

# On the Generalization Properties of Diffusion Models

Puheng Li\*, Zhong Li\*, Huishuai Zhang, Jiang Bian

Stanford & MSRA

Nov 10, 2023

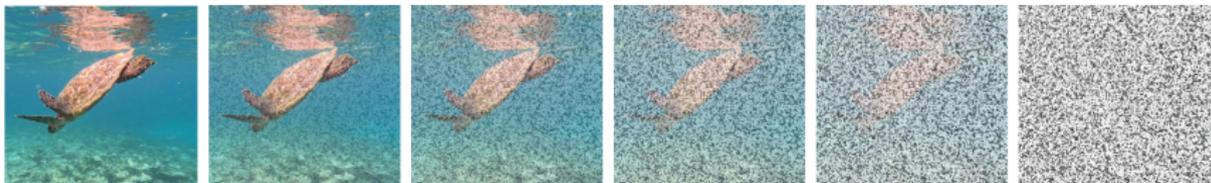
- **Diffusion model:** A generative model establishing a stochastic transport map between an empirically observed, yet unknown, target distribution and a known prior.

- **Diffusion model:** A generative model establishing a stochastic transport map between an empirically observed, yet unknown, target distribution and a known prior.
- Real-world applications: DALL·E, Imagen, Stable Diffusion...

- **Diffusion model:** A generative model establishing a stochastic transport map between an empirically observed, yet unknown, target distribution and a known prior.
- Real-world applications: DALL·E, Imagen, Stable Diffusion...
- Theoretical foundations of diffusion models remain under-explored, particularly, the fundamental **generalization** problem

# Formulation

$$\mathbf{x}(0) \xrightarrow{\quad\quad\quad} d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{W}_t \xrightarrow{\quad\quad\quad} \mathbf{x}(T)$$



$$\mathbf{x}(0) \xleftarrow{\quad\quad\quad} d\mathbf{x} = \left[ \mathbf{f}(\mathbf{x}, t) - \underbrace{g^2(t)\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}_{\approx s_{t,\theta}(\mathbf{x})} \right] dt + g(t)d\bar{\mathbf{W}}_t \xleftarrow{\quad\quad\quad} \mathbf{x}(T)$$
$$\approx s_{t,\theta}(\mathbf{x}) := \frac{1}{m} \mathbf{A}\sigma(\mathbf{W}\mathbf{x} + \mathbf{U}e(t)), \theta = \mathbf{A}$$

**Target:** Finitely-supported prob. & Gaussian mixtures    **Notations:**

**Loss:** Time-dependent score matching (Eq. (7))

**Algorithm:** Gradient flow

$t$ : SDE time     $T$ : maximal SDE time

$p_T \approx \pi$ : a known prior

$\tau$ : training time

# Main Result I

## Theorem (Data-Independent Generalization Gap)

Suppose  $p_0$  is continuously differentiable with a **compact support** set, and there exists a reproducing kernel Hilbert space (RKHS)  $\mathcal{H} (:=\mathcal{H}_{k_{\rho_0}})$  such that  $\bar{\mathbf{s}}_0, \bar{\theta}^* \in \mathcal{H}$ . Assume the initial loss, trainable parameters, the embedding function  $\mathbf{e}(t)$  and weighting function  $\lambda(t)$  are all bounded. Then with high probability, we have

$$D_{\text{KL}}\left(p_0 \| p_{0, \hat{\theta}_n(\tau)}\right) \lesssim \left[ \frac{\tau^4}{mn} + \frac{\tau^3}{m^2} + \frac{1}{\tau} \right] + \left[ \frac{1}{m} + \tilde{\mathcal{L}}(\bar{\theta}^*) + \tilde{\mathcal{L}}(\theta^*) \right] + D_{\text{KL}}(p_T \| \pi).$$

- Early-stopping generalization gap:

$$\tau_{\text{es}} = \Theta\left(n^{\frac{2}{5}}\right) \Rightarrow D_{\text{KL}}\left(p_0 \| p_{0, \hat{\theta}_n(\tau_{\text{es}})}\right) \lesssim (1/n)^{\frac{2}{5}} + (1/m)^{\frac{4}{5}}.$$

# Early-Stopping Generalization

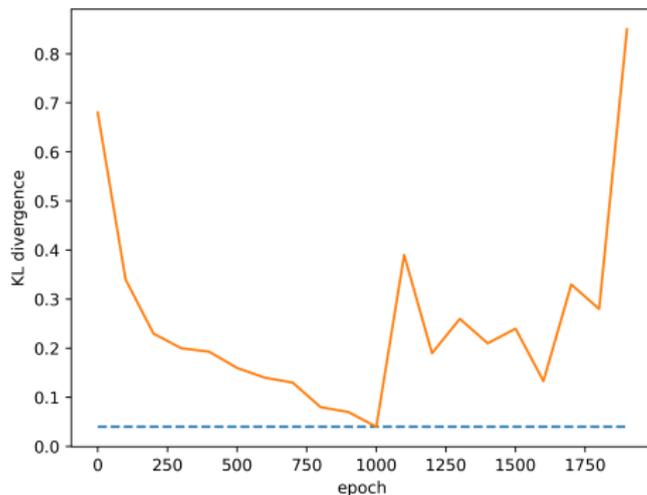


Figure: The KL divergence dynamics.

# Main Result II

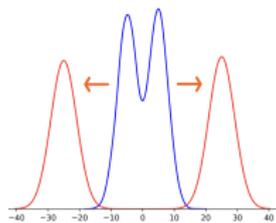


Figure: Illustration of modes shift.

## Theorem (Data-Dependent Generalization Gap)

Suppose  $p_0(x) = q_1 \mathcal{N}(x; -\mu, 1) + q_2 \mathcal{N}(x; \mu, 1)$ , where  $q_1, q_2 > 0$  with  $q_1 + q_2 = 1$ . Under the same conditions, with high probability we have

$$D_{\text{KL}}(p_0 \| p_{0, \hat{\theta}_n(\tau)}) \\ \lesssim \text{Poly}(\mu) \left[ \frac{\tau^4}{mn} + \frac{\tau^3}{m^2} \right] + \frac{1}{\tau} + \left[ \frac{\mu^2}{m} + \tilde{\mathcal{L}}(\bar{\theta}^*) + \tilde{\mathcal{L}}(\theta^*) \right] + D_{\text{KL}}(p_T \| \pi).$$

# Modes Shift Effect

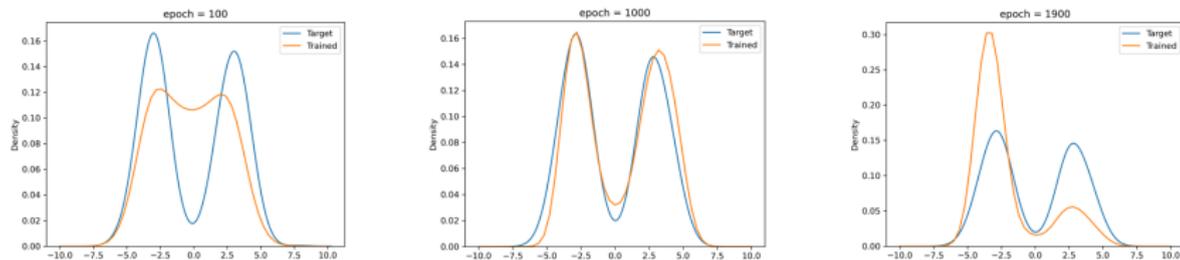


Figure: Training dynamics when the distance between two modes is 6 ( $\mu = 3$ ).

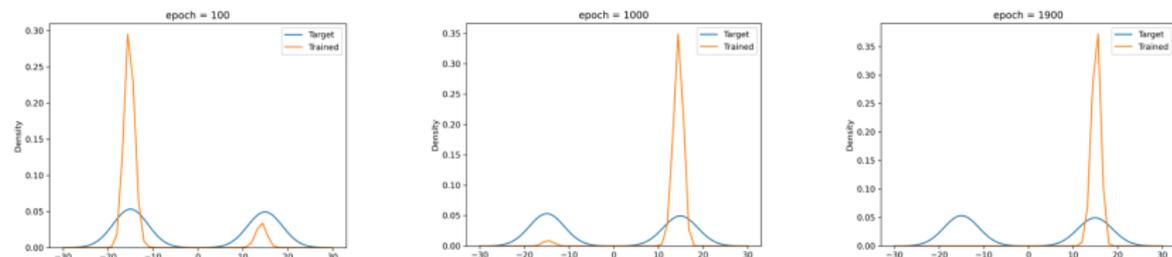


Figure: Training dynamics when the distance between two modes is 30 ( $\mu = 15$ ).

## U-net + MNIST:

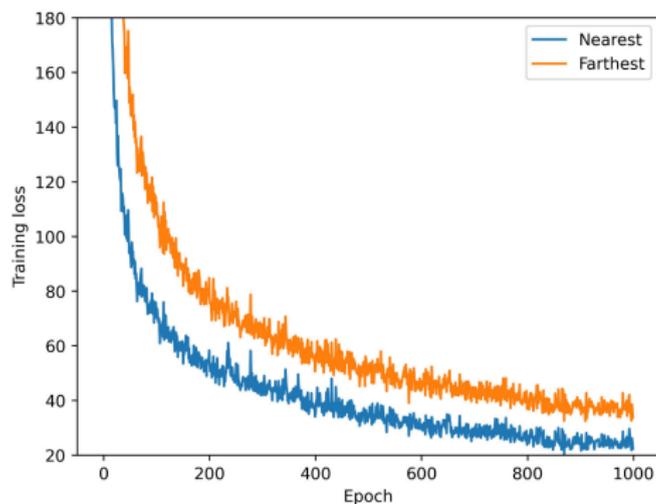


Figure: The training loss dynamics.

# Real-World Experiment



Figure: Sampling from the farthest (left) and nearest (right) clusters.

- Data-independent regime: For target distributions with finite supports, diffusion models have polynomially small generalization errors in  $n$  (sample size) and  $m$  (model capacity) with early-stopping.

# Conclusion

- Data-independent regime: For target distributions with finite supports, diffusion models have polynomially small generalization errors in  $n$  (sample size) and  $m$  (model capacity) with early-stopping.
- Data-dependent regime: For target distributions with increasing modes distances, the generalization performance of diffusion models becomes significantly worse.

**Thank you!**