

# Distilled Diffusion Models

July 3, 2024

## Contents

<b>1</b>	<b>Consistency Models</b>	<b>2</b>
1.1	Motivation . . . . .	2
1.2	Consistency function . . . . .	3
1.3	Consistency via distillation . . . . .	3
1.3.1	Training details . . . . .	5
1.4	The difference between DM and CM . . . . .	6
1.5	Consistency from scratch . . . . .	7
<b>2</b>	<b>Improved CT</b>	<b>8</b>
2.1	EMA . . . . .	8
2.2	Technical improvements . . . . .	11
2.2.1	Pseudo-Huber . . . . .	11
<b>3</b>	<b>Multistep CM</b>	<b>12</b>
3.1	Multistep CD . . . . .	12
3.1.1	VP case . . . . .	12
3.1.2	VE case . . . . .	12
<b>4</b>	<b>Pagoda</b>	<b>14</b>
<b>5</b>	<b>Score distillation</b>	<b>15</b>
5.1	DMD . . . . .	15
5.2	ADD . . . . .	16
5.3	Score identity distillation . . . . .	17
5.3.1	1D Gaussian . . . . .	17

# 1 Consistency Models

## 1.1 Motivation

Why Diffusion Models cannot generate plausible samples in the few-step (1-5) regime?

We consider the following reverse process to generate data using Diffusion Model:

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\mathbf{w}_t, \quad t \in [1, 0]. \quad (1)$$

If the forward process satisfies  $p_t(\mathbf{x}|\mathbf{x}_0) = \mathcal{N}(\mathbf{x} | \mathbf{x}_0, s_t^2)$ , where  $\mathbf{x}_0$  is a sample from data distribution. Then, the score function has the following form:

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = -\frac{1}{s_t^2} (\mathbf{x} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}]). \quad (2)$$

We can find the form of the conditional expectation assuming that

$$p_{\text{data}}(\mathbf{x}_0) = \frac{1}{N} \sum_{j=1}^N \delta(\mathbf{x}_0 - \mathbf{x}_0^j). \quad (3)$$

Then,

$$\begin{aligned} \mathbb{E}[\mathbf{x}_0 | \mathbf{x}] &= \int \mathbf{x}_0 p_t(\mathbf{x}_0 | \mathbf{x}) d\mathbf{x}_0, \\ p_t(\mathbf{x}_0 | \mathbf{x}) &= \frac{p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0)}{\int p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0}, \end{aligned} \quad (4)$$

$$\begin{aligned} \mathbb{E}[\mathbf{x}_0 | \mathbf{x}] &= \int \mathbf{x}_0 \frac{p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0)}{\int p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0} d\mathbf{x}_0 \\ &= \frac{1}{\int p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0} \int \mathbf{x}_0 p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0 \\ &= \frac{1}{\int p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0} \int \mathbf{x}_0 \mathcal{N}(\mathbf{x} | \mathbf{x}_0, s_t^2) \frac{1}{N} \sum_{j=1}^N \delta(\mathbf{x}_0 - \mathbf{x}_0^j) d\mathbf{x}_0 \\ &= \frac{1}{\int p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0} \frac{1}{N} \sum_{j=1}^N \mathbf{x}_0^j \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2) \end{aligned} \quad (5)$$

Finally,

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}] = \frac{\sum_{j=1}^N \mathbf{x}_0^j \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2)}{\sum_{j=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2)} \approx \mathbf{f}_{\theta}(\mathbf{x}, t) \quad (6)$$

For large  $t$  we obtain the mean of the dataset. Diffusion Models provides progressive denoising and thus cannot produce a sample in a few steps. We have to make samples less and less averaged.

The main idea behind consistency models:

How to train  $\mathbf{f}_\theta(\mathbf{x}, t)$  in such a way it produces a concrete sample,  $\mathbf{x}_0^j$ , instead of the mixture. That is,  $\mathbf{f}_\theta(\mathbf{x}, t) \approx \mathbf{x}_0^j$ .

## 1.2 Consistency function

So, we want  $\mathbf{f}_\theta(\mathbf{x}, t) \approx \mathbf{x}_0 \forall t$ , where  $\mathbf{x}$  - solution of the ODE at timestep  $t$ . A function that satisfies this property called **consistency function**. But how to obtain such function in practice?

Lets consider a discrete case,  $t_1, \dots, t_n$ . Note that  $\mathbf{f}_\theta(\mathbf{x}, t) \approx \mathbf{x}_0 \forall t$  leads to  $\mathbf{f}_\theta(\mathbf{x}_1, t_1) = \mathbf{f}_\theta(\mathbf{x}_2, t_2) = \dots = \mathbf{f}_\theta(\mathbf{x}_n, t_n)$ . We can train a consistency model by enforcing this property.

$$\|\mathbf{f}_\theta(\mathbf{x}_{n-1}, t_{n-1}) - \mathbf{f}_\theta(\mathbf{x}_n, t_n)\| \rightarrow \min_{\theta} \quad (7)$$

Note that this optimization problem may lead to the trivial solution:  $\mathbf{f}_\theta(\mathbf{x}, t) = 0$ . To avoid this, the repamaterization is introduced:

$$\mathbf{f}_\theta(\mathbf{x}, t_1) = \mathbf{x}, \quad (8)$$

where  $t_1$  is small. In practice, to make the consistency model to be continuous-time, we consider  $\hat{\mathbf{f}}_\theta(\mathbf{x}, t) = c_1(t)\mathbf{x} + c_2(t)\mathbf{f}_\theta(\mathbf{x}, t)$ , where  $c_1, c_2$  are differentiable functions,  $c_1(t_1) = 1, c_2(t_2) = 0$ .

Can we ensure that Eq. 27 lead to the desired consistency property?

**Proposition 1.** *The optimal solution to the consistency distillation problem Eq. (27) has the following form:*

$$\mathbf{f}_\theta(\mathbf{x}, t) = \mathbf{x}_0^k. \quad (9)$$

*Proof.*

$$\begin{aligned} \|\mathbf{f}_\theta(\mathbf{x}_{n-1}, t_{n-1}) - \mathbf{f}_\theta(\mathbf{x}_n, t_n)\| &= 0 \\ \rightarrow \mathbf{f}_\theta(\mathbf{x}_{n-1}, t_{n-1}) &= \mathbf{f}_\theta(\mathbf{x}_n, t_n) \end{aligned} \quad (10)$$

Let  $n = 2$ , then

$$\mathbf{f}_\theta(\mathbf{x}_1, t_1) = \mathbf{x}_1 = \mathbf{x}_0 + s_1 \cdot \mathbf{z} = \mathbf{f}_\theta(\mathbf{x}_2, t_2), \mathbf{z} \sim \mathcal{N}(0, 1) \quad (11)$$

Assuming that  $s_1 \approx 0$ , we obtain  $\mathbf{f}_\theta(\mathbf{x}_2, t_2) \approx \mathbf{x}_0$ .  $\square$

We have two strategies to train the consistency model.

## 1.3 Consistency via distillation

. In this case, we initialize the consistency model from the diffusion model and solve the following optimization problem.

$$\mathbb{E}_{\mathcal{U}(n|2,N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)p_{t_n}(\mathbf{x}|\mathbf{x}_0)} w(t_n) \|\mathbf{f}_{\theta^-}(\hat{\mathbf{x}}, t_{n-1}) - \mathbf{f}_\theta(\mathbf{x}, t_n)\|, \quad (12)$$

where  $\hat{\mathbf{x}}$  is one step of ODE solver with parameters  $\phi$ , from  $t_n$  to  $t_{n-1}$ , starting from  $\mathbf{x}$ ;  $\theta^- = \text{stopgrad}(\mu\theta^- + (1 - \mu)\theta)$ ;  $w(t)$  - weighting function.  $\mathbf{f}_{\theta^-}(\hat{\mathbf{x}}, t_{n-1})$  and  $\mathbf{f}_{\theta}(\mathbf{x}, t_n)$  are called as a teacher and student, respectively

### Questions.

**Q.** Why do we need a solver? Why using solver is better than simply  $\hat{\mathbf{x}} \sim p_{t_{n-1}}(\hat{\mathbf{x}}|\mathbf{x}_0)$

**A.** Intuition in Figure 1. If we do  $\hat{\mathbf{x}} \sim p_{t_{n-1}}(\hat{\mathbf{x}}|\mathbf{x}_0)$  and  $\mathbf{x} \sim p_{t_n}(\mathbf{x}|\mathbf{x}_0)$  then  $\mathbf{f}_{\theta}(\mathbf{x})$  and  $\mathbf{f}_{\theta}(\hat{\mathbf{x}})$  will look at slightly different modes,  $\mathbf{x}_0$ , because different trajectories intersect. However, we need that  $\mathbf{f}_{\theta}$  look at the same mode for different  $t$ . I think, it potentially leads to slower convergence and worse final performance.

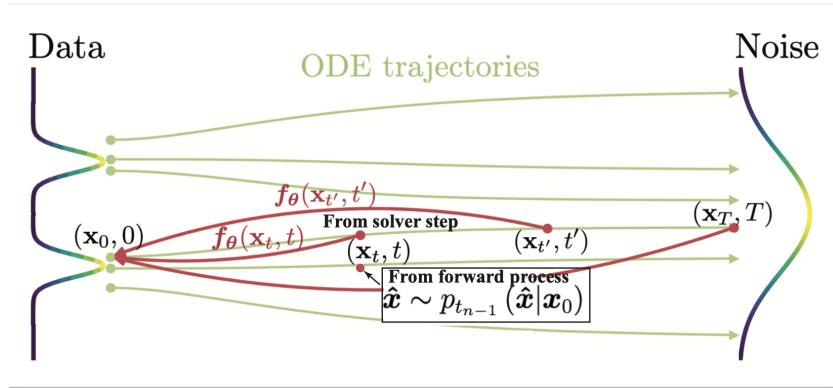


Figure 1

**Q.** Why do we need EMA?

**A.** Possible because we need the student and the teacher to be different models in order to avoid trivial solution.

**Q.** What is the right rule for distillation? Now we need much training steps and full dataset to perform distillation. Is it the right way?

**A.** Intuition: distillation is not a training, the model already knows the distribution. Thus, we do not need to learn a whole distribution. **Which subset of data we need?; How to change the model weights?**

**P.** Possible ablations to try:

- align not only neighboring  $t_n$  and  $t_{n-1}$  but also more distant steps. Concretely:

$$\|\mathbf{f}_{\theta}(\mathbf{x}_k, t_k) - \mathbf{f}_{\theta}(\mathbf{x}_j, t_j)\| \rightarrow \min_{\theta} \quad (13)$$

where  $k \neq j$ . Intuition: in the definition of the consistency function we said that the function should be the same for different time steps,  $t$ .

In practice we can do: **1)** make more ODE steps, not only a single step. However this is a time-consuming. But maybe will converge much faster; **2)** Combine distillation and training loss of consistency:

$$\begin{aligned} & \mathbb{E}_{\mathcal{U}(k,j|1,N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)p(\mathbf{z})} w(t_n) \|\mathbf{f}_{\theta}(\mathbf{x}_0 + t_k \mathbf{z}, t_k) - \mathbf{f}_{\theta}(\mathbf{x}_0 + t_j \mathbf{z}, t_j)\| \\ & + \mathbb{E}_{\mathcal{U}(n|2,N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)p_{t_n}(\mathbf{x}|\mathbf{x}_0)} w(t_n) \|\mathbf{f}_{\theta^-}(\hat{\mathbf{x}}, t_{n-1}) - \mathbf{f}_{\theta}(\mathbf{x}, t_n)\| \end{aligned} \quad (14)$$

Or maybe combine 1) and 2). For example, if  $k$  and  $j$  are close - make 1) otherwise 2).

- try solvers of different order (DPM, DDIM, EDM), does it really matter?
- ablate EMA, does we really need this?
- ablate different weighting schemes.

### 1.3.1 Training details

**Timesteps schedule,  $N$  (?).** Song writes the intuition for the consistency training: *For improved practical performance, we propose to progressively increase  $N$  during training according to a schedule function. The consistency training loss has less “variance” but more “bias” when  $N$  is small, which facilitates faster convergence at the beginning of training. On the contrary, it has more “variance” but less “bias” when  $N$  is large (i.e.,  $\Delta t$  is small), which is desirable when closer to the end of training. Is it true for CD?*

**Ablation for  $N$  and solver order.** The best performing configuration of  $N$  corresponds to the best sampling configuration of diffusion model (Figure 2c). That is, for EDM cifar is 18 steps. **Larger  $N$  better results does not work?;** As shown in Figure 2d, adaptive  $N$  and  $\mu$  works good for CT, **Is it true for CD?**

**Think about it:** As demonstrated in Fig. 3d, the convergence of CT is highly sensitive to  $N$  —smaller  $N$  leads to faster convergence but worse samples, whereas larger  $N$  leads to slower convergence but better samples upon convergence. This matches our analysis in Section 5, and motivates our practical choice of progressively growing  $N$  and  $\mu$  for CT to balance the trade-off between convergence speed and sample quality.

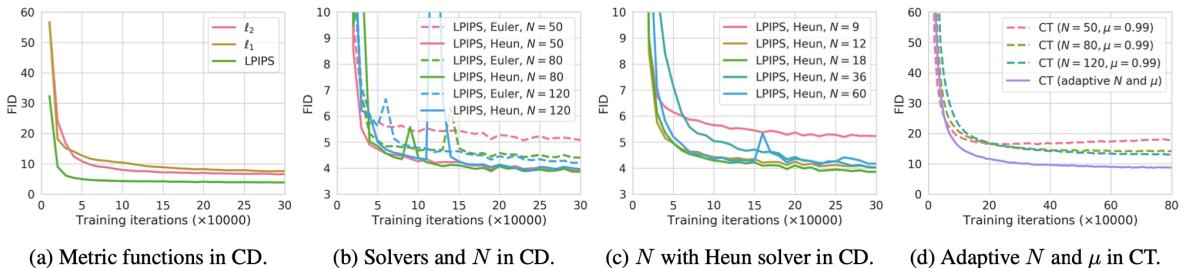


Figure 2

Table 3: Hyperparameters used for training CD and CT models

Hyperparameter	CIFAR-10		ImageNet 64 × 64		LSUN 256 × 256	
	CD	CT	CD	CT	CD	CT
Learning rate	4e-4	4e-4	8e-6	8e-6	1e-5	1e-5
Batch size	512	512	2048	2048	2048	2048
$\mu$	0		0.95		0.95	
$\mu_0$		0.9		0.95		0.95
$s_0$		2		2		2
$s_1$		150		200		150
$N$		18		40		40
ODE solver	Heun		Heun		Heun	
EMA decay rate	0.9999	0.9999	0.999943	0.999943	0.999943	0.999943
Training iterations	800k	800k	600k	800k	600k	1000k
Mixed-Precision (FP16)	No	No	Yes	Yes	Yes	Yes
Dropout probability	0.0	0.0	0.0	0.0	0.0	0.0
Number of GPUs	8	8	64	64	64	64

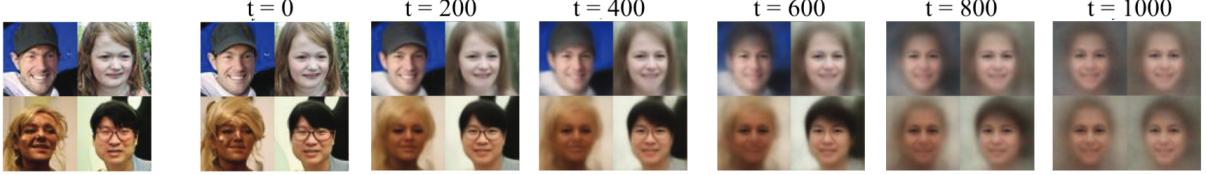
Figure 3

## 1.4 The difference between DM and CM

### 1. Approximation difference

DM:  $\|\mathbf{f}_\theta(\mathbf{x}_t, t) - \mathbf{x}_0\| \rightarrow [\mathbf{f}_\theta(\mathbf{x}_t, t) = \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]] \rightarrow \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = -\frac{1}{t}(\mathbf{x} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}])$

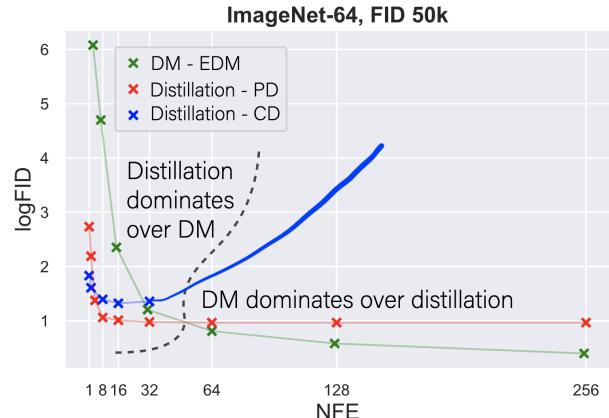
CM:  $\|\mathbf{f}_\theta(\mathbf{x}_{t-1}, t-1) - \mathbf{f}_\theta(\mathbf{x}_t, t)\| \rightarrow [\mathbf{f}_\theta(\mathbf{x}_t, t) = \mathbf{x}_0] \not\rightarrow \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$



### 2. Sampling difference

DM: Score function → Solvers

CM: Consistency function ↗ Solvers // → Stochastic sampler



### Q. Why CM degrades with the NFE increasing?

## 1.5 Consistency from scratch

Here, we also have to solve the same problem

$$\|\mathbf{f}_\theta(\mathbf{x}_{n-1}, t_{n-1}) - \mathbf{f}_\theta(\mathbf{x}_n, t_n)\| \rightarrow \min_\theta \quad (15)$$

We sample  $\mathbf{x}_n \sim p_{t_n}(\mathbf{x}_n | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_n | \mathbf{x}_0, t_n^2)$ , but how to obtain  $\mathbf{x}_{n-1}$  if we do not have a diffusion model that approximates the score,  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ . In other words, we cannot solve the equation (or even approximate it) without the known score. However, let's look at this in more detail.

$$\begin{aligned} d\mathbf{x} &= -t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) dt, \\ \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) &= -\frac{1}{t^2} (\mathbf{x} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}]). \end{aligned} \quad (16)$$

The unknown component, that we approximate via diffusion, it is  $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}]$ . This is an integral that we cannot solve for the experimental data only. The authors consider a Monte-Carlo approximation:

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}] = \int \mathbf{x}_0 p_t(\mathbf{x}_0 | \mathbf{x}) d\mathbf{x}_0 \approx \frac{1}{n} \sum_{i=1}^n \mathbf{x}_0^i | \mathbf{x}, \quad \mathbf{x}_0^i \sim p_t(\mathbf{x}_0^i | \mathbf{x}) \quad (17)$$

The authors consider  $n = 1$ . Thus,  $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}] \approx \mathbf{x}_0^i$  and we have the following.

$$d\mathbf{x} = -t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) dt \approx \frac{1}{t} (\mathbf{x} - \mathbf{x}_0) dt \quad (18)$$

Interestingly, now we can solve the equation analytically:

$$\begin{aligned} d\mathbf{x} &= \frac{1}{t} (\mathbf{x} - \mathbf{x}_0) dt \\ \frac{d\mathbf{x}}{(\mathbf{x} - \mathbf{x}_0)} &= \frac{1}{t} dt \\ \int \frac{d(\mathbf{x} - \mathbf{x}_0)}{(\mathbf{x} - \mathbf{x}_0)} &= \frac{dt}{t} \\ \ln |\mathbf{x}_{t_n} - \mathbf{x}_0| - \ln |\mathbf{x}_{t_{n+1}} - \mathbf{x}_0| &= \ln |t_n| - \ln |t_{n+1}| \\ \ln \left| \frac{\mathbf{x}_{t_n} - \mathbf{x}_0}{\mathbf{x}_{t_{n+1}} - \mathbf{x}_0} \right| &= \ln \left| \frac{t_n}{t_{n+1}} \right| \\ \mathbf{x}_{t_n} &= \mathbf{x}_0 + \frac{t_n}{t_{n+1}} (\mathbf{x}_{t_{n+1}} - \mathbf{x}_0) \\ \mathbf{x}_{t_n} &= \frac{t_n}{t_{n+1}} \mathbf{x}_{t_{n+1}} - \mathbf{x}_0 \left( \frac{t_n - t_{n+1}}{t_{n+1}} \right) \end{aligned} \quad (19)$$

Note that we can take arbitrary  $t_n$  and  $t_{n+1}$ . So, in practice we sample  $\mathbf{x}_{t_{n+1}} = \mathbf{x}_0 + t_{n+1} \cdot \mathbf{z}$ ,  $\mathbf{z} \sim \mathcal{N}(0, I)$ . And calculate  $\mathbf{x}_{t_n}$  according to the equation:

$$\begin{aligned} \mathbf{x}_{t_n} &= \frac{t_n}{t_{n+1}} \mathbf{x}_{t_{n+1}} - \mathbf{x}_0 \left( \frac{t_n - t_{n+1}}{t_{n+1}} \right) \\ &= \frac{t_n}{t_{n+1}} (\mathbf{x}_0 + t_{n+1} \cdot \mathbf{z}) - \mathbf{x}_0 \left( \frac{t_n - t_{n+1}}{t_{n+1}} \right) \\ &= \mathbf{x}_0 + t_n \mathbf{z} \end{aligned} \quad (20)$$

Finally, our optimization problem becomes:

$$\|\mathbf{f}_\theta(\mathbf{x}_0 + t_{n+1} \cdot \mathbf{z}, t_{n+1}) - \mathbf{f}_\theta(\mathbf{x}_0 + t_n \cdot \mathbf{z}, t_n)\| \rightarrow \min_\theta \quad (21)$$

### Questions.

**Q.** What happens if we take more accurate estimation of the conditional expectation? Or even analytical solution?

**A.** **Case 1.** More accurate estimation.

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] \approx \mathbf{x}_0^i \rightarrow \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] \approx \frac{1}{n} \sum_{i=1}^n \mathbf{x}_0^i | \mathbf{x}_t, \quad (22)$$

The number  $n$  depends on the current time  $t$ . For large  $t$ , e.g. 1, we have to take large  $n$  (approximately all dataset). For small  $t \approx 0$  one sample is enough. To sample from dataset, we can use k-nn to find corresponding  $x_0$ . Then,

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] \approx \frac{1}{n(t, \mathbf{x}_t)} \sum_{i=1}^{n(t, \mathbf{x}_t)} \mathbf{x}_0^i | \mathbf{x}_t, \quad (23)$$

The equation and its solution becomes.

$$\begin{aligned} d\mathbf{x} &= \frac{1}{t} \left( \mathbf{x} - \frac{1}{m(t, \mathbf{x}_t)} \sum_{i=1}^{m(t, \mathbf{x}_t)} \mathbf{x}_0^i \right) dt \\ \mathbf{x}_{t_n} &= \frac{t_n}{t_{n+1}} \mathbf{x}_{t_{n+1}} - \left( \frac{t_n - t_{n+1}}{t_{n+1}} \right) \frac{1}{m} \sum_{i=1}^m \mathbf{x}_0^i \end{aligned} \quad (24)$$

Taking into account  $\mathbf{x}_{t_{n+1}} = \mathbf{x}_0 + t_{n+1} \cdot \mathbf{z}$ , the consistency loss transforms to

$$\begin{aligned} \mathbf{x}_{t_n} &= \frac{t_n}{t_{n+1}} \mathbf{x}_0 + t_n \cdot \mathbf{z} - \left( \frac{t_n - t_{n+1}}{t_{n+1}} \right) \frac{1}{m} \sum_{i=1}^m \mathbf{x}_0^i \\ \mathbf{x}_{t_n} &= \mathbf{x}_0 \left( \frac{1 + t_n/t_{n+1} (m-1)}{m} \right) - \left( \frac{t_n - t_{n+1}}{t_{n+1}} \right) \frac{1}{m} \sum_{\mathbf{x}_0^i \setminus \mathbf{x}_0} \mathbf{x}_0^i + t_n \cdot \mathbf{z} \end{aligned} \quad (25)$$

Looks creepy.

**Case 2.** Analytical form of the expectation. No way :(, i'm full rubbish, biotrash.

**Q.** Now we take a single step of the ODE using approximation of the conditional expectation. What if we do more?

## 2 Improved CT

### 2.1 EMA

Previously, Song uses the following loss to train the Consistency Models

$$\mathcal{L}^N(\theta, \theta^-) = \|\mathbf{f}_\theta(\mathbf{x}_{t_{n+1}}, t_{n+1}) - \mathbf{f}_{\theta^-}(\mathbf{x}_{t_n}, t_n)\| \rightarrow \min_\theta, \quad (26)$$

where  $\theta^- = \text{stopgrad}(\mu\theta^- + (1 - \mu)\theta)$ . For Consistency Training, Song uses an Euler step, approximating the score as a 1-sample Monte Carlo estimation.

$$\mathcal{L}_{\text{CT}}^N(\theta, \theta^-) = \|\mathbf{f}_\theta(\mathbf{x}_0 + t_{n+1} \cdot \mathbf{z}, t_{n+1}) - \mathbf{f}_{\theta^-}(\mathbf{x}_0 + t_n \cdot \mathbf{z}, t_n)\| \rightarrow \min_{\theta}, \quad (27)$$

However, he states that the EMA degrades the quality of the learned network. To confirm this, lets consider the following toy example.

Let

$$\begin{aligned} p_{\text{data}}(\mathbf{x}) &= \delta(\mathbf{x} - \boldsymbol{\epsilon}), \\ \underset{\theta}{\operatorname{argmin}} \mathcal{L}^N(\theta, \theta^-) &=? \end{aligned} \quad (28)$$

The question is the following: what is the optimum of  $\mathcal{L}^N(\theta, \theta^-)$ . To this end, we solve the ODE, using the Variance Exploding scheme.

$$\begin{aligned} \frac{d\mathbf{x}}{d\sigma} &= -\sigma \nabla_x \log p_\sigma(\mathbf{x}), \quad \sigma \in [\sigma_{\min}, \sigma_{\max}], \\ \mathbf{x}_\sigma &\sim p_\sigma(\mathbf{x}), \quad p_{\sigma_{\min}}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x}) \end{aligned} \quad (29)$$

Firstly,

$$\begin{aligned} p_\sigma(\mathbf{x}) &= \int p_\sigma(\mathbf{x}|\mathbf{x}_0)p_{\text{data}}(\mathbf{x}_0)d\mathbf{x}_0 \\ &= \int \mathcal{N}(\mathbf{x}|\mathbf{x}_0, \sigma^2)\delta(\mathbf{x}_0 - \boldsymbol{\epsilon})d\mathbf{x}_0 \\ &= \mathcal{N}(\mathbf{x}|\boldsymbol{\epsilon}, \sigma^2) \end{aligned} \quad (30)$$

Secondly,

$$\nabla_x \log p_\sigma(\mathbf{x}) = -\frac{\mathbf{x} - \boldsymbol{\epsilon}}{\sigma^2} \quad (31)$$

Then, we can solve the equation.

$$\begin{aligned} \frac{d\mathbf{x}}{d\sigma} &= -\sigma \nabla_x \log p_\sigma(\mathbf{x}) = \frac{\mathbf{x} - \boldsymbol{\epsilon}}{\sigma} \\ \int_{\mathbf{x}_{\sigma_{i+1}}}^{\mathbf{x}_{\sigma_i}} \frac{d\mathbf{x}}{\mathbf{x} - \boldsymbol{\epsilon}} &= \int_{\sigma_{i+1}}^{\sigma_i} \frac{d\sigma}{\sigma} \\ \ln \frac{\mathbf{x}_{\sigma_{i+1}} - \boldsymbol{\epsilon}}{\mathbf{x}_{\sigma_i} - \boldsymbol{\epsilon}} &= \frac{\sigma_{i+1}}{\sigma_i} \end{aligned} \quad (32)$$

Finally, the solution is the following.

$$\mathbf{x}_{\sigma_i} = \boldsymbol{\epsilon} + \frac{\sigma_i}{\sigma_{i+1}} (\mathbf{x}_{\sigma_{i+1}} - \boldsymbol{\epsilon}) \quad (33)$$

One more thing that should be noted is that the real consistency function has the following form in this case.

$$\mathbf{f}(\mathbf{x}, \sigma) = \frac{\sigma_{\min}}{\sigma} \mathbf{x} + \left(1 - \frac{\sigma_{\min}}{\sigma}\right) \boldsymbol{\epsilon} \quad (34)$$

It is clear that this function is satisfied the requirements consistency function has.

- 1)  $\mathbf{f}(\mathbf{x}, \sigma_{\min}) = \mathbf{x}$
- 2) If  $\mathbf{x} \sim p_\sigma(\mathbf{x})$ , then  $\mathbf{f}(\mathbf{x}, \sigma) = \mathbf{x}_{\sigma_{\min}}$

Moreover, we can approximate the real consistency function using the following parameterization:  $\mathbf{f}_\theta(\mathbf{x}, \sigma) = \frac{\sigma_{\min}}{\sigma} \mathbf{x} + \left(1 - \frac{\sigma_{\min}}{\sigma}\right) \boldsymbol{\theta}$ .  
To sum it up, if  $p_{\text{data}}(\mathbf{x}) = \delta(\mathbf{x} - \boldsymbol{\epsilon})$ , then

$$\boxed{\begin{aligned} p_\sigma(\mathbf{x}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\epsilon}, \sigma^2) \\ \nabla_{\mathbf{x}} \log p_\sigma(\mathbf{x}) &= -\frac{\mathbf{x} - \boldsymbol{\epsilon}}{\sigma^2} \\ \mathbf{x}_{\sigma_i} &= \boldsymbol{\epsilon} + \frac{\sigma_i}{\sigma_{i+1}} (\mathbf{x}_{\sigma_{i+1}} - \boldsymbol{\epsilon}) \\ \mathbf{f}(\mathbf{x}, \sigma) &= \frac{\sigma_{\min}}{\sigma} \mathbf{x} + \left(1 - \frac{\sigma_{\min}}{\sigma}\right) \boldsymbol{\epsilon} \\ \mathbf{f}_\theta(\mathbf{x}, \sigma) &= \frac{\sigma_{\min}}{\sigma} \mathbf{x} + \left(1 - \frac{\sigma_{\min}}{\sigma}\right) \boldsymbol{\theta} \end{aligned}} \quad (36)$$

Then we can back to the main question about the optimum of the loss function.

$$\mathcal{L}^N(\theta, \theta^-) = \|\mathbf{f}_\theta(\mathbf{x}_{\sigma_{i+1}}, \sigma_{i+1}) - \mathbf{f}_{\theta^-}(\mathbf{x}_{\sigma_i}, \sigma_i)\| \quad (37)$$

We have that  $\mathbf{x}_{\sigma_{i+1}} = \boldsymbol{\epsilon} + \sigma_{i+1} \cdot \mathbf{z}$ . Using the solution of the ODE, we can find that  $\mathbf{x}_{\sigma_i} = \boldsymbol{\epsilon} + \sigma_i \cdot \mathbf{z}$ .

$$\mathcal{L}^N(\theta, \theta^-) = \|\mathbf{f}_\theta(\boldsymbol{\epsilon} + \sigma_{i+1} \cdot \mathbf{z}, \sigma_{i+1}) - \mathbf{f}_{\theta^-}(\boldsymbol{\epsilon} + \sigma_i \cdot \mathbf{z}, \sigma_i)\| \quad (38)$$

which means that  $\mathcal{L}^N(\theta, \theta^-) = \mathcal{L}_{\text{CT}}^N(\theta, \theta^-)$  in this case.

$$\begin{aligned} \mathcal{L}_{\text{CT}}^N(\theta, \theta^-) &= \mathbb{E}[(\mathbf{f}_\theta(\boldsymbol{\epsilon} + \sigma_{i+1} \cdot \mathbf{z}, \sigma_{i+1}) - \mathbf{f}_{\theta^-}(\boldsymbol{\epsilon} + \sigma_i \cdot \mathbf{z}, \sigma_i))] \\ &= \mathbb{E} \left[ \left( \frac{\sigma_{\min}}{\sigma_{i+1}} (\boldsymbol{\epsilon} + \sigma_{i+1} \cdot \mathbf{z}) + \left(1 - \frac{\sigma_{\min}}{\sigma_{i+1}}\right) \boldsymbol{\theta} - \frac{\sigma_{\min}}{\sigma_i} (\boldsymbol{\epsilon} + \sigma_i \cdot \mathbf{z}) + \left(1 - \frac{\sigma_{\min}}{\sigma_i}\right) \boldsymbol{\theta}^- \right)^2 \right] \\ &= \mathbb{E} \left[ \left( \frac{\sigma_{\min}}{\sigma_{i+1}} \boldsymbol{\epsilon} + \left(1 - \frac{\sigma_{\min}}{\sigma_{i+1}}\right) \boldsymbol{\theta} - \frac{\sigma_{\min}}{\sigma_i} \boldsymbol{\epsilon} - \left(1 - \frac{\sigma_{\min}}{\sigma_i}\right) \boldsymbol{\theta}^- \right)^2 \right] \end{aligned} \quad (39)$$

We know that  $\sigma_{i+1} = \sigma_i + \Delta\sigma$ . Assuming that  $\Delta\sigma$  is small, we can write the following decomposition

$$\frac{\sigma_{\min}}{\sigma_{i+1}} = \frac{\sigma_{\min}}{\sigma_i + \Delta\sigma} = \frac{\sigma_{\min}}{\sigma_i (1 + \Delta\sigma/\sigma_i)} = \frac{\sigma_{\min}}{\sigma_i} \left(1 - \frac{\Delta\sigma}{\sigma_i} + O(\Delta\sigma^2)\right) \quad (40)$$

Then,

$$\begin{aligned} \mathcal{L}_{\text{CT}}^N(\theta, \theta^-) &= \\ &= \mathbb{E} \left[ \left( \frac{\sigma_{\min}}{\sigma_i} \left(1 - \frac{\Delta\sigma}{\sigma_i} + O(\Delta\sigma^2)\right) \boldsymbol{\epsilon} + \left(1 - \frac{\sigma_{\min}}{\sigma_i} \left(1 - \frac{\Delta\sigma}{\sigma_i} + O(\Delta\sigma^2)\right)\right) \boldsymbol{\theta} - \frac{\sigma_{\min}}{\sigma_i} \boldsymbol{\epsilon} - \left(1 - \frac{\sigma_{\min}}{\sigma_i}\right) \boldsymbol{\theta}^- \right)^2 \right] \\ &= \mathbb{E} \left[ \left( \frac{\sigma_{\min}}{\sigma_i} \left(-\frac{\Delta\sigma}{\sigma_i} + O(\Delta\sigma^2)\right) \boldsymbol{\epsilon} + \left(1 - \frac{\sigma_{\min}}{\sigma_i} \left(1 - \frac{\Delta\sigma}{\sigma_i} + O(\Delta\sigma^2)\right)\right) \boldsymbol{\theta} - \left(1 - \frac{\sigma_{\min}}{\sigma_i}\right) \boldsymbol{\theta}^- \right)^2 \right] \\ &= \mathbb{E} \left[ \left( \frac{\sigma_{\min}}{\sigma_i} \left(-\frac{\Delta\sigma}{\sigma_i} + O(\Delta\sigma^2)\right) \boldsymbol{\epsilon} + \left(1 - \frac{\sigma_{\min}}{\sigma_i} \left(1 - \frac{\Delta\sigma}{\sigma_i} + O(\Delta\sigma^2)\right)\right) \boldsymbol{\theta} - \left(1 - \frac{\sigma_{\min}}{\sigma_i}\right) \boldsymbol{\theta}^- \right)^2 \right] \\ &= \mathbb{E} \left[ \left( \frac{\sigma_{\min} \Delta\sigma}{\sigma_i^2} (\boldsymbol{\theta} - \boldsymbol{\epsilon}) + \left(1 - \frac{\sigma_{\min}}{\sigma_i}\right) (\boldsymbol{\theta} - \boldsymbol{\theta}^-) + O(\Delta\sigma^2) \right)^2 \right] \end{aligned} \quad (41)$$

We can use simple formula,  $(a + b)^2 = a^2 + 2ab + b^2$ , to obtain the final result.

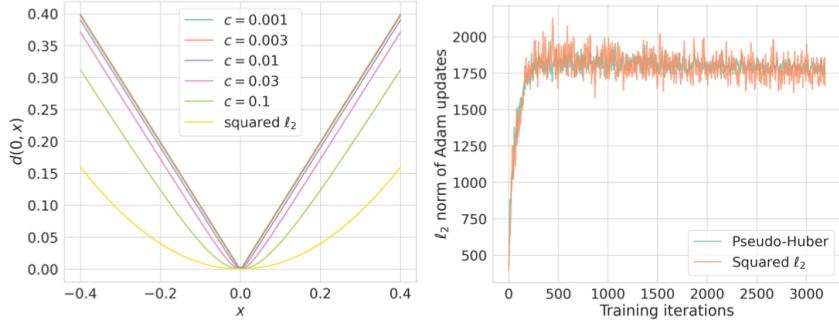
$$\mathcal{L}_{\text{CT}}^N(\theta, \theta^-) = \mathbb{E} \left[ \left( \frac{\sigma_{\min} \Delta \sigma}{\sigma_i^2} (\theta - \epsilon) + \left(1 - \frac{\sigma_{\min}}{\sigma_i}\right) (\theta - \theta^-) \right)^2 \right] + O(\Delta \sigma^2) \quad (42)$$

So, it can be seen that while minimizing this loss we have two components affecting our optimum:  $\frac{\sigma_{\min} \Delta \sigma}{\sigma_i^2} (\theta - \epsilon)$  and  $\left(1 - \frac{\sigma_{\min}}{\sigma_i}\right) (\theta - \theta^-)$ . The former is our goal since the optimum is  $\theta = \epsilon$ . However, if  $\theta \neq \theta^-$ , the latter moves us towards the wrong direction ( $\theta = \theta^-$ ). Moreover, the first component is multiplied by small value  $\Delta \sigma$  and thus has less influence than the second component.

Motivating by these arguments, Song propose to remove EMA, that is, use  $\theta = \theta^-$ .

## 2.2 Technical improvements

### 2.2.1 Pseudo-Huber



(a)  $d(\mathbf{0}, \mathbf{x})$  as a function of  $\mathbf{x}$ .      (b) Comparing Adam updates.

Figure 6: (a) The shapes of various metric functions. (b) The  $\ell_2$  norms of parameter updates in Adam optimizer. Curves are rescaled to have the same mean. The Pseudo-Huber metric has lower variance compared to the squared  $\ell_2$  metric.

Instead of L2 loss, Song propose to use the following function

$$d(x, y) = \sqrt{\|x - y\|^2 + c^2} - c \quad (43)$$

The intuition behind is that compared to the squared  $\ell_2$  metric, the Pseudo-Huber metric is more robust to outliers as it imposes a smaller penalty for large errors than the squared  $\ell_2$  metric does, yet behaves similarly for smaller errors. We posit that this added robustness can reduce variance during training. **Apparently, consistency loss has large variance during training.**

### 3 Multistep CM

#### 3.1 Multistep CD

##### 3.1.1 VP case

This work splits the solution interval into  $m$  segments and train the consistency model on each of these segments separately. They parameterize consistency function as

$$\mathbf{f}_{\theta}^m(\mathbf{x}_t, t, s_t^m, w) = \text{DDIM}(\mathbf{x}_t, t, s_t^m, w), \quad (44)$$

$$\mathbf{x}_s^w = \text{DDIM}(\mathbf{x}_t, t, s, \mathbf{c}, w) = \sqrt{\frac{\alpha_s}{\alpha_t}} \mathbf{x}_t + \boldsymbol{\epsilon}_{\theta}^w(\mathbf{x}_t, t, \mathbf{c}) \left( \sqrt{1 - \alpha_s} - \sqrt{\frac{\alpha_s}{\alpha_t}} \sqrt{1 - \alpha_t} \right), \quad (45)$$

The CD loss in this case

$$\mathcal{L}_{\text{CD}}(\theta) = \mathbb{E} [\|\mathbf{f}_{\theta}(\mathbf{x}_{t_{n-1}}^w, t_{n-1}, s_{t_{n-1}}^m) - \mathbf{f}_{\theta}^m(\mathbf{x}_{t_n}, t_n, s_{t_n}^m)\|] \rightarrow \min_{\theta}, \quad (46)$$

Note that  $s_{t_{n-1}} = s_{t_n}$ . We parameterize consistency function by  $\epsilon_{\theta}$ . Thus, we can rewrite this loss directly via  $\epsilon_{\theta}$ .

$$\mathcal{L}_{\text{CD}}(\theta) = \mathbb{E} \left[ \left\| \frac{\mathbf{x}_{s_{t_n}}^w - \sqrt{\frac{\alpha_{t_n}}{\alpha_{s_{t_n}}}} \mathbf{x}_{t_n}}{\sqrt{1 - \alpha_{t_n}} - \sqrt{\frac{\alpha_{t_n}}{\alpha_{s_{t_n}}}} \sqrt{1 - \alpha_{s_{t_n}}}} - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_{t_n}, t_n) \right\| \right] \rightarrow \min_{\theta}, \quad (47)$$

##### 3.1.2 VE case

Basic equations for the VE case:

$$\begin{aligned} dx &= -t \nabla_x \log p_t(x) dt \\ \nabla_x \log p_t(x) &= -\frac{1}{t^2} (x - \mathbb{E}[x_0|x]) \\ \mathbb{E}[x_0|x] &\approx x^{\phi}(x, t) - \text{teacher network} \end{aligned} \quad (48)$$

Euler integrator:

$$\begin{aligned} dx &= \frac{1}{t} (x - x^{\phi}(x, t)), \\ x_s^{\phi} &= x_t + \int_t^s (x_{\tau} - x^{\phi}(x_{\tau}, \tau)) \frac{d\tau}{\tau}, \\ x_s^{\phi} &\approx x_t + \frac{s-t}{t} (x_t - x^{\phi}(x_t, t)) \end{aligned} \quad (49)$$

Consistency loss:

$$\begin{aligned} \mathcal{L}_{\text{CD}}(\theta) &= \mathbb{E} [\|f_{\theta}(x_{t_{n-1}}^{\phi}, t_{n-1}) - f_{\theta}(x_{t_n}, t_n)\|] \rightarrow \min_{\theta}, \\ x_{t_{n-1}}^{\phi} &= x_{t_n} + \frac{t_{n-1} - t_n}{t_n} (x_{t_n} - x^{\phi}(x_{t_n}, t_n)) \end{aligned} \quad (50)$$

#### Multistep case

Parameterization:

$$\begin{aligned} f_{\theta}(x_{t_n}, t_n, b_{t_n}) &= x_{t_n} + \frac{b_{t_n} - t_n}{t_n} (x_{t_n} - x^{\theta}(x_{t_n}, t_n)), \\ x^{\theta} - \text{student}, b_{t_n} - \text{boundary} \end{aligned} \quad (51)$$

Updated consistency loss

$$\mathcal{L}_{\text{CD}}(\theta) = \mathbb{E} \left[ \|f_\theta(x_{t_{n-1}}^\phi, t_{n-1}, b_{t_{n-1}}) - f_\theta(x_{t_n}, t_n, b_{t_n})\| \right] \quad (52)$$

Note that,  $b_{t_n} = b_{t_{n-1}}$ . Lets consider in more detail the optimum

$$\begin{aligned} & f_\theta(x_{t_n}, t_n, b_{t_n}) = x_{t_n} + \frac{b_{t_n} - t_n}{t_n} (x_{t_n} - x^\theta(x_{t_n}, t_n)) \\ & f_\theta(x_{t_{n-1}}^\phi, t_{n-1}, b_{t_{n-1}}) = x_{t_{n-1}}^\phi + \frac{b_{t_n} - t_{n-1}}{t_{n-1}} (x_{t_{n-1}}^\phi - x^\theta(x_{t_{n-1}}^\phi, t_{n-1})) \\ & x_{t_{n-1}}^\phi = x_{t_n} + \frac{t_{n-1} - t_n}{t_n} (x_{t_n} - x^\phi(x_{t_n}, t_n)) \end{aligned} \quad (53)$$

Merge the second and the third:

$$\begin{aligned} & f_\theta(x_{t_{n-1}}^\phi, t_{n-1}, b_{t_{n-1}}) = x_{t_n} + \frac{t_{n-1} - t_n}{t_n} (x_{t_n} - x^\phi(x_{t_n}, t_n)) + \\ & + \frac{b_{t_n} - t_{n-1}}{t_{n-1}} \left( x_{t_n} + \frac{t_{n-1} - t_n}{t_n} (x_{t_n} - x^\phi(x_{t_n}, t_n)) - x^\theta(x_{t_{n-1}}^\phi, t_{n-1}) \right) \end{aligned} \quad (54)$$

**Q.** Can we achieve any insights from there reparameterization?

## 4 Pagoda

This is a combination of Knowledge Distillation (but in the different way, we collect noise and images using inversion) and GAN. Without CFG:

$$\begin{aligned}\mathcal{L}_{rec}(\theta) &= \mathbb{E}_{p_{\text{data}}(x,c)} [\|x - G_\theta(E(x,c), c)\|] \\ &= \mathbb{E}_{p_{\text{data}}(x|c)p(c)} [\|x - G_\theta(E(x,c), c)\|]\end{aligned}\tag{55}$$

$$\begin{aligned}\mathcal{L}_{adv}(\theta) &= \mathbb{E}_{p_{\text{data}}(x,c)} [\log D_\psi(x,c)] + \mathbb{E}_{p(z,c)} [1 - \log D_\psi(G_\theta(z,c), c)] \\ &= \mathbb{E}_{p_{\text{data}}(x|c)p(c)} [\log D_\psi(x,c)] + \mathbb{E}_{p(z)p(c)} [1 - \log D_\psi(G_\theta(z,c), c)]\end{aligned}$$

With CFG:

$$\begin{aligned}\mathcal{L}_{rec}(\theta) &= \mathbb{E}_{p_{\text{data}}(x,c,w)} [\|x - G_\theta(E(x,c,w), c, w)\|] \\ &= \mathbb{E}_{p_{\text{data}}(x|c,w)p(c)p(w)} [\|x - G_\theta(E(x,c,w), c, w)\|] \\ \mathcal{L}_{adv}(\theta) &= \mathbb{E}_{p_{\text{data}}(x,c,w)} [\log D_\psi(x,c,w)] + \mathbb{E}_{p(z,c,w)} [1 - \log D_\psi(G_\theta(z,c,w), c, w)] \\ &= \mathbb{E}_{p_{\text{data}}(x|c,w)p(c)p(w)} [\log D_\psi(x,c,w)] + \mathbb{E}_{p(z)p(c)p(w)} [1 - \log D_\psi(G_\theta(z,c,w), c, w)]\end{aligned}\tag{56}$$

Bayes rule:

$$\begin{aligned}p_{\text{data}}(x,c,w) &= p_{\text{data}}(x|c,w)p(c)p(w) \\ p_{\text{data}}(x,c,w) &= p(w|x,c)p_{\text{data}}(x|c)p_{\text{data}}(c)\end{aligned}\tag{57}$$

With CFG and Bayes rule:

$$\begin{aligned}\mathcal{L}_{rec}(\theta) &= \mathbb{E}_{p_{\text{data}}(x,c,w)} [\|x - G_\theta(E(x,c,w), c, w)\|] \\ &= \mathbb{E}_{p(w|x,c)p_{\text{data}}(x|c)p(c)} [\|x - G_\theta(E(x,c,w), c, w)\|] \\ \mathcal{L}_{adv}(\theta) &= \mathbb{E}_{p_{\text{data}}(x,c,w)} [\log D_\psi(x,c,w)] + \mathbb{E}_{p(z,c,w)} [1 - \log D_\psi(G_\theta(z,c,w), c, w)] \\ &= \mathbb{E}_{p(w|x,c)p_{\text{data}}(x|c)p(c)} [\log D_\psi(x,c,w)] + \mathbb{E}_{p(z)p(c)p(w)} [1 - \log D_\psi(G_\theta(z,c,w), c, w)]\end{aligned}\tag{58}$$

## 5 Score distillation

### 5.1 DMD

In DMD, the authors consider Variational Score Distillation:

$$\begin{aligned} \text{KL}(p_\theta || p_{real}) &= \mathbb{E}_{p_\theta} \log \left( \frac{p_\theta}{p_{real}} \right) = \int p_\theta(x) (\log p_\theta(x) - \log p_{real}(x)) dx \\ &= \mathbb{E}_{x=G_\theta(z), z \sim \mathcal{N}(0,1)} (\log p_\theta(x) - \log p_{real}(x)) \\ &= \int \mathcal{N}(z|0,1) [\log p_\theta(G_\theta(z)) - \log p_{real}(G_\theta(z))] dz \end{aligned} \quad (59)$$

We sample from the generator  $x \sim p_\theta(x)$  and calculate the difference between real and fake scores. We cannot calculate this reverse KL-divergence, since we do not know the densities. We can take the gradient only:

$$\begin{aligned} \nabla_\theta \text{KL}(p_\theta || p_{real}) &= \nabla_\theta \int \mathcal{N}(z|0,1) [\log p_\theta(G_\theta(z)) - \log p_{real}(G_\theta(z))] dz \\ &= \int \mathcal{N}(z|0,1) \nabla_\theta [\log p_\theta(G_\theta(z)) - \log p_{real}(G_\theta(z))] dz \\ \nabla_\theta [\log p_\theta(G_\theta(z)) - \log p_{real}(G_\theta(z))] &= \nabla_\theta \log p_\theta(G_\theta(z)) - \nabla_\theta \log p_{real}(G_\theta(z)) \end{aligned} \quad (60)$$

The second term is easy:

$$\nabla_\theta \log p_{real}(G_\theta(z)) = \frac{1}{p_{real}(x)} \nabla_x \log p_{real}(x) \cdot \nabla_\theta G_\theta(z) \quad (61)$$

However, the first one is more tricky. We need to consider it as a function of two arguments.

$$\begin{aligned} \nabla_\theta \log p_\theta(G_\theta(z)) &= \frac{1}{p_\theta(x)} \nabla_\theta [p_\theta(G_\theta(z))] \\ &= \frac{1}{p_\theta(x)} [\nabla_x p_\theta(x) \cdot \nabla_\theta G_\theta(z) + \nabla_\theta p_\theta(x) \nabla_\theta \theta] \\ &= \frac{1}{p_\theta(x)} [\nabla_x p_\theta(x) \cdot \nabla_\theta G_\theta(z) + \nabla_\theta p_\theta(x)] \\ &= \nabla_x \log p_\theta(x) \cdot \nabla_\theta G_\theta(z) + \color{blue}{\nabla_\theta \log p_\theta(x)} \end{aligned} \quad (62)$$

The first term is the same as in the previous formula but with the fake score  $\nabla_x \log p_\theta(x)$ . However, the blue one is intractable. To solve this issue, the authors consider the fake score as a function independent of  $\theta$ , that is  $p_\theta(x) = p_{fake}(x)$ , then  $\nabla_\theta \log p_\theta(x) = 0$ . So, we obtain the following:

$$\boxed{\nabla_\theta \text{KL}(p_\theta || p_{real}) = \mathbb{E}_{x=G_\theta(z), z \sim \mathcal{N}(0,1)} [\nabla_x \log p_{fake}(x) - \nabla_x \log p_{real}(x)] \nabla_\theta G_\theta(z)} \quad (63)$$

The authors note: *Computing this gradient is still challenging for two reasons: first, the scores diverge for samples with low probability — in particular preal vanishes for fake samples, and second, our intended tool for estimating score, namely the diffusion models, only provide scores of the diffused distribution. Score-SDE [71, 72] provides an answer to these two issues. By perturbing the data distribution with random Gaussian*

noise of varying standard deviations, we create a family of “blurred” distributions that are fully-supported over the ambient space, and therefore overlap, so that the gradient in Eq. (2) is well-defined (Figure 4).

$$\begin{aligned} \nabla_{\theta} \text{KL}(p_{\theta} || p_{\text{real}}) &= \mathbb{E}_{x=G_{\theta}(z), z \sim \mathcal{N}(0,1)} [\nabla_x \log p_{\text{fake}}(x) - \nabla_x \log p_{\text{real}}(x)] \nabla_{\theta} G_{\theta}(z) + \\ &\quad + \mathbb{E}_{x=G_{\theta}(z), z \sim \mathcal{N}(0,1)} \nabla_{\theta} \log p_{\theta}(x) \end{aligned} \quad (64)$$

Want but cannot:

$$\mathbb{E}_{x=G_{\theta}(z), z \sim \mathcal{N}(0,1)} \nabla_{\theta} \log p_{\theta}(x) = -\nabla_{\theta} \mathbb{E}_z \log \frac{p_{\text{real}}(x)}{p_{\theta}(x)} = \text{KL}(p_{\text{real}} || p_{\theta}) \quad (65)$$

The real score is a pretrained teacher diffusion. The fake score trains dynamically on generated samples using diffusion loss since the fake distribution changes during distillation:

$$\boxed{\mathcal{L}_{\text{diffusion}} = \mathbb{E}_{z \sim \mathcal{N}(0,1)} \|x_{\phi}(x_t, t) - G_{\theta}(z)\|^2, x_t \sim p_t(x_t | G_{\theta}(z), \sigma_t^2)} \quad (66)$$

$x_{\phi}(x_t, t)$  approximates the fake score. Moreover, authors consider the classical regression loss (knowledge distillation) calling it as a regression term.

**The authors note that without the regression term, the only reverse-KL does not converge, maybe because of the  $p_{\theta}(x) = p_{\text{fake}}(x)$ ?** This question is considered in the ‘score identity distillation’ paper.

## 5.2 ADD

In Adversarial Diffusion Distillation there is also a score distillation loss but in a slightly different form.

$$\mathcal{L}_{\text{ADD}} = \mathbb{E}_{z,t} \|G_{\theta}(z) - x_{\psi}(G_{\theta}^t(z), t)\|^2, G_{\theta}^t(z) = G_{\theta}(z) + \sigma_t^2 \epsilon \quad (67)$$

The provide ”the same” gradients:

$$\nabla_{\theta} \mathcal{L}_{\text{ADD}} = \mathbb{E}_{z,t} 2 [G_{\theta}(z) - x_{\psi}(G_{\theta}^t(z), t)] \nabla_{\theta} G_{\theta}(z), \nabla_x x_{\psi}(x, t) = 1??? \quad (68)$$

The loss is written in terms of  $x_0$  estimates,  $x_{\psi}$  plays the role of real score. We can rewrite in terms of scores:

$$\begin{aligned} \nabla_x \log p_{\text{real}}(G_{\theta}^t(z)) &= -\frac{G_{\theta}^t(z) - x_{\psi}(G_{\theta}^t(z), t)}{\sigma^2} \\ \nabla_x \log p_{\text{fake}}(G_{\theta}^t(z)) &= -\frac{G_{\theta}^t(z) - x_{\phi}(G_{\theta}^t(z), t)}{\sigma^2} \end{aligned} \quad (69)$$

If we set that  $x_{\phi}(G_{\theta}^t(z), t) = G_{\theta}(z)$ , then we obtain the same gradient as in the DMD:

$$\nabla_{\theta} \mathcal{L}_{\text{ADD}} = \mathbb{E}_{z,t} 2 [\nabla_x \log p_{\text{fake}}(x) - \nabla_x \log p_{\text{real}}(x)] \nabla_{\theta} G_{\theta}(z) \quad (70)$$

In other words, in ADD the authors does not learn the fake score while in the DMD it is happening dynamically during training. Instead they approximate the fake score using the distilled model.

### 5.3 Score identity distillation

In this work, the authors consider the score matching loss between fake and real losses.

$$\mathcal{L}_\theta = \mathbb{E}_{z \sim \mathcal{N}(0,1)} \|\nabla_{G_\theta^t(z)} \log p_{real}(G_\theta^t(z)) - \nabla_{G_\theta^t(z)} \log p_\theta(G_\theta^t(z))\|^2, G_\theta^t(z) = G_\theta(z) + \sigma_t^2 \epsilon \quad (71)$$

As it can be seen, this is a slightly different loss, we match scores instead of densities. The fake score can be trained similarly via diffusion loss using an additional network.

$$\mathcal{L}_{\text{diffusion}}(\phi) = \mathbb{E}_{z \sim \mathcal{N}(0,1)} \|x_\phi(x_t, t) - G_\theta(z)\|^2, x_t \sim p_t(x_t | G_\theta(z)) \quad (72)$$

In this work, the authors reason what happens if the real score depends on the generator  $\phi = \phi(\theta)$ . So, they firstly assume, as previously,  $\phi = \phi$ . Thus,

$$\begin{aligned} \mathcal{L}_\theta \approx \mathcal{L}_\theta^1 &= \mathbb{E}_{z \sim \mathcal{N}(0,1)} \|x_\phi(G_\theta^t(z), t) - x_\psi(G_\theta^t(z), t)\|^2 \\ &= \mathbb{E}_{z \sim \mathcal{N}(0,1)} \|\delta_{\phi,\psi}(G_\theta^t(z))\|^2 \end{aligned} \quad (73)$$

Lets consider the difference between optimal solution  $\phi^*(\theta)$  depending on the  $\theta$  and our approximation  $\phi$ .

$$\begin{aligned} \Delta_{\phi^*(\theta),\phi}(G_\theta^t(z)) &= x_{\phi^*(\theta)}(G_\theta^t(z), t) - x_\phi(G_\theta^t(z), t) \\ \delta_{\phi,\psi}(G_\theta^t(z)) &= \delta_{\phi^*(\theta),\psi}(G_\theta^t(z)) - \Delta_{\phi^*(\theta),\phi}(G_\theta^t(z)) \end{aligned} \quad (74)$$

Thus, we obtain

$$\begin{aligned} \mathcal{L}_\theta^1 &= \mathcal{L}_\theta + \mathbb{E}_{z \sim \mathcal{N}(0,1)} \|\Delta_{\phi^*(\theta),\phi}(G_\theta^t(z))\|^2 - \\ &\quad - 2\Delta_{\phi^*(\theta),\phi}(G_\theta^t(z))^T \cdot \delta_{\phi^*(\theta),\psi}(G_\theta^t(z)) \end{aligned} \quad (75)$$

The second term we minimize via our diffusion loss. However, the third one requires that the fake score,  $\phi^*(\theta)$ , equals to the real score,  $\psi$ , controlled by  $\delta_{\phi^*(\theta),\psi}(G_\theta^t(z))$ . This loss is large especially in the initial training steps, when  $p_\theta$  does not match  $p_{\text{data}}$ .

#### 5.3.1 1D Gaussian

Let  $p_{\text{data}}(x_0) = \mathcal{N}(x_0 | 0, 1)$ ,  $p_t(x_t | G_\theta(z)) = \mathcal{N}(x_t | G_\theta(z), \sigma_t^2)$ ,  $p_\theta(G_\theta(z)) = \mathcal{N}(G_\theta(z) | \theta, 1)$ . Find the optimal  $\theta$  for both  $\mathcal{L}_\theta^1$  and  $\mathcal{L}_\theta$ .

For  $\mathcal{L}_\theta^1$  and  $\mathcal{L}_\theta$  we need to know the conditional expectation. We know that (from the toy examples in SDE/ODE framework file)

$$\begin{aligned} p_t(x_t) &= \mathcal{N}(x_t | 0, \sigma_t^2 + 1) \rightarrow \nabla_{x_t} \log p_t(x_t) = -\frac{x_t}{1+\sigma_t^2} \rightarrow \mathbb{E} x_0 | x_t = \frac{x_t}{1+\sigma_t^2} = x_\psi(x_t, t) \\ p_t^\theta(x_t) &= \mathcal{N}(x_t | \theta, \sigma_t^2 + 1) \rightarrow \nabla_{x_t} \log p_t^\theta(x_t) = -\frac{x_t - \theta}{1+\sigma_t^2} \rightarrow \mathbb{E} G_\theta | x_t = \frac{x_t + \theta \sigma_t^2}{1+\sigma_t^2} = x_{\phi^*(\theta)}(x_t, t); \\ \text{Let's take } x_\phi(x_t, t) &= \frac{x_t + \phi \sigma_t^2}{1+\sigma_t^2}, \text{ then } \phi^*(\theta) = \theta. \text{ Finally, we have} \end{aligned}$$

$$\mathcal{L}_\theta^1 = \frac{\phi \sigma_t^2}{1 + \sigma_t^2}, \quad \mathcal{L}_\theta = \frac{\theta \sigma_t^2}{1 + \sigma_t^2}$$

(76)

We can see that  $\mathcal{L}_\theta^1$  does not provide meaningful gradients.

Let's also consider the DMD case from the toy example perspective.

$$\begin{aligned} \nabla_\theta \text{KL}(p_\phi || p_{\text{real}}) &= \mathbb{E}_{z \sim \mathcal{N}(0,1)} \left( \nabla_{G_\theta^t(z)} \log p_t^\phi(G_\theta^t(z)) - \nabla_{G_\theta^t(z)} \log p_t(G_\theta^t(z)) \right) \nabla_\theta G_\theta(z) \\ \nabla_{x_t} \log p_t^\phi(x_t) &= -\frac{x_t - \phi}{1 + \sigma_t^2}; \quad \nabla_{x_t} \log p_t(x_t) = -\frac{x_t}{1 + \sigma_t^2} \end{aligned} \quad (77)$$

So, we do not have the problem here(?)

$$\nabla_{\theta} \text{KL}(p_{\phi} || p_{real}) = \mathbb{E}_{z \sim \mathcal{N}(0,1)} \left( \frac{\phi}{1 + \sigma_t^2} \right) \nabla_{\theta} G_{\theta}(z) \quad (78)$$

**Authors method.**

**Note 1.**

$$\begin{aligned} \mathbb{E}_z \left[ u^T(G_{\theta}^t) \nabla_{G_{\theta}^t} \log p_{\theta}(G_{\theta}^t) \right] &= \mathbb{E}_z \left[ u^T(G_{\theta}^t) \frac{\nabla_{G_{\theta}^t} p_{\theta}(G_{\theta}^t)}{p_{\theta}(G_{\theta}^t)} \right] \\ &= \int u^T(G_{\theta}^t) \nabla_{G_{\theta}^t} p_{\theta}(G_{\theta}^t) dG_{\theta}^t \\ p_{\theta}(G_{\theta}^t) &= \int p_t(G_{\theta}^t | G_{\theta}) p_{\theta}(G_{\theta}) dz \\ \int u^T(G_{\theta}^t) \nabla_{G_{\theta}^t} p_{\theta}(G_{\theta}^t) dz &= \int u^T(G_{\theta}^t) \left[ \int p_t(G_{\theta}^t | G_{\theta}) \nabla_{G_{\theta}^t} \log p_t(G_{\theta}^t | G_{\theta}) p_{\theta}(G_{\theta}) dz \right] dG_{\theta}^t \\ \boxed{\mathbb{E}_z \left[ u^T(G_{\theta}^t) \nabla_{G_{\theta}^t} \log p_{\theta}(G_{\theta}^t) \right]} &= \mathbb{E}_{G_{\theta}^t | G_{\theta}, z} \left[ u^T(G_{\theta}^t) \nabla_{G_{\theta}^t} \log p_t(G_{\theta}^t | G_{\theta}) \right] \end{aligned} \quad (79)$$

**Note 2.**

$$\begin{aligned} \mathcal{L}_{\theta} &= \mathbb{E}_z \| \nabla_{G_{\theta}^t(z)} \log p_{real}(G_{\theta}^t) - \nabla_{G_{\theta}^t} \log p_{\theta}(G_{\theta}^t) \|^2 \\ &= \mathbb{E}_z \| \frac{1}{\sigma_t^2} (x_{\psi}(G_{\theta}^t, t)) - x_{\phi^*(\theta)}(G_{\theta}^t, t) \|^2 \\ &= \frac{1}{\sigma_t^2} \mathbb{E}_z \left[ (x_{\psi}(G_{\theta}^t, t)) - x_{\phi^*(\theta)}(G_{\theta}^t, t) \right]^T \left( \nabla_{G_{\theta}^t(z)} \log p_{real}(G_{\theta}^t) - \nabla_{G_{\theta}^t} \log p_{\theta}(G_{\theta}^t) \right] \end{aligned} \quad (80)$$

We need to expand the brackets and use note 1 (see app.)

$$\boxed{\mathcal{L}_{\theta} = \frac{1}{\sigma_t^2} \mathbb{E}_{G_{\theta}^t | G_{\theta}, z} \left[ \delta_{\psi, \phi^*(\theta)}(G_{\theta}^t)^T (x_{\psi}(G_{\theta}^t, t) - G_{\theta}) \right]} \quad (81)$$

Let's consider the same approximation,  $\phi^*(\theta) = \phi$

$$\begin{aligned} \mathcal{L}_{\theta}^2 &= \frac{1}{\sigma_t^2} \mathbb{E}_{G_{\theta}^t | G_{\theta}, z} \left[ \delta_{\psi, \phi}(G_{\theta}^t)^T (x_{\psi}(G_{\theta}^t, t) - G_{\theta}) \right] \\ &= \mathcal{L}_{\theta} + \frac{1}{\sigma_t^2} \mathbb{E}_{G_{\theta}^t | G_{\theta}, z} \left[ \Delta_{\phi^*(\theta), \phi}(G_{\theta}^t)^T (x_{\psi}(G_{\theta}^t, t) - G_{\theta}) \right] \end{aligned} \quad (82)$$

As we can see, there is no  $\delta_{\phi^*(\theta), \psi}$ .

**Final loss**

$$\boxed{\begin{aligned} \mathcal{L}_{\theta} &= \frac{w(t)}{\sigma_t^2} \mathbb{E}_{G_{\theta}^t | G_{\theta}, z} \left[ (x_{\phi}(G_{\theta}^t, t) - x_{\psi}(G_{\theta}^t, t))^T (x_{\psi}(G_{\theta}^t, t) - G_{\theta}) \right] \\ &\quad - \alpha w(t) \mathbb{E}_z \| x_{\phi}(G_{\theta}^t, t) - x_{\psi}(G_{\theta}^t, t) \|^2 \end{aligned}} \quad (83)$$