

# Generative Modeling by Estimating Gradients of the Data Distribution

This blog post focuses on a promising new direction for generative modeling. We can learn score functions (gradients of log probability density functions) on a large number of noise-perturbed data distributions, then generate samples with Langevin-type sampling. The resulting generative models, often called *score-based generative models*, has several important advantages over existing model families: GAN-level sample quality without adversarial training, flexible model architectures, exact log-likelihood computation, and inverse problem solving without re-training models. In this blog post, we will show you in more detail the intuition, basic concepts, and potential applications of score-based generative models.

AUTHORS  
Yang Song

PUBLISHED  
May 5, 2021

AFFILIATIONS  
Stanford University

## Contents

### Introduction

The score function, score-based models, and score matching

Langevin dynamics

Naive score-based generative modeling and its pitfalls

Score-based generative modeling with multiple noise perturbations

Score-based generative modeling with stochastic differential equations (SDEs)

Perturbing data with an SDE

Reversing the SDE for sample generation

Estimating the reverse SDE with score-based models and score matching

How to solve the reverse SDE

Probability flow ODE

Controllable generation for inverse problem solving

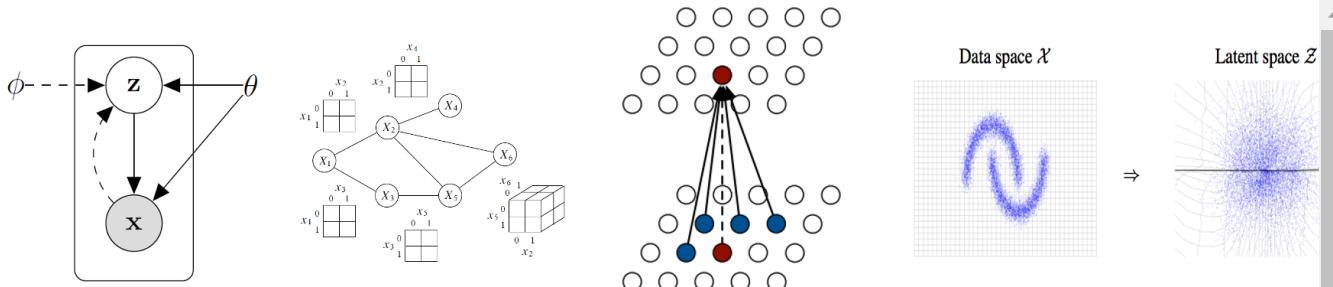
Connection to diffusion models and others

Concluding remarks

## Introduction

Existing generative modeling techniques can largely be grouped into two categories based on how they represent probability distributions.

1. **likelihood-based models**, which directly learn the distribution's probability density (or mass) function via (approximate) maximum likelihood. Typical likelihood-based models include autoregressive models [1, 2, 3], normalizing flow models [4, 5], energy-based models (EBMs) [6, 7], and variational auto-encoders (VAEs) [8, 9].
2. **implicit generative models** [10], where the probability distribution is implicitly represented by a model of its sampling process. The most prominent example is generative adversarial networks (GANs) [11], where new samples from the data distribution are synthesized by transforming a random Gaussian vector with a neural network.



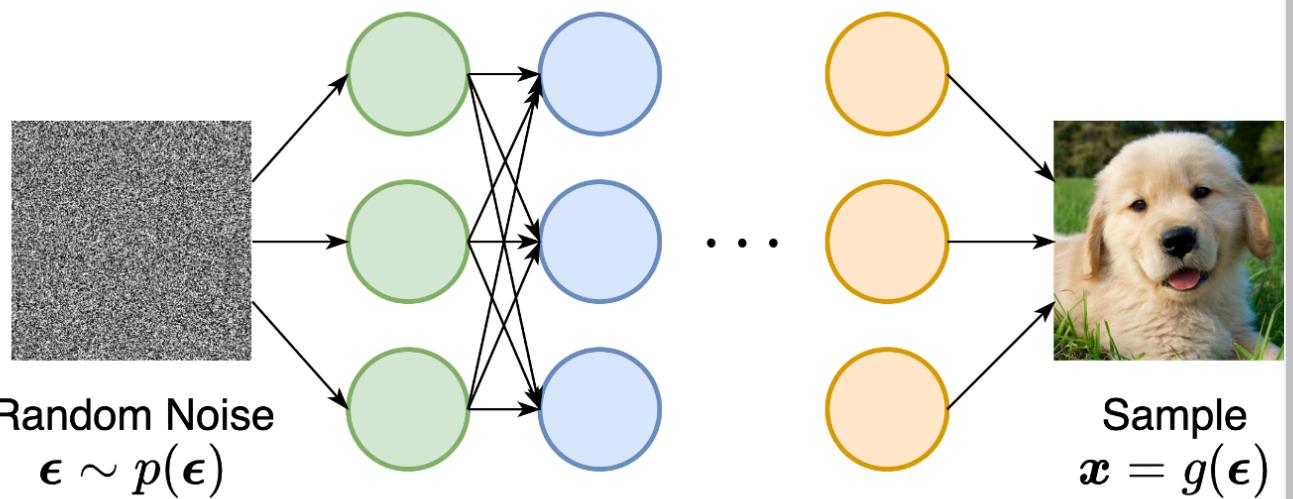
## Bayesian networks (e.g., VAEs)

## MRF

## Autoregressive models

## Flow models

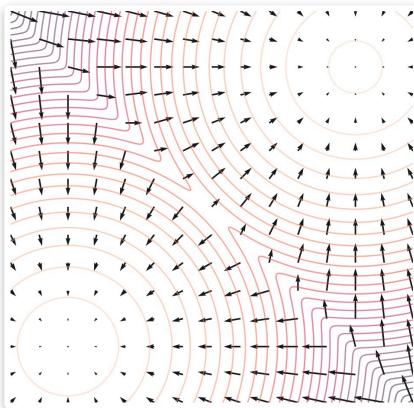
Bayesian networks, Markov random fields (MRF), autoregressive models, and normalizing flow models are all examples of likelihood-based models. All these models represent the probability density or mass function of a distribution.



GAN is an example of implicit models. It implicitly represents a distribution over all objects that can be produced by the generator network.

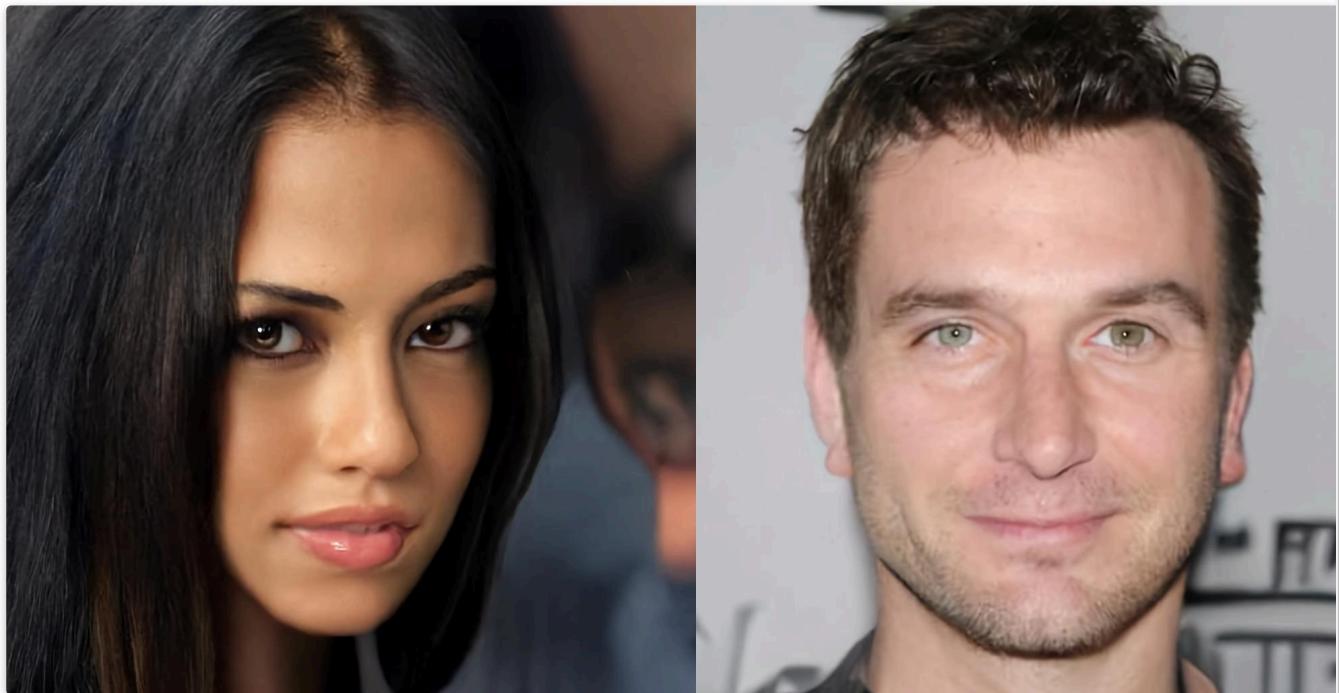
Likelihood-based models and implicit generative models, however, both have significant limitations. Likelihood-based models either require strong restrictions on the model architecture to ensure a tractable normalizing constant for likelihood computation, or must rely on surrogate objectives to approximate maximum likelihood training. Implicit generative models, on the other hand, often require adversarial training, which is notoriously unstable [12] and can lead to mode collapse [13].

In this blog post, I will introduce another way to represent probability distributions that may circumvent several of these limitations. The key idea is to model *the gradient of the log probability density function*, a quantity often known as the (Stein) **score function** [14, 15]. Such **score-based models** are not required to have a tractable normalizing constant, and can be directly learned by **score matching** [16, 17].



Score function (the vector field) and density function (contours) of a mixture of two Gaussians.

Score-based models have achieved state-of-the-art performance on many downstream tasks and applications. These tasks include, among others, image generation [18, 19, 20, 21, 22, 23] (Yes, better than GANs!), audio synthesis [24, 25, 26], shape generation [27], and music generation [28]. Moreover, score-based model have connections to normalizing flow models, therefore allowing exact likelihood computation and representation learning. Additionally, modeling and estimat scores facilitates inverse problem solving, with applications such as image inpainting [18, 21], image colorization [21], compressive sensing, and medical image reconstruction (e.g., CT, MRI) [29].



1024 x 1024 samples generated from score-based models [21]

This post aims to show you the motivation and intuition of score-based generative modeling, as well as its basic concepts, properties and applications.

## The score function, score-based models, and score matching

Suppose we are given a dataset  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , where each point is drawn independently from an underlying data distribution  $p(\mathbf{x})$ . Given this dataset, the goal of generative modeling is to fit a model to the data distribution such that we can synthesize new data points at will by sampling from the distribution.

In order to build such a generative model, we first need a way to represent a probability distribution. One such way, as in likelihood-based models, is to directly model the probability density function (p.d.f.) or probability mass function (p.m.f.). Let  $f_\theta(\mathbf{x}) \in \mathbb{R}$  be a real-valued function parameterized by a learnable parameter  $\theta$ . We can define a p.d.f. <sup>1</sup> via

$$p_\theta(\mathbf{x}) = \frac{e^{-f_\theta(\mathbf{x})}}{Z_\theta},$$

where  $Z_\theta > 0$  is a normalizing constant dependent on  $\theta$ , such that  $\int p_\theta(\mathbf{x}) d\mathbf{x} = 1$ . Here the function  $f_\theta(\mathbf{x})$  is often called an unnormalized probabilistic model or energy-based model [7].

We can train  $p_\theta(\mathbf{x})$  by maximizing the log-likelihood of the data

$$\max_{\theta} \sum_{i=1}^N \log p_\theta(\mathbf{x}_i).$$

However, equation (2) requires  $p_\theta(\mathbf{x})$  to be a normalized probability density function. This is undesirable because in order to compute  $p_\theta(\mathbf{x})$ , we must evaluate the normalizing constant  $Z_\theta$ —a typically intractable quantity for any general  $f_\theta(\mathbf{x})$ . Thus to make maximum likelihood training feasible, likelihood-based models must either restrict their model architectures (e.g., causal convolutions in autoregressive models, invertible networks in normalizing flow models) to make  $Z_\theta$  tractable, or approximate the normalizing constant (e.g., variational inference in VAEs, or MCMC sampling used in contrastive divergence [30]) which may be computationally expensive.

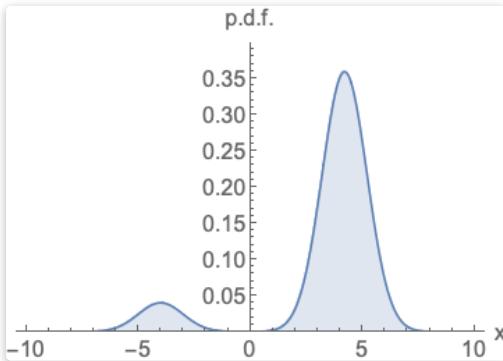
By modeling the score function instead of the density function, we can sidestep the difficulty of intractable normalizing constants. The **score function** of a distribution  $p(\mathbf{x})$  is defined as

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}),$$

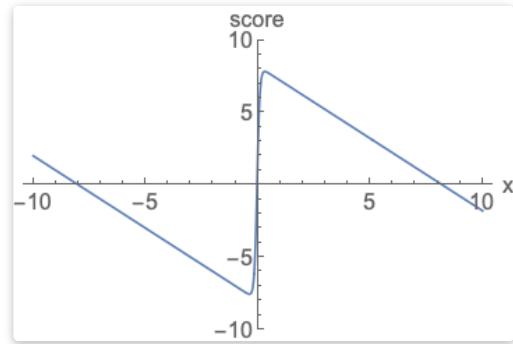
and a model for the score function is called a **score-based model** [18], which we denote as  $\mathbf{s}_\theta(\mathbf{x})$ . The score-based model is learned such that  $\mathbf{s}_\theta(\mathbf{x}) \sim \nabla_{\mathbf{x}} \log p(\mathbf{x})$ , and can be parameterized without worrying about the normalizing constant. For example, we can easily parameterize a score-based model

$$\mathbf{s}_\theta(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}) = -\nabla_{\mathbf{x}} f_\theta(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log Z_\theta}_{=0} = -\nabla_{\mathbf{x}} f_\theta(\mathbf{x}).$$

Note that the score-based model  $\mathbf{s}_\theta(\mathbf{x})$  is independent of the normalizing constant  $Z_\theta$ ! This significantly expands the family of models that we can tractably since we don't need any special architectures to make the normalizing constant tractable.



Parameterizing probability density functions. No matter how you change the model family and parameters, it has to be normalized (area under the curve must integrate to one).



Parameterizing score functions. No need to worry about normalization.

Similar to likelihood-based models, we can train score-based models by minimizing the **Fisher divergence**<sup>2</sup> between the model and the data distributions, def as

$$\mathbb{E}_{p(\mathbf{x})}[\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x})\|_2^2]$$

Intuitively, the Fisher divergence compares the squared  $\ell_2$  distance between the ground-truth data score and the score-based model. Directly computing this divergence, however, is infeasible because it requires access to the unknown data score  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ . Fortunately, there exists a family of methods called **score matching**<sup>3</sup> [16, 17, 31] that minimize the Fisher divergence without knowledge of the ground-truth data score. Score matching objectives can directly be estimated on a dataset and optimized with stochastic gradient descent, analogous to the log-likelihood objective for training likelihood-based models (with known normalizing constants). We can train the score-based model by minimizing a score matching objective, **without requiring adversarial optimization**.

Additionally, using the score matching objective gives us a considerable amount of modeling flexibility. The Fisher divergence itself does not require  $\mathbf{s}_\theta(\mathbf{x})$  to be the actual score function of any normalized distribution—it simply compares the  $\ell_2$  distance between the ground-truth data score and the score-based model, without additional assumptions on the form of  $\mathbf{s}_\theta(\mathbf{x})$ . In fact, the only requirement on the score-based model is that it should be a vector-valued function with the same input and output dimensionality, which is easy to satisfy in practice.

As a brief summary, we can represent a distribution by modeling its score function, which can be estimated by training a score-based model of free-form architectures with score matching.

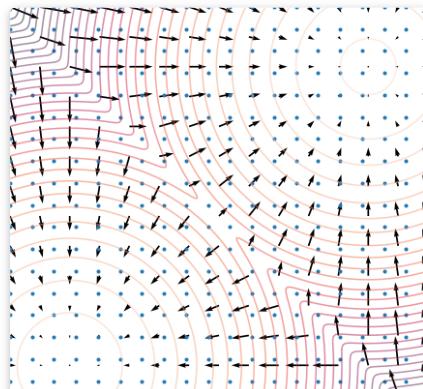
## Langevin dynamics

Once we have trained a score-based model  $\mathbf{s}_\theta(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$ , we can use an iterative procedure called **Langevin dynamics** [32, 33] to draw samples from

Langevin dynamics provides an MCMC procedure to sample from a distribution  $p(\mathbf{x})$  using only its score function  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ . Specifically, it initializes the chain from an arbitrary prior distribution  $\mathbf{x}_0 \sim \pi(\mathbf{x})$ , and then iterates the following

$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \epsilon \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \sqrt{2\epsilon} \mathbf{z}_i, \quad i = 0, 1, \dots, K,$$

where  $\mathbf{z}_i \sim \mathcal{N}(0, I)$ . When  $\epsilon \rightarrow 0$  and  $K \rightarrow \infty$ ,  $\mathbf{x}_K$  obtained from the procedure in (6) converges to a sample from  $p(\mathbf{x})$  under some regularity conditions. In practice, the error is negligible when  $\epsilon$  is sufficiently small and  $K$  is sufficiently large.

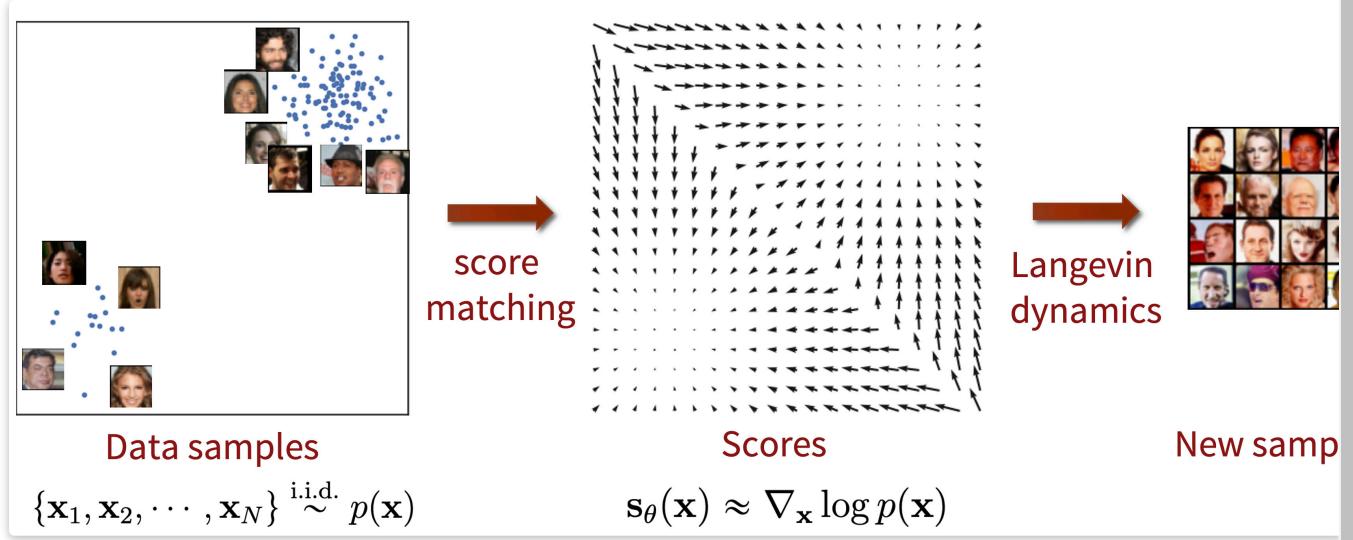


Using Langevin dynamics to sample from a mixture of two Gaussians.

Note that Langevin dynamics accesses  $p(\mathbf{x})$  only through  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ . Since  $\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$ , we can produce samples from our score-based model  $\mathbf{s}_{\theta}(\mathbf{x})$  by plugging it into equation (6).

## Naive score-based generative modeling and its pitfalls

So far, we've discussed how to train a score-based model with score matching, and then produce samples via Langevin dynamics. However, this naive approach has had limited success in practice—we'll talk about some pitfalls of score matching that received little attention in prior works [18].

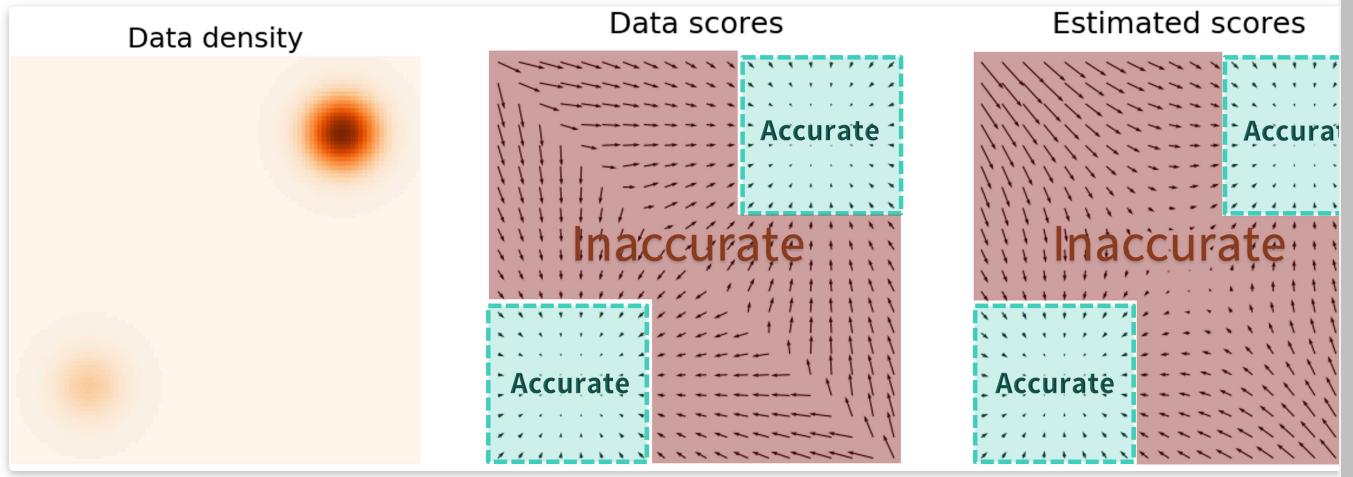


Score-based generative modeling with score matching + Langevin dynamics.

The key challenge is the fact that the estimated score functions are inaccurate in low density regions, where few data points are available for computing the score matching objective. This is expected as score matching minimizes the Fisher divergence

$$\mathbb{E}_{p(\mathbf{x})}[\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x})\|_2^2] = \int p(\mathbf{x}) \|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x})\|_2^2 d\mathbf{x}.$$

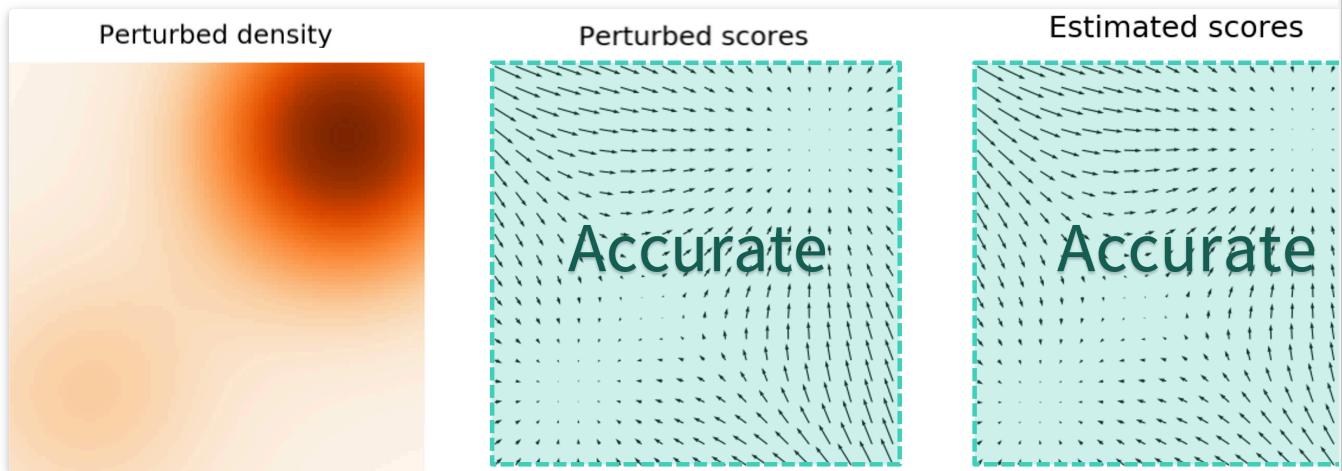
Since the  $\ell_2$  differences between the true data score function and score-based model are weighted by  $p(\mathbf{x})$ , they are largely ignored in low density regions where  $p(\mathbf{x})$  is small. This behavior can lead to subpar results, as illustrated by the figure below:



When sampling with Langevin dynamics, our initial sample is highly likely in low density regions when data reside in a high dimensional space. Therefore, having an inaccurate score-based model will derail Langevin dynamics from the very beginning of the procedure, preventing it from generating high quality samples that are representative of the data.

## Score-based generative modeling with multiple noise perturbations

based models on the noisy data points instead. When the noise magnitude is sufficiently large, it can populate low data density regions to improve the accuracy of estimated scores. For example, here is what happens when we perturb a mixture of two Gaussians perturbed by additional Gaussian noise.



Estimated scores are accurate everywhere for the noise-perturbed data distribution due to reduced low data density regions.

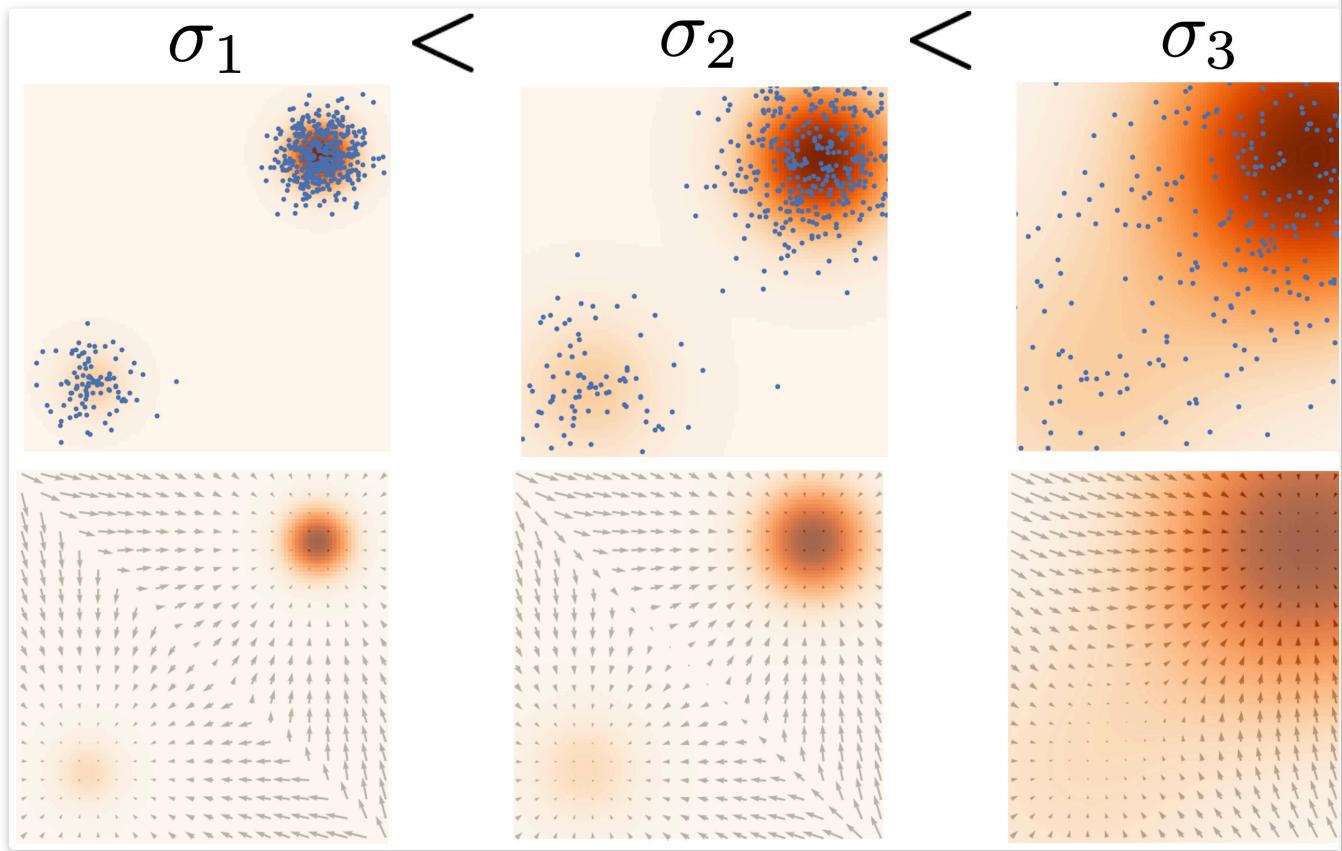
Yet another question remains: how do we choose an appropriate noise scale for the perturbation process? Larger noise can obviously cover more low density regions for better score estimation, but it over-corrupts the data and alters it significantly from the original distribution. Smaller noise, on the other hand, causes less corruption of the original data distribution, but does not cover the low density regions as well as we would like.

To achieve the best of both worlds, we use multiple scales of noise perturbations simultaneously [18, 19]. Suppose we always perturb the data with isotropic Gaussian noise, and let there be a total of  $L$  increasing standard deviations  $\sigma_1 < \sigma_2 < \dots < \sigma_L$ . We first perturb the data distribution  $p(\mathbf{x})$  with each of the Gaussian noise  $\mathcal{N}(0, \sigma_i^2 I)$ ,  $i = 1, 2, \dots, L$  to obtain a noise-perturbed distribution

$$p_{\sigma_i}(\mathbf{x}) = \int p(\mathbf{y}) \mathcal{N}(\mathbf{x}; \mathbf{y}, \sigma_i^2 I) d\mathbf{y}.$$

Note that we can easily draw samples from  $p_{\sigma_i}(\mathbf{x})$  by sampling  $\mathbf{x} \sim p(\mathbf{x})$  and computing  $\mathbf{x} + \sigma_i \mathbf{z}$ , with  $\mathbf{z} \sim \mathcal{N}(0, I)$ .

Next, we estimate the score function of each noise-perturbed distribution,  $\nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x})$ , by training a **Noise Conditional Score-Based Model**  $s_{\theta}(\mathbf{x}, i)$  (also a Noise Conditional Score Network, or NCSN [18, 19, 21], when parameterized with a neural network) with score matching, such that  $s_{\theta}(\mathbf{x}, i) \approx \nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x})$  all  $i = 1, 2, \dots, L$ .





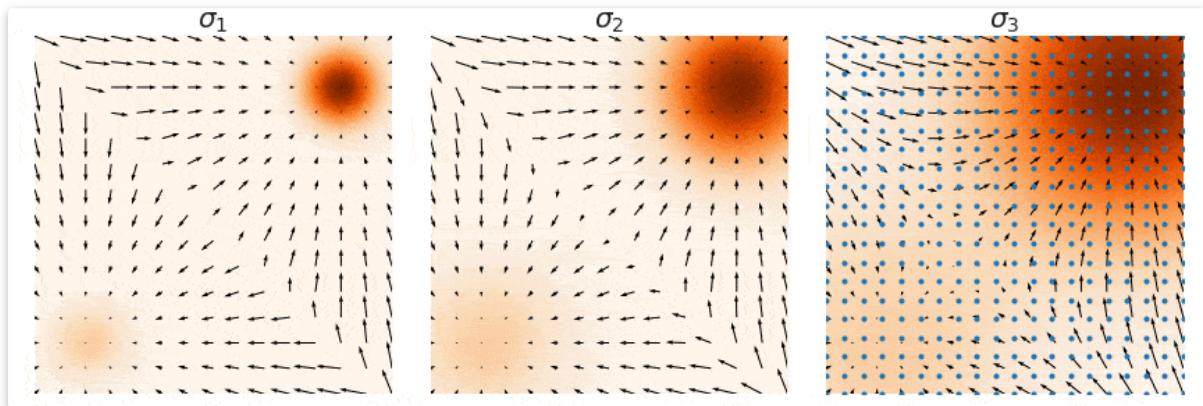
Perturbing an image with multiple scales of Gaussian noise.

The training objective for  $\mathbf{s}_\theta(\mathbf{x}, i)$  is a weighted sum of Fisher divergences for all noise scales. In particular, we use the objective below:

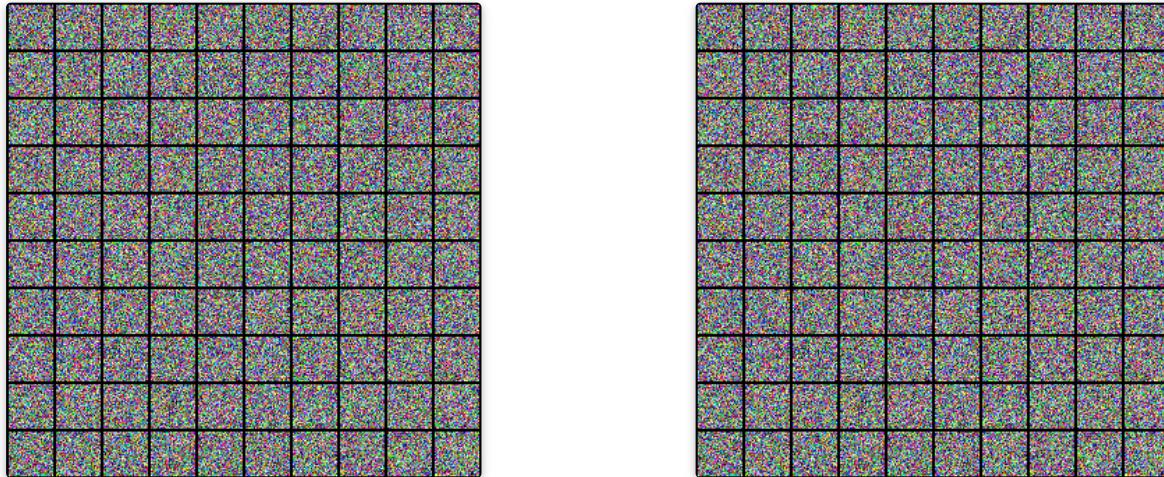
$$\sum_{i=1}^L \lambda(i) \mathbb{E}_{p_{\sigma_i}(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, i)\|_2^2],$$

where  $\lambda(i) \in \mathbb{R}_{>0}$  is a positive weighting function, often chosen to be  $\lambda(i) = \sigma_i^2$ . The objective (7) can be optimized with score matching, exactly as in optimizing the naive (unconditional) score-based model  $\mathbf{s}_\theta(\mathbf{x})$ .

After training our noise-conditional score-based model  $\mathbf{s}_\theta(\mathbf{x}, i)$ , we can produce samples from it by running Langevin dynamics for  $i = L, L-1, \dots, 1$  in sequence. This method is called **annealed Langevin dynamics** (defined by Algorithm 1 in [18], and improved by [19, 34]), since the noise scale  $\sigma_i$  decreases (anneals) gradually over time.



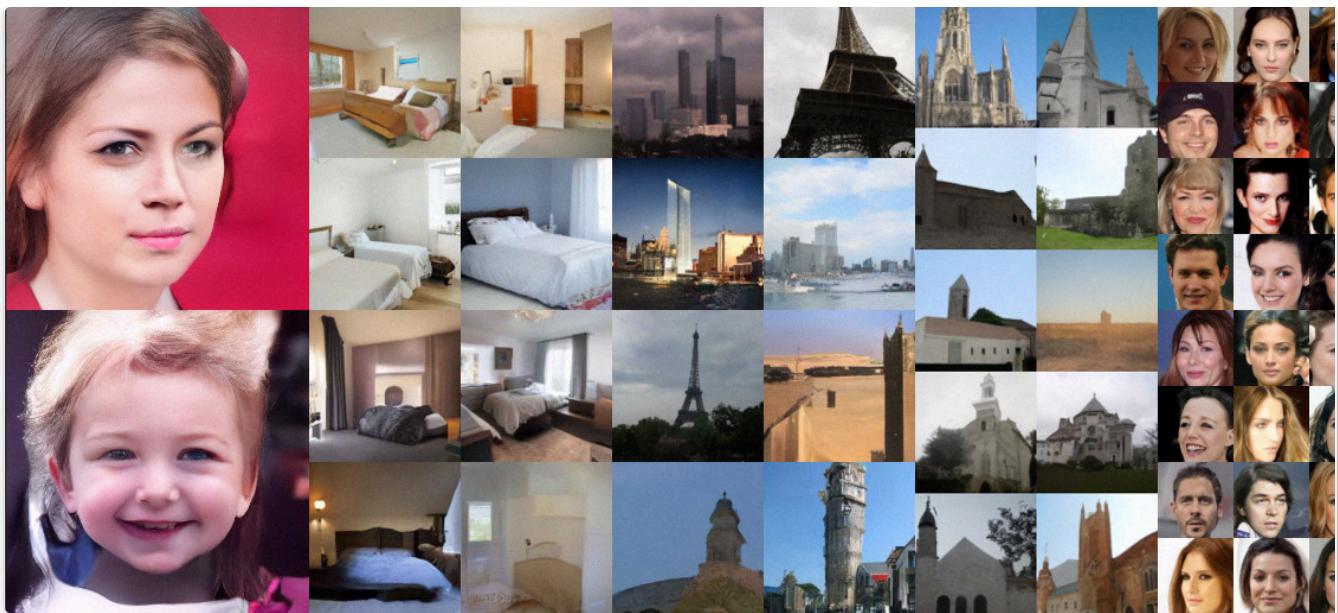
Annealed Langevin dynamics combine a sequence of Langevin chains with gradually decreasing noise scales.

Annealed Langevin dynamics for the Noise Conditional Score Network (NCSN) model (from ref. [18]) trained on CelebA (**left**) and CIFAR-10 (**right**). We can start from unstructured noise and modify images according to the scores, and generate nice samples. The method achieved state-of-the-art Inception score on CIFAR-10 at its time.

Here are some practical recommendations for tuning score-based generative models with multiple noise scales:

- Choose  $\sigma_1 < \sigma_2 < \dots < \sigma_L$  as a geometric progression, with  $\sigma_1$  being sufficiently small and  $\sigma_L$  comparable to the maximum pairwise distance between training data points [19].  $L$  is typically on the order of hundreds or thousands.
- Parameterize the score-based model  $\mathbf{s}_\theta(\mathbf{x}, i)$  with U-Net skip connections [18, 20].
- Apply exponential moving average on the weights of the score-based model when used at test time [19, 20].

With such best practices, we are able to generate high quality image samples with comparable quality to GANs on various datasets, such as below:



Samples from the NCSNv2 [19] model. From left to right: FFHQ 256x256, LSUN bedroom 128x128, LSUN tower 128x128, LSUN church\_outdoor 96x96, and CelebA 64x64.

## Score-based generative modeling with stochastic differential equations (SDEs)

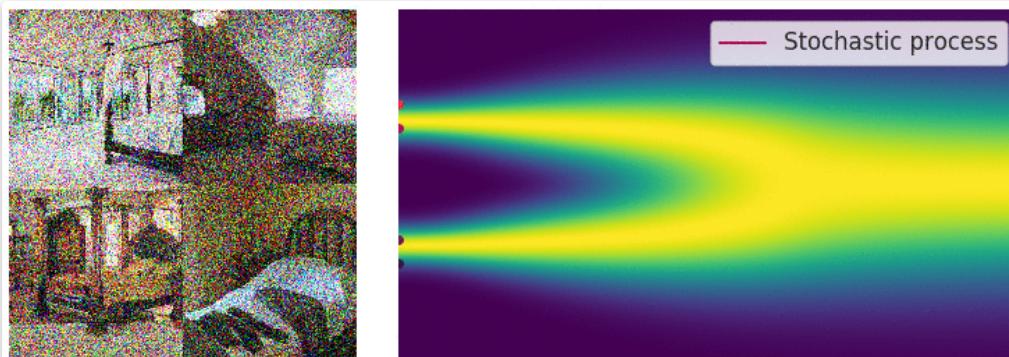
As we already discussed, adding multiple noise scales is critical to the success of score-based generative models. By generalizing the number of noise scales to infinity [21], we obtain not only **higher quality samples**, but also, among others, **exact log-likelihood computation**, and **controllable generation for inverse problem solving**.

In addition to this introduction, we have tutorials written in [Google Colab](#) to provide a step-by-step guide for training a toy model on MNIST. We also have more advanced code repositories that provide full-fledged implementations for large scale applications.

Link	Description
<a href="#"> Open in Colab</a>	Tutorial of score-based generative modeling with SDEs in JAX + FLAX
<a href="#"> Open in Colab</a>	Load our pretrained checkpoints and play with sampling, likelihood computation, and controllable synthesis (JAX + FLAX)
<a href="#"> Open in Colab</a>	Tutorial of score-based generative modeling with SDEs in PyTorch
<a href="#"> Open in Colab</a>	Load our pretrained checkpoints and play with sampling, likelihood computation, and controllable synthesis (PyTorch)
<a href="#">Code in JAX</a>	Score SDE codebase in JAX + FLAX
<a href="#">Code in PyTorch</a>	Score SDE codebase in PyTorch

## Perturbing data with an SDE

When the number of noise scales approaches infinity, we essentially perturb the data distribution with continuously growing levels of noise. In this case, the noise perturbation procedure is a continuous-time stochastic process, as demonstrated below



Perturbing data to noise with a continuous-time stochastic process.

How can we represent a stochastic process in a concise way? Many stochastic processes (diffusion processes in particular) are solutions of stochastic differential equations (SDEs). In this section, we will introduce the basic concepts of SDEs and how they relate to generative modeling.

© Copyright 2024 Yang Song. Powered by Jekyll with al-folio theme.

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w},$$

where  $\mathbf{f}(\cdot, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a vector-valued function called the drift coefficient,  $g(t) \in \mathbb{R}$  is a real-valued function called the diffusion coefficient,  $d\mathbf{w}$  denotes standard Brownian motion, and  $d\mathbf{w}$  can be viewed as infinitesimal white noise. The solution of a stochastic differential equation is a continuous collection of random variables  $\{\mathbf{x}(t)\}_{t \in [0, T]}$ . These random variables trace stochastic trajectories as the time index  $t$  grows from the start time  $0$  to the end time  $T$ . Let  $p$  denote the (marginal) probability density function of  $\mathbf{x}(t)$ . Here  $t \in [0, T]$  is analogous to  $i = 1, 2, \dots, L$  when we had a finite number of noise scales, and  $p_t(\mathbf{x})$  is analogous to  $p_{\sigma_i}(\mathbf{x})$ . Clearly,  $p_0(\mathbf{x}) = p(\mathbf{x})$  is the data distribution since no perturbation is applied to data at  $t = 0$ . After perturbing  $p(\mathbf{x})$  with the stochastic process for a sufficiently long time  $T$ ,  $p_T(\mathbf{x})$  becomes close to a tractable noise distribution  $\pi(\mathbf{x})$ , called a **prior distribution**. We note that  $p_T(\mathbf{x})$  is analogous to  $p_{\sigma_L}(\mathbf{x})$  in the case of finite noise scales, which corresponds to applying the largest noise perturbation  $\sigma_L$  to the data.

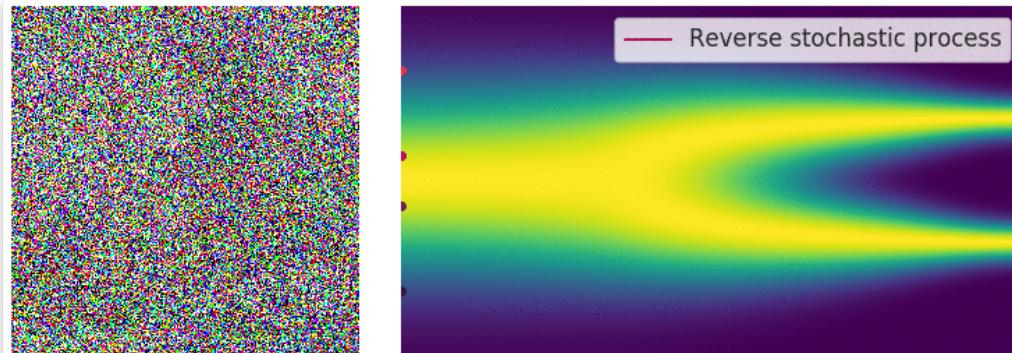
The SDE in (8) is **hand designed**, similarly to how we hand-designed  $\sigma_1 < \sigma_2 < \dots < \sigma_L$  in the case of finite noise scales. There are numerous ways to add noise perturbations, and the choice of SDEs is not unique. For example, the following SDE

$$d\mathbf{x} = e^t d\mathbf{w}$$

perturbs data with a Gaussian noise of mean zero and exponentially growing variance, which is analogous to perturbing data with  $\mathcal{N}(0, \sigma_1^2 I), \mathcal{N}(0, \sigma_2^2 I), \dots, \mathcal{N}(0, \sigma_L^2 I)$  when  $\sigma_1 < \sigma_2 < \dots < \sigma_L$  is a **geometric progression**. Therefore, the SDE should be viewed as part of the mode much like  $\{\sigma_1, \sigma_2, \dots, \sigma_L\}$ . In [21], we provide three SDEs that generally work well for images: the Variance Exploding SDE (VE SDE), the Variance Preserving (VP SDE), and the sub-VP SDE.

## Reversing the SDE for sample generation

Recall that with a finite number of noise scales, we can generate samples by reversing the perturbation process with **annealed Langevin dynamics**, i.e., sequentially sampling from each noise-perturbed distribution using Langevin dynamics. For infinite noise scales, we can analogously reverse the perturbation process for sample generation by using the reverse SDE.

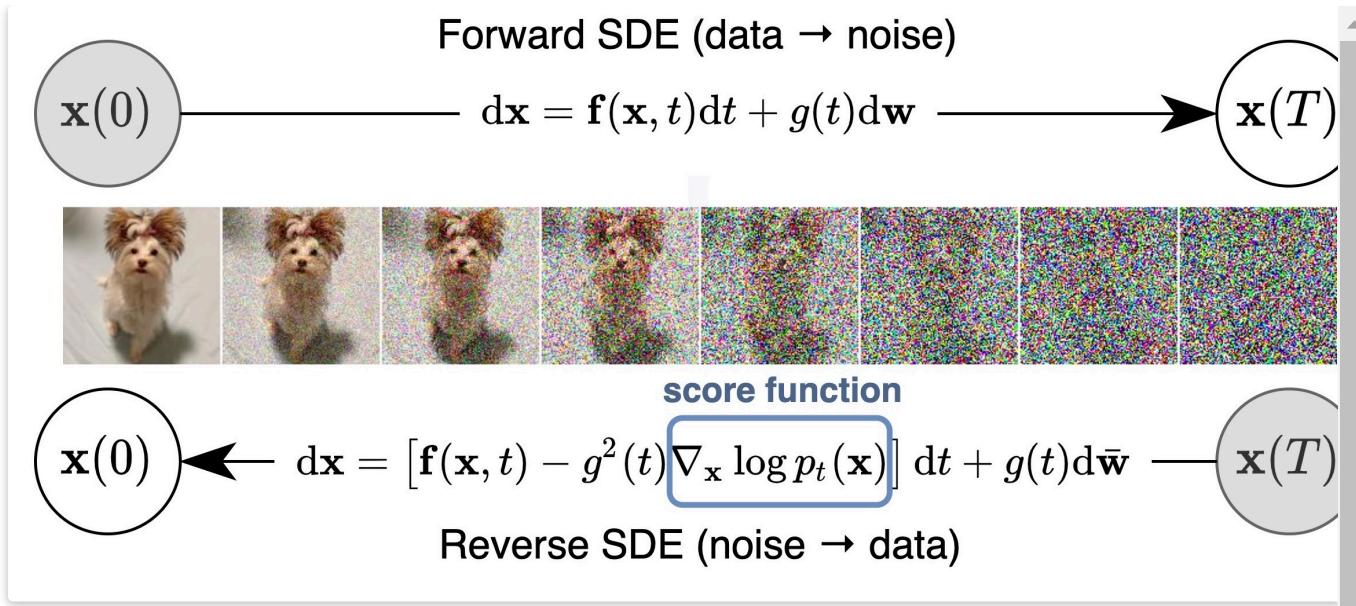


Generate data from noise by reversing the perturbation procedure.

Importantly, any SDE has a corresponding reverse SDE [35], whose closed form is given by

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g^2(t) \nabla_{\mathbf{x}} \log p_t(\mathbf{x})]dt + g(t)d\mathbf{w}.$$

Here  $dt$  represents a negative infinitesimal time step, since the SDE (10) needs to be solved backwards in time (from  $t = T$  to  $t = 0$ ). In order to compute the reverse SDE, we need to estimate  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ , which is exactly the **score function** of  $p_t(\mathbf{x})$ .



Solving a reverse SDE yields a score-based generative model. Transforming data to a simple noise distribution can be accomplished with an SDE. It can be reversed to generate sample noise if we know the score of the distribution at each intermediate time step.

### Estimating the reverse SDE with score-based models and score matching

Solving the reverse SDE requires us to know the terminal distribution  $p_T(\mathbf{x})$ , and the score function  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ . By design, the former is close to the prior distribution  $\pi(\mathbf{x})$  which is fully tractable. In order to estimate  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ , we train a **Time-Dependent Score-Based Model**  $\mathbf{s}_\theta(\mathbf{x}, t)$ , such that  $\mathbf{s}_\theta(\mathbf{x}, t) \approx \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ . This is analogous to the noise-conditional score-based model  $\mathbf{s}_\theta(\mathbf{x}, i)$  used for finite noise scales, trained such that  $\mathbf{s}_\theta(\mathbf{x}, i) \approx \nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x})$ .

Our training objective for  $\mathbf{s}_\theta(\mathbf{x}, t)$  is a continuous weighted combination of Fisher divergences, given by

$$\mathbb{E}_{t \in \mathcal{U}(0, T)} \mathbb{E}_{p_t(\mathbf{x})} [\lambda(t) \|\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, t)\|_2^2],$$

where  $\mathcal{U}(0, T)$  denotes a uniform distribution over the time interval  $[0, T]$ , and  $\lambda : \mathbb{R} \rightarrow \mathbb{R}_{>0}$  is a positive weighting function. Typically we use  $\lambda(t) \propto 1/\mathbb{E}[\|\nabla_{\mathbf{x}(t)} \log p(\mathbf{x}(t) | \mathbf{x}(0))\|_2^2]$  to balance the magnitude of different score matching losses across time.

As before, our weighted combination of Fisher divergences can be efficiently optimized with score matching methods, such as denoising score matching [17] sliced score matching [31]. Once our score-based model  $\mathbf{s}_\theta(\mathbf{x}, t)$  is trained to optimality, we can plug it into the expression of the reverse SDE in (10) to obtain estimated reverse SDE.

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g^2(t)\mathbf{s}_\theta(\mathbf{x}, t)]dt + g(t)d\mathbf{w}.$$

We can start with  $\mathbf{x}(T) \sim \pi$ , and solve the above reverse SDE to obtain a sample  $\mathbf{x}(0)$ . Let us denote the distribution of  $\mathbf{x}(0)$  obtained in such way as  $p_\theta$ . When the score-based model  $\mathbf{s}_\theta(\mathbf{x}, t)$  is well-trained, we have  $p_\theta \approx p_0$ , in which case  $\mathbf{x}(0)$  is an approximate sample from the data distribution  $p_0$ .

When  $\lambda(t) = g^2(t)$ , we have an important connection between our weighted combination of Fisher divergences and the KL divergence from  $p_0$  to  $p_\theta$  under some regularity conditions [36]:

$$\begin{aligned} \text{KL}(p_0(\mathbf{x}) \| p_\theta(\mathbf{x})) &\leq \frac{T}{2} \mathbb{E}_{t \in \mathcal{U}(0, T)} \mathbb{E}_{p_t(\mathbf{x})} [\lambda(t) \|\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, t)\|_2^2] \\ &\quad + \text{KL}(p_T \| \pi). \end{aligned}$$

Due to this special connection to the KL divergence and the equivalence between minimizing KL divergences and maximizing likelihood for model training, we have  $\lambda(t) = g(t)^2$  the **likelihood weighting function**. Using this likelihood weighting function, we can train score-based generative models to achieve very high likelihoods, comparable or even superior to state-of-the-art autoregressive models [36].

### How to solve the reverse SDE

By solving the estimated reverse SDE with numerical SDE solvers, we can simulate the reverse stochastic process for sample generation. Perhaps the simplest numerical SDE solver is the Euler-Maruyama method. When applied to our estimated reverse SDE, it discretizes the SDE using finite time steps and small Gaussian noise. Specifically, it chooses a small negative time step  $\Delta t \approx 0$ , initializes  $t \leftarrow T$ , and iterates the following procedure until  $t \approx 0$ :

$$\begin{aligned}\Delta \mathbf{x} &\leftarrow [\mathbf{f}(\mathbf{x}, t) - g^2(t)\mathbf{s}_\theta(\mathbf{x}, t)]\Delta t + g(t)\sqrt{|\Delta t|}\mathbf{z}_t \\ \mathbf{x} &\leftarrow \mathbf{x} + \Delta \mathbf{x} \\ t &\leftarrow t + \Delta t,\end{aligned}$$

Here  $\mathbf{z}_t \sim \mathcal{N}(0, I)$ . The Euler-Maruyama method is qualitatively similar to Langevin dynamics—both update  $\mathbf{x}$  by following score functions perturbed with Gaussian noise.

Aside from the Euler-Maruyama method, other numerical SDE solvers can be directly employed to solve the reverse SDE for sample generation, including, for example, [Milstein method](#) and [stochastic Runge-Kutta methods](#). In [21], we provided a reverse diffusion solver similar to Euler-Maruyama, but more tailored for solving reverse-time SDEs. More recently, authors in [37] introduced adaptive step-size SDE solvers that can generate samples faster with better quality.

In addition, there are two special properties of our reverse SDE that allow for even more flexible sampling methods:

- We have an estimate of  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$  via our time-dependent score-based model  $\mathbf{s}_\theta(\mathbf{x}, t)$ .
- We only care about sampling from each marginal distribution  $p_t(\mathbf{x})$ . Samples obtained at different time steps can have arbitrary correlations and do not have to form a particular trajectory sampled from the reverse SDE.

As a consequence of these two properties, we can apply MCMC approaches to fine-tune the trajectories obtained from numerical SDE solvers. Specifically, we propose **Predictor-Corrector samplers**. The **predictor** can be any numerical SDE solver that predicts  $\mathbf{x}(t + \Delta t) \sim p_{t+\Delta t}(\mathbf{x})$  from an existing sample  $\mathbf{x}(t) \sim p_t(\mathbf{x})$ . The **corrector** can be any MCMC procedure that solely relies on the score function, such as Langevin dynamics and Hamiltonian Monte Carlo.

At each step of the Predictor-Corrector sampler, we first use the predictor to choose a proper step size  $\Delta t < 0$ , and then predict  $\mathbf{x}(t + \Delta t)$  based on the current sample  $\mathbf{x}(t)$ . Next, we run several corrector steps to improve the sample  $\mathbf{x}(t + \Delta t)$  according to our score-based model  $\mathbf{s}_\theta(\mathbf{x}, t + \Delta t)$ , so that  $\mathbf{x}(t + \Delta t)$  becomes a higher-quality sample from  $p_{t+\Delta t}(\mathbf{x})$ .

With Predictor-Corrector methods and better architectures of score-based models, we can achieve **state-of-the-art** sample quality on CIFAR-10 (measured in [38]) and Inception scores [12], outperforming the best GAN model to date (StyleGAN2 + ADA [39]).

Method	FID ↓	Inception score ↑
StyleGAN2 + ADA [39]	2.92	9.83
Ours [21]	<b>2.20</b>	<b>9.89</b>

The sampling methods are also scalable for extremely high dimensional data. For example, it can successfully generate high fidelity images of resolution  $1024 \times 1024$ .



1024 x 1024 samples from a score-based model trained on the FFHQ dataset.

Some additional (uncurated) samples for other datasets (taken from this [GitHub repo](#)):



256 x 256 samples on LSUN bedroom.



256 x 256 samples on CelebA-HQ.

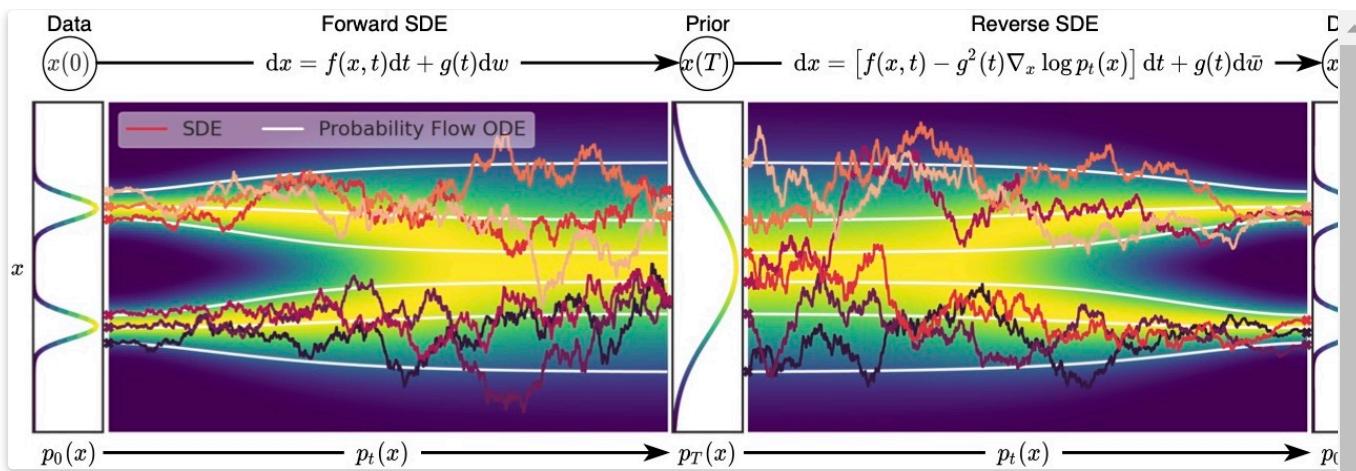
## Probability flow ODE

Despite capable of generating high-quality samples, samplers based on Langevin MCMC and SDE solvers do not provide a way to compute the exact log-likelihood of score-based generative models. Below, we introduce a sampler based on ordinary differential equations (ODEs) that allow for exact likelihood computation.

In [21], we show it is possible to convert any SDE into an ordinary differential equation (ODE) without changing its marginal distributions  $\{p_t(\mathbf{x})\}_{t \in [0, T]}$ . Thus by solving this ODE, we can sample from the same distributions as the reverse SDE. The corresponding ODE of an SDE is named **probability flow ODE** [21], given by

$$d\mathbf{x} = \left[ \mathbf{f}(\mathbf{x}, t) - \frac{1}{2} g^2(t) \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt.$$

The following figure depicts trajectories of both SDEs and probability flow ODEs. Although ODE trajectories are noticeably smoother than SDE trajectories, they convert the same data distribution to the same prior distribution and vice versa, sharing the same set of marginal distributions  $\{p_t(\mathbf{x})\}_{t \in [0, T]}$ . In other words,



We can map data to a noise distribution (the prior) with an SDE, and reverse this SDE for generative modeling. We can also reverse the associated probability flow ODE, which yields a deterministic process that samples from the same distribution as the SDE. Both the reverse-time SDE and probability flow ODE can be obtained by estimating score functions.

This probability flow ODE formulation has several unique advantages.

When  $\nabla_x \log p_t(\mathbf{x})$  is replaced by its approximation  $s_\theta(\mathbf{x}, t)$ , the probability flow ODE becomes a special case of a neural ODE [40]. In particular, it is an example of continuous normalizing flows [41], since the probability flow ODE converts a data distribution  $p_0(\mathbf{x})$  to a prior noise distribution  $p_T(\mathbf{x})$  (since it shares the same marginal distributions as the SDE) and is fully invertible.

As such, the probability flow ODE inherits all properties of neural ODEs or continuous normalizing flows, including exact log-likelihood computation. Specifically, one can leverage the instantaneous change-of-variable formula (Theorem 1 in [40], Equation (4) in [41]) to compute the unknown data density  $p_0$  from the known prior density  $p_T$  with numerical ODE solvers.

In fact, our model achieves the **state-of-the-art** log-likelihoods on uniformly dequantized <sup>4</sup> CIFAR-10 images [21], even without maximum likelihood training.

Method	Negative log-likelihood (bits/dim) ↓
RealNVP	3.49
iResNet	3.45
Glow	3.35
FFJORD	3.40
Flow++	3.29
Ours	<b>2.99</b>

When training score-based models with the **likelihood weighting** we discussed before, and using **variational dequantization** to obtain likelihoods on discrete images, we can achieve comparable or even superior likelihood to the state-of-the-art autoregressive models (all without any data augmentation) [36].

Method	Negative log-likelihood (bits/dim) ↓ on CIFAR-10	Negative log-likelihood (bits/dim) ↓ on ImageNet 32x32
Sparse Transformer	<b>2.80</b>	-
Image Transformer	2.90	3.77
Ours	2.83	<b>3.76</b>

## Controllable generation for inverse problem solving

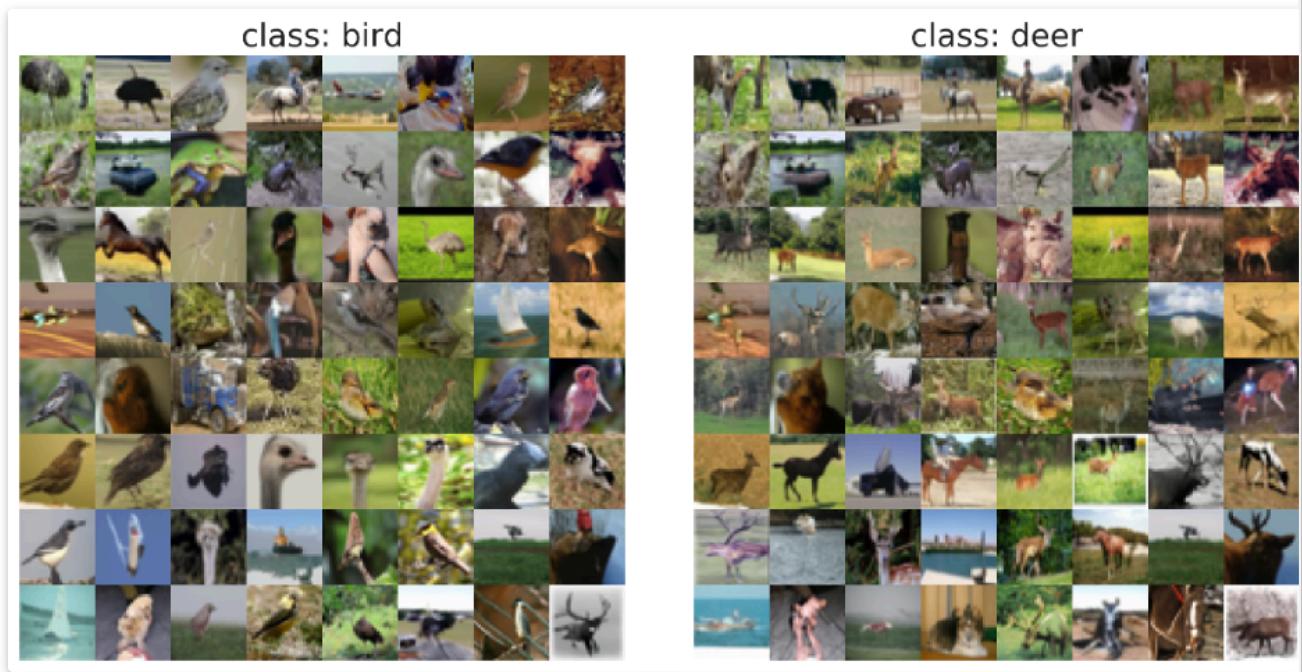
Score-based generative models are particularly suitable for solving inverse problems. At its core, inverse problems are same as Bayesian inference problems. Let  $\mathbf{x}$  and  $\mathbf{y}$  be two random variables, and suppose we know the forward process of generating  $\mathbf{y}$  from  $\mathbf{x}$ , represented by the transition probability distribution  $p(\mathbf{y} | \mathbf{x})$ . The inverse problem is to compute  $p(\mathbf{x} | \mathbf{y})$ . From Bayes' rule, we have  $p(\mathbf{x} | \mathbf{y}) = p(\mathbf{x})p(\mathbf{y} | \mathbf{x}) / \int p(\mathbf{x})p(\mathbf{y} | \mathbf{x})d\mathbf{x}$ . This expression can be greatly simplified by taking gradients with respect to  $\mathbf{x}$  on both sides, leading to the following Bayes' rule for score functions:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{y}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}).$$

Through score matching, we can train a model to estimate the score function of the unconditional data distribution, i.e.,  $s_\theta(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$ . This will allow us to easily compute the posterior score function  $\nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{y})$  from the known forward process  $p(\mathbf{y} | \mathbf{x})$  via equation (15), and sample from it with Langevin-type sampling [21].

A recent work from UT Austin [29] has demonstrated that score-based generative models can be applied to solving inverse problems in medical imaging, such as accelerating magnetic resonance imaging (MRI). Concurrently in [42], we demonstrated superior performance of score-based generative models not only on accelerated MRI, but also sparse-view computed tomography (CT). We were able to achieve comparable or even better performance than supervised or unrolled methods.

Below we show some examples on solving inverse problems for computer vision.



Class-conditional generation with an unconditional time-dependent score-based model, and a pre-trained noise-conditional image classifier on CIFAR-10.



Image inpainting with a time-dependent score-based model trained on LSUN bedroom. The leftmost column is ground-truth. The second column shows masked images (y in our framework). The rest columns show different inpainted images, generated by solving the conditional reverse-time SDE.



Image colorization with a time-dependent score-based model trained on LSUN church\_outdoor and bedroom. The leftmost column is ground-truth. The second column shows gray-scale images (y in our framework). The rest columns show different colorized images, generated by solving the conditional reverse-time SDE.



We can even colorize gray-scale portraits of famous people in history (Abraham Lincoln) with a time-dependent score-based model trained on FFHQ. The image resolution is  $1024 \times 1024$ .

## Connection to diffusion models and others

I started working on score-based generative modeling since 2019, when I was trying hard to make score matching scalable for training deep energy-based models on high-dimensional datasets. My first attempt at this led to the method sliced score matching [31]. Despite the scalability of sliced score matching for training energy-based models, I found to my surprise that Langevin sampling from those models fails to produce reasonable samples even on the MNIST dataset. I started investigating this issue and discovered three crucial improvements that can lead to extremely good samples: (1) perturbing data with multiple scales of noise, training score-based models for each noise scale; (2) using a U-Net architecture (we used RefineNet since it is a modern version of U-Nets) for the score-based model; (3) applying Langevin MCMC to each noise scale and chaining them together. With those methods, I was able to obtain the state-of-the-art Inception Score on CIFAR-10 in [18] (even better than the best GANs!), and generate high-fidelity image samples of resolution up to  $256 \times 256$  in [19].

The idea of perturbing data with multiple scales of noise is by no means unique to score-based generative models though. It has been previously used in, for example, simulated annealing, annealed importance sampling [43], diffusion probabilistic models [44], infusion training [45], and variational walkback [46] for generative stochastic networks [47]. Out of all these works, diffusion probabilistic modeling is perhaps the closest to score-based generative modeling. Diffusion probabilistic models are hierarchical latent variable models first proposed by Jascha and his colleagues [44] in 2015, which generate samples by learning a variational decoder to reverse a discrete diffusion process that perturbs data to noise. Without awareness of this work, score-based generative modeling was proposed and motivated independently from a very different perspective. Despite both perturbing data with multiple scales of noise, the connection between score-based generative modeling and diffusion probabilistic modeling seemed superficial at that time, since the former is trained by score matching and sampled by Langevin dynamics, while the latter is trained by the evidence lower bound (ELBO) and sampled with a learned decoder.

In 2020, Jonathan Ho and colleagues [20] significantly improved the empirical performance of diffusion probabilistic models and first unveiled a deeper connection to score-based generative modeling. They showed that the ELBO used for training diffusion probabilistic models is essentially equivalent to the weighted combination of score matching objectives used in score-based generative modeling. Moreover, by parameterizing the decoder as a sequence of score-based models with a U-Net architecture, they demonstrated for the first time that diffusion probabilistic models can also generate high quality image samples comparable or superior to GANs.

Inspired by their work, we further investigated the relationship between diffusion models and score-based generative models in an ICLR 2021 paper [21]. We found that the sampling method of diffusion probabilistic models can be integrated with annealed Langevin dynamics of score-based models to create a unified and more powerful sampler (the Predictor-Corrector sampler). By generalizing the number of noise scales to infinity, we further proved that score-based generative models and diffusion probabilistic models can both be viewed as discretizations to stochastic differential equations determined by score functions. This work bridges both score-based generative modeling and diffusion probabilistic modeling into a unified framework.

Collectively, these latest developments seem to indicate that both score-based generative modeling with multiple noise perturbations and diffusion probabilistic models are different perspectives of the same model family, much like how wave mechanics and matrix mechanics are equivalent formulations of quantum mechanics in the history of physics <sup>5</sup>. The perspective of score matching and score-based models allows one to calculate log-likelihoods exactly, solve inverse problems naturally, and is directly connected to energy-based models, Schrödinger bridges and optimal transport [48]. The perspective of diffusion models is naturally connected to VAEs, lossy compression, and can be directly incorporated with variational probabilistic inference. This blog post focuses on the first perspective, but I highly recommend interested readers to learn about the alternative perspective of diffusion models as well (see a great blog by Lilian Weng)

Many recent works on score-based generative models or diffusion probabilistic models have been deeply influenced by knowledge from both sides of research (see a [website](#) curated by researchers at the University of Oxford). Despite this deep connection between score-based generative models and diffusion model, it is hard to come up with an umbrella term for the model family that they both belong to. Some colleagues in DeepMind propose to call them "Generative Diffusion Processes". It remains to be seen if this will be adopted by the community in the future.

## Concluding remarks

This blog post gives a detailed introduction to score-based generative models. We demonstrate that this new paradigm of generative modeling is able to produce high quality samples, compute exact log-likelihoods, and perform controllable generation for inverse problem solving. It is a compilation of several papers we published in the past few years. Please visit them if you are interested in more details:

- [Yang Song\\*, Sahaj Garg\\*, Jiaxin Shi, and Stefano Ermon. Sliced Score Matching: A Scalable Approach to Density and Score Estimation. UAI 2019 \(Oral\)](#).
- [Yang Song, and Stefano Ermon. Generative Modeling by Estimating Gradients of the Data Distribution. NeurIPS 2019 \(Oral\)](#).
- [Yang Song, and Stefano Ermon. Improved Techniques for Training Score-Based Generative Models. NeurIPS 2020](#)
- [Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. ICLR 2021 \(Outstanding Paper Award\)](#).
- [Yang Song\\*, Conor Durkan\\*, Iain Murray, and Stefano Ermon. Maximum Likelihood Training of Score-Based Diffusion Models. NeurIPS 2021 \(Spotlight\)](#).
- [Yang Song\\*, Liyue Shen\\*, Lei Xing, and Stefano Ermon. Solving Inverse Problems in Medical Imaging with Score-Based Generative Models. ICLR 2022](#)

For a list of works that have been influenced by score-based generative modeling, researchers at the University of Oxford have built a very useful (but necessarily incomplete) website: <https://scorebasedgenerativemodeling.github.io/>.

There are two major challenges of score-based generative models. First, the sampling speed is slow since it involves a large number of Langevin-type iterations. Second, it is inconvenient to work with discrete data distributions since scores are only defined on continuous distributions.

The first challenge can be partially solved by using numerical ODE solvers for the probability flow ODE with lower precision (a similar method, denoising diffusion implicit modeling, has been proposed in [49]). It is also possible to learn a direct mapping from the latent space of probability flow ODEs to the image space, as shown in [50]. However, all such methods to date result in worse sample quality.

The second challenge can be addressed by learning an autoencoder on discrete data and performing score-based generative modeling on its continuous latent space [28, 51]. Jascha's original work on diffusion models [44] also provides a discrete diffusion process for discrete data distributions, but its potential for large-scale applications remains yet to be proven.

It is my conviction that these challenges will soon be solved with the joint efforts of the research community, and score-based generative models/ diffusion-based models will become one of the most useful tools for data generation, density estimation, inverse problem solving, and many other downstream tasks in machine learning.

## Footnotes

1. Hereafter we only consider probability density functions. Probability mass functions are similar. [↪]

2. Fisher divergence is typically between two distributions p and q, defined as

$$\mathbb{E}_{p(\mathbf{x})}[\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \nabla_{\mathbf{x}} \log q(\mathbf{x})\|_2^2]. \quad (4)$$

Here we slightly abuse the term as the name of a closely related expression for score-based models. [↩]

3. Commonly used score matching methods include denoising score matching [17] and sliced score matching [31]. Here is an introduction to score matching and sliced score matching. [↪]

4. It is typical for normalizing flow models to convert discrete images to continuous ones by adding small uniform noise to them. [↪]

5. Goes without saying that the significance of score-based generative models/diffusion probabilistic models is in no way comparable to quantum mechanics. [↪]

## References

1. [The neural autoregressive distribution estimator](#)  
Larochelle, H. and Murray, I., 2011. International Conference on Artificial Intelligence and Statistics, pp. 29–37.

2. [Made: Masked autoencoder for distribution estimation](#)  
Germain, M., Gregor, K., Murray, I. and Larochelle, H., 2015. International Conference on Machine Learning, pp. 881–889.

**3. Pixel recurrent neural networks**

Van Oord, A., Kalchbrenner, N. and Kavukcuoglu, K., 2016. International Conference on Machine Learning, pp. 1747–1756.

**4. NICE: Non-linear independent components estimation**

Dinh, L., Krueger, D. and Bengio, Y., 2014. arXiv preprint arXiv:1410.8516.

**5. Density estimation using Real NVP**

Dinh, L., Sohl-Dickstein, J. and Bengio, S., 2017. International Conference on Learning Representations.

**6. A tutorial on energy-based learning**

LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M. and Huang, F., 2006. Predicting structured data, Vol 1(0).

**7. How to Train Your Energy-Based Models**

Song, Y. and Kingma, D.P., 2021. arXiv preprint arXiv:2101.03288.

**8. Auto-encoding variational bayes**

Kingma, D.P. and Welling, M., 2014. International Conference on Learning Representations.

**9. Stochastic backpropagation and approximate inference in deep generative models**

Rezende, D.J., Mohamed, S. and Wierstra, D., 2014. International conference on machine learning, pp. 1278–1286.

**10. Learning in implicit generative models**

Mohamed, S. and Lakshminarayanan, B., 2016. arXiv preprint arXiv:1610.03483.

**11. Generative adversarial nets**

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. and Bengio, Y., 2014. Advances in neural information processing systems, pp. 2672–2680.

**12. Improved techniques for training gans**

Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A. and Chen, X., 2016. Advances in Neural Information Processing Systems, pp. 2226–2234.

**13. Unrolled Generative Adversarial Networks [\[link\]](#)**

Metz, L., Poole, B., Pfau, D. and Sohl-Dickstein, J., 2017. 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.

**14. A Kernel Test of Goodness of Fit [\[HTML\]](#)**

Chwialkowski, K., Strathmann, H. and Gretton, A., 2016. Proceedings of The 33rd International Conference on Machine Learning, Vol 48, pp. 2606–2615. PMLR.

**15. A kernelized Stein discrepancy for goodness-of-fit tests**

Liu, Q., Lee, J. and Jordan, M., 2016. International conference on machine learning, pp. 276–284.

**16. Estimation of non-normalized statistical models by score matching**

Hyvärinen, A., 2005. Journal of Machine Learning Research, Vol 6(Apr), pp. 695–709.

**17. A connection between score matching and denoising autoencoders**

Vincent, P., 2011. Neural computation, Vol 23(7), pp. 1661–1674. MIT Press.

**18. Generative Modeling by Estimating Gradients of the Data Distribution [\[PDF\]](#)**

Song, Y. and Ermon, S., 2019. Advances in Neural Information Processing Systems, pp. 11895–11907.

**19. Improved Techniques for Training Score-Based Generative Models [\[PDF\]](#)**

Song, Y. and Ermon, S., 2020. Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.

**20. Denoising diffusion probabilistic models**

Ho, J., Jain, A. and Abbeel, P., 2020. arXiv preprint arXiv:2006.11239.

**21. Score-Based Generative Modeling through Stochastic Differential Equations [\[link\]](#)**

Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Ermon, S. and Poole, B., 2021. International Conference on Learning Representations.

**22. Diffusion models beat gans on image synthesis**

Dhariwal, P. and Nichol, A., 2021. arXiv preprint arXiv:2105.05233.

**23. Cascaded Diffusion Models for High Fidelity Image Generation**

Ho, J., Saharia, C., Chan, W., Fleet, D.J., Norouzi, M. and Salimans, T., 2021.

**24. WaveGrad: Estimating Gradients for Waveform Generation [\[link\]](#)**

Chen, N., Zhang, Y., Zen, H., Weiss, R.J., Norouzi, M. and Chan, W., 2021. International Conference on Learning Representations.

**25. DiffWave: A Versatile Diffusion Model for Audio Synthesis [\[link\]](#)**

Kong, Z., Ping, W., Huang, J., Zhao, K. and Catanzaro, B., 2021. International Conference on Learning Representations.

**26. Grad-tts: A diffusion probabilistic model for text-to-speech**

Popov, V., Vovk, I., Gogoryan, V., Sadekova, T. and Kudinov, M., 2021. arXiv preprint arXiv:2105.06337.

**27. Learning Gradient Fields for Shape Generation**

Cai, R., Yang, G., Averbuch-Elor, H., Hao, Z., Belongie, S., Snavely, N. and Hariharan, B., 2020. Proceedings of the European Conference on Computer Vision (ECCV).

**28. Symbolic Music Generation with Diffusion Models**

Mittal, G., Engel, J., Hawthorne, C. and Simon, I., 2021. arXiv preprint arXiv:2103.16091.

**29. Robust Compressed Sensing MRI with Deep Generative Priors**

Jalal, A., Arvinte, M., Daras, G., Price, E., Dimakis, A.G. and Tamir, J.I., 2021. Advances in neural information processing systems.

31. **Sliced score matching: A scalable approach to density and score estimation** [\[PDF\]](#)  
 Song, Y., Garg, S., Shi, J. and Ermon, S., 2020. Uncertainty in Artificial Intelligence, pp. 574–584.
32. **Correlation functions and computer simulations**  
 Parisi, G., 1981. Nuclear Physics B, Vol 180(3), pp. 378–384. Elsevier.
33. **Representations of knowledge in complex systems**  
 Grenander, U. and Miller, M.I., 1994. Journal of the Royal Statistical Society: Series B (Methodological), Vol 56(4), pp. 549–581. Wiley Online Library.
34. **Adversarial score matching and improved sampling for image generation** [\[link\]](#)  
 Jolicoeur-Martineau, A., Piche-Taillefer, R., Mitiagkas, I. and Combes, R.T.d., 2021. International Conference on Learning Representations.
35. **Reverse-time diffusion equation models**  
 Anderson, B.D., 1982. Stochastic Processes and their Applications, Vol 12(3), pp. 313–326. Elsevier.
36. **Maximum Likelihood Training of Score-Based Diffusion Models**  
 Song, Y., Durkan, C., Murray, I. and Ermon, S., 2021. Advances in Neural Information Processing Systems (NeurIPS).
37. **Gotta Go Fast When Generating Data with Score-Based Models**  
 Jolicoeur-Martineau, A., Li, K., Piche-Taillefer, R., Kachman, T. and Mitiagkas, I., 2021. arXiv preprint arXiv:2105.14080.
38. **GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium**  
 Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B. and Hochreiter, S., 2017. Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, {USA}, pp. 6626–6637.
39. **Training Generative Adversarial Networks with Limited Data**  
 Karras, T., Aittala, M., Hellsten, J., Laine, S., Lehtinen, J. and Aila, T., 2020. Proc. NeurIPS.
40. **Neural Ordinary Differential Equations**  
 Chen, T.Q., Rubanova, Y., Bettencourt, J. and Duvenaud, D., 2018. Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pp. 6572–6583.
41. **Scalable Reversible Generative Models with Free-form Continuous Dynamics** [\[link\]](#)  
 Grathwohl, W., Chen, R.T.Q., Bettencourt, J. and Duvenaud, D., 2019. International Conference on Learning Representations.
42. **Solving Inverse Problems in Medical Imaging with Score-Based Generative Models** [\[PDF\]](#)  
 Song, Y., Shen, L., Xing, L. and Ermon, S., 2022. International Conference on Learning Representations.
43. **Annealed importance sampling**  
 Neal, R.M., 2001. Statistics and computing, Vol 11(2), pp. 125–139. Springer.
44. **Deep unsupervised learning using nonequilibrium thermodynamics**  
 Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N. and Ganguli, S., 2015. International Conference on Machine Learning, pp. 2256–2265.
45. **Learning to generate samples from noise through infusion training**  
 Bordes, F., Honari, S. and Vincent, P., 2017. arXiv preprint arXiv:1703.06975.
46. **Variational walkback: Learning a transition operator as a stochastic recurrent net**  
 Goyal, A., Ke, N.R., Ganguli, S. and Bengio, Y., 2017. arXiv preprint arXiv:1711.02282.
47. **GSNs: generative stochastic networks**  
 Alain, G., Bengio, Y., Yao, L., Yosinski, J., Thibodeau-Laufer, E., Zhang, S. and Vincent, P., 2016. Information and Inference: A Journal of the IMA, Vol 5(2), pp. 210–249. Oxford University Press.
48. **Diffusion Schrödinger Bridge with Applications to Score-Based Generative Modeling**  
 De Bortoli, V., Thornton, J., Heng, J. and Doucet, A., 2021. Advances in Neural Information Processing Systems (NeurIPS).
49. **Denoising Diffusion Implicit Models** [\[link\]](#)  
 Song, J., Meng, C. and Ermon, S., 2021. International Conference on Learning Representations.
50. **Knowledge Distillation in Iterative Generative Models for Improved Sampling Speed**  
 Luhman, E. and Luhman, T., 2021. arXiv e-prints, pp. arXiv–2101.
51. **Score-based Generative Modeling in Latent Space**  
 Vahdat, A., Kreis, K. and Kautz, J., 2021. Advances in Neural Information Processing Systems (NeurIPS).

**What do you think?**

531 Responses



Upvote



Funny



Love



Surprised



Angry



Sad

**17 Comments****Login ▾****G**

Join the discussion...

LOG IN WITH

OR SIGN UP WITH DISQUS [?](#)

Name

**♡** 62**Share****Best** **Newest** **Oldest****Muhammad Firmansyah Kasim**

2 years ago

Thank you for this great blog post! Why did you model the score  $s$  directly? Why not model the unnormalized  $\log(p)$ , and then take its derivative with respect to  $x$  to compute  $s$ ? The second one produces the score profile that is automatically a differential 1-form (i.e. its Jacobian is symmetric), while modelling  $s$  directly does not necessarily produce a differential 1-form (its Jacobian is not necessarily symmetric).

4 1 Reply Share ›

**C****Carl** → Muhammad Firmansyah Kasim

2 years ago

For very good reasons: 1) score can be directly used in Langevin sampling & reverse SDE in diffusion models 2) Taking derivative is extremely expensive for large neural networks. 3) Most realistic data distribution has compact support, the log of 0 is minus infinity, modeling  $\log(p)$  can lead to numerical instability .

1 0 Reply Share ›

**Z****Zheneng Shi**

3 months ago

您好，您在文中提到“to make  $Z$  tractable, or approximate the normalizing constant (e.g., variational inference in VAEs) ”，请问VAE的变分推理具体是怎么最大似然归一化常数的呢？

2 0 Reply Share ›

**J****jianfeng cai**

a year ago

Thanks for the post. I have some questions. In your notebook example (Pytorch version), it seems that you do not use score matching to transform the loss function, but directly derive the data distribution from the SDE formula and obtain  $\mathbf{p}_{\theta}(x(t))$ . The calculation is then performed using the simple Fisher divergence. Are all SDE formulas directly inferable to derive their data distributions, or are there any restrictions that require score matching to be used?

1 0 Reply Share ›

**R****Renjie Wang**

a year ago

为什么朗之万采样最后会收敛到目标分布？

1 0 Reply Share ›

**Z****Zheneng Shi** → Renjie Wang

3 months ago

可以看下这篇文章里有讲：arXiv:2403.18103

1 0 Reply Share ›

**jiahao Lu**

3 years ago

Such an inspiring and revealing blog ! It really helps understanding. Many thanks !

1 0 Reply Share ↗

**jiahao Lu**

3 years ago

Actually I have one question here .. Why score-based models have exact log-likelihood computation ? From my understanding, matching the gradients of log-likelihood is kind of approximating the exact log-likelihood computation ...

0 0 Reply Share ↗

**Carl Thomé**

3 years ago

Amazing blog post! Thanks so much for this.

1 0 Reply Share ↗

**Cheng Shen**

7 months ago

Thanks for the great blog post that clarifies the thinking process behind score-based model and its' connection to diffusion model. I was a student in CS 236 and this post has been very helpful. Should've read it earlier!

0 0 Reply Share ↗

**Renjie Wang**

a year ago

受益匪浅

0 0 Reply Share ↗

**Jaechul Lee**

2 years ago

It makes me much easier to figure out your great job. I look forward to your next publication. Thank you so much for the awesome post!

0 0 Reply Share ↗

**ivan**

2 years ago

Hi, thanks for the post. I am looking at your notebook example (Pytorch version). You rescale your ScoreNet output by multiplying  $1/\text{std}[:, \text{None}, \text{None}, \text{None}]$ , then in the loss\_fn, you multiply  $\text{std}[:, \text{None}, \text{None}, \text{None}]$  back again, what's the point? It seems to do nothing with the output if you don't normalize and rescale your output.

0 0 Reply Share ↗

**Sukjun Hwang**

↗ ivan

a year ago edited

$\frac{1}{\sigma}$  gets multiplied to the ScoreNet due to the characteristic that the output of optimally trained  $s_\theta(x)$  is proportional to  $\frac{1}{\sigma}$ .

Say it is not multiplied at the end. Then, parameters of the ScoreNet has to control L2 norm of its output with respect to the given time step  $t$ .

Because  $\frac{1}{\sigma}$  gets multiplied at the end, the ScoreNet can be better trained as now it does not have to control L2 norm of its output, getting independent to  $t$ .

The latter multiplication of  $\sigma$  in the loss function is driven from the denoising score matching, which you can refer to Section 4.2. of the NCSN paper.

2 0 Reply Share ↗

**Zweig Shao**

2 years ago

Why does the pitfall of naive score-based generative (SBG) modeling (mentioned in the post) exist and can be solved by noise perturbing? I mean, when training SBG model, we use the denoising score matching technique, the noise perturbing is already involved. So does the SBG model trained by denoising score matching still have that pitfall? If yes, is there a conflict that we still use noise perturbing to eliminate this problem?

感谢您提供了可读性如此高的post，为了表达清晰，我重新用中文表达一遍我的困惑。文章中提到的naive score-based generative

0 0 Reply Share ›



**Miao Song** → Zweig Shao  
2 years ago

Use the sliced score matching to estimate scores without noise, which is written on P4 of the paper

0 0 Reply Share ›



**舒展** → Miao Song  
2 months ago

为什么 likelihood based 方法一定要同时估计归一化系数，这个归一化系数生成采样是什么关系？

0 0 Reply Share ›

---

Subscribe

Privacy

Do Not Sell My Data