log p_t(x)$$

In practice, learn

$$s_\theta(x, t) \approx \nabla \log p_t(x)$$

$$\pi_{\text{data}} \approx \pi_\infty Q_T^\theta$$

3. Discretization

The reverse SDE should be discretized. Choose

$$h = \frac{T}{N}, \quad t_k = kh.$$

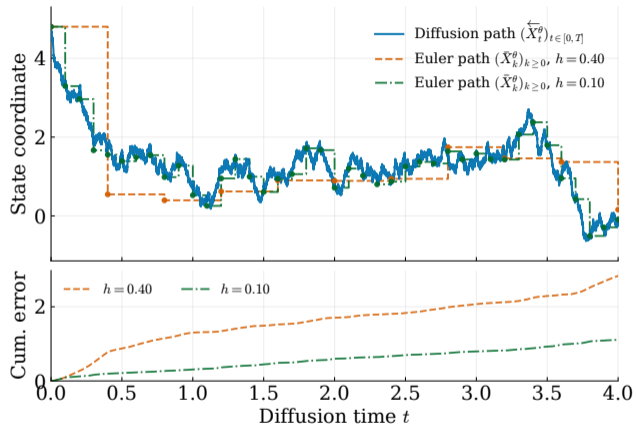
and apply Euler-Maruyama.

$$\pi_{\text{data}} \approx \pi_\infty \bar{Q}_{T,N}^\theta$$

$$\bar{X}_{k+1}^{\theta,N} = \bar{X}_k^{\theta,N} + h \left(\bar{X}_k^{\theta,N} + 2s_\theta(\bar{X}_k^{\theta,N}, T - t_k) \right) + \sqrt{2h} Z_{k+1}, \quad \bar{X}_0^{\theta,N} \sim \pi_\infty, \quad Z_{k+1} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$$

$$\hat{\pi}_{\infty,N}^\theta := \pi_\infty \bar{Q}_{T,N}^\theta$$

1-d pathwise intuition



All discretizations are coupled with the same Brownian path and $h = T/N$.

2-Wasserstein Stability Theory

The \mathcal{W}_2 distance

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

Score approximation assumption

Assume that,

$$\sup_{0 \leq k \leq N-1} \left\| \nabla \log p_{T-t_k}(\bar{X}_k^{N, \theta}) - s_{\theta}(T-t_k, \bar{X}_k^{N, \theta}) \right\|_{L^2} \leq M.$$

Typical \mathcal{W}_2 bound

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

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

Extending the Phase Space: CLD Noising

Kinetic idea

Augment the data space with a velocity:

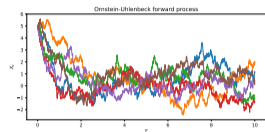
$$\vec{U}_t := \begin{pmatrix} \vec{X}_t \\ \vec{V}_t \end{pmatrix} \in \mathbb{R}^d \times \mathbb{R}^d, \quad (\vec{X}_0, \vec{V}_0) \sim \pi_{\text{data}} \otimes \mathcal{N}(0, I_d).$$

- Noise is injected only through the velocity.
- X_t and V_t are coupled by the drift.
- Hypocoelliptic structure !

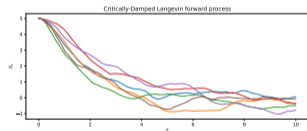
Critically damped Langevin noising

$$d \begin{pmatrix} \vec{X}_t \\ \vec{V}_t \end{pmatrix} = \underbrace{\begin{pmatrix} 0 & I_d \\ -I_d & -2I_d \end{pmatrix}}_A \begin{pmatrix} \vec{X}_t \\ \vec{V}_t \end{pmatrix} dt + \underbrace{\begin{pmatrix} 0 & 0 \\ 0 & 2I_d \end{pmatrix}}_\Sigma dB_t.$$

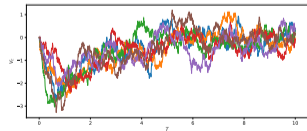
Noising process comparison



Standard OU noising



CLD position marginal X_t



CLD velocity marginal V_t

What makes a forward SDE a generative model?

1. Interpolation

It transforms data into an easy-to-sample prior:

$$\pi_{\text{data}} \longrightarrow \pi_{\infty}.$$

2. Learnability

The reverse drift should be learnable from samples *i.e.*

$$\nabla \log p_t$$

should be learnable.

3. Sampling

The reverse dynamics should be numerically simulable.

CLD forward process

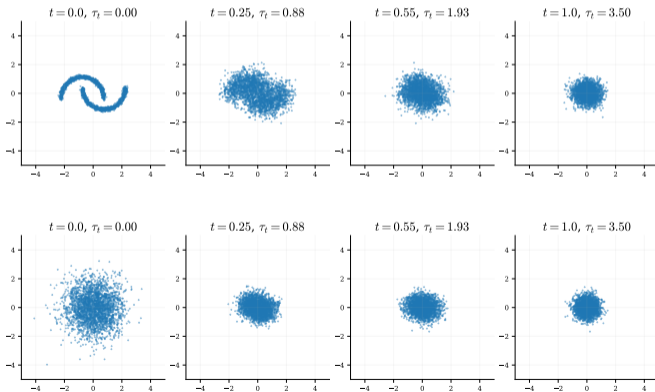
$$d\vec{\mathbf{U}}_t = A\vec{\mathbf{U}}_t dt + \Sigma dB_t, \quad \vec{\mathbf{U}}_0 \sim \pi_{\text{data}} \otimes \mathcal{N}(0, I_d).$$

$$\vec{\mathbf{U}}_t = \begin{pmatrix} \vec{X}_t \\ \vec{V}_t \end{pmatrix}, \quad A = \begin{pmatrix} 0 & I_d \\ -I_d & -2I_d \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 0 & 0 \\ 0 & 2I_d \end{pmatrix}.$$

1. Interpolate between the data distribution and a prior

- The dynamics is linear on the extended space:

$$\vec{U}_t = e^{tA} \vec{U}_0 + \int_0^t e^{(t-s)A} \Sigma dB_s, \quad p_T \approx \pi_\infty = \mathcal{N}(0, I_{2d}) \quad \text{for large } T.$$



2. The reverse process can be learned

Time reversal on phase space

Let p_t be the density of $\vec{\mathbf{U}}_t = (\vec{\mathbf{X}}_t, \vec{\mathbf{V}}_t) \in \mathbb{R}^{2d}$. Time reversal gives

$$(\overleftarrow{\mathbf{U}}_t)_{t \in [0, T]} \stackrel{\mathcal{L}}{=} (\vec{\mathbf{U}}_{T-t})_{t \in [0, T]},$$

and the reverse dynamics is

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

Learning the phase-space score

As in standard SGMs, train a score network on corrupted samples:

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

Since $V_0 \sim \mathcal{N}(0, I_d)$ is known, we may marginalize over V_0 and use the closed-form target

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

3. The reverse dynamics can be sampled

Numerical simulation

- Euler–Maruyama baseline.
- Splitting / kinetic integrators.

Takeaway

CLD keeps the SGM recipe, but changes the geometry of the reverse sampler. Empirically, this can improve generation quality.

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	
	Score	Adversarial DSM (Jolicoeur-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 (Jolicoeur-Martineau et al., 2021a)	-	2.44
		DDPM Distillation (Luhman & Luhman, 2021)	-	9.36

Empirical comparison from (Dockhorn, 2022).

Classical \mathcal{W}_2 contraction: reverse O.U. I

For the O.U. forward noising process, the reverse dynamics is

$$d\overleftarrow{X}_t = \left(\overleftarrow{X}_t + 2\nabla \log p_{T-t}(\overleftarrow{X}_t) \right) dt + \sqrt{2} dB_t.$$

Synchronous coupling

Fix $x, y \in \mathbb{R}^d$. Let \overleftarrow{X}_t^x and \overleftarrow{X}_t^y solve the same reverse SDE, driven by the same Brownian motion, with

$$\overleftarrow{X}_0^x = x, \quad \overleftarrow{X}_0^y = y.$$

Set

$$Z_t := \overleftarrow{X}_t^x - \overleftarrow{X}_t^y.$$

The Brownian noises cancel, hence

$$\frac{d}{dt} Z_t = Z_t + 2\nabla \log p_{T-t}(\overleftarrow{X}_t^x) - \nabla \log p_{T-t}(\overleftarrow{X}_t^y).$$

Therefore

$$\frac{d}{dt} \|Z_t\|^2 = 2\|Z_t\|^2 + 4 \left\langle Z_t, \nabla \log p_{T-t}(\overleftarrow{X}_t^x) - \nabla \log p_{T-t}(\overleftarrow{X}_t^y) \right\rangle.$$

Classical \mathcal{W}_2 contraction: reverse O.U. II

Uniform log-concavity assumption

Assume that p_{T-t} is uniformly λ -log-concave:

$$\langle x - y, \nabla \log p_{T-t}(x) - \nabla \log p_{T-t}(y) \rangle \leq -\lambda \|x - y\|^2.$$

Equivalently, if p_{T-t} is sufficiently regular,

$$\nabla^2 \log p_{T-t} \preceq -\lambda I_d.$$

Therefore,

$$\frac{d}{dt} \|Z_t\|^2 \leq 2(1 - 2\lambda) \|Z_t\|^2.$$

By Grönwall,

$$\|Z_t\|^2 \leq \exp\{-2(2\lambda - 1)t\} \|Z_0\|^2.$$

Consequence

If $\lambda > 1/2$, the reverse O.U. dynamics is contractive:

$$\mathcal{W}_2(\mu Q_t, \nu Q_t) \leq \exp\{-(2\lambda - 1)t\} \mathcal{W}_2(\mu, \nu),$$

where Q_t is the reverse Markov kernel.

Why the naive Euclidean contraction proof breaks for CLD

Consider two synchronously coupled reverse CLD processes and set $\mathbf{z}_t = \overleftarrow{\mathbf{U}}_t^x - \overleftarrow{\mathbf{U}}_t^y$. For the reverse dynamics

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

the Brownian noises cancel. By the mean-value theorem,

$$\nabla \log p_{T-t}(\overleftarrow{\mathbf{U}}_t^x) - \nabla \log p_{T-t}(\overleftarrow{\mathbf{U}}_t^y) = \bar{H}_t \mathbf{z}_t,$$

with \bar{H}_t an averaged Hessian of $\log p_{T-t}$. Hence


$$\frac{d}{dt} \|\mathbf{z}_t\|^2 = -2\mathbf{z}_t^\top A \mathbf{z}_t + 2\mathbf{z}_t^\top \Sigma \Sigma^\top \bar{H}_t \mathbf{z}_t.$$

Hypoelliptic obstruction

Even if p_{T-t} is log-concave, so that $\bar{H}_t \preceq 0$, the term $\mathbf{z}_t^\top \Sigma \Sigma^\top \bar{H}_t \mathbf{z}_t$, is not necessarily nonpositive. Indeed,

$$\Sigma \Sigma^\top \bar{H}_t = \begin{pmatrix} 0 & 0 \\ 0 & 4\mathbf{I}_d \end{pmatrix} \begin{pmatrix} H_{xx} & H_{xv} \\ H_{vx} & H_{vv} \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 4H_{vx} & 4H_{vv} \end{pmatrix}$$

only acts on the velocity directions.

 Euclidean synchronous contraction is not automatic in hypoelliptic phase space.

Solution: Long-term regularity of the renormalized score

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

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

Key properties.

1. $\tilde{\mathbf{A}}$ is negative definite.
2. \tilde{p}_t "quantifies" **deviation from equilibrium** p_∞ .
3. Its curvature $\nabla^2 \log \tilde{p}_t$ characterizes the **regularity of the score**, for all $t \in (0, T]$,

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

4. Recover a bound of the type, as for general SGMs

$$\mathcal{W}_2(\pi_{\text{data}}, \hat{\pi}_{\infty, N}^\theta) \leq e^{-T} c_1 + M c_2 + \sqrt{h} c_3,$$

Kinetic contraction mechanism

Regularized score estimate

Under structural assumptions on the data distribution weaker than log-concavity, we proved

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

The singularity at $t = 0$ is integrable and the large-time behavior decays exponentially.

Weighted contraction

There exist a positive definite matrix \mathfrak{M} and $\eta > 0$ such that, for the synchronous difference,

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

Hence, by Grönwall, there exists $C \geq 0$ such that

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

Selected references

- Anderson (1982). *Reverse-time diffusion equation models*.
- Hyvärinen (2005). *Estimation of non-normalized statistical models by score matching*.
- Vincent (2011). *A connection between score matching and denoising autoencoders*.
- Song and Ermon (2019). *Generative modeling by estimating gradients of the data distribution*.
- Song et al. (2021). *Score-based generative modeling through stochastic differential equations*.
- Dockhorn et al. (2022). *Score-based generative modeling with critically-damped Langevin diffusion*.
- Strasman, Surendran, Boyer, Le Corff, Lemaire and Ocello (2025). *Wasserstein stability of kinetic score-based generative models*.

Solution 2: restore ellipticity

Idea

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

$$\Sigma_\varepsilon = \begin{pmatrix} \varepsilon & 0 \\ 0 & \sigma \end{pmatrix} \otimes \mathbf{I}_d, \quad \varepsilon > 0.$$

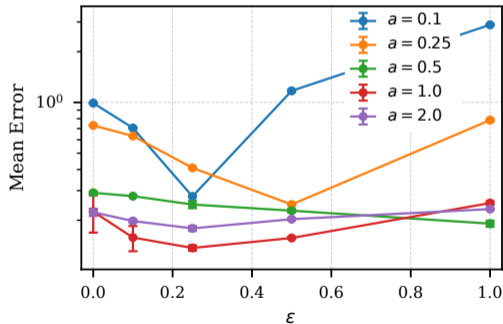
Consequences

- **Uniform ellipticity:** full phase-space noise.
- **More quantitative bounds:** standard log-concave tools apply.
- **Practice:** ε is an extra parameter controlling path regularity.

Numerical aspects

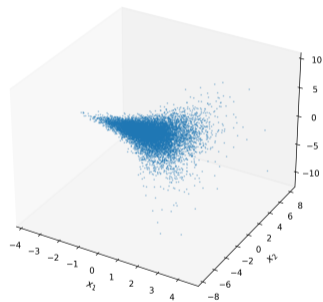
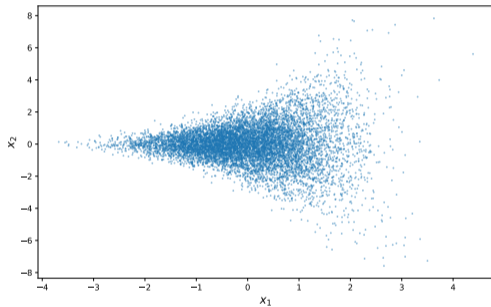
Empirics: Funnel dataset, $d = 100$

- Small $\varepsilon > 0$ can improve sliced- \mathcal{W}_2 compared with the CLD baseline ($\varepsilon = 0$).
- There is a trade-off with the other model hyperparameters.



Mean sliced- \mathcal{W}_2 over 5 runs; error bars are ± 1 standard deviation.

Funnel distribution



Samples from a funnel distribution. Left: first two coordinates. Right: first three coordinates.