

Reviewing the Paper “Generative modeling by estimating gradients of the data distribution” and Extending the Framework for A Multi-Resolution Diffusion Generative Model

Dominic Padova and Zixuan Liu

05/15/2023

1 Introduction

Machine learning has benefited from having high quality data. But collecting and cleaning data to get high quality can be complicated and expensive in terms of time and monetary cost. Furthermore, when the data generation setting changes, often new data must be collected and cleaned to accommodate the new setting and new constraints; this compounds issues stemming from the data collection process. However, estimating the data generation process from existing data, then simulating new data based on the estimated data generation process could alleviate some of the costs of high quality data collection. Estimating the data generation process often amounts to estimating an unknown and complicated data distribution, say of images of dogs, using a statistical model, called a generative model. The statistical model assigns a probability to data points, like an image of a chihuahua or an image of a blueberry muffin, and builds the model from these assignments. Once this model is estimated, new data can be produced by sampling from the generative model. This is the idea behind generative modeling (see Figure 1).

2 Paper Review: Generative modeling by estimating gradients of the data distribution

2.1 What is generative modeling?

Generative modeling is an unsupervised learning task in machine learning that involves automatically discovering and learning regularities or patterns in input data so that the model can be used to generate or output new examples that may be drawn from the original dataset.

2.2 A key challenge for building complex generative models

One main challenge of building a generative model from data is that our data distribution can be extremely complicated, especially for data with high dimensions. So consider how complicated it might be for the distribution of images, video, and audio. It might have millions of dimensions. So in our result, we need to build a complex model to fit the data

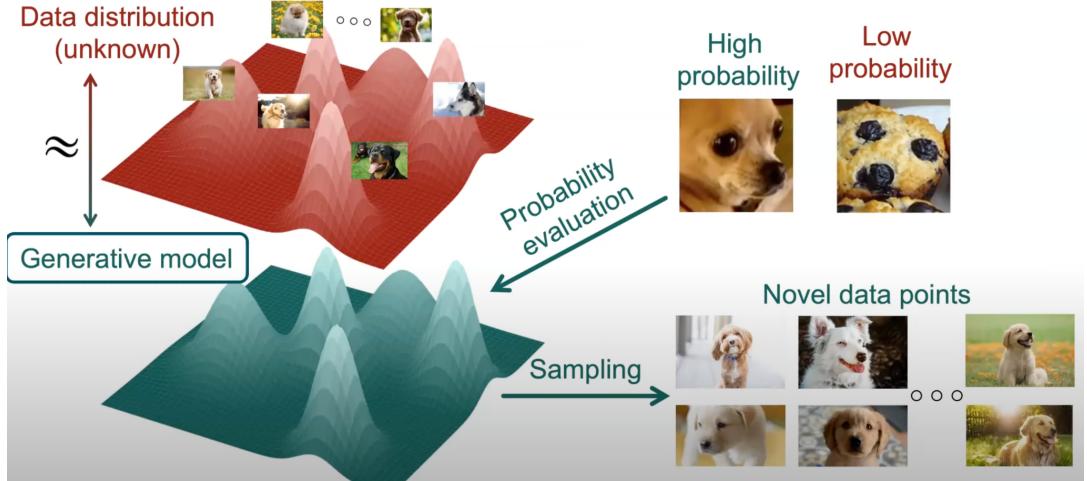


Figure 1: **Generative modeling.** Generative modeling is performed by estimating the data distribution using a statistical model that assigns probabilities to data to build the model, and then produces new data by sampling from the model. Image credit: Song, Y. <https://www.youtube.com/watch?v=nv-WTeKRL10&t=3174s> (oral presentation).

distribution. Suppose we have a continuous probability distribution where we use $p(x)$ to represent the probability density function, we define the score function as the gradient of $\log p(x)$. Given the density function, we can compute the score function very easily because we can just take the derivative of the logarithm:

$$p(x) = \frac{1}{Z(\theta)} e^{-E(x, \theta)}$$

$$\log(p(x)) = -E(x, \theta) - \log(Z(\theta))$$

$$\underbrace{\nabla_x \log(p(x))}_{\text{Stein score function}} = \underbrace{-\nabla_x E(x, \theta)}_{\substack{\text{Neg. gradient of energy} \\ \text{function w.r.t. data}}}$$

Conversely, with the score function, we can also recover the density function in principle by computing integrals. So mathematically, this function preserves all the information in the density function. But computationally, this score function is much easier to work with compared to the density function. Intuitively, the Stein score function is the vector field that points in the direction where the (log) probability grows the fastest. Therefore, it captures information about the data distribution.

2.3 Score-matching

Mathematically, we are given a bunch of data points which are assumed to be i.i.d. sampled from the data distribution $p(x)$, and our goal is to estimate this score function of the data density. So how can we train this model to be close to our ground truth data score function?

One approach is to minimize an objective function that seeks to bring two vector fields close³ together: the ground truth score function and an estimate of the score function. This approach is called score matching. How can we compute the score matching objective? Mathematically, we can capture this with the Fisher divergence objective. However, the fisher divergence cannot be directly computed because we do not know the ground truth value of the score function. Still, there is a way to address this challenge, using an two approaches called denoising score matching and sliced score matching:

- Explicit (Fisher Divergence-based) score matching

$$\frac{1}{2} \mathbb{E}_{p_{data}(x)} [\| \underbrace{\nabla_x \log(p_{data}(x))}_{\text{Fisher divergence: unknown}} - s_\theta(x) \|_2^2] \quad (1)$$

- Denoising score matching

$$\frac{1}{2} \mathbb{E}_{\substack{q_\sigma(\tilde{x} | x) p_{data}(x) \\ q_\sigma(\tilde{x}): \text{noised data dist.}}} [\| s_\theta(\tilde{x}) - \underbrace{\nabla_{\tilde{x}} \log(q_\sigma(\tilde{x} | x))}_{\text{Gradient of Gaussian noise kernel, conditioned on original data}} \|_2^2] \quad (2)$$

- Sliced score matching

$$\mathbb{E}_{p_v} \mathbb{E}_{p_{data}} [\underbrace{v^T \nabla_x s_\theta(x) v}_{\text{Random projections onto Jacobian of score model}} + \frac{1}{2} \|s_\theta(x)\|_2^2] \quad (3)$$

where the sliced score matching objective is achieved by applying the Divergence theorem to the Explicit (Fisher divergence-based) score matching objective to achieve the Implicit score matching objective

$$\mathbb{E}_{p_{data}(x)} [\frac{1}{2} \|s_\theta(x)\|_2^2 + \underbrace{\text{trace}(\nabla_x s_\theta(x))}_{\text{divergence of } s_\theta(x)}]$$

and using an efficient random projection estimator to estimate the $\text{trace}(\nabla_x s_\theta(x))$ divergence term.

2.4 Langevin dynamics

Langevin dynamics can produce samples from a probability density $p(x)$ using only the score function $\nabla_x \log p(x)$ (See Figure 2). Given a fixed step size $\epsilon > 0$, and an initial value $\tilde{x}_0 \sim \pi(x)$ with π being a prior distribution (e.g. a Uniform distribution), the Langevin method recursively computes the following

$$\tilde{x}_t = \tilde{x}_{t-1} + \frac{\epsilon}{2} \nabla_x \log(p(\tilde{x}_{t-1})) + \sqrt{\epsilon} z_t, z_t \sim N(0, Id), \epsilon > 0 \quad (4)$$

This procedure can be interpreted as a noisy gradient ascent. The noise here allows the samples to be distributed around the peaks instead of directly on the peaks. This allows better recovery of the true distribution.

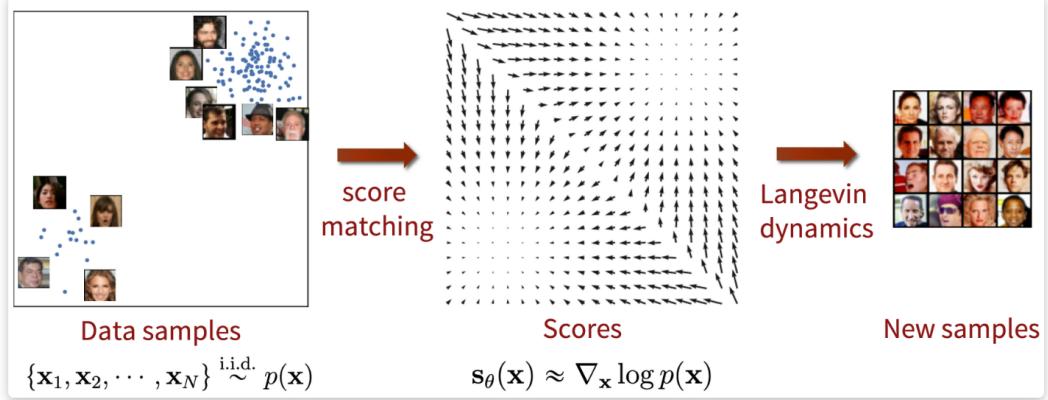


Figure 2: **Langevin dynamics.** New samples are produced using the estimated score model and Langevin dynamics. Image credit: [1].

2.5 Challenges of score-based generative modeling

- The manifold hypothesis

Since the score function is a gradient taken in the ambient space, it is undefined when x is confined to a low dimensional manifold.

The score-matching objective provides a consistent score estimator only when the support of the data distribution is the whole space (See Figure 3).

- Inaccurate score estimation with score matching

In regions of low data density, score matching may not have enough evidence to estimate score functions accurately, due to a lack of data samples (See Figure 4).

- Slow mixing of Langevin dynamics When two modes of the data distribution are separated by low-density regions, Langevin dynamics will not be able to correctly recover the relative weights of these two modes in a reasonable time, and therefore might not converge to the true distribution (See Figure 5).

2.6 Adding noise to the data and annealing the Langevin sampling solves the challenges

Perturbing data with random Gaussian noise makes the data distribution more amenable to score-based generative modeling. When using Langevin dynamics to generate samples, initially use scores corresponding to large noise models, and gradually anneal down the step size based on the noise level. These two additions allow score matching models to estimate the true score function, effectively passing the benefits from estimating large noise models down to the low noise models, and allow the annealed Langevin dynamics process to correctly recover the relative weights of the modes separated by low data density regions (See Figure 6).

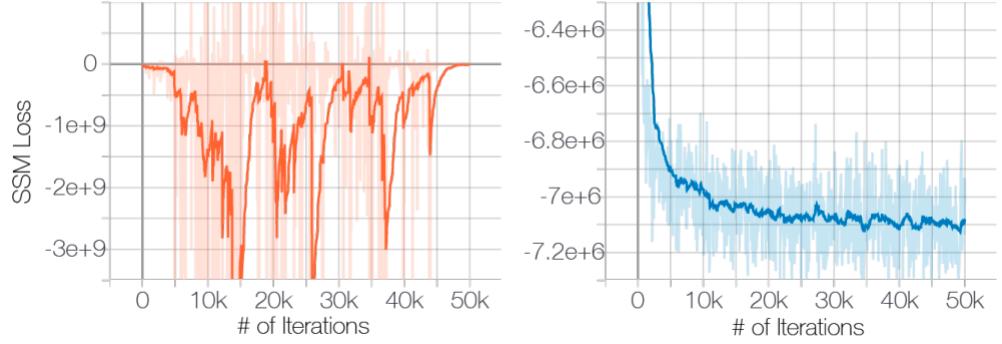


Figure 3: Manifold Hypothesis Challenge. Assuming the data lies on a low-dimensional manifold (the Manifold Hypothesis) presents a challenge for a score-based model to estimate the score function, since the score depends on the gradient of the full ambient space, not a gradient along a low-dimensional manifold. Whereas the sliced score matching (SSM) loss function under the Manifold Hypothesis (left) shows the model has trouble learning the score function, the SSM loss after injecting noise into the data to support the data on the full ambient space shows a marked improvement in learning ability. Image credit: [1].

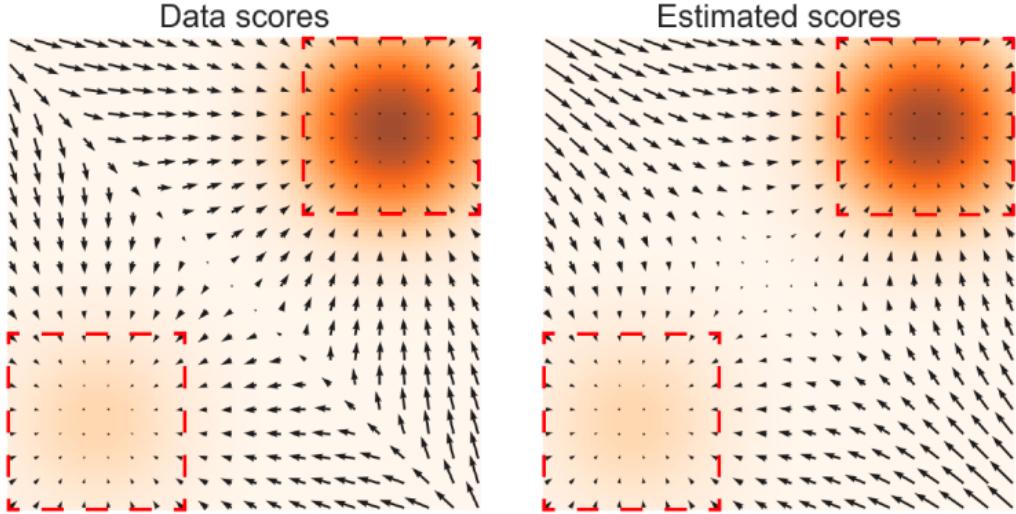


Figure 4: Inaccurate Score Estimation Challenge. Score matching suffers in low density regions of the data distribution, due to a dearth of samples. The true data scores (left) in the low data density region (bottom red square) are not recapitulated by the estimated scores (right). Image credit: [1].

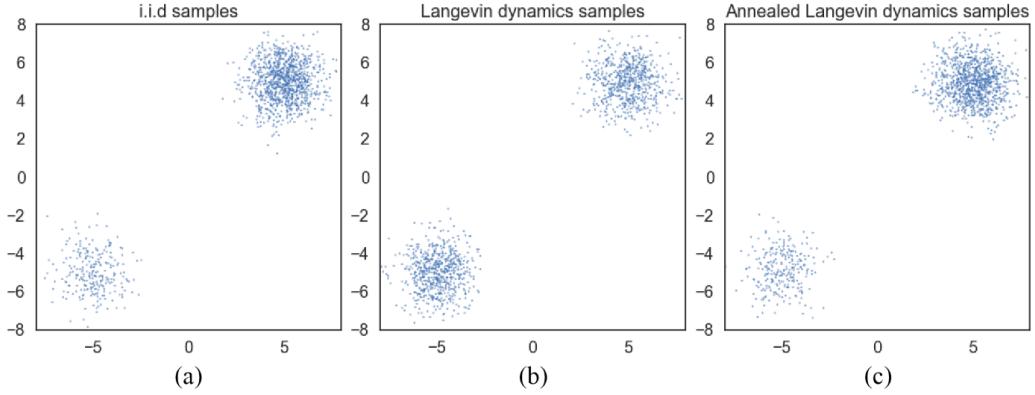
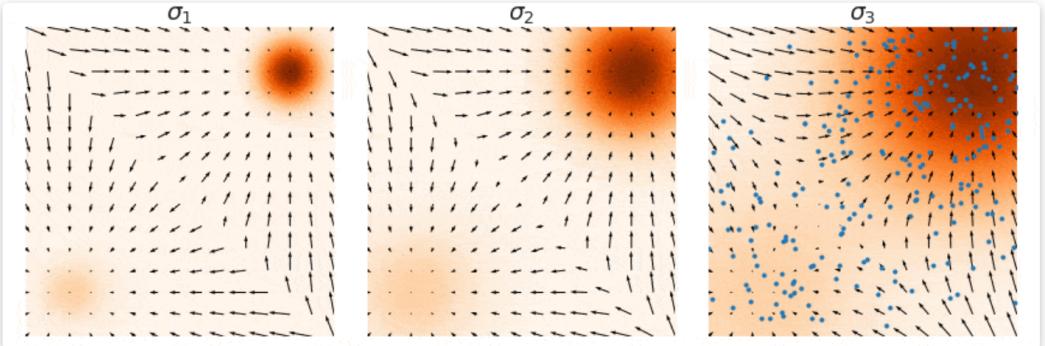


Figure 5: Slow Mixing of Langevin Dynamics Challenge. Langevin Dynamics suffers from the slow mixing problem, where the process may not recover the relative weights of modes are separated by low density regions of the data distribution. Sample i.i.d. from an arbitrary distribution (a), use Langevin dynamics to push samples towards modes, where the relative weights are the same, thus not converging to the true data distribution (b), and annealing the Langevin dynamics recovers the proper weights (c). Image credit: [1].



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

Figure 6: Annealed Langevin Dynamics. Learning noise-conditional score models and then using a sequence of Langevin dynamics processes, starting with the large noise score models and sequentially using the lower noise score models solves the three challenges of score based modeling. [1].

Model	Inception	FID
CIFAR-10 Unconditional		
PixelCNN [59]	4.60	65.93
PixelIQN [42]	5.29	49.46
EBM [12]	6.02	40.58
WGAN-GP [18]	7.86 ± .07	36.4
MoLM [45]	7.90 ± .10	18.9
SNGAN [36]	8.22 ± .05	21.7
ProgressiveGAN [25]	8.80 ± .05	-
NCSN (Ours)	8.87 ± .12	25.32
CIFAR-10 Conditional		
EBM [12]	8.30	37.9
SNGAN [36]	8.60 ± .08	25.5
BigGAN [6]	9.22	14.73

Table 1: Inception and FID scores for CIFAR-10

Figure 7: **Inception and FID Scores for CIFAR-10.** The noise conditional score network pipeline achieves the highest Inception score on the unconditional CIFAR-10 image generation task and a low FID score; these metrics are competitive with the best performing GANs and method-of-moments-trained networks (MoLM). Image credit: [1].

2.7 Results

Table 7 clearly shows the comparison results with different models evaluated by Inception and FID two metrics. Inception is a metric to evaluate the model performance and compare the accuracy and training time. FID is a measure of the similarity between two sets of images. It is commonly used as a metric to evaluate the performance of generative models. A lower FID score indicates that the generated images are more similar to the real images, and thus the generative model is considered to perform better.

This paper built models on three different datasets, which are MNIST, CelebA, and CIFAR-10, and Figure 8 shows the new generating samples generated from the generative model. Through the image results, we can clearly see that the generated model has a good effect, which is not much different from the data points in the original data set.

3 Conclusion of Paper Review

Estimating the data distribution, which is unknown and high-dimensional, is hard. It requires calculating the partition function is hard or intractable in general. Rather, estimating a score-based neural network model $s_\theta(x) \approx \nabla_x \log p(x)$ from data avoids this issue and is tractable via score-matching. Samples can be produced using Langevin dynamics, an MCMC procedure that samples from the data distribution $p(x)$ using only its score function with added noise. This process that can be interpreted as noisy gradient ascent. The noise here allows the samples to be distributed around the peaks instead of directly on the peaks. This allows better recovery of the true distribution. Noising the data and annealing (i.e. gradually lowering the learning rate of) the Langevin dynamics greatly assists the score estimation and sample generation process.

Thus, score-based generative modeling avoids calculation of partition function by relying on the Stein score function, benefits from injecting noise into the data, avoids the need to sample from a Markov chain during training, is performed in two, decoupled stages: (1) noising, (2)



(a) MNIST

(b) CelebA

(c) CIFAR-10

Figure 5: Uncurated samples on MNIST, CelebA, and CIFAR-10 datasets.

Figure 8: New samples on MNIST, CelebA, and CIFAR-10. The noise conditional score network and annealed Langevin dynamics pipeline produces new samples that resemble the data on which they were trained. Image credit: [1].

sample generation, and can be used to train energy-based models by using the gradient of an energy-based model as the score model.

4 Extending the Score-Based Generative Modeling Framework with a Multi-Resolution Diffusion Generative Model

Since adding noise improves score-based generative modeling performance, and adding more noise from the same family improves performance even more, then adding an infinite number of noise perturbations from that same family may improve the model performance even more. This idea of appealing to a continuous noising process places score-based generative models into the continuous diffusion with drift regime [2]. These models are referred to as generative diffusion models.

Generative diffusion models perform well, but training them and sampling from them can be slow. One way to improve the speed of training and performance is to introduce a multi-resolution inductive bias into the model. This multi-resolution inductive bias incorporates into the model the idea that images have objects of different scales and have locally varying statistics. Rather than train a single diffusion model on a single resolution of images, which is slow and can get stuck in local optima, we aim to train a single diffusion model for each resolution, where the coarser models inform the learning process of the higher resolution models, speeding up the learning process for higher resolutions and avoiding getting caught in inferior local optima. Furthermore, using a multi-resolution representation based on residual images could further improve the efficiency of the model by reducing redundancy of information at neighboring resolutions. A multi-resolution model could therefore benefit from the simplicity and efficiency of the low-resolution model as well as the efficiency and detailed accuracy of the higher-resolution residual models that provide insight into how the image information changes between resolutions.

5 Methods

9

5.1 Laplacian Pyramid Representation of Images

Let $\mathcal{X}_i := \mathbb{R}^{H_i \times W_i \times C}$ denote the space of C -channel images at resolution $H_i \times W_i$, where typically $H_{i+1} = H_i/2$ and $W_{i+1} = W_i/2$ (integer-valued) for $i = 0, \dots, k-1$. We define, for each level $i = 0, \dots, k-1$:

$$d_i : \mathcal{X}_i \rightarrow \mathcal{X}_{i+1} \quad (\text{downsampling: blur + decimate}), \quad u_i : \mathcal{X}_{i+1} \rightarrow \mathcal{X}_i \quad (\text{upsampling: interpolate + smooth})$$

(For example, d_i could be average pooling after blurring; u_i could be bilinear interpolation.)

Given a full-resolution image $I \in \mathcal{X}_0$, define the associated *Gaussian pyramid* $\{I_i\}_{i=0}^k$ and *Laplacian residuals* $\{r_i\}_{i=0}^{k-1}$ by

$$I_0 := I,$$

$$I_{i+1} := d_i(I_i), \quad i = 0, \dots, k-1,$$

$$r_i := I_i - u_i(I_{i+1}), \quad i = 0, \dots, k-1. \quad (8)$$

Each residual $r_i \in \mathcal{X}_i$ is a band-pass (high-frequency) image at scale i .

Laplacian pyramid transform and its inverse. Define the (*minimal*) *Laplacian pyramid transform*

$$L : \mathcal{X}_0 \rightarrow \mathcal{X}_k \times \prod_{i=0}^{k-1} \mathcal{X}_i,$$

$$L(I) := (I_k, r_{k-1}, r_{k-2}, \dots, r_0). \quad (9)$$

This representation is *not redundant*: it stores the coarsest image I_k and one residual per finer level.

Define the reconstruction (inverse) map R by the backward recurrence

$$\hat{I}_k := I_k,$$

$$\hat{I}_i := u_i(\hat{I}_{i+1}) + r_i, \quad i = k-1, k-2, \dots, 0,$$

$$R(I_k, r_{k-1}, \dots, r_0) := \hat{I}_0. \quad (10)$$

By construction, $R(L(I)) = I$ for all $I \in \mathcal{X}_0$ (exact inversion in the discrete setting where u_i, d_i are fixed operators).

Remark (relation to the redundant list notation). Sometimes one writes a redundant list $\{I_0, r_0, \dots, r_{k-1}, I_k\}$; however I_0 is recoverable from $\{r_0, \dots, r_{k-1}, I_k\}$ via (10).

There are many advantages to this representation: (i) the process is linear, (ii) it is invertible, (iii) it reduces redundancy since we work not only with the low frequency shape information but also with the high frequency texture information not present in low resolution images through the residual images, (iv) the residual images are sparse and should be easier to forward and backwards transform.

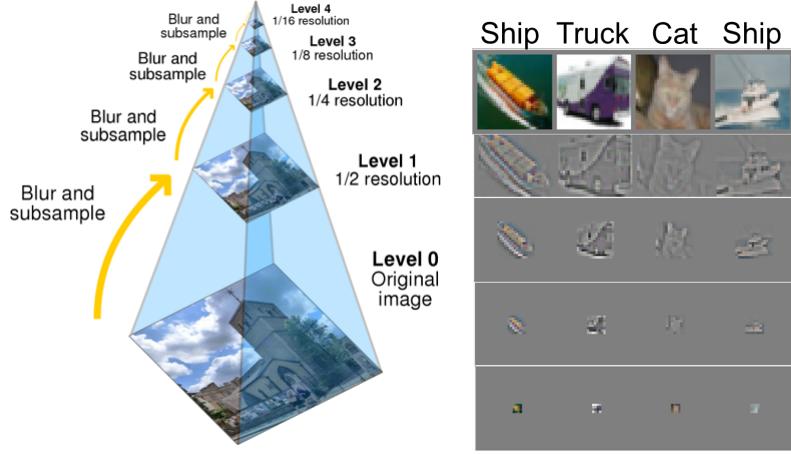


Figure 9: **Laplacian Pyramid Image Representation.** An illustration of the Laplacian pyramid (left) and a grid of Laplacian pyramid images representation of four images sampled from CIFAR-10 (right). Cartoon image credit: [https://en.wikipedia.org/wiki/Pyramid_\(image_processing\)#/media/File:Image_pyramid.svg](https://en.wikipedia.org/wiki/Pyramid_(image_processing)#/media/File:Image_pyramid.svg)

5.2 Laplacian Pyramid Diffusion

For each training sample $I \sim p_{\text{data}}$ we deterministically compute $(I_k, r_{k-1}, \dots, r_0) = L(I)$. A generative model should capture (i) a prior over the coarsest image I_k and (ii) the *conditional* distribution of each residual given the next-coarser content.

Coarse-to-fine probabilistic factorization. Let $\ell_i := u_i(I_{i+1}) \in \mathcal{X}_i$ denote the upsampled coarse image at level i (same resolution as r_i). We model

$$p_\theta(I_k, r_{k-1}, \dots, r_0) := p_{\theta,k}(I_k) \prod_{i=0}^{k-1} p_{\theta,i}(r_i | \ell_i), \quad \ell_i = u_i(I_{i+1}), \quad I_{i+1} = d_i(I_i). \quad (11a)$$

The induced distribution over full-resolution images is the *pushforward* of this joint through the deterministic reconstruction map R in (10):

$$I_0 := R(I_k, r_{k-1}, \dots, r_0), \quad I_0 \sim (R)_\# p_\theta. \quad (11b)$$

Importantly, we do *not* identify a marginal density $p(I_0)$ with a joint density $p(I_k, r_{k-1}, \dots, r_0)$; instead we define a coherent joint model on $(I_k, r_{k-1:0})$ and obtain I_0 by reconstruction.

Diffusion models at each level (with correct domains). For each level we define a diffusion on a variable z_i living in the appropriate space:

$$z_k := I_k \in \mathcal{X}_k, \quad z_i := r_i \in \mathcal{X}_i \text{ for } i = 0, \dots, k-1.$$

Let $d_i := H_i W_i C$ and identify $\mathcal{X}_i \cong \mathbb{R}^{d_i}$ by vectorization when writing SDEs.

For each $i \in \{0, \dots, k\}$ we choose drift and diffusion coefficients

$$f_i : \mathbb{R}^{d_i} \times [0, T] \rightarrow \mathbb{R}^{d_i}, \quad g_i : [0, T] \rightarrow \mathbb{R}_{>0},$$

and define the forward Itô SDE (unconditional in the state, but the *score network* may be conditional):

$$dz_i(t) = f_i(z_i(t), t) dt + g_i(t) dw_i(t).$$

Its reverse-time SDE takes the form

$$dz_i(t) = \left(f_i(z_i(t), t) - g_i(t)^2 \nabla_z \log p_{i,t}(z_i(t) | c_i) \right) dt + g_i(t) d\bar{w}_i(t),$$

where c_k is empty (no conditioning) and $c_i := \ell_i$ for $i = 0, \dots, k-1$.

Accordingly, we learn a time-dependent (conditional) score model

$$s_{\theta_i}(z, t; c_i) \approx \nabla_z \log p_{i,t}(z | c_i),$$

by minimizing the (conditional) denoising score matching objective

$$\theta_i^* = \arg \min_{\theta_i} \mathbb{E}_t \left[\lambda(t) \mathbb{E}_{(z_i(0), c_i)} \mathbb{E}_{z_i(t) | z_i(0)} \| s_{\theta_i}(z_i(t), t; c_i) - \nabla_{z_i(t)} \log p_{0,t}(z_i(t) | z_i(0)) \|_2^2 \right]. \quad (12)$$

Coarse-to-fine sampling (cascaded, but internally accelerated if desired). To sample a full-resolution image:

1. Sample I_k by solving the reverse-time sampler for the level- k diffusion.
2. For $i = k-1, \dots, 0$:
 - (a) Set $\ell_i := u_i(I_{i+1})$.
 - (b) Sample r_i by solving the reverse-time sampler for the level- i *conditional* diffusion with conditioning ℓ_i .
 - (c) Set $I_i := \ell_i + r_i$.
3. Return I_0 .

This procedure guarantees cross-scale consistency by construction (each residual is generated conditional on the coarser content that will be used in reconstruction).

5.3 Likelihood relations and score relations across pyramid variables

Let $x_0 \sim p_{\text{data}}$ be a random image at full resolution. Using the (deterministic) Laplacian pyramid operators, define for $i = 0, \dots, k-1$

$$x_{i+1} := d_i(x_i), \quad r_i := x_i - u_i(x_{i+1}).$$

Define the pyramid code as the tuple

$$z := (x_k, r_{k-1}, \dots, r_0) \in \mathcal{X}_k \times \prod_{i=0}^{k-1} \mathcal{X}_i, \quad z = L(x_0). \quad (13)$$

Important note (what we do and do not claim). Because L is deterministic and (for Laplacian pyramids) typically *not* a bijection between equal-dimensional Euclidean spaces, it is generally *not correct* to identify a marginal density $p(x_0)$ with a joint density $p(x_k, r_{k-1}, \dots, r_0)$. Instead, we define a coherent *joint generative model* over the pyramid variables z and obtain x_0 by deterministic reconstruction.

$$p_\theta(z) := p_{\theta,k}(x_k) \prod_{i=0}^{k-1} p_{\theta,i}(r_i \mid x_{i+1}), \quad x_{i+1} = d_i(x_i), \quad x_i = u_i(x_{i+1}) + r_i. \quad (14)$$

Equivalently, the log-likelihood of the pyramid code decomposes additively as

$$\log p_\theta(z) = \log p_{\theta,k}(x_k) + \sum_{i=0}^{k-1} \log p_{\theta,i}(r_i \mid x_{i+1}). \quad (15)$$

Score relations (where additivity is valid). Taking gradients of (15) yields *score decompositions with respect to the variables in z* . In particular, for each residual variable r_i we have the clean identity

$$\nabla_{r_i} \log p_\theta(z) = \nabla_{r_i} \log p_{\theta,i}(r_i \mid x_{i+1}), \quad i = 0, \dots, k-1, \quad (16)$$

and similarly $\nabla_{x_k} \log p_\theta(z) = \nabla_{x_k} \log p_{\theta,k}(x_k)$ up to additional terms if one allows the conditionals $p_{\theta,i}(r_i \mid x_{i+1})$ to depend on x_k through the deterministic recursion for x_{i+1} .

Induced distribution on full-resolution images (pushforward). Given a sample $z \sim p_\theta(z)$, define the reconstruction map R (as in Eq. (10)) and set

$$x_0 := R(z), \quad x_0 \sim (R)_\# p_\theta. \quad (17)$$

We emphasize that (15)–(16) provide rigorous likelihood/score relations for the modeled pyramid variables z . Obtaining an explicit closed-form density $p_\theta(x_0)$ or score $\nabla_{x_0} \log p_\theta(x_0)$ generally requires either: (i) a separate model directly on x_0 , or (ii) replacing L by an *invertible* (dimension-preserving) transform so that a change-of-variables formula applies.

Remark (Bayes rule view of forward vs. reverse diffusion.). For any diffusion variable z (e.g. $z = x_k$ or $z = r_i$) with forward Markov kernel $q(z_t \mid z_{t-\Delta})$, Bayes rule gives

$$q(z_{t-\Delta} \mid z_t, c) \propto q(z_t \mid z_{t-\Delta}) q(z_{t-\Delta} \mid c), \quad (18)$$

so the reverse-time dynamics can be written in terms of the *score* $\nabla_z \log q_t(z \mid c)$ of the perturbed (conditional) density. This is the continuous-time origin of the reverse-time SDE/ODE formulations used in score-based diffusion models.

5.4 Model Intuition

Intuitively, one can appeal to classical mechanical models to understand the behavior of diffusion under drift. For instance, imagine that you are throwing a ball to a friend when there is no wind. The ball is given an initial velocity and it shoots through the air (following a parabola) and lands in your friend's glove. This trajectory is the marginal ODE: because there is no wind or random forces acting on the ball to change its trajectory, the trajectory is determined completely by the initial velocity.

Now imagine that you throw the ball on a windy day. When there is a steady gust of wind acting on the ball while it travels through the air, the trajectory changes. The wind modifies the velocity of the ball by a constant magnitude and by a change in the direction. This is the marginal ODE with drift.

Now imagine that you throw the ball on a really windy day. Wind is gusting with different strengths in different directions, randomly perturbing the trajectory of the ball as it sails through the air. This is the stochastic ODE: the trajectory of the ball is no longer determined by the initial velocity alone but also on the random forces that act on the ball, changing its trajectory. How will you know where the ball will go if it is subject to randomness? The best you can say is the expected terminal point.¹³

To reverse the process, one must follow the negative trajectory backwards in time.

5.5 Training

Training can be performed in separate stages in parallel or end-to-end using a cascade model. In the separate stage model, the lowest resolution model is trained at the same time as the next highest level model, and so on, so each level is trained for its specific resolution and residual, while the weights of the previous level are fixed and separate from the other models. In the cascade model, the weights are not fixed but update with full passes through the network.

There are three experiments that can achieve these training pipelines:

- Train only the lowest resolution and compare.
- Train each level of the pyramid separately and compare.
- Train each level of the pyramid simultaneously in one model and compare.

However, in the work, we only accomplish the separate training of each level and generate images unconditionally (i.e. not conditioned on the previous levels of the pyramid).

5.6 Image Generation

Given a k -level Laplacian pyramid representation of the data, we can sample from the Gaussian noise distribution at each level and use the reverse Ito SDE to flow from the noise distribution to produce new samples at each level of the pyramid. This approach was done using a black-box ODE sampler. However, the conditional sampling approach could be achieved using a annealed Langevin dynamics conditioned on the previous levels of the pyramid. This approach in theory would sequentially build up a high resolution image by updating the sample state according the pyramid-level score model.

Algorithm 1 Annealed Conditional Langevin MCMC

```

Require:  $\sigma, \epsilon, T, k$ 
    Initialize  $\tilde{x}_0 \sim U(0, 1)$ 
    for  $i \leftarrow 1$  to  $k$  do                                 $\triangleright$  Loop over pyramid levels
         $\alpha_i \leftarrow \epsilon \cdot \sigma^2 / 2^i$            $\triangleright$  Anneal step size based on pyramid level
        for  $t \leftarrow 1$  to  $T$  do
            draw  $z_t \sim N(0, Id)$ 
             $\tilde{x}_t \leftarrow \tilde{x}_{t-1} + \frac{\alpha_i}{2} s_\theta(\tilde{x}_{t-1}, t - 1) + \sqrt{\alpha_i} z_t$ 
        end for
         $\tilde{x}_0 \leftarrow \tilde{x}_T$ 
    end for
    return  $\tilde{x}_T$ 

```

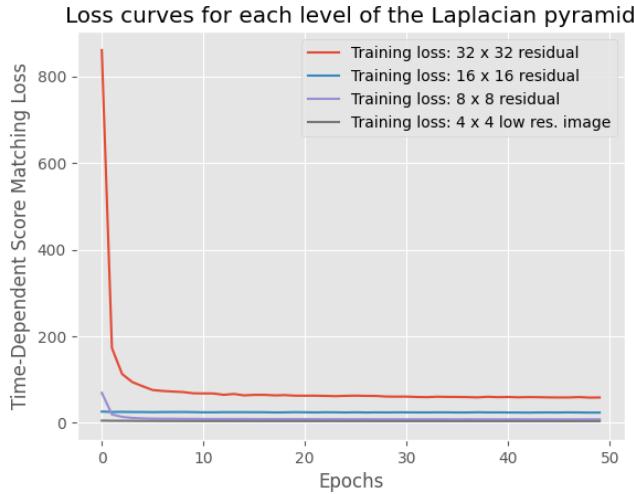


Figure 10: **Training loss curves for the Laplacian pyramid diffusion models.** Four distinct diffusion models, one model per level of the pyramid, were trained separately and the average time dependent score matching loss during training is plotted against the training time (50 epochs).

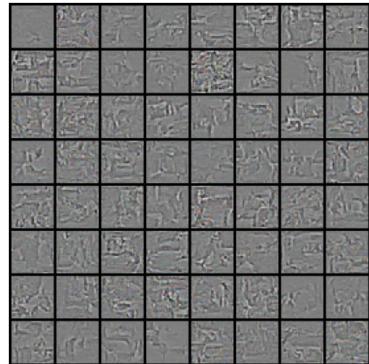
The sample produced at the final time should have the properties of the last level of the pyramid. In our case, this should be the full resolution image. This should produce high resolution samples quickly, since each iteration over the levels of the pyramid is conditioned on the final state of the sample at the last level of the pyramid.

5.7 Data

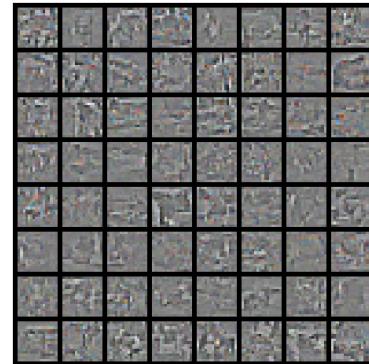
- Dataset: CIFAR-10 ($3 \times 32 \times 32$, 10 classes, 60k RGB images)
- Denoising: unconditional; score-matching
- Sampling: unconditional
- Task: image generation

5.8 Results on CIFAR-10: Image Generation

The results of training experiments number 1 and number 2 can be seen in Figures 10 and 11. The loss curves show that the models do learn the score function at each level of the pyramid, as indicated by the rapidly decreasing average loss values. The new samples produced are uncurated images of each level of the pyramid. However, the images are sampled unconditionally to one another. The Google Collab Notebook which modifies Dr. Yang Song’s tutorial to achieve these results can be found here: https://colab.research.google.com/drive/1EK8SpL60IRbU64iyP47tmARXhmrvS0Eq#scrollTo=zX1_hSXpK09R.



(a) 32x32 Residual Image Samples



(b) 16x16 Residual Image Samples



(c) 8x8 Residual Image Samples



(d) 4x4 Low Res. Image Samples

Figure 11: **New Laplacian pyramid samples on CIFAR-10.** New samples generated at each level of the Laplacian pyramid after training 4 distinct diffusion models in parallel and sampling using a black-box ODE solver.

5.9 Discussion and Future Work

16

Preliminary experiments are inconclusive, as only one distinct model was fit for each Laplacian pyramid image distribution, and sampling was performed unconditionally, i.e. distinct from the other levels of the pyramid. Future work will involve conditional sampling, conditioned on the previous level of the pyramid, and training a cascaded model, where training is conditionally performed as well as conditional sampling. Training the cascaded model may be accomplished by using a sum of denoising score matching terms, summed over the levels of the Laplacian pyramid, and conditional training may be performed by embedding a new layer into the score network for each level of the pyramid. Furthermore, the authors believe this model can be extended for class conditional sampling, by embedding class labels as a layer in the score network. Finally, the authors would like to use optimization and sampling methods from optimal control (e.g. shooting method for training, Hamiltonian Monte Carlo for sampling). These ideas are expounded below:

- Biased and conditional sampling

Each class should get its own drift which biases the trajectory, which could be considered analogous to the slope term in linear regression, and each class gets its own diffusion term. Also, the diffusion models share a background space and so the drifts and diffusions should be added.

$$\begin{cases} dx = \sum_{j=1}^C f_j(x, t)dt + \sum_{j=1}^C g_j(t)dw & (\text{Biased Forward Ito SDE}) \\ dx = \sum_{j=1}^C (f_j(x, t) - g_j^2(t)\nabla_x \log p_t(y | x))dt + \sum_{j=1}^C g_j(t)d\bar{w} & (\text{Biased Reverse Ito SDE}) \end{cases} \quad (5)$$

- Use Optimization Methods from Optimal Control

Examples that could be used at least in the marginal setting (i.e. drift only) are the single/multiple shooting methods or the adjoint sensitivity methods. The shooting methods could reduce the dimensionality of the problem by formulating the problem in terms of the initial parameters. This means samples can be generated by sampling from the distribution of initial momenta and shooting the initial sample through time and space according to the probability flow ODE.

- Hamiltonian Monte Carlo sampling conditioned on class and on images

- Introduce time-optimality condition in the cost function.

The diffusion problem has a solution on the infinite horizon in time. In practice, one can only approximate this with finite terminal time T . By leaving T unfixed, one can solve for the T that satisfies some time optimality condition in the optimization procedure. This provides a principled tuning knob balancing the trade-off between the accuracy of the solution and the time it takes to train the model.

- Introduce penalty terms on the initial and terminal conditions for tunable, approximate matching.

References

- [1] Yang Song and Stefano Ermon. “Generative modeling by estimating gradients of the data distribution”. In: *Advances in neural information processing systems* 32 (2019).

- [2] Yang Song et al. “Score-based generative modeling through stochastic differential equations”. In: *arXiv preprint arXiv:2011.13456* (2020).