# Data-free Distillation of Diffusion Models with Bootstrapping

Jiatao Gu, Chen Wang, Shuangfei Zhai, Yizhe Zhang, Lingjie Liu, Josh Susskind ICML 2024  $\cdot$  Apple Inc., UPenn

# Introduction and Motivation

TLDR: We introduce a distillation method that can reduce diffusion models into one-step inference without using any training data.



## **Problems with Existing Distillation Methods**

1). Direct distillation approaches need to generate noiseimage targets through multi-step DDIM sampling, which is expensive.

 $\mathcal{L}_{\theta}^{\text{Direct}} = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0,I)} \| \boldsymbol{g}_{\theta}(\boldsymbol{\epsilon}) - \texttt{ODE-Solver}(\boldsymbol{f}_{\phi}, \boldsymbol{\epsilon}, T \rightarrow 0) \|_{2}^{2}$ 

2). Consistency models removed the requirements of synthetic dataset by learning from "bootstrapping", but they need to access the real data and requires EMA to

avoid trivial solutions. Data



## **Our Method**

#### Overview

We introduce singal-ODE, a modeling technique focused exclusively on signals and its corresponding distillation process. We also adopt our method on text-to-image diffusion models.

#### Signal-ODE

In this work, we propose to learn from a "signal-ODE" derived from the original diffusion model by predict the "signal" part of the original DDIM path:

$$\boldsymbol{y}_t = (\boldsymbol{x}_t - \sigma_t \boldsymbol{\epsilon})/\alpha_t$$

We can rewrite DDIM sampling as:

$$\boldsymbol{y}_{s} = \left(1 - e^{\lambda_{s} - \lambda_{t}}\right) \boldsymbol{f}_{\phi}(\boldsymbol{x}_{t}, t) + e^{\lambda_{s} - \lambda_{t}} \boldsymbol{y}_{t}$$
$$\frac{\mathrm{d}\boldsymbol{y}_{t}}{\mathrm{d}t} = -\lambda_{t}' \cdot \left(\boldsymbol{f}_{\phi}(\boldsymbol{x}_{t}, t) - \boldsymbol{y}_{t}\right) \quad \text{Signal-ODE}$$

#### Learning with Bootstrapping

The training object is to match the student output at two different timesteps with the help of teacher ODE solvers. Through bootstrapping, the student model can learn to generate samples from noise directly.



Matching two sides of the signal ODE:

$$\mathcal{L}_{\theta}^{\text{DE}} = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0,I), t \sim [0,T]} \left\| \left| \frac{\mathrm{d} \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},t)}{\mathrm{d} t} + \lambda'_{t} \cdot (\boldsymbol{f}_{\phi}(\hat{\boldsymbol{x}}_{t},t) - \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},t)) \right\|_{2}^{2}$$

$$\begin{aligned} \mathcal{L}_{\theta} &= \mathbb{E}_{\boldsymbol{\epsilon},t} \left[ \tilde{\omega}_{t} \left\| \left\| \frac{\mathrm{d}\boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},t)}{\mathrm{d}t} + \lambda_{t}' \cdot \left(\boldsymbol{f}_{\phi}(\hat{\boldsymbol{x}}_{t},t) - \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},t)\right) \right\|_{2}^{2} \right] \\ &\approx \mathbb{E}_{\boldsymbol{\epsilon},t} \left[ \tilde{\omega}_{t} \left\| \frac{\boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},s) - \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},t)}{\delta} - \lambda_{t}' \left(\boldsymbol{f}_{\phi}(\hat{\boldsymbol{x}}_{t},t) - \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},t)\right) \right\|_{2}^{2} \right] \\ &= \mathbb{E}_{\boldsymbol{\epsilon},t} \left[ \frac{\tilde{\omega}_{t}}{\delta^{2}} \left\| \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},s) - \mathrm{SG} \left[ \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},t) + \underbrace{\delta\lambda_{t}' \left(\boldsymbol{f}_{\phi}(\hat{\boldsymbol{x}}_{t},t) - \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon},t)\right)}_{\text{incremental improvement}} \right] \right\|_{2}^{2} \right] \end{aligned}$$

No ema needed

#### **Error Accumulation**

Imperfect predictions at large time steps might propagate to subsequent timesteps during bootstrapping. We employ two methods: (1) we sample timesteps uniformly throughout training; (2) We use high-order solvers such as Heun's Method.

#### **Boundary Condition**

In theory, the boundary can be arbitrary values since  $\alpha_T = 0$ . We make the student model learn in a truncated range  $t \in [t_{\min}, t_{\max}]$  and add an auxiliary boundary loss:

 $\mathcal{L}_{\theta}^{\text{BC}} = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0,I)} \left[ \| \boldsymbol{f}_{\phi}(\boldsymbol{\epsilon}, t_{\max}) - \boldsymbol{y}_{\theta}(\boldsymbol{\epsilon}, t_{\max}) \|_{2}^{2} \right]$ 



(a) without boundary loss



## **Distillation of Text-to-Image Models**

Our method can be readily applied for distilling conditional diffusion models in either pixel space or latent space.

 $\tilde{\boldsymbol{f}}_{\phi}(\boldsymbol{x}_t, t, \boldsymbol{c}) = \boldsymbol{f}_{\phi}(\boldsymbol{x}_t, t, \boldsymbol{n}) + w \cdot (\boldsymbol{f}_{\phi}(\boldsymbol{x}_t, t, \boldsymbol{c}) - \boldsymbol{f}_{\phi}(\boldsymbol{x}_t, t, \boldsymbol{n}))$ 

## **Our Algorithm**

Algorithm 1 Distillation using BOOT for Conditional Diffusion Models. **Require:** pretrained diffusion model  $f_{\phi}$ , initial student parameter from the teacher  $\theta \leftarrow \phi$ , step size  $\delta$ , learning rate  $\eta$ , CFG weight w, context dataset  $\mathcal{D}$ , negative condition  $\boldsymbol{n} = \emptyset$ ,  $t_{\min}$ ,  $t_{\max}$ ,  $\beta$ . : while not converged do Sample noise input  $\epsilon \sim \mathcal{N}(0, I)$ Sample context input  $\boldsymbol{c} \sim \mathcal{D}$ Sample  $t \sim (t_{\min}, t_{\max}), s = \min(t - \delta, t_{\min})$ Compute noise schedule  $\alpha_t, \sigma_t, \alpha_s, \sigma_s$ Compute  $\lambda'_t \approx (1 - \frac{\alpha_t \sigma_s}{\delta})/\delta$ Generate the model predictions:  $oldsymbol{y}_t = oldsymbol{y}_ heta(oldsymbol{\epsilon},t,oldsymbol{c}), \quad oldsymbol{y}_s = oldsymbol{y}_ heta(oldsymbol{\epsilon},s,oldsymbol{c}), \quad oldsymbol{y}_{t_{ ext{max}}} = oldsymbol{y}_ heta(oldsymbol{\epsilon},t_{ ext{max}},oldsymbol{c})$ Generate the noisy sample  $\hat{x}_t = \alpha_t y_t + \sigma_t \epsilon$ Compute the denoised target:  $\tilde{\boldsymbol{f}}_t = \boldsymbol{f}_{\phi}(\hat{\boldsymbol{x}}_t, t, \boldsymbol{n}) + w \cdot (\boldsymbol{f}_{\phi}(\hat{\boldsymbol{x}}_t, t, \boldsymbol{c}) - \boldsymbol{f}_{\phi}(\hat{\boldsymbol{x}}_t, t, \boldsymbol{n}))$  $oldsymbol{f}_{t_{\max}} = oldsymbol{f}_{\phi}(oldsymbol{\epsilon}, t_{\max}, oldsymbol{n}) + w \cdot (oldsymbol{f}_{\phi}(oldsymbol{\epsilon}, t_{\max}, oldsymbol{c}) - oldsymbol{f}_{\phi}(oldsymbol{\epsilon}, t_{\max}, oldsymbol{n}))$ 13: Compute the bootstrapping loss  $\mathcal{L}_{\theta}^{BS} = \frac{1}{(\delta \lambda'_{t})^{2}} \| \boldsymbol{y}_{s} - SG(\boldsymbol{y}_{t} + \delta \lambda'_{t}(\tilde{\boldsymbol{f}}_{t} - \boldsymbol{y}_{t})) \|_{2}^{2}$ 14: Compute the boundary loss  $\mathcal{L}_{\theta}^{\text{BC}} = \| \boldsymbol{y}_{t_{\text{max}}} - \tilde{\boldsymbol{f}}_{t_{\text{max}}} \|_2^2$ 15: Update model parameters  $\theta \leftarrow \theta - \eta \cdot \nabla_{\theta} \left( \mathcal{L}_{\theta}^{BS} + \beta \mathcal{L}_{\theta}^{BC} \right)$ 16: **end while** 17: **return** Trained model parameters  $\theta$ 



# **Quantitative Results**

Method	NFE	FID
Diffusion Model+Solver		
DDPM (Ho et al., 2020)	1000	3.17
DDIM (Song et al., 2021)	50	4.67
DPM-solver-2 (Lu et al., 2022)	10	5.94
DEIS (Zhang & Chen, 2022)	10	4.17
EDM (Karras et al., 2022)	35	2.04
Distillation		
Direct* (Luhman & Luhman, 2021)	1	9.36
DFNO* (Zheng et al., 2023)	1	4.12
ReFlow* (Liu et al., 2022)	1	6.18
PD (Salimans & Ho, 2022)	1	8.34
CM (Song et al., 2023)	1	3.55
Data-free Distillation		
Ours (L2 loss)	1	6.88
Ours (LPIPS loss)	1	4.38

	Steps	FFHQ FID	64  imes 64 fps	ImageN FID	et 64 × 64 fps
DDPM	250	5.4	0.2	11.0	0.1
DDIM	50	7.6	1.2	13.7	0.6
	10	18.3	5.3	18.3	3.3
	1	225	54	237	34
Ours	1	9.0	54	12.3	34

Results on FFHQ and ImageNet

Results on CIFAR-10

# **Qualitative Results**



DDIM 8 steps (16 NFEs)





DDIM1 step (2 NFEs) LCM1 step (1 NFEs) Comparison with other baselines



BOOT 1 step (1 NFEs, Ours)



Random Samples from our single-step student model distilled from DeepFloyd-IF

# Links



## https://machinelearning.apple.com/research/boot

TM and © 2024 Apple Inc. All rights reserved.