

Diffusion Study Group #8

11/5/2022

Tanishq Abraham

Perception Prioritized Training of Diffusion Models

Elucidating the Design Space of Diffusion-Based Generative Models

Motivation

Diffusion models are a powerful framework for neural image synthesis and even useful in other domains. How can we further improve the image/distribution quality, training cost, and generation speed?

The literature is dense on theory, design choices are made to ensure the models are on a solid theoretical footing.

However, this approach has a danger of obscuring the available design space — a proposed model may appear as a tightly coupled package where no individual component can be modified without breaking the entire system

Expressing diffusion models in a common framework

Data distribution: $p_{\text{data}}(\mathbf{x})$ with standard deviation σ_{data}

Noisy data distributions: $p(\mathbf{x}; \sigma)$ (add Gaussian noise of standard deviation σ to the data)

For $\sigma_{\text{max}} \gg \sigma_{\text{data}}$, indistinguishable from noise

Start at $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{max}}^2 \mathbf{I})$ and denoise into images $\mathbf{x}_i \sim p(\mathbf{x}_i; \sigma_i)$ with noise levels $\sigma_0 = \sigma_{\text{max}} > \sigma_1 > \dots > \sigma_N = 0$

The endpoint \mathbf{x}_N is distributed according to the data

The ODE formulation

To treat the process as continuous, we define $\mathbf{x}_a \sim p(\mathbf{x}_a; \sigma(t_a))$ which is a sample at timestep t_a which can be evolved to $\mathbf{x}_b \sim p(\mathbf{x}_b; \sigma(t_b))$

The noise levels are now a continuous function $\sigma(t)$

This evolution can be performed by integrating the ODE:

$$d\mathbf{x} = -\dot{\sigma}(t)\sigma(t)\nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma(t))dt$$

where $\dot{\sigma}(t) = \frac{d\sigma(t)}{dt}$ and $\nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma(t))$ is the score function

A forward step ($t_a \rightarrow t_b$) nudges the sample away from the data, at a rate that depends on the change in noise level. The reverse is true for a backward step ($t_b \rightarrow t_a$).

Denoising score matching

We define $D(\mathbf{x}; \sigma)$ to be a denoiser function that is trained with the following objective:

$$\mathbb{E}_{\mathbf{y} \sim p_{\text{data}}} \mathbb{E}_{\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} \|D(\mathbf{y} + \mathbf{n}; \sigma) - \mathbf{y}\|_2^2$$

where \mathbf{y} is a training image, and \mathbf{n} is Gaussian noise. As we've discussed (ex: Tweedie's formula), the following is true:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma) = (D(\mathbf{x}; \sigma) - \mathbf{x}) / \sigma^2$$

The denoiser function is implemented as a neural network $D_{\theta}(\mathbf{x}; \sigma)$

Ideal Denoiser

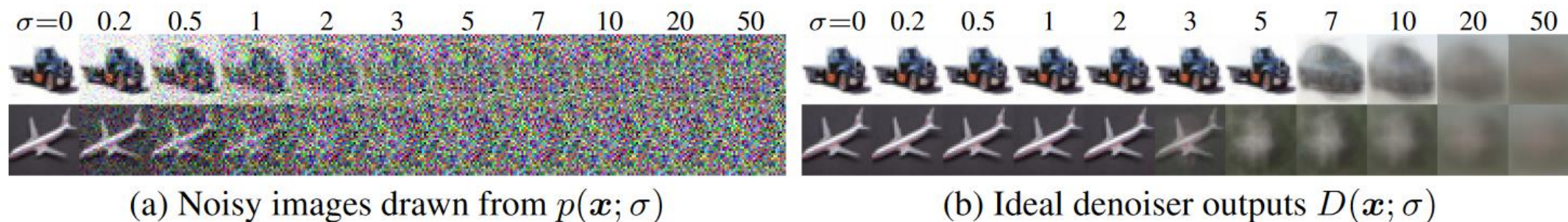


Figure 1: Denoising score matching on CIFAR-10. **(a)** Images from the training set corrupted with varying levels of additive Gaussian noise. High levels of noise lead to oversaturated colors; we normalize the images for cleaner visualization. **(b)** Optimal denoising result from minimizing Eq. 2 analytically (see Appendix B.3). With increasing noise level, the result approaches dataset mean.

Time-dependent scaling

Some methods (such as VP ODEs) introduce an additional scale $s(t)$ such that $\mathbf{x} = s(t)\hat{\mathbf{x}}$ is a scaled down version of the original non-scaled sample $\hat{\mathbf{x}}$. The more generalized ODE taking this into account is:

$$d\mathbf{x} = \left[\frac{\dot{s}(t)}{s(t)} \mathbf{x} - s(t)^2 \dot{\sigma}(t) \sigma(t) \nabla_{\mathbf{x}} \log p \left(\frac{\mathbf{x}}{s(t)}; \sigma(t) \right) dt \right]$$

Solution by discretization

Plugging the learned score function into the ODE, we can numerically solve the ODE in order to sample from our distribution.

Different ODE solvers (also called integration schemes) and discrete time steps (where $t_N = 0$) can be chosen

Usually Euler's method is used:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + (t_{i+1} - t_i)d_i$$

where d_i is the right-hand side of the ODE evaluated at timestep t_i

Expressing diffusion models in a common framework - sampling

		VP [49]	VE [49]	iDDPM [37] + DDIM [47]	Ours (“EDM”)
Sampling (Section 3)					
ODE solver		Euler	Euler	Euler	2 nd order Heun
Time steps	$t_{i < N}$	$1 + \frac{i}{N-1}(\epsilon_s - 1)$	$\sigma_{\max}^2 (\sigma_{\min}^2 / \sigma_{\max}^2)^{\frac{i}{N-1}}$	$u_{\lfloor j_0 + \frac{M-1-j_0}{N-1}i + \frac{1}{2} \rfloor}$, where $u_M = 0$ $u_{j-1} = \sqrt{\frac{u_j^2 + 1}{\max(\bar{\alpha}_{j-1}/\bar{\alpha}_j, C_1)}} - 1$	$\left(\sigma_{\max}^{\frac{1}{\rho}} + \frac{i}{N-1}(\sigma_{\min}^{\frac{1}{\rho}} - \sigma_{\max}^{\frac{1}{\rho}}) \right)^\rho$
Schedule	$\sigma(t)$	$\sqrt{e^{\frac{1}{2}\beta_d t^2 + \beta_{\min} t} - 1}$	\sqrt{t}	t	t
Scaling	$s(t)$	$1/\sqrt{e^{\frac{1}{2}\beta_d t^2 + \beta_{\min} t}}$	1	1	1

Improvements to deterministic sampling

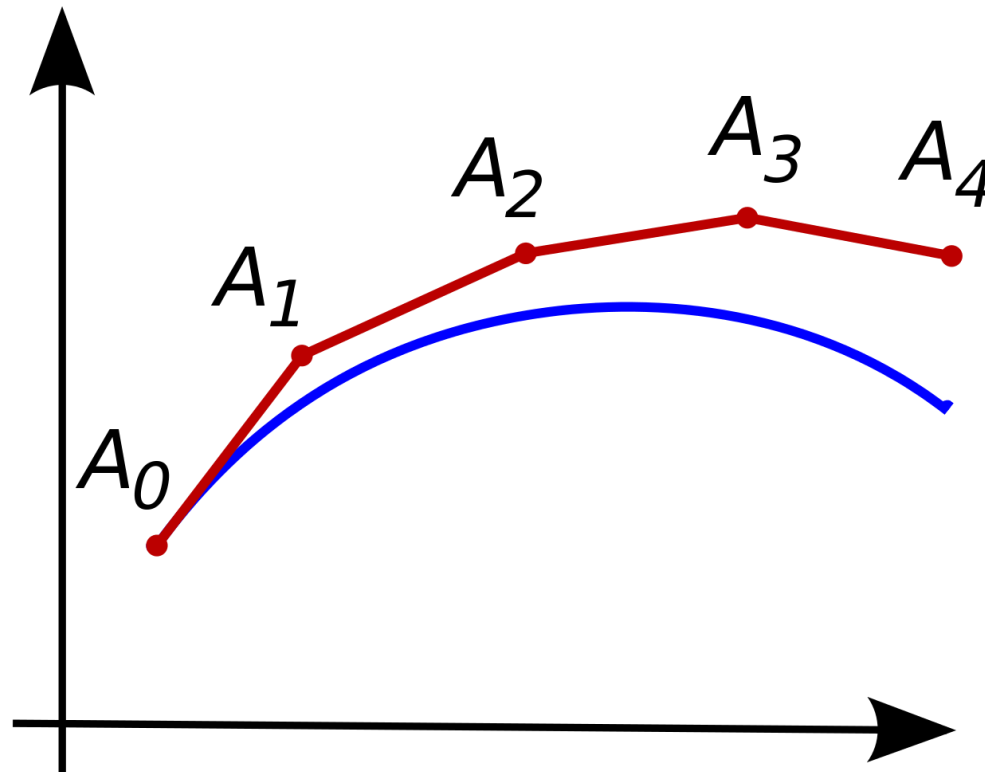
Choices related to the sampling process are largely independent of the other components, such as network architecture and training details.

From the viewpoint of the sampler, D_θ is a black box!

To demonstrate this, authors tried the original sampler, a re-implemented sampler, and a new sampler with a model trained under the VP ODE (Score SDE paper) VE ODE (Score SDE paper), and iDDPM (ADM paper) frameworks.

Discretization and higher-order integrators

Truncation error by ODEs (accumulates over timesteps):



Discretization and higher-order integrators

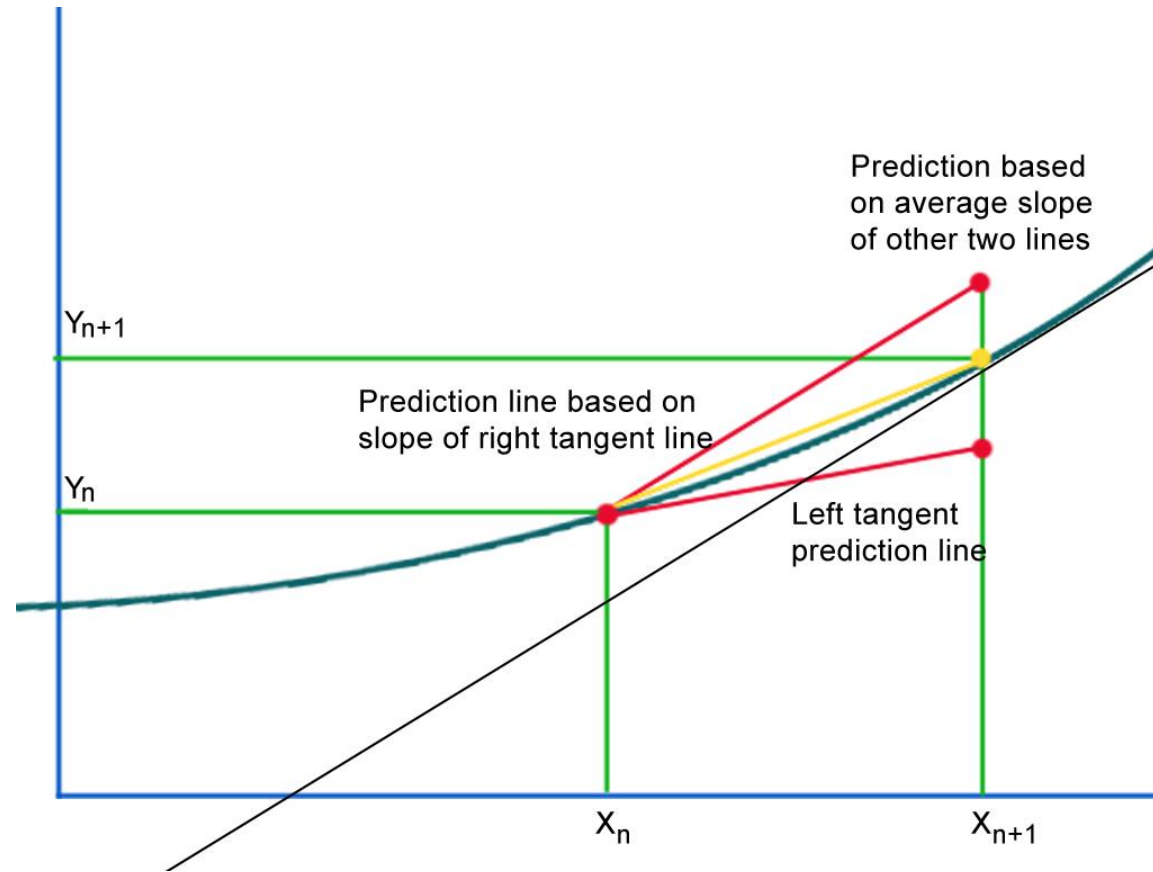
Euler's method - local truncation error scales with $\mathcal{O}(h^2)$

Heun's method (second-order method) – local truncation error scales with $\mathcal{O}(h^3)$

$$\begin{aligned}\tilde{\mathbf{x}}_{i+1} &= \mathbf{x}_i + (t_{i+1} - t_i) d_i(\mathbf{x}_i, t_i) \\ d'_i &= d_i(\tilde{\mathbf{x}}_{i+1}, t_{i+1}) \\ x_{i+1} &= x_i + \frac{(t_{i+1} - t_i)}{2} (d_i + d'_i)\end{aligned}$$

Adds an additional evaluation of D_θ per step

Heun's method



Heun's method

Algorithm 1 Deterministic sampling using Heun's 2nd order method with arbitrary $\sigma(t)$ and $s(t)$.

```
1: procedure HEUNSAMPLER( $D_\theta(\mathbf{x}; \sigma)$ ,  $\sigma(t)$ ,  $s(t)$ ,  $t_{i \in \{0, \dots, N\}}$ )
2:   sample  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \sigma^2(t_0) s^2(t_0) \mathbf{I})$  ▷ Generate initial sample at  $t_0$ 
3:   for  $i \in \{0, \dots, N-1\}$  do ▷ Solve Eq. 4 over  $N$  time steps
4:      $\mathbf{d}_i \leftarrow \left( \frac{\dot{\sigma}(t_i)}{\sigma(t_i)} + \frac{\dot{s}(t_i)}{s(t_i)} \right) \mathbf{x}_i - \frac{\dot{\sigma}(t_i)s(t_i)}{\sigma(t_i)} D_\theta \left( \frac{\mathbf{x}_i}{s(t_i)}; \sigma(t_i) \right)$  ▷ Evaluate  $d\mathbf{x}/dt$  at  $t_i$ 
5:      $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + (t_{i+1} - t_i) \mathbf{d}_i$  ▷ Take Euler step from  $t_i$  to  $t_{i+1}$ 
6:     if  $\sigma(t_{i+1}) \neq 0$  then ▷ Apply 2nd order correction unless  $\sigma$  goes to zero
7:        $\mathbf{d}'_i \leftarrow \left( \frac{\dot{\sigma}(t_{i+1})}{\sigma(t_{i+1})} + \frac{\dot{s}(t_{i+1})}{s(t_{i+1})} \right) \mathbf{x}_{i+1} - \frac{\dot{\sigma}(t_{i+1})s(t_{i+1})}{\sigma(t_{i+1})} D_\theta \left( \frac{\mathbf{x}_{i+1}}{s(t_{i+1})}; \sigma(t_{i+1}) \right)$  ▷ Eval.  $d\mathbf{x}/dt$  at  $t_{i+1}$ 
8:        $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + (t_{i+1} - t_i) \left( \frac{1}{2} \mathbf{d}_i + \frac{1}{2} \mathbf{d}'_i \right)$  ▷ Explicit trapezoidal rule at  $t_{i+1}$ 
9:   return  $\mathbf{x}_N$  ▷ Return noise-free sample at  $t_N$ 
```

Choosing time steps

Truncation error is very large for low noise levels and smaller for high noise levels (confirmed through empirical experiments in the appendix)

This indicates the best selection of step sizes should be if it decreased monotonically with decreasing σ . Therefore, $t_i = \sigma^{-1}(\sigma_i)$ and we set:

$$\sigma_{i < N} = \left(\sigma_{\max}^{\frac{1}{\rho}} + \frac{i}{N-1} (\sigma_{\min}^{\frac{1}{\rho}} - \sigma_{\max}^{\frac{1}{\rho}}) \right)^{\rho} \quad \text{and} \quad \sigma_N = 0.$$

ρ controls how much the steps near σ_{\min} are shortened at the expense of longer steps near σ_{\max} .

$\rho = 3$ equalizes the truncation error over all steps, but $5 < \rho = 7 < 10$ is better for sampling images. This suggests that errors near σ_{\min} have a large impact.

Trajectory curvature and noise schedule

The shape of the ODE solution trajectories is defined by $s(t), \sigma(t)$

The best choice is $\sigma(t) = t, s(t) = 1$, which is also used in DDIM. The ODE becomes:

$$\frac{d\mathbf{x}}{dt} = \frac{(\mathbf{x} - D(\mathbf{x}; t))}{t} \rightarrow d\mathbf{x} = \frac{(\mathbf{x} - D(\mathbf{x}; t))}{t} dt$$

A single Euler step to $t = 0$ is the denoised image. The tangent of the solution trajectory always points to the denoiser output.

Trajectory curvature and noise schedule

This ODE provides largely linear trajectories

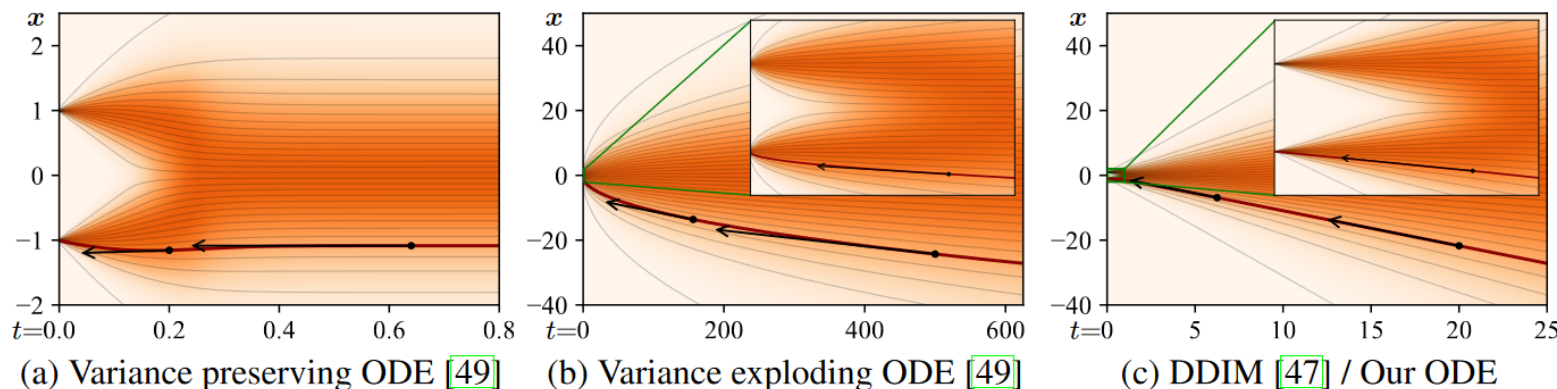


Figure 3: A sketch of ODE curvature in 1D where p_{data} is two Dirac peaks at $x = \pm 1$. Horizontal t axis is chosen to show $\sigma \in [0, 25]$ in each plot, with insets showing $\sigma \in [0, 1]$ near the data. Example local gradients are shown with black arrows. **(a)** Variance preserving ODE of Song et al. [49] has solution trajectories that flatten out to horizontal lines at large σ . Local gradients start pointing towards data only at small σ . **(b)** Variance exploding variant has extreme curvature near data and the solution trajectories are curved everywhere. **(c)** With the schedule used by DDIM [47] and us, as σ increases the solution trajectories approach straight lines that point towards the mean of data. As $\sigma \rightarrow 0$, the trajectories become linear and point towards the data manifold.

Results – improvements to deterministic sampling

The re-implementations of the deterministic samplers in the unified framework resulted in consistently better results.

“The differences are explained by certain oversights in the original implementations as well as our more careful treatment of discrete noise levels in the case of DDIM”

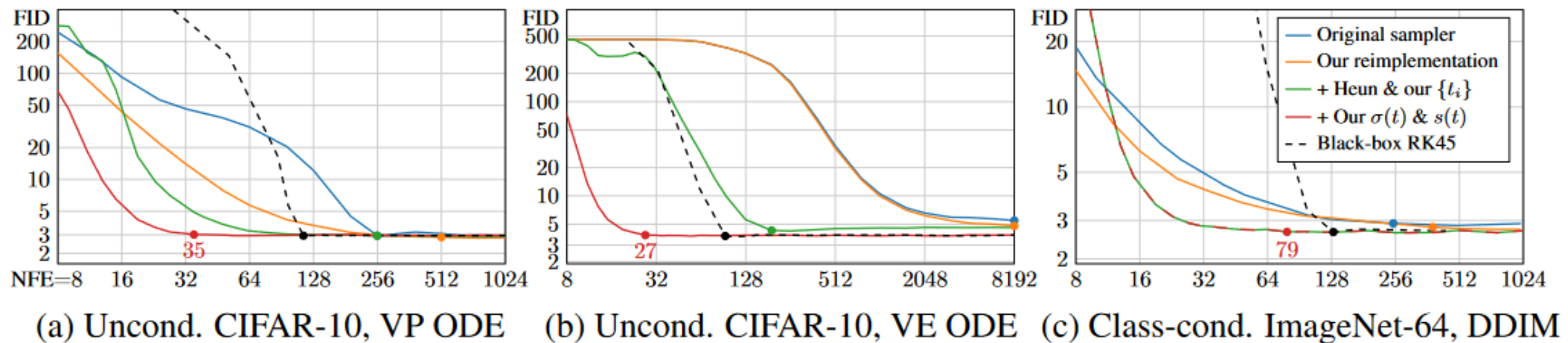


Figure 2: Comparison of deterministic sampling methods using three pre-trained models. For each curve, the dot indicates the lowest NFE whose FID is within 3% of the lowest observed FID.

Stochastic sampling

While they describe the same distributions, ODEs result in lower output quality compared to SDEs, so what is the role of stochasticity in practice?

Let's examine the SDE where $d\mathbf{x}_+$ is for the forward process and $d\mathbf{x}_-$ is for the reverse process

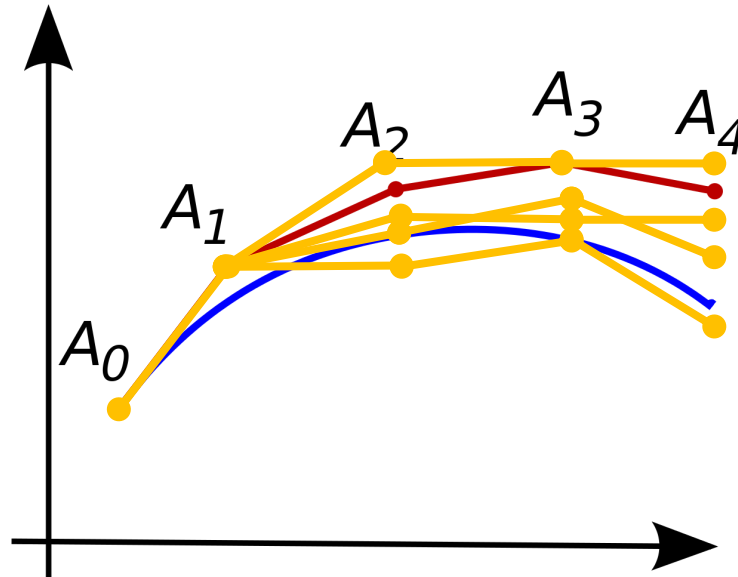
$$d\mathbf{x}_{\pm} = \underbrace{-\dot{\sigma}(t)\sigma(t)\nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma(t)) dt}_{\text{probability flow ODE (Eq. 1)}} \pm \underbrace{\beta(t)\sigma(t)^2 \nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma(t)) dt}_{\text{deterministic noise decay}} + \underbrace{\sqrt{2\beta(t)}\sigma(t) d\omega_t}_{\text{noise injection}},$$

Langevin diffusion SDE

Song et al. uses $\beta(t) = \dot{\sigma}(t)/\sigma(t)$

Stochastic sampling

The implicit Langevin diffusion drives the sample towards the desired marginal distribution at a given time, actively correcting for any errors made in earlier sampling steps. On the other hand, approximating the Langevin term with discrete SDE solver steps introduces error in itself.



Proposed stochastic sampler

Go to higher noise level $\hat{t}_i = t_i + \gamma_i t_i$ and get a noisier sample

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \sqrt{\hat{t}_i^2 - t_i^2} \boldsymbol{\epsilon}_i$$

Numerically solve the ODE back to t_{i+1}

Algorithm 2 Our stochastic sampler with $\sigma(t) = t$ and $s(t) = 1$.

```

1: procedure STOCHASTICSAMPLER( $D_\theta(\mathbf{x}; \sigma)$ ,  $t_{i \in \{0, \dots, N\}}$ ,  $\gamma_{i \in \{0, \dots, N-1\}}$ ,  $S_{\text{noise}}$ )
2:   sample  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, t_0^2 \mathbf{I})$ 
3:   for  $i \in \{0, \dots, N-1\}$  do
4:     sample  $\boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, S_{\text{noise}}^2 \mathbf{I})$ 
5:      $\hat{t}_i \leftarrow t_i + \gamma_i t_i$ 
6:      $\hat{\mathbf{x}}_i \leftarrow \mathbf{x}_i + \sqrt{\hat{t}_i^2 - t_i^2} \boldsymbol{\epsilon}_i$ 
7:      $\mathbf{d}_i \leftarrow (\hat{\mathbf{x}}_i - D_\theta(\hat{\mathbf{x}}_i; \hat{t}_i)) / \hat{t}_i$ 
8:      $\mathbf{x}_{i+1} \leftarrow \hat{\mathbf{x}}_i + (t_{i+1} - \hat{t}_i) \mathbf{d}_i$ 
9:     if  $t_{i+1} \neq 0$  then
10:       $\mathbf{d}'_i \leftarrow (\mathbf{x}_{i+1} - D_\theta(\mathbf{x}_{i+1}; t_{i+1})) / t_{i+1}$ 
11:       $\mathbf{x}_{i+1} \leftarrow \hat{\mathbf{x}}_i + (t_{i+1} - \hat{t}_i) (\frac{1}{2} \mathbf{d}_i + \frac{1}{2} \mathbf{d}'_i)$ 
12:   return  $\mathbf{x}_N$ 

```

$\triangleright \gamma_i = \begin{cases} \min\left(\frac{S_{\text{churn}}}{N}, \sqrt{2}-1\right) & \text{if } t_i \in [S_{\text{tmin}}, S_{\text{tmax}}] \\ 0 & \text{otherwise} \end{cases}$

\triangleright Select temporarily increased noise level \hat{t}_i
 \triangleright Add new noise to move from t_i to \hat{t}_i
 \triangleright Evaluate $d\mathbf{x}/dt$ at \hat{t}_i
 \triangleright Take Euler step from \hat{t}_i to t_{i+1}

\triangleright Apply 2nd order correction

Proposed stochastic sampler

Similar to the predictor-corrector sampler of Song et al.

Comparison to Euler-Maruyama (?)

Increasing the amount of stochasticity is effective in correcting errors made by earlier sampling steps but has drawbacks.

Excessive stochasticity results in gradual loss of detail in generations

There is also a drift toward oversaturated colors at very low and high noise levels.

Degradation is likely due to flaws in D_θ that violate the premises of Langevin diffusion

Proposed stochastic sampler

Empirical hacks used like:

- Only enabling stochasticity for $t_i \in [S_{\text{tmin}}, S_{\text{tmax}}]$
- Define $\gamma_i = S_{\text{churn}}/N$ where S_{churn} controls the overall amount of stochasticity
- Clamp γ_i to never introduce more new noise than what is already present in the image (?)
- Set S_{noise} slightly above 1 (inflate standard deviation of noise added)

The optimal value for these parameters were found by grid search on a case-by-case basis.

“[...]using stochastic sampling as the primary means of evaluating model improvements may inadvertently end up influencing the design choices related to model architecture and training.”

Proposed stochastic sampler

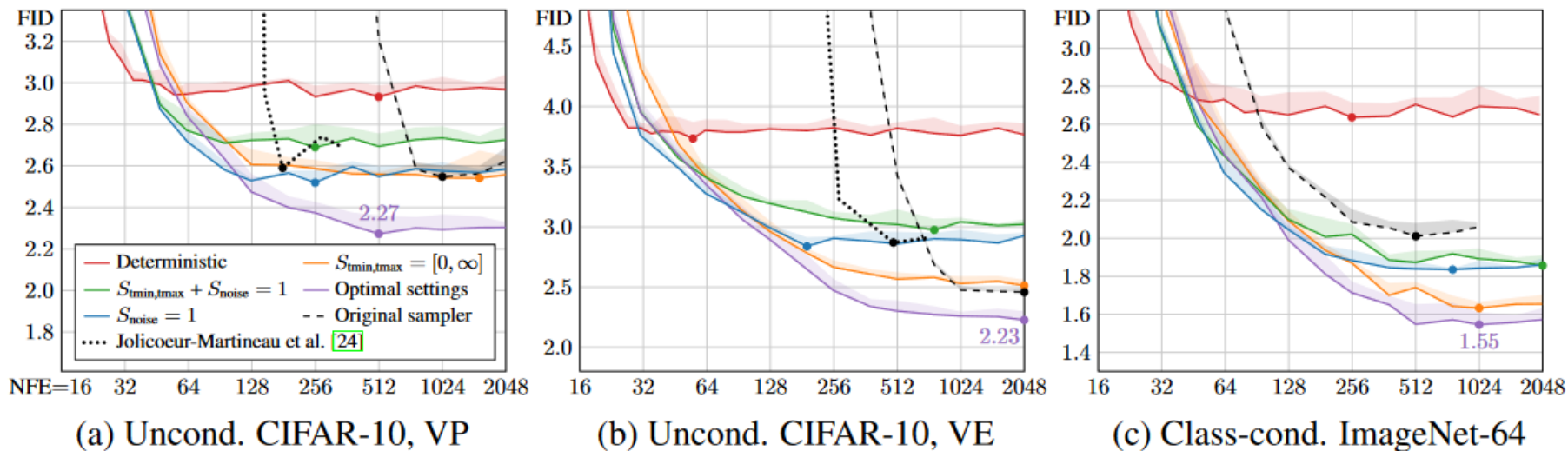


Figure 4: Evaluation of our stochastic sampler (Algorithm 2). The purple curve corresponds to optimal choices for $\{S_{\text{churn}}, S_{\text{tmin}}, S_{\text{tmax}}, S_{\text{noise}}\}$; orange, blue, and green correspond to disabling the effects of $S_{\text{tmin}, \text{tmax}}$ and/or S_{noise} . The red curves show reference results for our deterministic sampler (Algorithm 1), equivalent to setting $S_{\text{churn}} = 0$. The dashed black curves correspond to the original stochastic samplers from previous work: Euler–Maruyama [49] for VP, predictor–corrector [49] for VE, and iDDPM [37] for ImageNet-64. The dots indicate lowest observed FID.

Preconditioning and Training Objective

It is advisable to keep input and output signal magnitudes fixed e.g. unit variance

Modeling D directly is therefore not ideal as the input $\mathbf{y} = \mathbf{x} + \mathbf{n}$ and output \mathbf{x} magnitudes depend on the noise level σ

Most previous approaches train F_θ to predict \mathbf{n} scaled to unit variance, giving us $D_\theta(\mathbf{x}; \sigma) = \mathbf{x} - \sigma F_\theta(\mathbf{x}; \sigma)$

However, at large σ , the predicted scaled \mathbf{n} needs to be predicted exactly, as any minor errors will be amplified by a factor of σ . In this case, directly doing $D_\theta(\mathbf{x}; \sigma)$ might be optimal.

Preconditioning and Training Objective

An alternative denoiser that preconditions and post-processes the neural network $F_\theta(\cdot)$ is proposed in this work:

$$D_\theta(\mathbf{x}; \sigma) = c_{\text{skip}}(\sigma) \mathbf{x} + c_{\text{out}}(\sigma) F_\theta(c_{\text{in}}(\sigma) \mathbf{x}; c_{\text{noise}}(\sigma))$$

This results in the following training objective:

$$\mathbb{E}_{\sigma, \mathbf{y}, \mathbf{n}} \left[\underbrace{\lambda(\sigma) c_{\text{out}}(\sigma)^2}_{\text{effective weight}} \left\| \underbrace{F_\theta(c_{\text{in}}(\sigma) \cdot (\mathbf{y} + \mathbf{n}); c_{\text{noise}}(\sigma))}_{\text{network output}} - \underbrace{\frac{1}{c_{\text{out}}(\sigma)} (\mathbf{y} - c_{\text{skip}}(\sigma) \cdot (\mathbf{y} + \mathbf{n}))}_{\text{effective training target}} \right\|_2^2 \right].$$

Main difference: $\sigma \sim p_{\text{train}}(\sigma)$ and $\lambda(\sigma)$ weights different noise levels in the loss

$c_{\text{in}}(\sigma)$ is selected such that the networks inputs have unit variance, $c_{\text{out}}(\sigma)$ is selected such that the network outputs have unit variance, and $c_{\text{skip}}(\sigma)$ is selected to minimize error amplification

Preconditioning and Training Objective

Comparing over different models:

Network and preconditioning (Section 5)				
Architecture of F_θ	DDPM++	NCSN++	DDPM	(any)
Skip scaling $c_{\text{skip}}(\sigma)$	1	1	1	$\sigma_{\text{data}}^2 / (\sigma^2 + \sigma_{\text{data}}^2)$
Output scaling $c_{\text{out}}(\sigma)$	$-\sigma$	σ	$-\sigma$	$\sigma \cdot \sigma_{\text{data}} / \sqrt{\sigma_{\text{data}}^2 + \sigma^2}$
Input scaling $c_{\text{in}}(\sigma)$	$1/\sqrt{\sigma^2 + 1}$	1	$1/\sqrt{\sigma^2 + 1}$	$1/\sqrt{\sigma^2 + \sigma_{\text{data}}^2}$
Noise cond. $c_{\text{noise}}(\sigma)$	$(M - 1) \sigma^{-1}(\sigma)$	$\ln(\frac{1}{2}\sigma)$	$M - 1 - \arg \min_j u_j - \sigma $	$\frac{1}{4} \ln(\sigma)$
Training (Section 5)				
Noise distribution	$\sigma^{-1}(\sigma) \sim \mathcal{U}(\epsilon_t, 1)$	$\ln(\sigma) \sim \mathcal{U}(\ln(\sigma_{\min}), \ln(\sigma_{\max}))$	$\sigma = u_j, j \sim \mathcal{U}\{0, M - 1\}$	$\ln(\sigma) \sim \mathcal{N}(P_{\text{mean}}, P_{\text{std}}^2)$
Loss weighting $\lambda(\sigma)$	$1/\sigma^2$	$1/\sigma^2$	$1/\sigma^2$ (note: *)	$(\sigma^2 + \sigma_{\text{data}}^2) / (\sigma \cdot \sigma_{\text{data}})^2$
Parameters	$\beta_d = 19.9, \beta_{\min} = 0.1$	$\sigma_{\min} = 0.02$	$\bar{\alpha}_j = \sin^2(\frac{\pi}{2} \frac{j}{M(C_2+1)})$	$\sigma_{\min} = 0.002, \sigma_{\max} = 80$
	$\epsilon_s = 10^{-3}, \epsilon_t = 10^{-5}$	$\sigma_{\max} = 100$	$C_1 = 0.001, C_2 = 0.008$	$\sigma_{\text{data}} = 0.5, \rho = 7$
	$M = 1000$		$M = 1000, j_0 = 8^\dagger$	$P_{\text{mean}} = -1.2, P_{\text{std}} = 1.2$

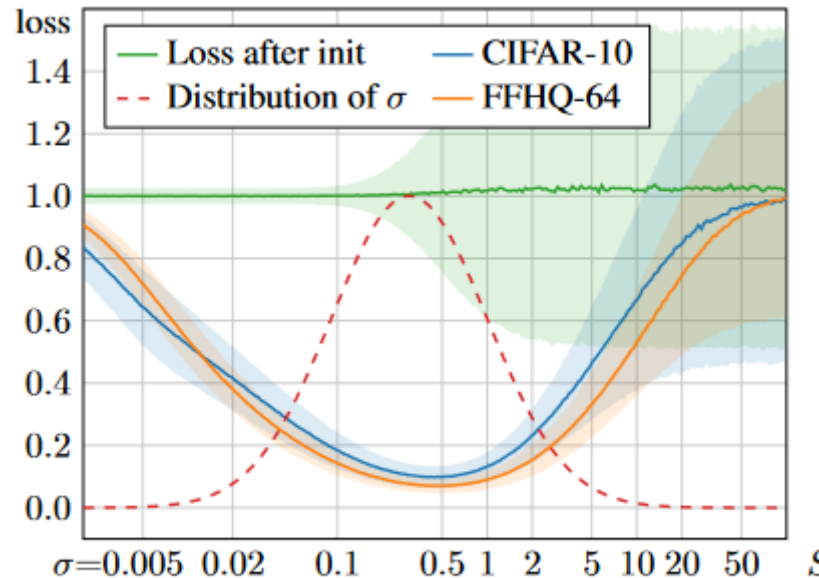
* iDDPM also employs a second loss term L_{vlb}

† In our tests, $j_0 = 8$ yielded better FID than $j_0 = 0$ used by iDDPM

Loss weighting and sampling

Balance the magnitude of the loss over different noise levels: $\lambda(\sigma) = 1/c_{\text{out}}(\sigma)^2$

The intermediate noise levels matter most so a log-normal distribution is chosen



(a) Loss & noise distribution

Augmentation

Augmentation is often an issue for smaller dataset

Borrow GAN augmentation pipeline (StyleGAN-ADA)

Use 6 geometric transformations, other transformations seem to be harmful for diffusion models.

The model is now conditioned on the augmentation during training and set to no augmentation conditioning during inference

Table 6: Our augmentation pipeline. Each training image undergoes a combined geometric transformation based on 8 random parameters that receive non-zero values with a certain probability. The model is conditioned with an additional 9-dimensional input vector derived from these parameters.

Augmentation	Transformation	Parameters	Prob.	Conditioning	Constants
x -flip	SCALE2D($1 - 2a_0$, 1)	$a_0 \sim \mathcal{U}\{0, 1\}$	100%	a_0	$A_{\text{prob}} = 12\%$
y -flip	SCALE2D(1, $1 - 2a_1$)	$a_1 \sim \mathcal{U}\{0, 1\}$	A_{prob}	a_1	or 15%
Scaling	SCALE2D($(A_{\text{scale}})^{a_2}$, $(A_{\text{scale}})^{a_2}$)	$a_2 \sim \mathcal{N}(0, 1)$	A_{prob}	a_2	$A_{\text{scale}} = 2^{0.2}$
Rotation	ROTATE2D($-a_3$)	$a_3 \sim \mathcal{U}(-\pi, \pi)$	A_{prob}	$\cos a_3 - 1$ $\sin a_3$	
Anisotropy	ROTATE2D(a_4) SCALE2D($(A_{\text{aniso}})^{a_5}$, $1/(A_{\text{aniso}})^{a_5}$) ROTATE2D($-a_4$)	$a_4 \sim \mathcal{U}(-\pi, \pi)$ $a_5 \sim \mathcal{N}(0, 1)$	A_{prob}	$a_5 \cos a_4$ $a_5 \sin a_4$	$A_{\text{aniso}} = 2^{0.2}$
Translation	TRANSLATE2D($(A_{\text{trans}})a_6$, $(A_{\text{trans}})a_7$)	$a_6 \sim \mathcal{N}(0, 1)$ $a_7 \sim \mathcal{N}(0, 1)$	A_{prob}	a_6 a_7	$A_{\text{trans}} = 1/8$

Results – training improvements

Preconditioning doesn't improve FID per se, but makes training more robust

P2 weighting provides improvement but only improves $\lambda(\sigma)$ and not $p_{\text{train}}(\sigma)$

New SOTA on CIFAR10 (VP – DDPM++, VE – NCSN++)

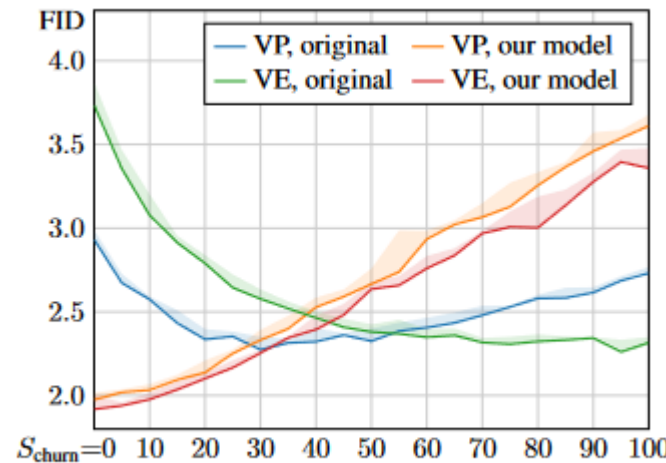
Table 2: Evaluation of our training improvements. The starting point (config A) is VP & VE using our **deterministic** sampler. At the end (configs E,F), VP & VE only differ in the architecture of F_θ .

Training configuration	CIFAR-10 [29] at 32×32				FFHQ [27] 64×64		AFHQv2 [7] 64×64	
	Conditional		Unconditional		Unconditional		Unconditional	
	VP	VE	VP	VE	VP	VE	VP	VE
A Baseline [49] (*pre-trained)	2.48	3.11	3.01*	3.77*	3.39	25.95	2.58	18.52
B + Adjust hyperparameters	2.18	2.48	2.51	2.94	3.13	22.53	2.43	23.12
C + Redistribute capacity	2.08	2.52	2.31	2.83	2.78	41.62	2.54	15.04
D + Our preconditioning	2.09	2.64	2.29	3.10	2.94	3.39	2.79	3.81
E + Our loss function	1.88	1.86	2.05	1.99	2.60	2.81	2.29	2.28
F + Non-leaky augmentation	1.79	1.79	1.97	1.98	2.39	2.53	1.96	2.16
NFE	35	35	35	35	79	79	79	79

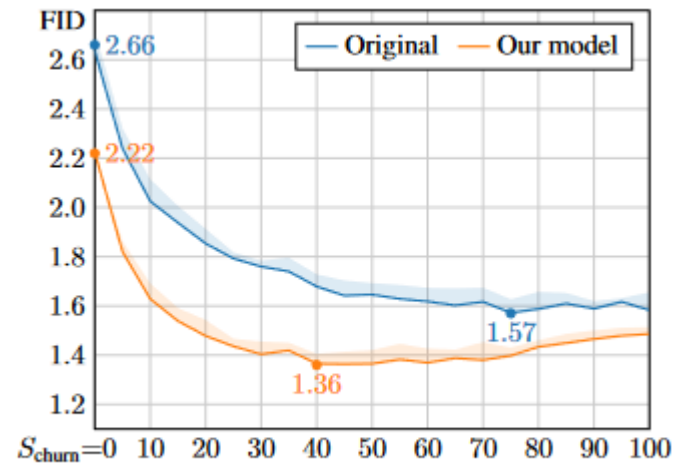
Results – training improvements

On CIFAR10, stochastic sampling worsens results! On ImageNet 64x64 stochastic sampling helps, indicating benefit of stochastic sampling for more diverse datasets

New SOTA on ImageNet 64x64 (ADM architecture, no aug reg)



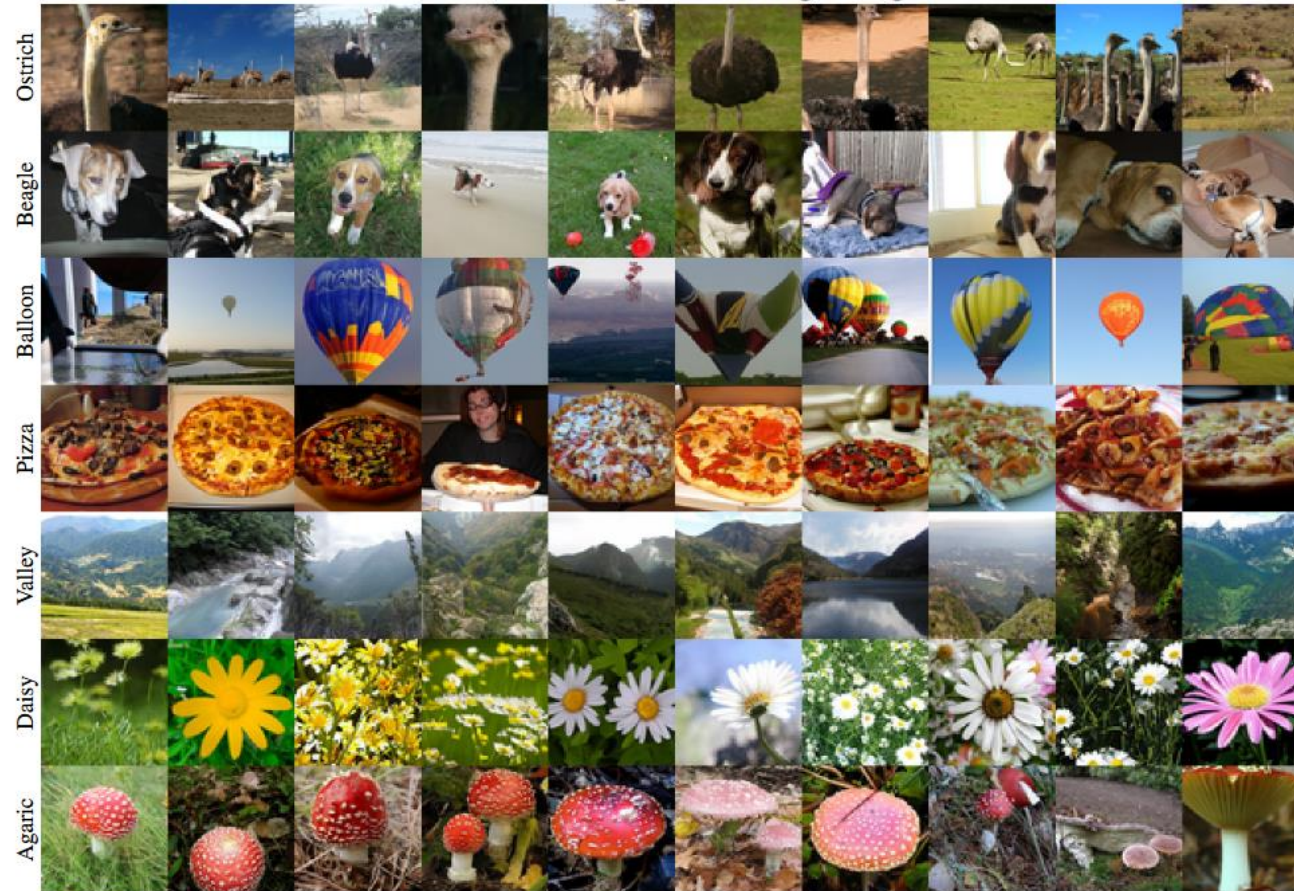
(b) Stochasticity on CIFAR-10



(c) Stochasticity on ImageNet-64

Generations

Stochastic, Our sampler & training configuration



FID 1.36 NFE 511