



How well does **diffusion model** generate?  
- a **training** and **sampling** combined **quantification**

arXiv:2406.12839

*Yuqing Wang*<sup>1,2</sup>

*Ye He*<sup>1</sup>

*Molei Tao*<sup>1</sup>

1 Georgia Tech, USA



2 Simons Institute, UC Berkeley



Sep 10, 2024 Emerging Generalization Settings Workshop @ Simons

90%

quantification of diffusion model **generation**

90%      quantification of diffusion model **generation**

$\leq 10\%$       **generalization**

why am I here?

90%          quantification of diffusion model **generation**

$\leq 10\%$           **generalization**

why am I here?

- one interesting generalization setting: understood

90%          quantification of diffusion model **generation**

$\leq 10\%$           **generalization**

why am I here?

- one interesting generalization setting: understood
- another: not

90%          quantification of diffusion model **generation**

$\leq 10\%$           **generalization**

why am I here?

- one interesting generalization setting: understood
- another: not
- the quantifications are interesting too (hopefully)

## Generative Modeling

Given **samples** of an unknown probability distribution (possibly in very high dim.), generate **more samples** of the same distribution.

## Generative Modeling

Given **samples** of an unknown probability distribution (possibly in very high dim.), generate **more samples** of the same distribution.





denoising diffusion model

(*Sohl-Dickstein*+ 15,  
*Ho*+ 20, *Song*+ 21 ...)

denoising diffusion model

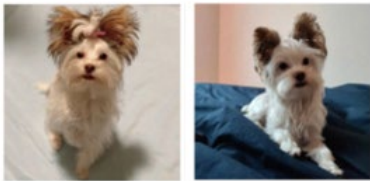
(*Sohl-Dickstein*+ 15,  
*Ho*+ 20, *Song*+ 21 ...)

Stable Diffusion, DALL·E, Midjourney; Sora;

(Chat)GPT, Gemini, Llama, Claude, ...

## denoising diffusion model

(*Sohl-Dickstein*+ 15,  
*Ho*+ 20, *Song*+ 21 ...)

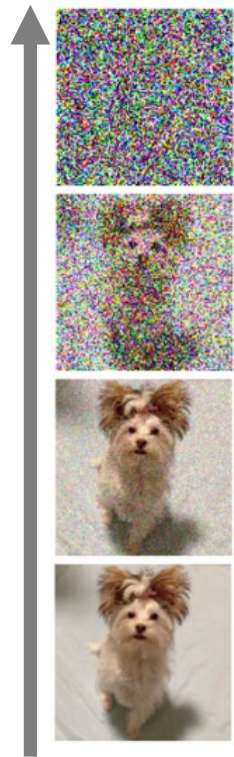


*Yang*+ 22

# denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

forward  
noising  
process:  
learn  
“score”  
(~  
evolution  
of  
data  
density)

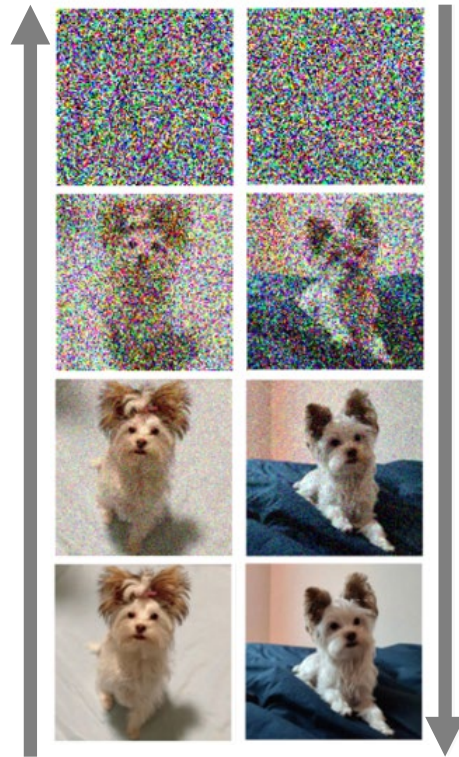


Yang+ 22

# denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

forward  
noising  
process:  
learn  
“score”  
(~  
evolution  
of  
data  
density)



backward  
denoising  
process:  
use  
“score”  
to  
generate  
data  
from  
noise

Yang+ 22

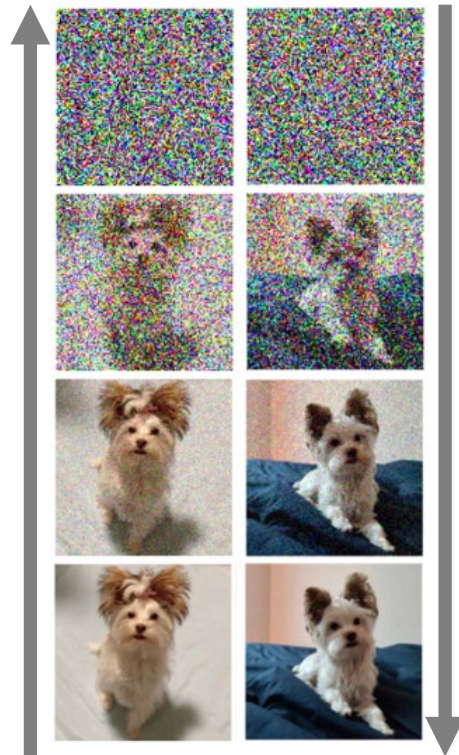
# denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

$$dX = -X dt + \sqrt{2}dW_t$$



forward  
noising  
process:  
learn  
“score”  
(~  
evolution  
of  
data  
density)

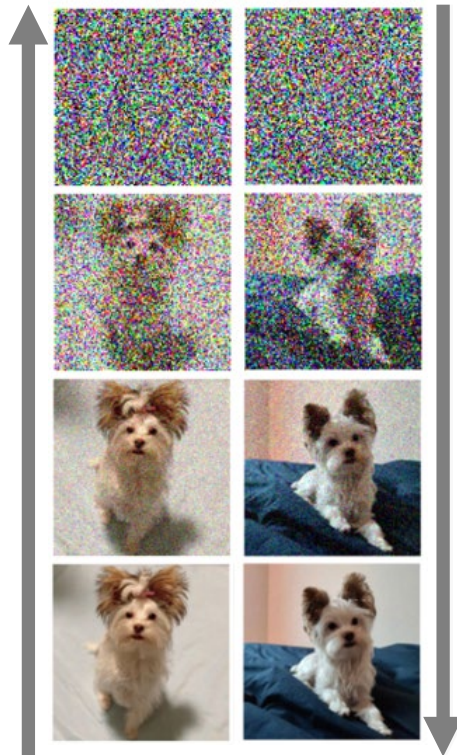


backward  
denoising  
process:  
use  
“score”  
to  
generate  
data  
from  
noise

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

forward  
noising  
process:  
learn  
“score”  
(~  
evolution  
of  
data  
density)



backward  
denoising  
process:  
use  
“score”  
to  
generate  
data  
from  
noise

Yang+ 22

$$dX = -X dt + \sqrt{2}dW_t$$

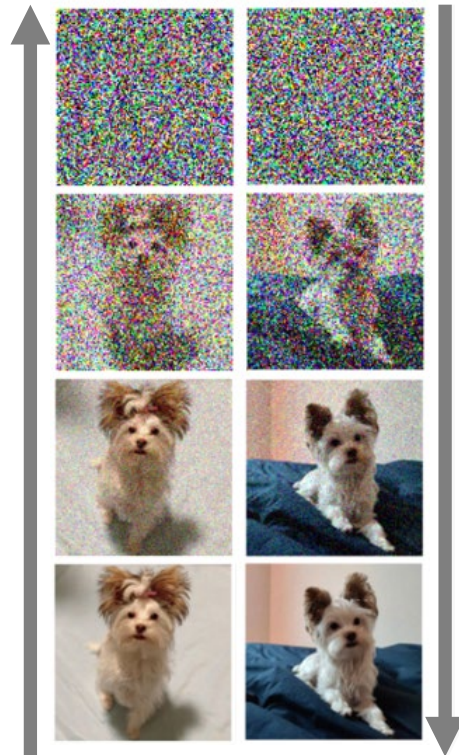
$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

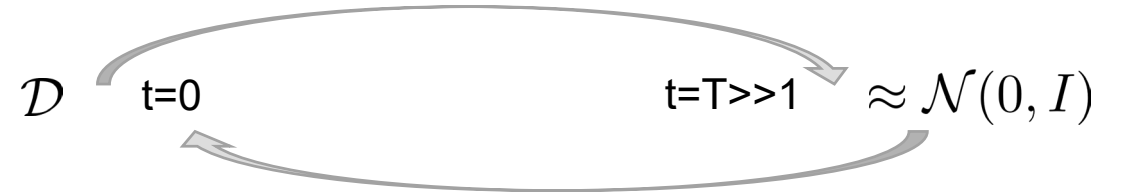
forward  
noising  
process:  
learn  
“score”  
(~  
evolution  
of  
data  
density)



backward  
denoising  
process:  
use  
“score”  
to  
generate  
data  
from  
noise

Yang+ 22

$$dX = -X dt + \sqrt{2}dW_t$$



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

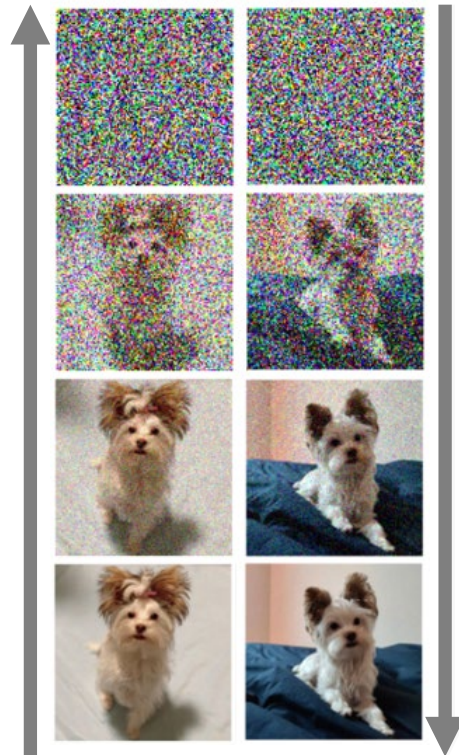
$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

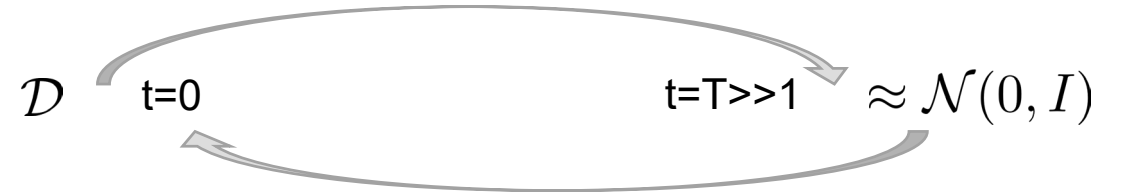


forward  
noising  
process:  
learn  
“score”  
(~  
evolution  
of  
data  
density)

backward  
denoising  
process:  
use  
“score”  
to  
generate  
data  
from  
noise

Yang+ 22

$$dX = -X dt + \sqrt{2}dW_t$$



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

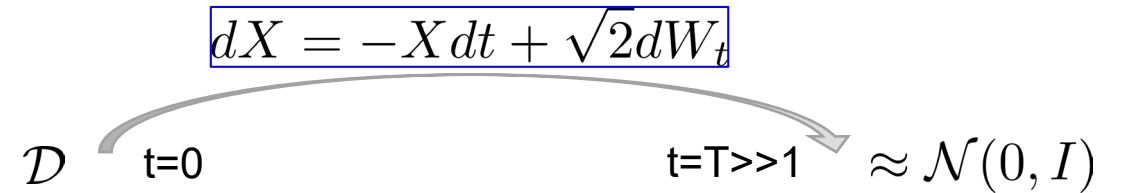
$$X(t) \stackrel{d}{=} Y(T - t), \forall t$$

$$X(t) \sim p(x, t)$$

## 0.2 Introduction: Denoising Diffusion Generative Model

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score ??

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 0.2 Introduction: Denoising Diffusion Generative Model

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

$$\mathcal{D} \quad t=0 \quad \xrightarrow{dX = -X dt + \sqrt{2}dW_t} \quad t=T \gg 1 \quad \approx \mathcal{N}(0, I)$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score ??

learned during forward process

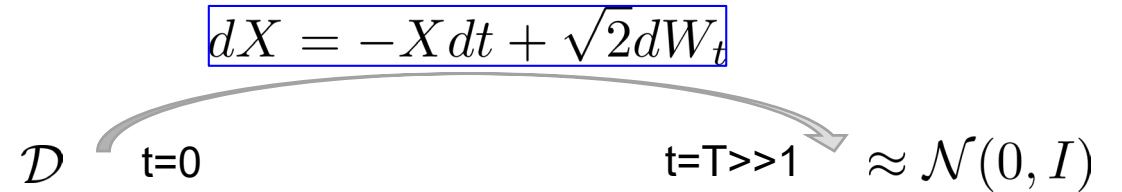


$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score ??

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



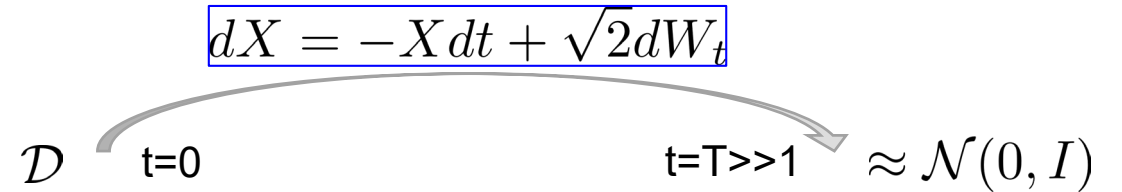
learned during forward process

score matching

$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score ??

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



learned during forward process

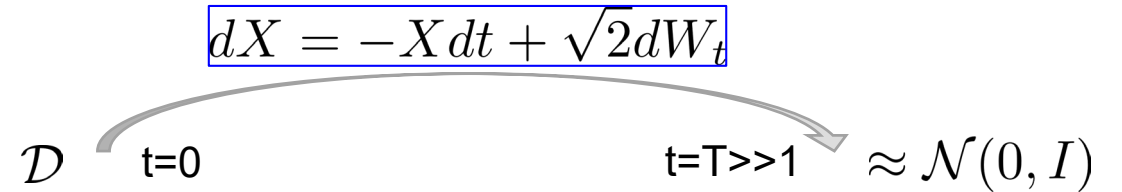
score matching

$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \underbrace{\|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2}_{\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2} dt$$

$$X_0^i \sim \mathcal{D}$$

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score ??

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



learned during forward process

score matching **unimplementable**

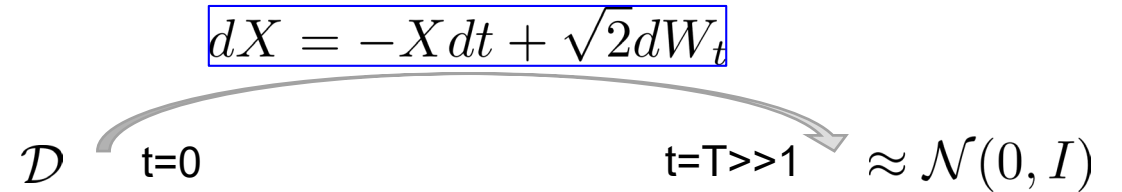
$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \left\| \nabla_x \log p(X_t, t) - s_{\theta}(X_t, t) \right\|^2 dt$$

$$\approx \frac{1}{n} \sum_{i=1}^n \left\| \nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t) \right\|^2$$

$$X_0^i \sim \mathcal{D}$$

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score ??

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

learned during forward process



score matching **unimplementable**

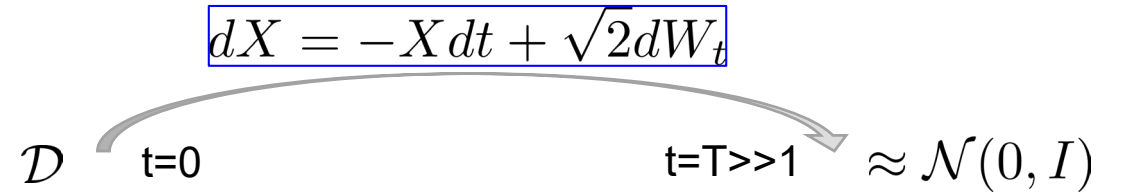
$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

denoising score matching

$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - s_{\theta}(X_t, t)\|^2 dt$$

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score ??

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

learned during forward process



score matching **unimplementable**

$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

denoising score matching

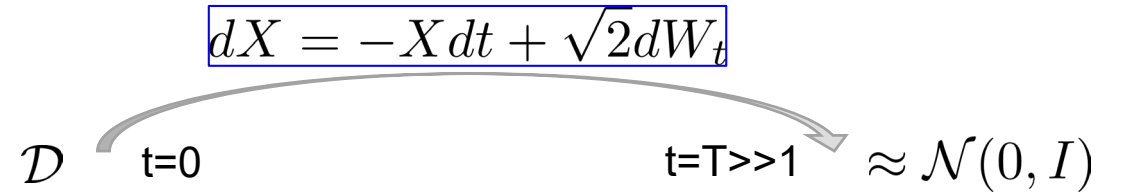
$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - s_{\theta}(X_t, t)\|^2 dt$$

analytically **available**



denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)



$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$

score ??

$s(x, t) := \nabla_x \log p(x, t)$

$X(t) \sim p(x, t)$



learned during forward process

score matching **unimplementable**

$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$

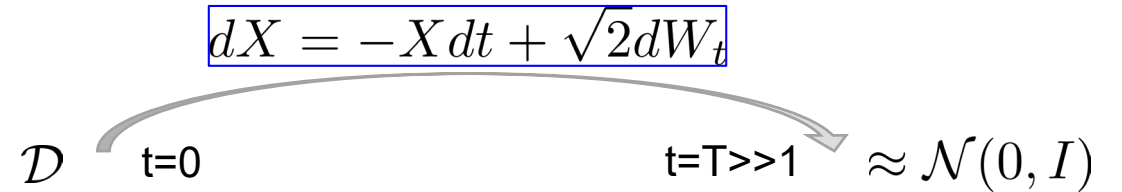
denoising score matching

$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - s_{\theta}(X_t, t)\|^2 dt$

analytically **available**

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)



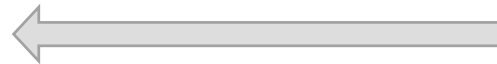
$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score ??

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

learned during forward process



score matching unimplementable

$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

denoising score matching less trivial than appeared, but provable

$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - s_{\theta}(X_t, t)\|^2 dt$$

analytically available

## 1.1 Quantification of Diffusion Model's Generation Quality: Overview

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

1. **forward** training/learning process:  
optimization  $\rightarrow$  score  $s$

$$\mathcal{D} \begin{array}{c} \xrightarrow{t=0} \\ \xrightarrow{t=T \gg 1} \end{array} \approx \mathcal{N}(0, I)$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

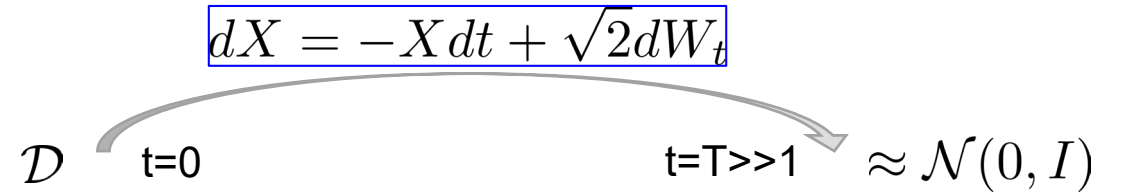
$$X(t) \sim p(x, t)$$

## 1.1 Quantification of Diffusion Model's Generation Quality: Overview

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

1. **forward** training/learning process:  
optimization  $\rightarrow$  score  $s$
2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

**score**

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.1 Quantification of Diffusion Model's Generation Quality: Overview

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

1. **forward** training/learning process:  
optimization  $\rightarrow$  score  $s$
2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

**?** quality of generated samples

$$\text{KL}(\text{Law}(X_0) | \text{Law}(Y_T)) \leq \dots$$

$$\mathcal{D} \xrightarrow[t=0]{dX = -X dt + \sqrt{2}dW_t} \xrightarrow[t=T \gg 1]{\approx \mathcal{N}(0, I)}$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

**score**

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

1. **forward** training/learning process:  
optimization  $\rightarrow$  score  $s$
2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

? quality of generated samples

$$\text{KL}(\text{Law}(X_0) | \text{Law}(Y_T)) \leq \dots$$

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

“main stream”

if score is approximated with error  $\leq \varepsilon$  in the sense of \_\_\_\_\_, then generated and training samples have statistical distance/divergence  $\leq$  \_\_\_\_, under assumptions \_\_\_\_\_.

Lee+ 22, de Bortoli 22, Yang & Wibisono 22, S Chen+ 23,  
H Chen+ 23, Benton+ 23, Conforti+ 23, ...

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

denoising diffusion model

(Sohl-Dickstein+ 15,  
Ho+ 20, Song+ 21 ...)

1. **forward** training/learning process:  
optimization  $\rightarrow$  score  $s$

only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

? quality of generated samples

$$\text{KL}(\text{Law}(X_0) | \text{Law}(Y_T)) \leq \dots$$

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

“main stream”

if score is approximated with error  $\leq \varepsilon$  in the sense of \_\_\_\_\_, then generated and training samples have statistical distance/divergence  $\leq$  \_\_\_\_, under assumptions \_\_\_\_\_.

Lee+ 22, de Bortoli 22, Yang & Wibisono 22, S Chen+ 23,  
H Chen+ 23, Benton+ 23, Conforti+ 23, ...

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$



## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1 \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

- score error

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1 \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error

$dX = -X dt + \sqrt{2}dW_t$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$dY = Y dt + 2s(Y, T-t)dt + \sqrt{2}dB_t$

$s_\theta$  score

$s(x, t) := \nabla_x \log p(x, t)$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1 \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error

common: exponential integrator

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

- score error
- integration error

common: exponential integrator

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

- score error
- integration error

common: **exponential integrator**

$Y_t$ : exact solution, N/A due to nonlinearity

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$dX = -X dt + \sqrt{2}dW_t$$

The diagram shows a forward process starting from a distribution  $\mathcal{D}$  at  $t=0$  and ending at a distribution  $\approx \mathcal{N}(0, I)$  at  $t=T \gg 1$ . A curved arrow points from  $t=0$  to  $t=T \gg 1$ .

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 numerical simulation  $\rightarrow$  sample  $Y$   
**only**

already highly nontrivial

### sources of error

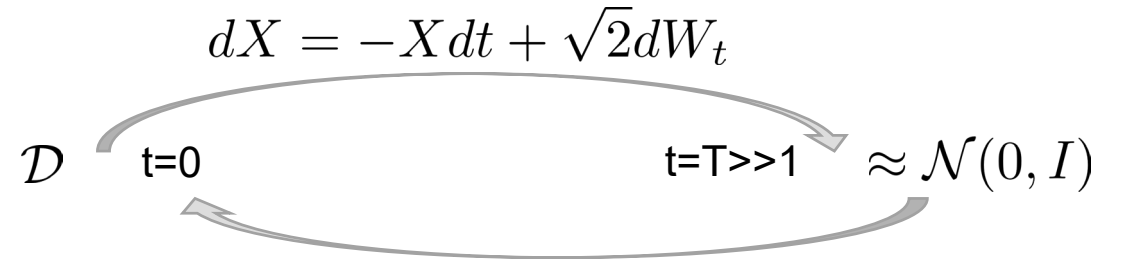
- score error
- integration error

common: **exponential integrator**

$Y_t$ : exact solution, N/A due to nonlinearity

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

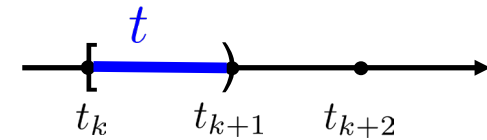
$$d\hat{Y}_t = \hat{Y}_t dt + 2s(\hat{Y}_{t_k}, T - t_k)dt + \sqrt{2}dB_t$$



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$





## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

### sources of error

- score error
- integration error

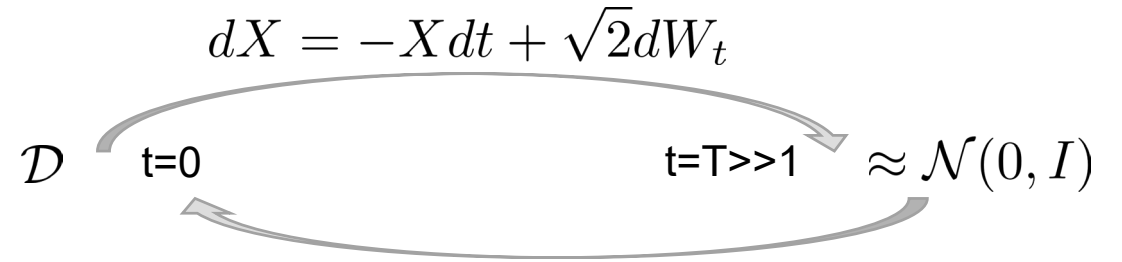
common: exponential integrator

$Y_t$ : exact solution, N/A due to nonlinearity

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$d\hat{Y}_t = \hat{Y}_t dt + 2s(\hat{Y}_{t_k}, T - t_k)dt + \sqrt{2}dB_t$$

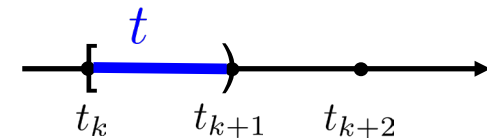
$\hat{Y}_t$ : exactly solvable, a numerical approximation of  $Y_t$



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$



## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

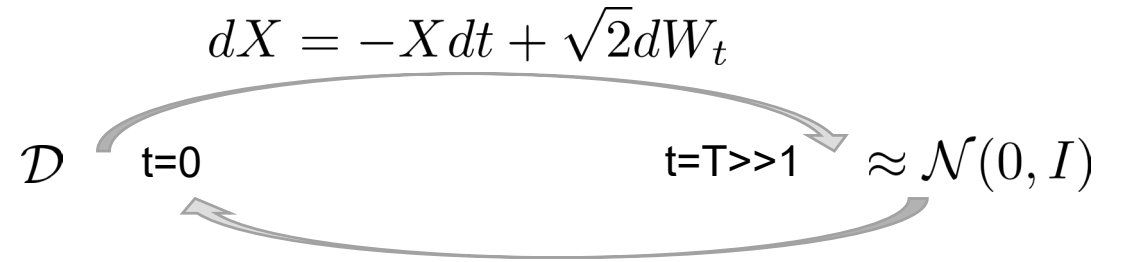
- score error
- integration error

efficient integration:  
lots of research

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$d\hat{Y}_t = \hat{Y}_t dt + 2s(\hat{Y}_{t_k}, T - t_k)dt + \sqrt{2}dB_t$$

$\hat{Y}_t$ : exactly solvable, a numerical approximation of  $Y_t$



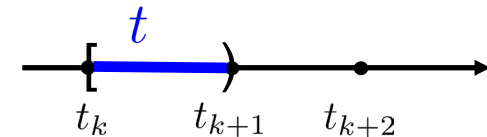
$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

common: exponential integrator

$Y_t$ : exact solution, N/A due to nonlinearity



## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error

efficient integration:  
lots of research

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

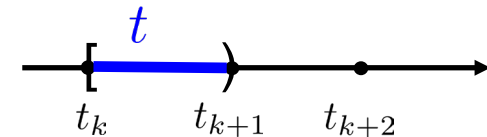
$$s(x, t) := \nabla_x \log p(x, t)$$

common: exponential integrator

$Y_t$ : exact solution, N/A due to nonlinearity

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$d\hat{Y}_t = \hat{Y}_t dt + 2s(\hat{Y}_{t_k}, T - t_k)dt + \sqrt{2}dB_t$$



$\hat{Y}_t$ : exactly solvable, a numerical approximation of  $Y_t$

$\hat{Y}_{t_k} \mapsto \hat{Y}_{t_{k+1}}$ : one function evaluation (of  $s$ )

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error

efficient integration:

lots of research

e.g., DEIS, gDDIM

(Zhang & Chen 2022)

(Zhang, Tao & Chen 2022)

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

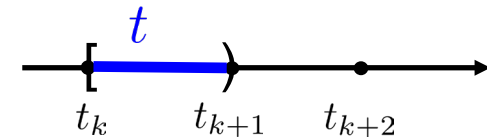
$$s(x, t) := \nabla_x \log p(x, t)$$

common: exponential integrator

$Y_t$ : exact solution, N/A due to nonlinearity

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$d\hat{Y}_t = \hat{Y}_t dt + 2s(\hat{Y}_{t_k}, T - t_k)dt + \sqrt{2}dB_t$$



$\hat{Y}_t$ : exactly solvable, a numerical approximation of  $Y_t$

$\hat{Y}_{t_k} \mapsto \hat{Y}_{t_{k+1}}$ : one function evaluation (of  $s$ )

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error

efficient integration:

lots of research

e.g., DEIS, gDDIM

(Zhang & Chen 2022)

(Zhang, Tao & Chen 2022)

NFE~10-20

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

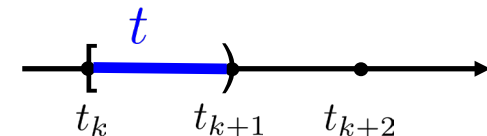
$$s(x, t) := \nabla_x \log p(x, t)$$

common: exponential integrator

$Y_t$ : exact solution, N/A due to nonlinearity

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$d\hat{Y}_t = \hat{Y}_t dt + 2s(\hat{Y}_{t_k}, T - t_k)dt + \sqrt{2}dB_t$$



$\hat{Y}_t$ : exactly solvable, a numerical approximation of  $Y_t$

$\hat{Y}_{t_k} \mapsto \hat{Y}_{t_{k+1}}$ : one function evaluation (of  $s$ )

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

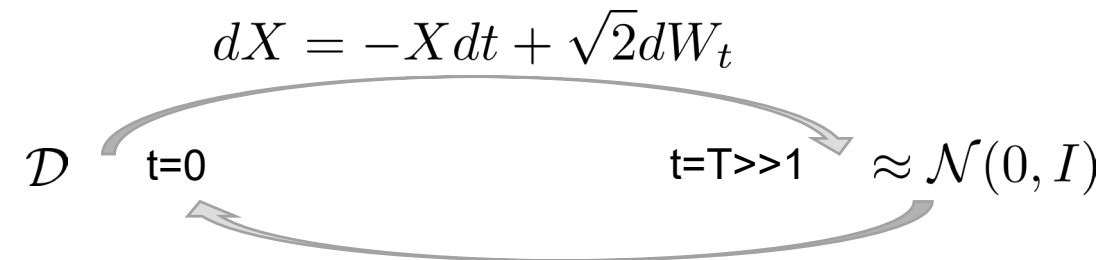
already highly nontrivial

sources of error

- score error
- integration error

efficient integration:  
 lots of research  
 e.g., DEIS, gDDIM  
 (Zhang & Chen 2022)  
 (Zhang, Tao & Chen 2022)

NFE~10-20  
 $\rightarrow$  distillation, etc.



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

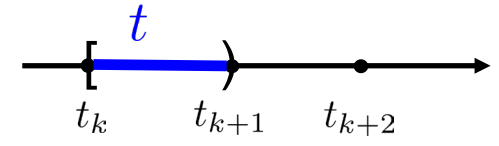
$$s(x, t) := \nabla_x \log p(x, t)$$

common: exponential integrator

$Y_t$ : exact solution, N/A due to nonlinearity

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$d\hat{Y}_t = \hat{Y}_t dt + 2s(\hat{Y}_{t_k}, T - t_k)dt + \sqrt{2}dB_t$$



$\hat{Y}_t$ : exactly solvable, a numerical approximation of  $Y_t$

$\hat{Y}_{t_k} \mapsto \hat{Y}_{t_{k+1}}$ : one function evaluation (of  $s$ )

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error

efficient integration:

lots of research

e.g., DEIS, gDDIM

(Zhang & Chen 2022)

(Zhang, Tao & Chen 2022)

NFE~10-20

$\rightarrow$  **distillation**, etc.

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

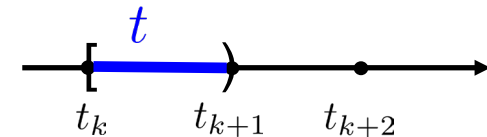
$$s(x, t) := \nabla_x \log p(x, t)$$

common: **exponential integrator**

$Y_t$ : exact solution, N/A due to nonlinearity

$$dY_t = Y_t dt + 2s(Y_t, T - t)dt + \sqrt{2}dB_t$$

$$d\hat{Y}_t = \hat{Y}_t dt + 2s(\hat{Y}_{t_k}, T - t_k)dt + \sqrt{2}dB_t$$



$\hat{Y}_t$ : exactly solvable, a numerical approximation of  $Y_t$

$\hat{Y}_{t_k} \mapsto \hat{Y}_{t_{k+1}}$ : one function evaluation (of  $s$ )

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

- score error
- integration error
- initialization error

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

- score error
- integration error
- initialization error

ideally

$$Y(0) \stackrel{d}{=} X(\infty) \sim \mathcal{N}(0, I)$$

$$dX = -X dt + \sqrt{2}dW_t$$

The diagram shows a forward process starting at  $\mathcal{D}$  at  $t=0$  and ending at  $\approx \mathcal{N}(0, I)$  at  $t=T \gg 1$ . A curved arrow points from  $\mathcal{D}$  to the Gaussian distribution. A second curved arrow points from the Gaussian distribution back to  $\mathcal{D}$ , representing the backward process.

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

- score error
- integration error
- initialization error

ideally

$$Y(0) \stackrel{d}{=} X(\infty) \sim \mathcal{N}(0, I)$$

in reality

$T$  large but finite

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error
- initialization error

ideally

$$Y(0) \stackrel{d}{=} X(\infty) \sim \mathcal{N}(0, I)$$

in reality

$T$  large but finite

then, ideally

$$Y(0) \stackrel{d}{=} X(T) \Rightarrow Y(T) \stackrel{d}{=} X(0)$$

$$dX = -X dt + \sqrt{2} dW_t$$

$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t) dt + \sqrt{2} dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error
- initialization error

ideally

$$Y(0) \stackrel{d}{=} X(\infty) \sim \mathcal{N}(0, I)$$

in reality

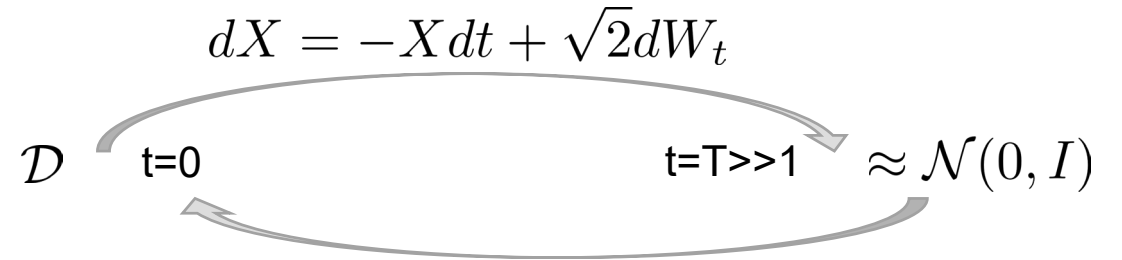
$T$  large but finite

then, ideally

$$Y(0) \stackrel{d}{=} X(T) \Rightarrow Y(T) \stackrel{d}{=} X(0)$$

but, in reality

$$\text{Law } X(T) \text{ unknown} \Rightarrow Y(0) \stackrel{d}{=} X(\infty) \sim \mathcal{N}(0, I)$$



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

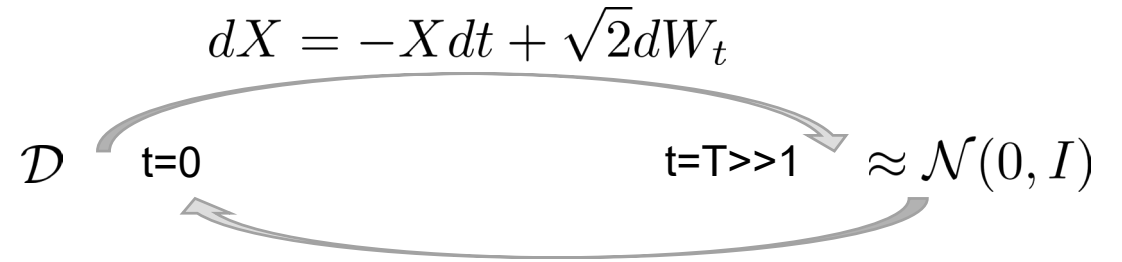
# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. backward sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

## sources of error

- score error
- integration error
- initialization error



ideally

$$Y(0) \stackrel{d}{=} X(\infty) \sim \mathcal{N}(0, I)$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

in reality

$T$  large but finite

$$X(t) \sim p(x, t)$$

then, ideally

$$Y(0) \stackrel{d}{=} X(T)$$

discrepancy propagated thru Y dynamics → error

but, in reality

$$\text{Law } X(T) \text{ unknown} \Rightarrow Y(0) \stackrel{d}{=} X(\infty) \sim \mathcal{N}(0, I)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

- score error
- integration error
- initialization error
- early stopping

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1 \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

- score error
- integration error
- initialization error
- early stopping

common:  
low dim. data manifold assumption

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1 \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

- score error
- integration error
- initialization error
- early stopping

common:  
low dim. data manifold assumption

$p(\cdot, 0)$  is supported on  $\mathcal{M}$

$x \in \mathbb{R}^D, \quad D > \dim \mathcal{M}$

$$dX = -X dt + \sqrt{2} dW_t$$

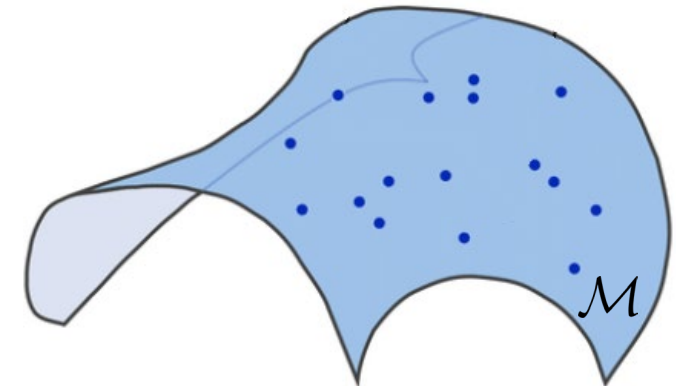
$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t) dt + \sqrt{2} dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$





## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

### sources of error

- score error
- integration error
- initialization error
- early stopping

common:  
 low dim. data manifold assumption  
 $p(\cdot, 0)$  is supported on  $\mathcal{M}$   
 $x \in \mathbb{R}^D, \quad D > \dim \mathcal{M}$   
 $\Rightarrow p(\cdot, 0)$  has a jump in  $n$  direction

$$dX = -X dt + \sqrt{2} dW_t$$

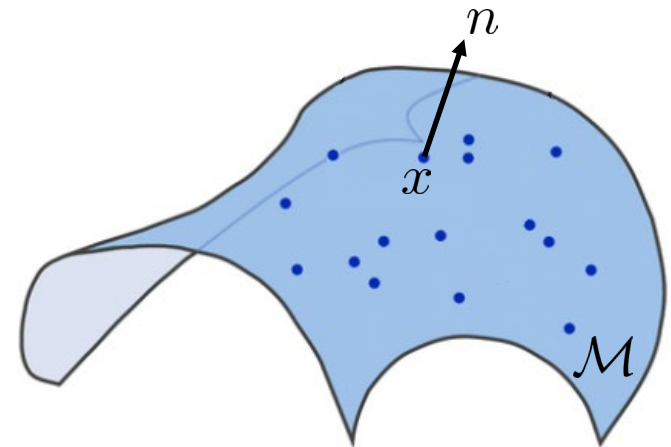
$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t) dt + \sqrt{2} dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

### sources of error

- score error
- integration error
- initialization error
- early stopping

common:  
 low dim. data manifold assumption

$p(\cdot, 0)$  is supported on  $\mathcal{M}$

$x \in \mathbb{R}^D, \quad D > \dim \mathcal{M}$

$\Rightarrow p(\cdot, 0)$  has a jump in  $n$  direction

$\Rightarrow \langle \nabla_x \log p(x, 0), n \rangle = \infty$

i.e. score at  $t=0$  is ill-defined

$$dX = -X dt + \sqrt{2} dW_t$$

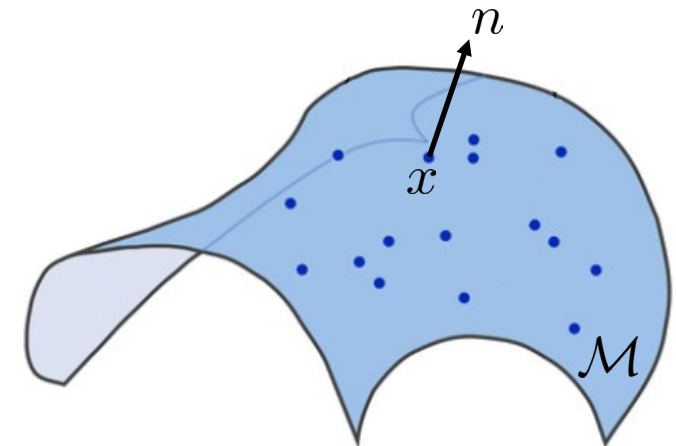
$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T-t) dt + \sqrt{2} dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

### sources of error

- score error
- integration error
- initialization error
- early stopping

common:  
 low dim. data manifold assumption

$p(\cdot, 0)$  is supported on  $\mathcal{M}$

$x \in \mathbb{R}^D$ ,  $D > \dim \mathcal{M}$

$\Rightarrow p(\cdot, 0)$  has a jump in  $n$  direction

$\Rightarrow \langle \nabla_x \log p(x, 0), n \rangle = \infty$

i.e. score at  $t=0$  is ill-defined

$$dX = -X dt + \sqrt{2} dW_t$$

$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T-t) dt + \sqrt{2} dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



$Y(T)$  not usable

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

### sources of error

- score error
- integration error
- initialization error
- early stopping

common:  
 low dim. data manifold assumption

$p(\cdot, 0)$  is supported on  $\mathcal{M}$

$x \in \mathbb{R}^D$ ,  $D > \dim \mathcal{M}$

$\Rightarrow p(\cdot, 0)$  has a jump in  $n$  direction

$\Rightarrow \langle \nabla_x \log p(x, 0), n \rangle = \infty$

i.e. score at  $t=0$  is ill-defined

$$dX = -X dt + \sqrt{2} dW_t$$

$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T-t) dt + \sqrt{2} dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



$Y(T)$  not usable

use  $Y(T-\delta)$  instead

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

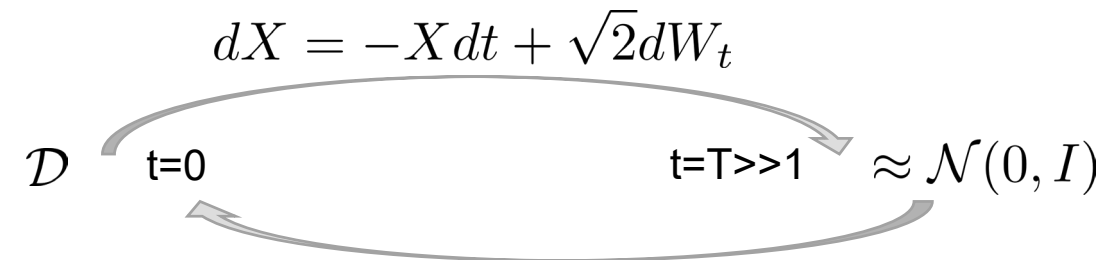
already highly nontrivial

### sources of error

- score error
- integration error
- initialization error
- early stopping

common:  
 low dim. data manifold assumption  
 $p(\cdot, 0)$  is supported on  $\mathcal{M}$   
 $x \in \mathbb{R}^D, \quad D > \dim \mathcal{M}$   
 $\Rightarrow p(\cdot, 0)$  has a jump in  $n$  direction  
 $\Rightarrow \langle \nabla_x \log p(x, 0), n \rangle = \infty$

i.e. score at  $t=0$  is ill-defined



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

$\rightarrow$   $Y(T)$  not usable  
 use  $Y(T-\delta)$  instead

### accuracy quantifier

$$KL(\text{Law}(X_0) | \text{Law}(Y_T))$$

# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. backward sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

## sources of error

- score error
- integration error
- initialization error
- early stopping

common:  
low dim. data manifold assumption  
 $p(\cdot, 0)$  is supported on  $\mathcal{M}$   
 $x \in \mathbb{R}^D, \quad D > \dim \mathcal{M}$   
 $\Rightarrow p(\cdot, 0)$  has a jump in  $n$  direction  
 $\Rightarrow \langle \nabla_x \log p(x, 0), n \rangle = \infty$

## accuracy quantifier

i.e. score at  $t=0$  is ill-defined

~~$\text{KL}(\text{Law}(X_0) | \text{Law}(Y_T))$~~  →

$$dX = -X dt + \sqrt{2} dW_t$$

$\mathcal{D} \xrightarrow{t=0} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T-t) dt + \sqrt{2} dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

→ Y(T) not usable

use  $Y(T-\delta)$  instead

$$\text{KL}(\text{Law}(X_\delta) | \text{Law}(Y_{T-\delta})) + W_2(\text{Law}(X_0), \text{Law}(X_\delta))^2$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

assumptions on data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1 \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
**only** numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1 \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$



## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

(overdamped) Langevin dynamics:

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$\mathcal{D} \xrightarrow[t=0]{} \xrightarrow[t=T \gg 1]{} \approx \mathcal{N}(0, I)$$

$dX = -X dt + \sqrt{2}dW_t$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

(overdamped) Langevin dynamics:

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$Z(\infty) \sim p(\cdot, 0)$$

need **long** time to conv.

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

(overdamped) Langevin dynamics:

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

How long?

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

$$Z(\infty) \sim p(\cdot, 0)$$

need **long** time to conv.

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation  $\rightarrow$  sample  $Y$

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

(overdamped) Langevin dynamics:

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

How long? Depends on  $V$  and thus data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $\approx \mathcal{N}(0, I)$   $t=T \gg 1$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

$Z(\infty) \sim p(\cdot, 0)$   
 need **long** time to conv.

# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

(overdamped) Langevin dynamics:

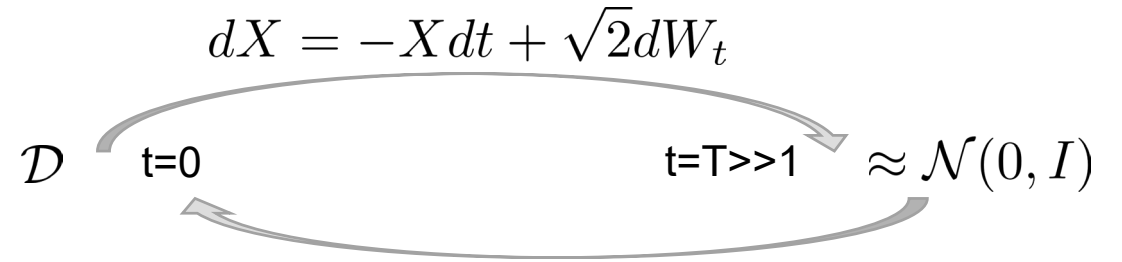
$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

How long? Depends on V and thus data distribution

~1/m

m-convex

unimodal



$$dY = Y dt + 2s(Y, T-t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

$$Z(\infty) \sim p(\cdot, 0)$$

need **long** time to conv.

# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

(overdamped) Langevin dynamics:

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$Z(\infty) \sim p(\cdot, 0)$$

need **long** time to conv.

How long? Depends on V and thus data distribution

~1/m

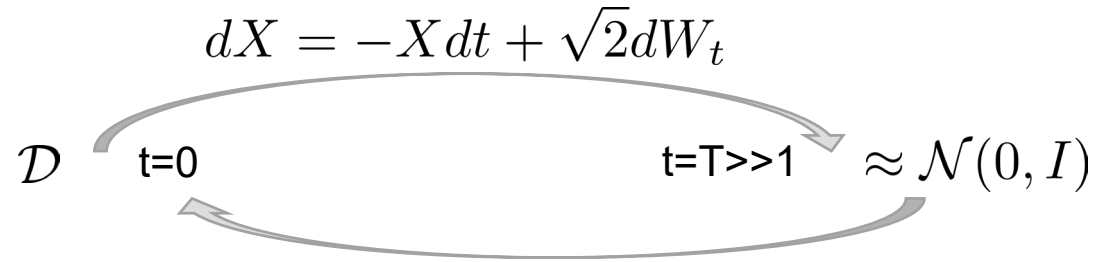
m-convex

unimodal

can be exp. longer

nonconvex

multimodal



score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

(overdamped) Langevin dynamics:

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

How long? Depends on V and thus data distribution

~1/m

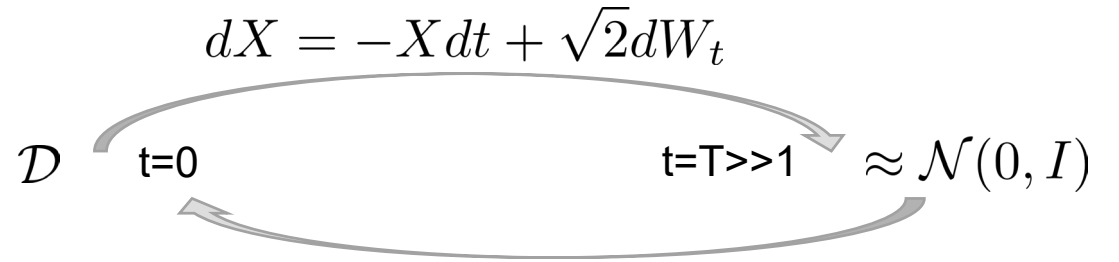
m-convex

unimodal

can be exp. longer

nonconvex

multimodal



need T time to conv.  
 reach  $p(\cdot, 0)$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

$$Z(\infty) \sim p(\cdot, 0)$$

need long time to conv.

# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution

Ex (prior to diffusion model & its analysis)

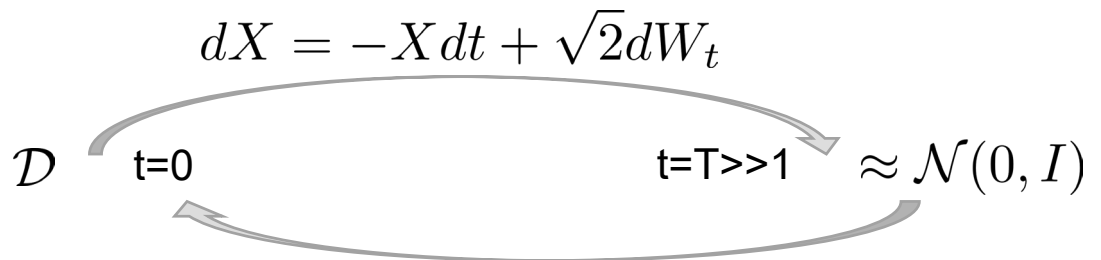
(overdamped) Langevin dynamics:

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

How long? Depends on V and thus data distribution

~1/m                      m-convex                      unimodal

can be exp. longer                      nonconvex                      multimodal



need **T** time to conv.  
 reach  $p(\cdot, 0)$

**agnostic** to multimodality!

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

$$Z(\infty) \sim p(\cdot, 0)$$

need **long** time to conv.



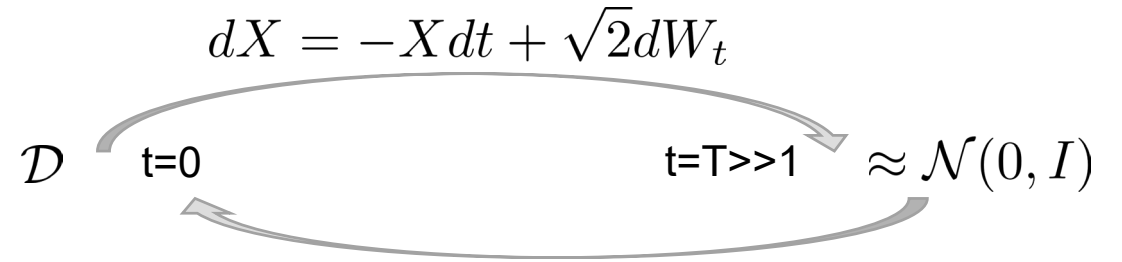
# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
 only numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

need  $T$  time to conv.  
 reach  $p(\cdot, 0)$

agnostic to multimodality!

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$Z(\infty) \sim p(\cdot, 0)$$

suffer from multimodality etc.

# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample Y

already highly nontrivial

sources of error

assumptions on data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $\xrightarrow{\hspace{10em}}$   $t=T \gg 1 \approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

need  $T$  time to conv.  
reach  $p(\cdot, 0)$

**agnostic** to multimodality!

score

$$s(x, t) := \nabla_x \log p(x, t)$$

$$X(t) \sim p(x, t)$$

difference?

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$Z(\infty) \sim p(\cdot, 0)$$

suffer from multimodality etc.

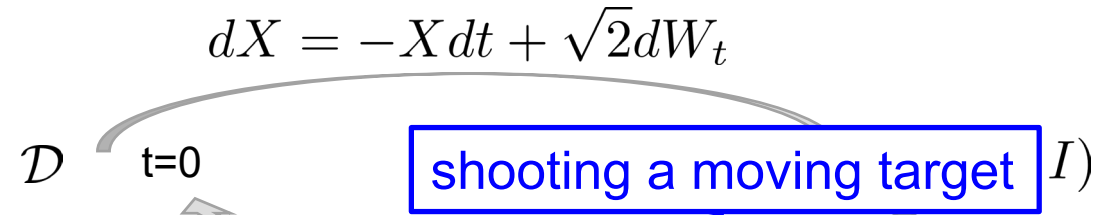
# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. backward sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution



$$dY = Y dt + 2s(Y, T-t)dt + \sqrt{2}dB_t$$

score

need  $T$  time to conv.  
reach  $p(\cdot, 0)$

$$s(x, t) := \nabla_x \log p(x, t)$$

agnostic to multimodality!

$$X(t) \sim p(x, t)$$

difference?

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$Z(\infty) \sim p(\cdot, 0)$$

suffer from multimodality etc.

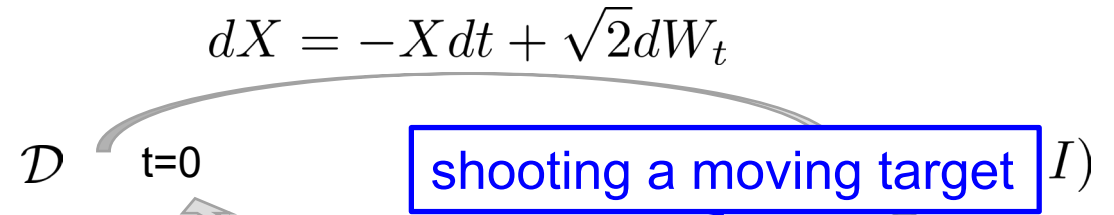
# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. backward sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution



$$dY = Y dt + 2s(Y, T-t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

need  $T$  time to conv.  
reach  $p(\cdot, 0)$

agnostic to multimodality!

$$X(t) \sim p(x, t)$$

difference?

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$Z(\infty) \sim p(\cdot, 0)$$

suffer from multimodality etc.

a specific annealing scheme made multimodal sampling effectively

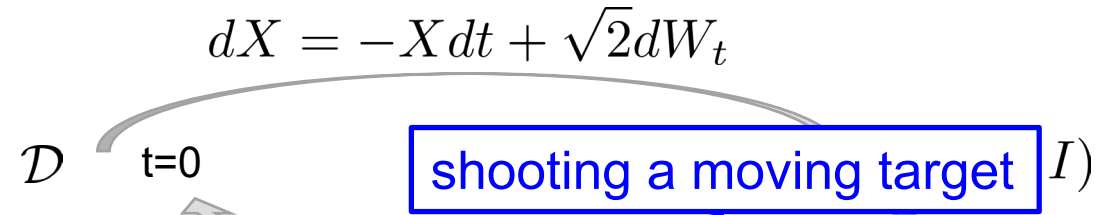
# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. backward sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution



$$dY = Y dt + 2s(Y, T-t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

need  $T$  time to conv.  
reach  $p(\cdot, 0)$

agnostic to multimodality!

$$X(t) \sim p(x, t)$$

difference?

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$Z(\infty) \sim p(\cdot, 0)$$

suffer from multimodality etc.

a specific annealing scheme made multimodal sampling effectively

Lee, Risteski, Ge 18

Chehab, Hyvarinen, Risteski 23

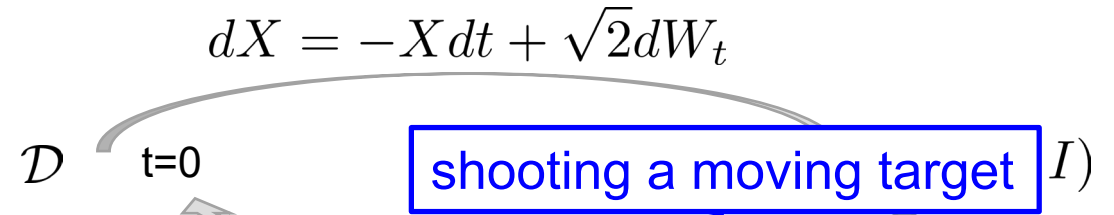
# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. backward sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

sources of error

assumptions on data distribution



$$dY = Y dt + 2s(Y, T-t)dt + \sqrt{2}dB_t$$

score

$$s(x, t) := \nabla_x \log p(x, t)$$

need  $T$  time to conv.  
reach  $p(\cdot, 0)$

agnostic to multimodality!

$$X(t) \sim p(x, t)$$

difference?

$$dZ_t = -\nabla V(Z_t)dt + \sqrt{2}dB_t \stackrel{V := -\log p(\cdot, 0)}{=} s(Z_t, 0)dt + \sqrt{2}dB_t$$

$$Z(\infty) \sim p(\cdot, 0)$$

suffer from multimodality etc.

a specific annealing scheme made multimodal sampling effectively

Lee, Risteski, Ge 18

Chehab, Hyvarinen, Risteski 23

Guo, Tao, Chen 24

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error

assumptions on data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error


assumptions on data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

ideally: 



## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error    score, integration, initialization...


assumptions on data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

ideally: 

practically: **nontrivial**

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$   
**only**

already highly nontrivial

sources of error    score, integration, initialization...


assumptions on data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

ideally: 

practically: **nontrivial**

errors propagate thru  $Y$  dynamics  
in controlled fashion?

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$   
**only**

already highly nontrivial

sources of error    score, integration, initialization...

assumptions on data distribution

- strongly log concave: contraction, relatively easy

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

ideally: 

practically: **nontrivial**

errors propagate thru  $Y$  dynamics  
in controlled fashion?

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

only

already highly nontrivial

sources of error    score, integration, initialization...

assumptions on data distribution


- strongly log concave: contraction, relatively easy
- isoperimetric ineq. (LSI, PI, ...): ~2-3 years ago

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

ideally: 

practically: **nontrivial**

errors propagate thru  $Y$  dynamics  
in controlled fashion?

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$   
**only**

already highly nontrivial

sources of error score, integration, initialization...

assumptions on data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

ideally: 

practically: **nontrivial**

errors propagate thru  $Y$  dynamics  
in controlled fashion?

- strongly log concave: contraction, relatively easy
- isoperimetric ineq. (LSI, PI, ...): ~2-3 years ago
- bounded 2<sup>nd</sup> moment + Lipschitz score: ~1-2 years ago

## 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$   
**only**

already highly nontrivial

sources of error score, integration, initialization...

assumptions on data distribution

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

ideally: 

practically: **nontrivial**

errors propagate thru  $Y$  dynamics  
in controlled fashion?

- strongly log concave: contraction, relatively easy
- isoperimetric ineq. (LSI, PI, ...): ~2-3 years ago
- bounded 2<sup>nd</sup> moment + Lipschitz score: ~1-2 years ago
- bounded 2<sup>nd</sup> moment: ~0-1 years ago



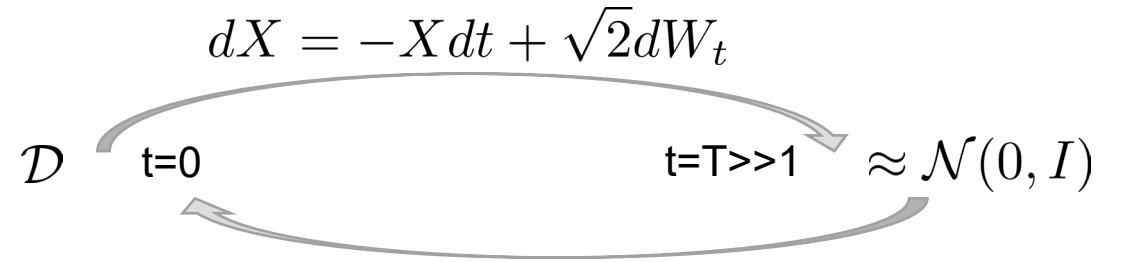
# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. **backward** sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

sources of error score, integration, initialization...

assumptions on data distribution



$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$

denoising diffusion annealing: **agnostic** to multimodality

ideally: ✓

practically: **nontrivial**

errors propagate thru Y dynamics in controlled fashion?

- strongly log concave: contraction, relatively easy
- isoperimetric ineq. (LSI, PI, ...): ~2-3 years ago
- bounded 2<sup>nd</sup> moment + Lipschitz score: ~1-2 years ago
- bounded 2<sup>nd</sup> moment: ~0-1 years ago [Benton+ 23] [Conforti+ 23]

time schedule

exponential  
constant step



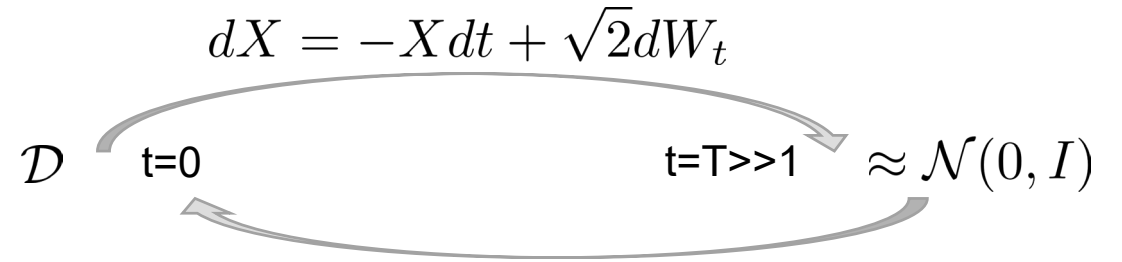
# 1.2 Quantification of Diffusion Model's Generation Quality: Existing Results - Sampling Only

**only** 2. **backward** sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

sources of error score, integration, initialization...

assumptions on data distribution



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

denoising diffusion annealing: **agnostic** to multimodality

ideally: ✓

practically: **nontrivial**

errors propagate thru Y dynamics  
in controlled fashion?

- strongly log concave: contraction, relatively easy
- isoperimetric ineq. (LSI, PI, ...): ~2-3 years ago
- bounded 2<sup>nd</sup> moment + Lipschitz score: ~1-2 years ago
- bounded 2<sup>nd</sup> moment: ~0-1 years ago [Benton+ 23]  
[Conforti+ 23]

time schedule

exponential  
constant step

Variance Preserving  
SDE

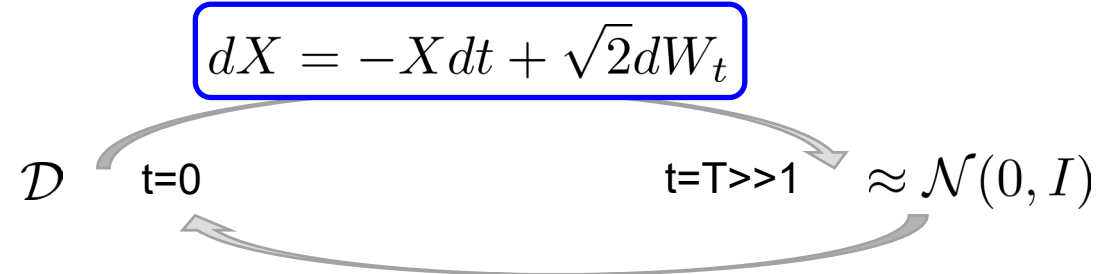
# 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

**only** 2. **backward** sampling/inference process:  
numerical simulation → sample Y

already highly nontrivial

sources of error score, integration, initialization...

assumptions on data distribution



denoising diffusion annealing: **agnostic** to multimodality

ideally: ✓      practically: **nontrivial**

errors propagate thru Y dynamics  
in controlled fashion?

- strongly log concave: contraction, relatively easy
- isoperimetric ineq. (LSI, PI, ...): ~2-3 years ago
- bounded 2<sup>nd</sup> moment + Lipschitz score: ~1-2 years ago
- bounded 2<sup>nd</sup> moment: ~0-1 years ago [Benton+ 23]  
[Conforti+ 23]

time schedule  
exponential  
constant step

Variance Preserving  
SDE

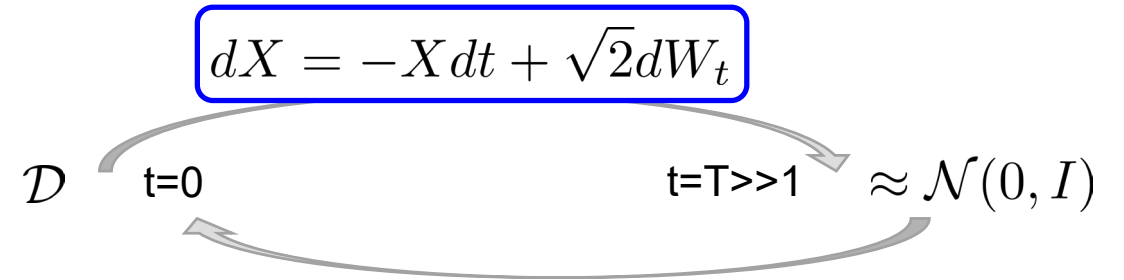
## 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

$$\mathcal{D} \xrightarrow[t=0]{} \boxed{dX = -X dt + \sqrt{2}dW_t} \xrightarrow[t=T \gg 1]{} \approx \mathcal{N}(0, I)$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

Variance Preserving SDE

### 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part



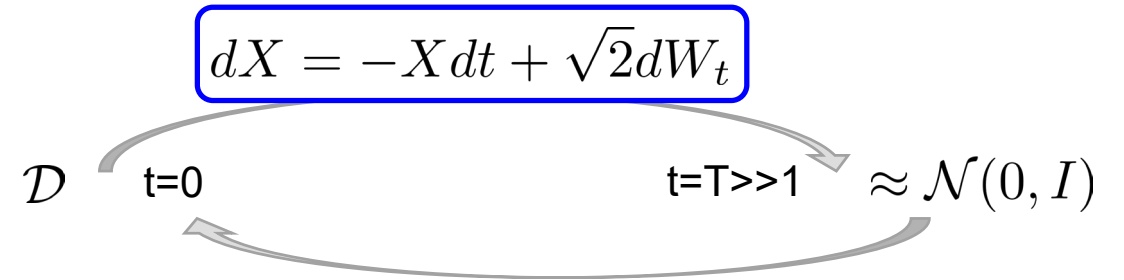
$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

Variance Preserving SDE

more general

$$dX = -f(X, t)dt + g(t)dW_t$$

### 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

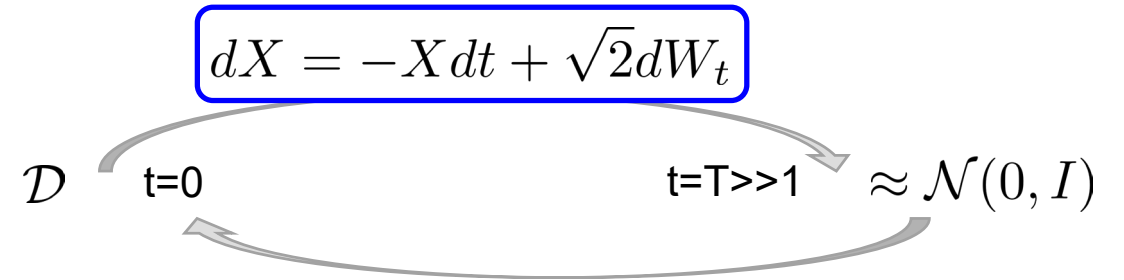
Variance Preserving SDE

more general

$$dX = -f(X, t)dt + g(t)dW_t$$

$$dY = f(Y, T - t)dt + gg^\top s(Y, T - t)dt + g(T - t)dB_t$$

### 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

Variance Preserving SDE

popular version (e.g., EDM [Karras+ 22])

more general

$$dX = -f(X, t)dt + g(t)dW_t$$

$$dY = f(Y, T - t)dt + gg^\top s(Y, T - t)dt + g(T - t)dB_t$$

## 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

$$\mathcal{D} \begin{array}{c} \xrightarrow{t=0} \boxed{dX = -X dt + \sqrt{2}dW_t} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I) \\ \xleftarrow{\hspace{10em}} \end{array}$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

more general

$$dX = -f(X, t)dt + g(t)dW_t$$

$$dY = f(Y, T - t)dt + gg^\top s(Y, T - t)dt + g(T - t)dB_t$$

Variance Preserving SDE

popular version (e.g., EDM [Karras+ 22])

(generalized) Variance Exploding SDE

$$dX = g(t)dW_t$$

$$dY = g^2(T - t)s(Y, T - t)dt + g(T - t)dB_t$$

### 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

minimal data assumption

$$\mathcal{D} \xrightarrow{t=0} \boxed{dX = -X dt + \sqrt{2}dW_t} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

[Benton+ 23]

[Conforti+ 23]

+

Variance Preserving SDE

popular version (e.g., EDM [Karras+ 22])

(generalized) Variance Exploding SDE

$$dX = g(t)dW_t$$

$$dY = g^2(T - t)s(Y, T - t)dt + g(T - t)dB_t$$



### 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

minimal data assumption

$$\mathcal{D} \xrightarrow{t=0} \boxed{dX = -X dt + \sqrt{2}dW_t} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

[Benton+ 23]

[Conforti+ 23]

+

Variance Preserving SDE

new

[Wang+ 24]

+

popular version (e.g., EDM [Karras+ 22])

(generalized) Variance Exploding SDE

$$dX = g(t)dW_t$$

$$dY = g^2(T - t)s(Y, T - t)dt + g(T - t)dB_t$$

### 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

minimal data assumption

$$\mathcal{D} \xrightarrow{t=0} \boxed{dX = -X dt + \sqrt{2}dW_t} \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I)$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

[Benton+ 23]

[Conforti+ 23]

+

Variance Preserving SDE

new

[Wang+ 24]

+

popular version (e.g., EDM [Karras+ 22])

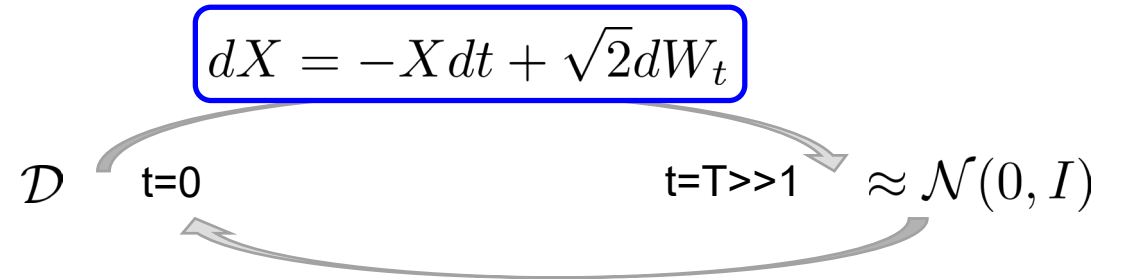
(generalized) Variance Exploding SDE

$$dX = g(t)dW_t$$

$$dY = g^2(T - t)s(Y, T - t)dt + g(T - t)dB_t$$

# 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

minimal data assumption



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

[Benton+ 23]

time schedule

exponential

+

Variance Preserving SDE

[Conforti+ 23]

constant step

new

popular version (e.g., EDM [Karras+ 22])

[Wang+ 24]

arbitrary

+

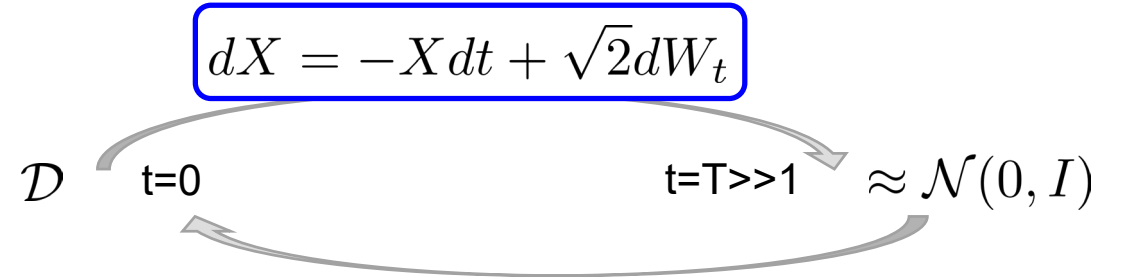
(generalized) Variance Exploding SDE

$$dX = g(t)dW_t$$

$$dY = g^2(T - t)s(Y, T - t)dt + g(T - t)dB_t$$

# 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

minimal data assumption



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

[Benton+ 23]

time schedule

exponential

+

Variance Preserving SDE

[Conforti+ 23]

constant step

new

popular version (e.g., EDM [Karras+ 22])

[Wang+ 24]

arbitrary

+

(generalized) Variance Exploding SDE

motivation:

how to choose  $t_k$  (and  $g$ ?)

$$dX = g(t)dW_t$$

$$dY = g^2(T - t)s(Y, T - t)dt + g(T - t)dB_t$$

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

**new**

[Wang+ 24]

arbitrary

+

popular version (e.g., EDM [Karras+ 22])

**(generalized) Variance Exploding SDE**

$$dX = g(t)dW_t$$

$$dY = g^2(T - t)s(Y, T - t)dt + g(T - t)dB_t$$

### 1.3 Quantification of Diffusion Model's Generation Quality: New Results - Sampling Part

2. **backward** sampling/inference process:  
numerical simulation  $\rightarrow$  sample  $Y$

+

1. **forward** training/learning process:  
optimization  $\rightarrow$  score  $s$

★ bigger contribution

**new**

[Wang+ 24]

arbitrary

+

popular version (e.g., EDM [Karras+ 22])

**(generalized) Variance Exploding SDE**

$$dX = g(t)dW_t$$

$$dY = g^2(T - t)s(Y, T - t)dt + g(T - t)dB_t$$

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

approximating score is essential

$$\mathcal{D} \quad \begin{array}{ccc} & \boxed{dX = -X dt + \sqrt{2}dW_t} & \\ \xrightarrow{t=0} & & \xrightarrow{t=T \gg 1} \approx \mathcal{N}(0, I) \end{array}$$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

approximating score is essential

“main stream” generation quality bound

if  $\|s_\theta - s\| \leq \epsilon$ , then  $\dots$

$$\mathcal{D} \xrightarrow[t=0]{t=T \gg 1} \approx \mathcal{N}(0, I)$$
$$dX = -X dt + \sqrt{2}dW_t$$
$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$



## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

approximating score is essential

“main stream” generation quality bound

if  $\|s_\theta - s\| \leq \epsilon$ , then ...



where does it come from ?

$$dX = -X dt + \sqrt{2}dW_t$$

$\mathcal{D}$   $t=0$   $t=T \gg 1$   $\approx \mathcal{N}(0, I)$

$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

approximating score is essential

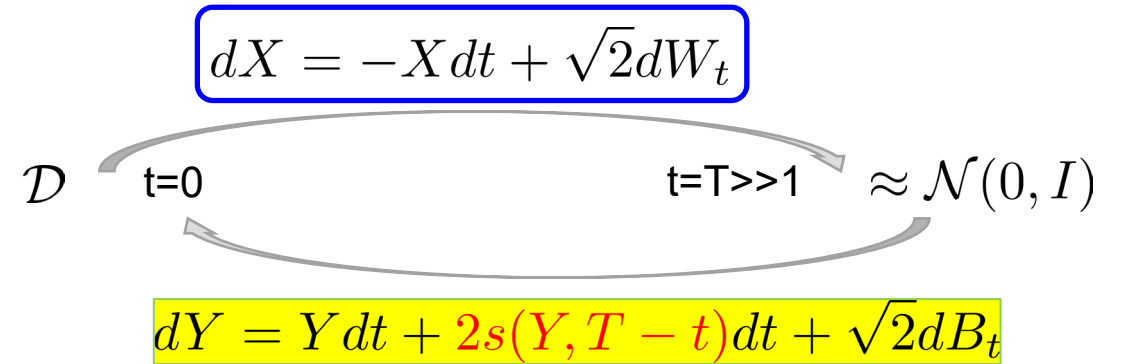
“main stream” generation quality bound

if  $\|s_\theta - s\| \leq \epsilon$ , then ...



where does it come from ?

- nontrivial even when the target density is known (e.g., [Huang+ 24], [He, Rojas, Tao 24], [Gupta+ 24])



approximating score is essential

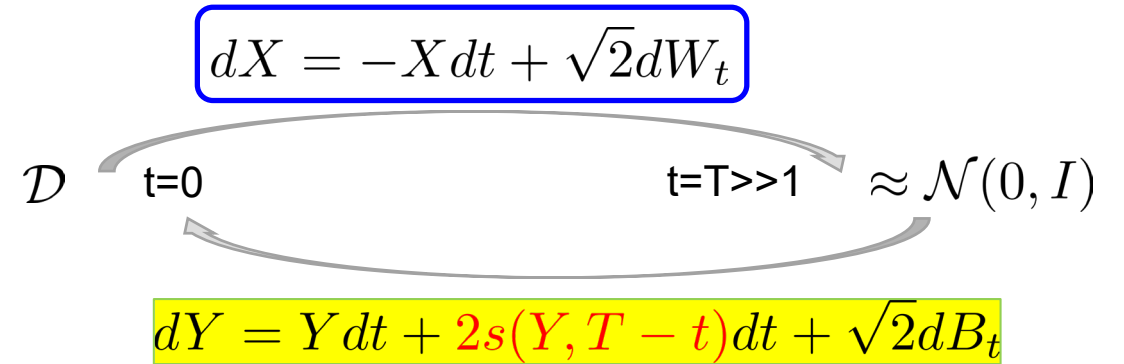
“main stream” generation quality bound

if  $\|s_\theta - s\| \leq \epsilon$ , then ...



where does it come from ?

- nontrivial even when the target density is known (e.g., [Huang+ 24], [He, Rojas, Tao 24], [Gupta+ 24])
- it also matters to innovation / generalization (an open question, to be discussed in the end)



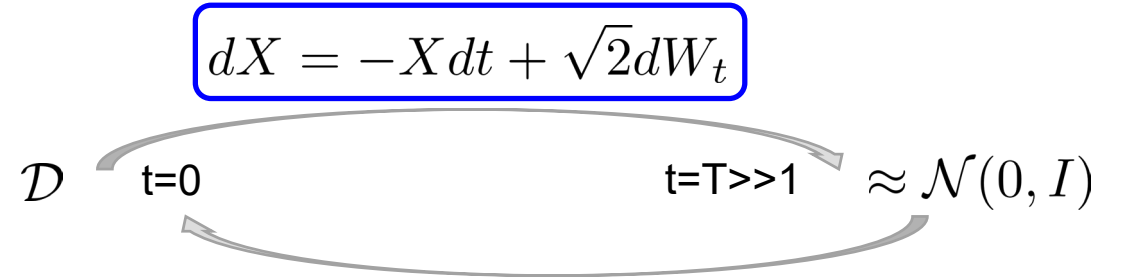
approximating score is essential

“main stream” generation quality bound

if  $\|s_\theta - s\| \leq \epsilon$ , then ...



where does it come from ?



$$dY = Y dt + 2s(Y, T - t)dt + \sqrt{2}dB_t$$

- nontrivial even when the target density is known (e.g., [Huang+ 24], [He, Rojas, Tao 24], [Gupta+ 24])
- it also matters to innovation / generalization (an open question, to be discussed in the end)

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - s_\theta(X_t, t)\|^2 dt$$

approximating score is essential

“main stream” generation quality bound

if  $\|s_\theta - s\| \leq \epsilon$ , then ...



where does it come from ?

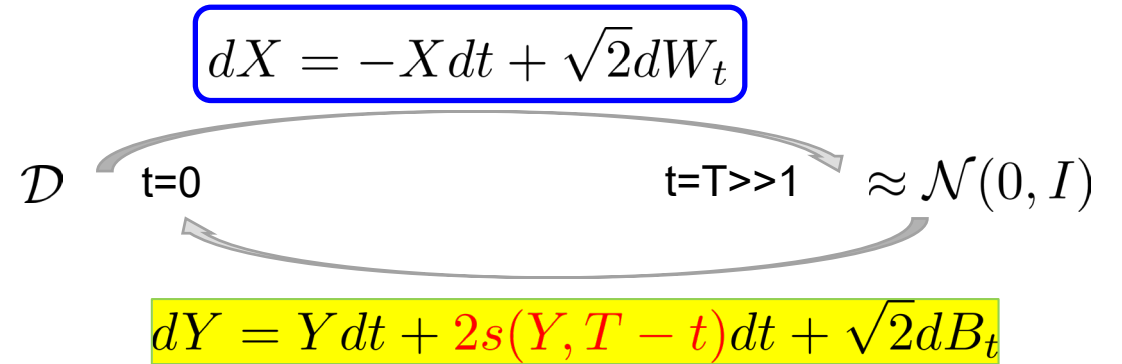
- nontrivial even when the target density is known (e.g., [Huang+ 24], [He, Rojas, Tao 24], [Gupta+ 24])
- it also matters to innovation / generalization (an open question, to be discussed in the end)

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - s_\theta(X_t, t)\|^2 dt$$

practice

time discretization/sampling → empirical approximation → optimization



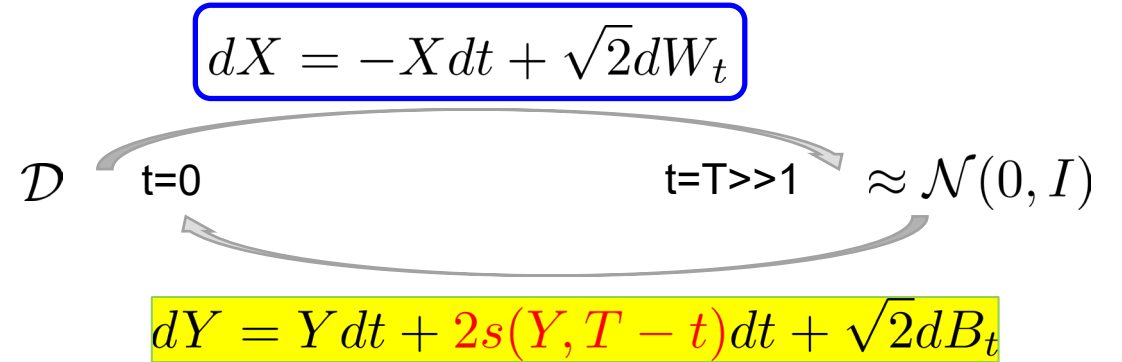
approximating score is essential

“main stream” generation quality bound

if  $\|s_\theta - s\| \leq \epsilon$ , then ...



where does it come from ?



- nontrivial even when the target density is known (e.g., [Huang+ 24], [He, Rojas, Tao 24], [Gupta+ 24])
- it also matters to innovation / generalization (an open question, to be discussed in the end)

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - s_\theta(X_t, t)\|^2 dt$$

$S(\theta; t, X_t)$

practice

time discretization/sampling → empirical approximation → optimization

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

$$\text{forward dynamics} \Rightarrow X_t = e^{-\mu t} X_0 + \bar{\sigma}_t \xi$$



## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

$$\text{forward dynamics} \Rightarrow X_t = e^{-\mu_t} X_0 + \bar{\sigma}_t \xi$$

0 for Variance Exploding SDE

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

forward dynamics  $\Rightarrow X_t = e^{-\mu_t} X_0 + \bar{\sigma}_t \xi$

$\bar{\sigma}_t$   $\rightarrow$  variance schedule

$\mu_t$   $\rightarrow$  0 for Variance Exploding SDE

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

forward dynamics  $\Rightarrow$



$$X_t = e^{-\mu_t} X_0 + \bar{\sigma}_t \xi$$

variance schedule

0 for Variance Exploding SDE

$$\min_{\theta} \frac{1}{2} \int_{t_0}^T w(t) \frac{1}{\bar{\sigma}_t} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_t S(\theta; t, X_0 + \bar{\sigma}_t \xi) + \xi\|^2 dt$$

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

forward dynamics  $\Rightarrow$

$$X_t = e^{-\mu t} X_0 + \bar{\sigma}_t \xi$$

variance schedule

0 for Variance Exploding SDE

$$\min_{\theta} \frac{1}{2} \int_{t_0}^T w(t) \frac{1}{\bar{\sigma}_t} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_t S(\theta; t, X_0 + \bar{\sigma}_t \xi) + \xi\|^2 dt$$

time discretization

$$\min_{\theta} \frac{1}{2} \underbrace{\sum_{j=1}^N w(t_j) (t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_0 + \bar{\sigma}_{t_j} \xi) + \xi\|^2}_{\bar{\mathcal{L}}(\theta)}$$

# 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

forward dynamics  $\Rightarrow$

$$X_t = e^{-\mu t} X_0 + \bar{\sigma}_t \xi$$

→ variance schedule  
→ 0 for Variance Exploding SDE

$$\min_{\theta} \frac{1}{2} \int_{t_0}^T w(t) \frac{1}{\bar{\sigma}_t} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_t S(\theta; t, X_0 + \bar{\sigma}_t \xi) + \xi\|^2 dt$$

time discretization

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^N w(t_j) (t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_0 + \bar{\sigma}_{t_j} \xi) + \xi\|^2$$

$\underbrace{\hspace{15em}}_{\bar{\mathcal{L}}(\theta)}$

→ time schedule

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

forward dynamics  $\Rightarrow$

$$X_t = e^{-\mu t} X_0 + \bar{\sigma}_t \xi$$

→ variance schedule  
→ 0 for Variance Exploding SDE

$$\min_{\theta} \frac{1}{2} \int_{t_0}^T w(t) \frac{1}{\bar{\sigma}_t} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_t S(\theta; t, X_0 + \bar{\sigma}_t \xi) + \xi\|^2 dt$$

time discretization

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^N w(t_j) (t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_0 + \bar{\sigma}_{t_j} \xi) + \xi\|^2$$

$\bar{\mathcal{L}}(\theta)$

→ time schedule

empirical version

$$\min_{\theta} \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

$\bar{\mathcal{L}}_{em}(\theta)$

# 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

forward dynamics  $\Rightarrow$

$$X_t = e^{-\mu t} X_0 + \bar{\sigma}_t \xi$$

variance schedule

0 for Variance Exploding SDE

$$\min_{\theta} \frac{1}{2} \int_{t_0}^T w(t) \frac{1}{\bar{\sigma}_t} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_t S(\theta; t, X_0 + \bar{\sigma}_t \xi) + \xi\|^2 dt$$

time discretization

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^N \underbrace{w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_0 + \bar{\sigma}_{t_j} \xi) + \xi\|^2}_{\bar{\mathcal{L}}(\theta)}$$

time schedule

empirical version

$$\min_{\theta} \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

$\bar{\mathcal{L}}_{em}(\theta)$

total weighting

# 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

ideal

$$\min_{\theta} \frac{1}{2} \int_0^T w(t) \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla_x \log p_{t|0}(X_t|X_0, t) - S(\theta; t, X_t)\|^2 dt$$

forward dynamics  $\Rightarrow$

$$X_t = e^{-\mu t} X_0 + \bar{\sigma}_t \xi$$

→ variance schedule  
→ 0 for Variance Exploding SDE

$$\min_{\theta} \frac{1}{2} \int_{t_0}^T w(t) \frac{1}{\bar{\sigma}_t} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_t S(\theta; t, X_0 + \bar{\sigma}_t \xi) + \xi\|^2 dt$$

time discretization

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^N \underbrace{w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_0 + \bar{\sigma}_{t_j} \xi) + \xi\|^2}_{\bar{\mathcal{L}}(\theta)}$$

→ time schedule

empirical version

$$\min_{\theta} \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \underbrace{\beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2}_{\bar{\mathcal{L}}_{em}(\theta)}$$

→ total weighting

training

$$\theta^{(k+1)} = \theta^{(k)} - h \nabla \bar{\mathcal{L}}_{em}(\theta^{(k)})$$



## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – **Training + Sampling**

GD training

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP

GD training

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

total weighting

time schedule

variance schedule

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

total weighting

time schedule

variance schedule

roadmap

non-asymptotic bound of GD  
optimization of  $\bar{\mathcal{L}}_{em}$

## 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

total weighting

time schedule

variance schedule

roadmap

$C > 0$

non-asymptotic bound of GD  
optimization of  $\bar{\mathcal{L}}_{em} \gtrsim C$

# 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

total weighting

time schedule

variance schedule

roadmap

$C > 0 \rightarrow$   
interesting  
generalization  
setting

non-asymptotic bound of GD  
optimization of  $\bar{\mathcal{L}}_{em} \gtrsim C$

# 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

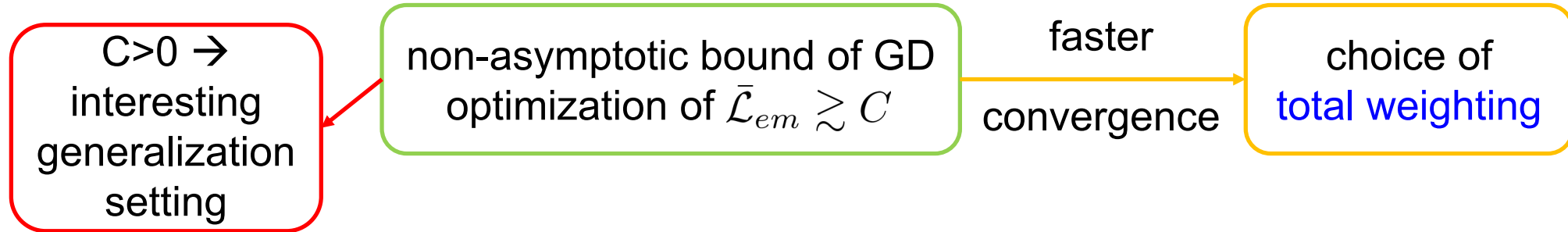
$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

total weighting

time schedule

variance schedule

roadmap





# 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

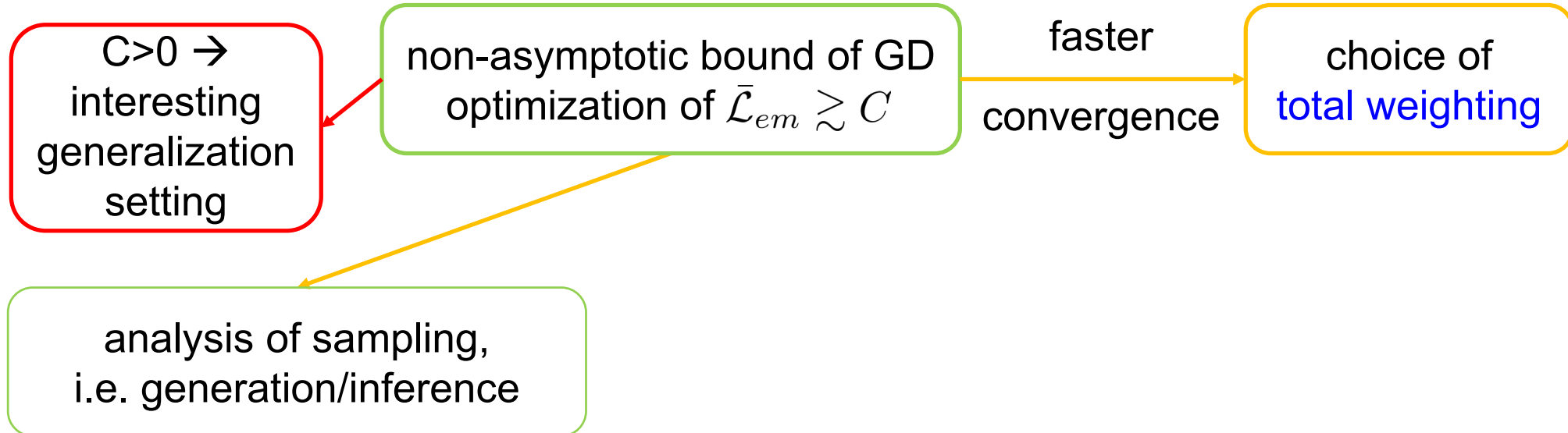
$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

total weighting

time schedule

variance schedule

roadmap



# 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

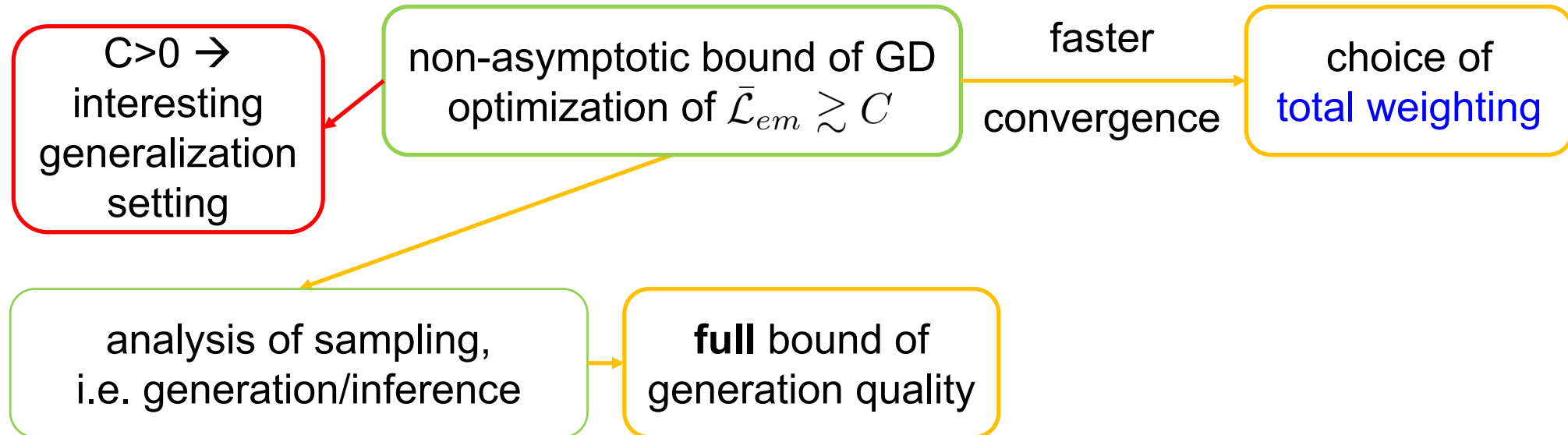
$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

total weighting

time schedule

variance schedule

roadmap



# 1.4 Quantification of Diffusion Model's Generation Quality: New Results – Training + Sampling

score parameterization

$S(\theta; \cdot)$  wide (& deep) ReLU MLP (challenge: U-Net? DiT?)

GD training

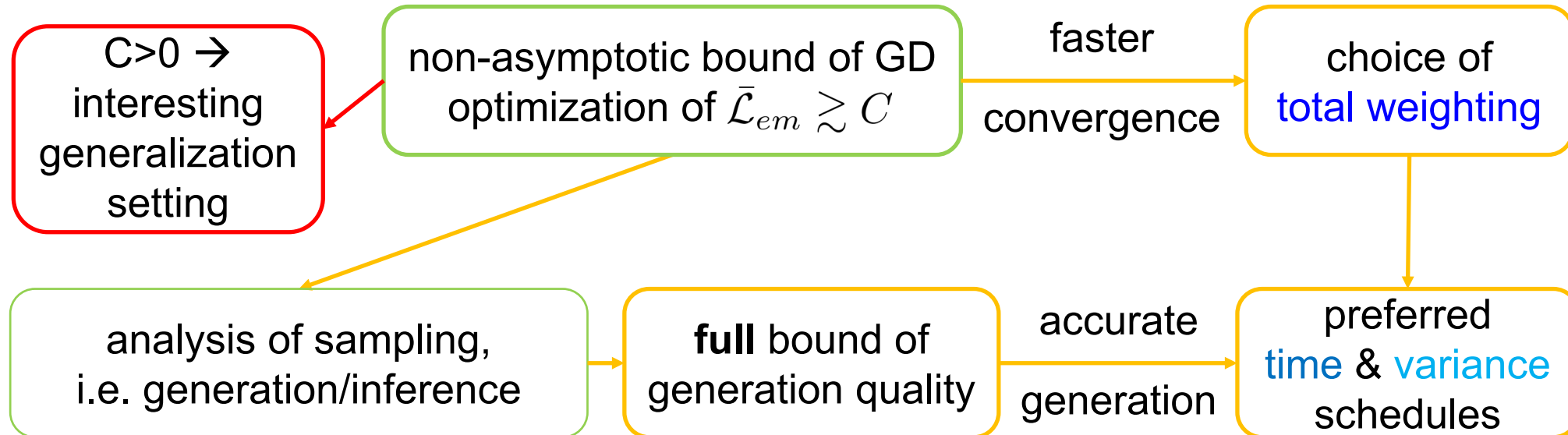
$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, x_i + \bar{\sigma}_{t_j} \xi_{ij}) + \xi_{ij}\|^2$$

total weighting

time schedule

variance schedule

roadmap



# Setup

- ▶ Objective function (minimized by GD):

# Setup

- ▶ Objective function (minimized by GD):

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \underbrace{\frac{w(t_j)}{\bar{\sigma}_{t_j}} (t_j - t_{j-1})}_{\beta_j} \underbrace{\|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2}_{f(\theta; i, j)}$$

# Setup

- ▶ Objective function (minimized by GD):

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \underbrace{\frac{w(t_j)}{\bar{\sigma}_{t_j}} (t_j - t_{j-1})}_{\beta_j} \underbrace{\|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2}_{f(\theta; i, j)}$$

- ▶ Architecture: deep ReLU network

## Setup

- ▶ Objective function (minimized by GD):

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \underbrace{\frac{w(t_j)}{\bar{\sigma}_{t_j}} (t_j - t_{j-1})}_{\beta_j} \underbrace{\|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2}_{f(\theta; i, j)}$$

- ▶ Architecture: deep ReLU network

$$S(\theta; X_{ij}) = W_{L+1} \sigma(W_L \cdots W_1 \sigma(W_0 [X_{ij}, t_j])),$$

# Setup

- ▶ Objective function (minimized by GD):

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \underbrace{\frac{w(t_j)}{\bar{\sigma}_{t_j}} (t_j - t_{j-1})}_{\beta_j} \underbrace{\|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2}_{f(\theta; i, j)}$$

- ▶ Architecture: deep ReLU network

$$S(\theta; X_{ij}) = W_{L+1} \sigma(W_L \cdots W_1 \sigma(W_0 [X_{ij}, t_j])),$$

- ▶  $W_{L+1} \in \mathbb{R}^{d \times m}$ ,  $W_\ell \in \mathbb{R}^{m \times m}$ ,  $W_0 \in \mathbb{R}^{m \times d}$



# Setup

- ▶ Objective function (minimized by GD):

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^N \underbrace{\frac{w(t_j)}{\bar{\sigma}_{t_j}} (t_j - t_{j-1})}_{\beta_j} \underbrace{\|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2}_{f(\theta; i, j)}$$

- ▶ Architecture: deep ReLU network

$$S(\theta; X_{ij}) = W_{L+1} \sigma(W_L \cdots W_1 \sigma(W_0 [X_{ij}, t_j])),$$

- ▶  $W_{L+1} \in \mathbb{R}^{d \times m}$ ,  $W_\ell \in \mathbb{R}^{m \times m}$ ,  $W_0 \in \mathbb{R}^{m \times d}$
- ▶  $\theta := (W_0, W_1, \dots, W_L, W_{L+1})$

# Setup

- ▶ Input data  $(t_j, X_{ij})$ :

# Setup

- ▶ Input data  $(t_j, X_{ij})$ :

$$X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij} \sim P_{t_j}(X|X(0) = x_i),$$

# Setup

- ▶ Input data  $(t_j, X_{ij})$ :

$$X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij} \sim P_{t_j}(X|X(0) = x_i),$$

where  $\bar{\sigma}_{t_j}^2$  is the variance of  $X_{t_j}|X_0$ , and  $\xi_{ij} \sim \mathcal{N}(0, I)$

# Setup

- ▶ Input data  $(t_j, X_{ij})$ :

$$X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij} \sim P_{t_j}(X|X(0) = x_i),$$

where  $\bar{\sigma}_{t_j}^2$  is the variance of  $X_{t_j}|X_0$ , and  $\xi_{ij} \sim \mathcal{N}(0, 1)$

- ▶  $\bar{\sigma}_t$ : monotonically increasing functions of  $t$ ;  $\bar{\sigma}_0 = 0$

# Setup

- ▶ Input data  $(t_j, X_{ij})$ :

$$X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij} \sim P_{t_j}(X|X(0) = x_i),$$

where  $\bar{\sigma}_{t_j}^2$  is the variance of  $X_{t_j}|X_0$ , and  $\xi_{ij} \sim \mathcal{N}(0, 1)$

- ▶  $\bar{\sigma}_t$ : monotonically increasing functions of  $t$ ;  $\bar{\sigma}_0 = 0$
- ▶ Output data:

$$\frac{-\xi_{ij}}{\bar{\sigma}_{t_j}}$$

# Setup

- ▶ Input data  $(t_j, X_{ij})$ :

$$X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij} \sim P_{t_j}(X|X(0) = x_i),$$

where  $\bar{\sigma}_{t_j}^2$  is the variance of  $X_{t_j}|X_0$ , and  $\xi_{ij} \sim \mathcal{N}(0, 1)$

- ▶  $\bar{\sigma}_t$ : monotonically increasing functions of  $t$ ;  $\bar{\sigma}_0 = 0$
- ▶ Output data:

$$\frac{-\xi_{ij}}{\bar{\sigma}_{t_j}}$$

- ▶ Very large if  $\bar{\sigma}_{t_j} \approx 0$

# Setup

- ▶ Number of samples:

$$X_{ij} = \underbrace{x_i}_{\substack{n \text{ samples} \\ \text{from initial} \\ \text{distribution } P_0}} + \underbrace{\bar{\sigma}_{t_j}}_{\substack{N \text{ time points} \\ \text{for forward/backward SDEs}} \xi_{ij}}$$



## Theory: assumptions

- ▶ Assumptions: mild + preserve the nature of diffusion models

## Theory: assumptions

- ▶ Assumptions: mild + preserve the nature of diffusion models
- ▶ For example,

## Theory: assumptions

- ▶ Assumptions: mild + preserve the nature of diffusion models
- ▶ For example,

### Assumption

- ▶ *Data scaling:  $\|x_i\| = \Theta(\sqrt{d})$  for all  $i$ .*

## Theory: assumptions

- ▶ Assumptions: mild + preserve the nature of diffusion models
- ▶ For example,

### Assumption

- ▶ *Data scaling:*  $\|x_i\| = \Theta(\sqrt{d})$  for all  $i$ .

Interpretation:

## Theory: assumptions

- ▶ Assumptions: mild + preserve the nature of diffusion models
- ▶ For example,

### Assumption

- ▶ *Data scaling:*  $\|x_i\| = \Theta(\sqrt{d})$  for all  $i$ .

Interpretation:

- ▶ Recall input data:

## Theory: assumptions

- ▶ Assumptions: mild + preserve the nature of diffusion models
- ▶ For example,

### Assumption

- ▶ *Data scaling*:  $\|x_i\| = \Theta(\sqrt{d})$  for all  $i$ .

Interpretation:

- ▶ Recall input data:

$$X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$$

# Theory: assumptions

- ▶ Assumptions: mild + preserve the nature of diffusion models
- ▶ For example,

## Assumption

- ▶ *Data scaling:*  $\|x_i\| = \Theta(\sqrt{d})$  for all  $i$ .

Interpretation:

- ▶ Recall input data:

$$X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$$

- ▶  $\xi_{ij} \sim \mathcal{N}(0, I) \Rightarrow \|\xi_{ij}\| \approx \sqrt{d}$

# Theory: assumptions

- ▶ Assumptions: mild + preserve the nature of diffusion models
- ▶ For example,

## Assumption

- ▶ *Data scaling*:  $\|x_i\| = \Theta(\sqrt{d})$  for all  $i$ .

Interpretation:

- ▶ Recall input data:

$$X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$$

- ▶  $\xi_{ij} \sim \mathcal{N}(0, I) \Rightarrow \|\xi_{ij}\| \approx \sqrt{d}$





# Theory: training

## Theorem

For any  $\epsilon_{\text{train}} > 0$ , consider  $m \geq M(\epsilon_{\text{train}})$ . With high probability,

$$\begin{aligned} & \bar{\mathcal{L}}_{em}(\theta^{(k)}) \\ & \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)})(t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)}) \end{aligned}$$

Moreover, when  $K = \Theta\left(d^{\frac{1-a_0}{2}} n^2 N \log\left(\frac{d}{\epsilon_{\text{train}}}\right)\right)$ ,

$$\bar{\mathcal{L}}_{em}(\theta^{(K)}) \leq \epsilon_{\text{train}}.$$

# Theory: training

## Theorem

For any  $\epsilon_{\text{train}} > 0$ , consider  $m \geq M(\epsilon_{\text{train}})$ . With high probability,

$$\begin{aligned} & \bar{\mathcal{L}}_{em}(\theta^{(k)}) \\ & \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)})(t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)}) \end{aligned}$$

Moreover, when  $K = \Theta\left(d^{\frac{1-a_0}{2}} n^2 N \log\left(\frac{d}{\epsilon_{\text{train}}}\right)\right)$ ,

$$\bar{\mathcal{L}}_{em}(\theta^{(K)}) \leq \epsilon_{\text{train}}.$$

- ▶ Exponential decay

# Theory: training

## Theorem

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)})(t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)})$$

# Theory: training

## Theorem

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)}) (t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)})$$

► Recall:

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1, j=1}^{n, N} f(\theta; i, j)$$

# Theory: training

## Theorem

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)}) (t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)})$$

► Recall:

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1, j=1}^{n, N} f(\theta; i, j)$$

►  $(i^*(s), j^*(s))$  = the index of the largest loss  $f(\theta^{(s)}; i, j)$

# Theory: training

## Theorem

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)}) (t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)})$$

► Recall:

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1, j=1}^{n, N} f(\theta; i, j)$$

- $(i^*(s), j^*(s))$  = the index of the largest loss  $f(\theta^{(s)}; i, j)$
- Faster convergence:

# Theory: training

## Theorem

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)}) (t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)})$$

► Recall:

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1, j=1}^{n, N} f(\theta; i, j)$$

- $(i^*(s), j^*(s))$  = the index of the largest loss  $f(\theta^{(s)}; i, j)$
- Faster convergence:

want to maximize over all the indices

# Theory: training

## Theorem

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)}) (t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)})$$

► Recall:

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2n} \sum_{i=1, j=1}^{n, N} f(\theta; i, j)$$

- $(i^*(s), j^*(s))$  = the index of the largest loss  $f(\theta^{(s)}; i, j)$
- Faster convergence:

want to maximize over all the indices  
⇒ want all  $f(\theta^{(s)}; i, j)$  to be the largest loss



# Total weighting: theory vs practice

## Corollary

When  $f(\theta^{(k)}; i, j) \approx f(\theta^{(k)}; i', j')$  for all  $(i, j), (i', j'), k$ , GD obtains the optimal rate of convergence

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \left( 1 - C_7 h \max_{j=1, \dots, N} w(t_j) (t_j - t_{j-1}) \bar{\sigma}_{t_j} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right)^k \bar{\mathcal{L}}_{em}(\theta^{(0)}).$$

# Total weighting: theory vs practice

## Corollary

When  $f(\theta^{(k)}; i, j) \approx f(\theta^{(k)}; i', j')$  for all  $(i, j), (i', j'), k$ , GD obtains the optimal rate of convergence

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \left( 1 - C_7 h \max_{j=1, \dots, N} w(t_j) (t_j - t_{j-1}) \bar{\sigma}_{t_j} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right)^k \bar{\mathcal{L}}_{em}(\theta^{(0)}).$$

►  $f(\theta^{(k)}; i, j) \approx f(\theta^{(k)}; i', j')$ :

# Total weighting: theory vs practice

## Corollary

When  $f(\theta^{(k)}; i, j) \approx f(\theta^{(k)}; i', j')$  for all  $(i, j), (i', j'), k$ , GD obtains the optimal rate of convergence

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \left( 1 - C_7 h \max_{j=1, \dots, N} w(t_j) (t_j - t_{j-1}) \bar{\sigma}_{t_j} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right)^k \bar{\mathcal{L}}_{em}(\theta^{(0)}).$$

►  $f(\theta^{(k)}; i, j) \approx f(\theta^{(k)}; i', j')$ :

$$f(\theta; i, j) = \underbrace{\beta_j}_{\text{total weighting}} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

# Total weighting: theory vs practice

## Corollary

When  $f(\theta^{(k)}; i, j) \approx f(\theta^{(k)}; i', j')$  for all  $(i, j), (i', j'), k$ , GD obtains the optimal rate of convergence

$$\bar{\mathcal{L}}_{em}(\theta^{(k)}) \leq \left( 1 - C_7 h \max_{j=1, \dots, N} w(t_j) (t_j - t_{j-1}) \bar{\sigma}_{t_j} \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right)^k \bar{\mathcal{L}}_{em}(\theta^{(0)}).$$

►  $f(\theta^{(k)}; i, j) \approx f(\theta^{(k)}; i', j')$ :

$$f(\theta; i, j) = \underbrace{\beta_j}_{\text{total weighting}} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

► Claim: this implies how to choose the total weighting  $\beta_j$

## Total weighting: theory vs practice

$$f(\theta; i, j) = \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

## Total weighting: theory vs practice

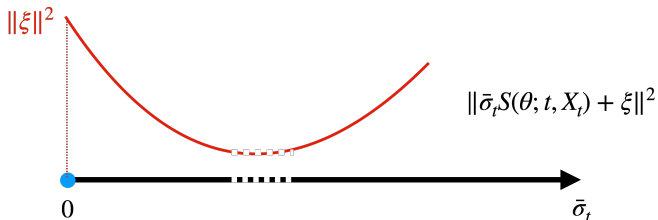
$$f(\theta; i, j) = \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

▶  $\|\bar{\sigma}_t S(\theta; t, x_0 + \bar{\sigma}_t \xi) + \xi\|^2$ :

## Total weighting: theory vs practice

$$f(\theta; i, j) = \beta_j \|\bar{\sigma}_t S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

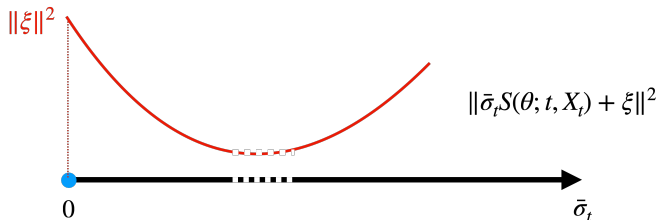
- ▶  $\|\bar{\sigma}_t S(\theta; t, x_0 + \bar{\sigma}_t \xi) + \xi\|^2$ :



## Total weighting: theory vs practice

$$f(\theta; i, j) = \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

- ▶  $\|\bar{\sigma}_t S(\theta; t, x_0 + \bar{\sigma}_t \xi) + \xi\|^2$ :



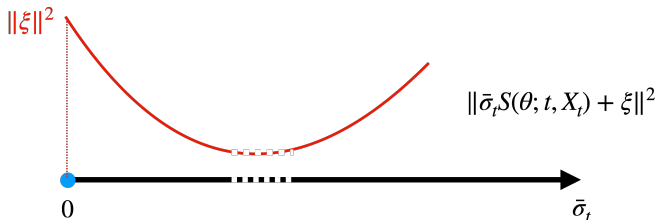
- ▶ Thus,  $\beta_j \propto \frac{1}{\|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2}$



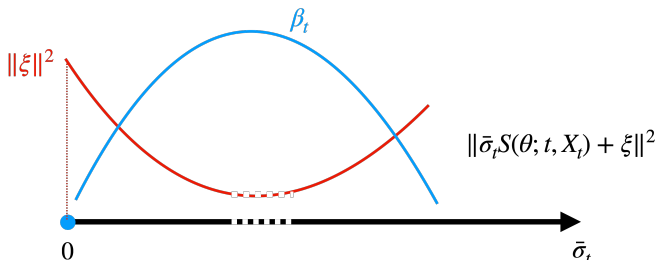
## Total weighting: theory vs practice

$$f(\theta; i, j) = \beta_j \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

- ▶  $\|\bar{\sigma}_t S(\theta; t, x_0 + \bar{\sigma}_t \xi) + \xi\|^2$ :



- ▶ Thus,  $\beta_j \propto \frac{1}{\|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2}$



# Total weighting: theory vs practice

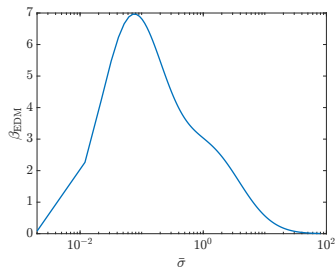
In practice:

- ▶ EDM [Karras et al., 2022]

# Total weighting: theory vs practice

In practice:

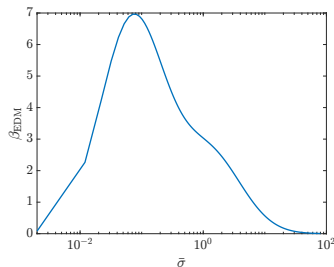
- ▶ EDM [Karras et al., 2022]



# Total weighting: theory vs practice

In practice:

- ▶ EDM [Karras et al., 2022]

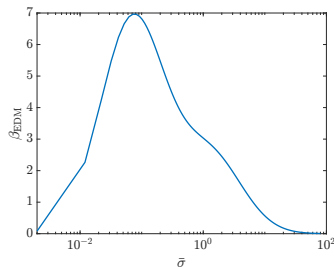


- ▶ Other total weighting functions used in practice (mostly monotone): e.g.,  $\beta_{\bar{\sigma}} = \frac{1}{\bar{\sigma}}$  [Song et al., 2021]

# Total weighting: theory vs practice

In practice:

- ▶ EDM [Karras et al., 2022]

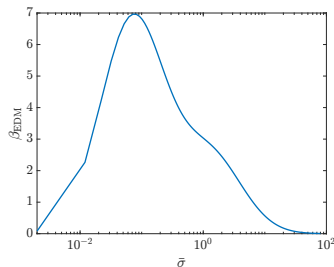


- ▶ Other total weighting functions used in practice (mostly monotone): e.g.,  $\beta_{\bar{\sigma}} = \frac{1}{\bar{\sigma}}$  [Song et al., 2021]
- ▶ Performance: EDM  $>$  other models

# Total weighting: theory vs practice

In practice:

- ▶ EDM [Karras et al., 2022]

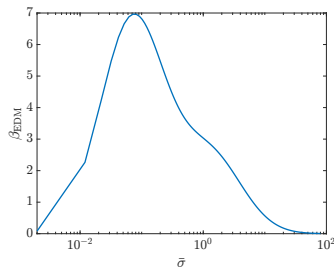


- ▶ Other total weighting functions used in practice (mostly monotone): e.g.,  $\beta_{\bar{\sigma}} = \frac{1}{\bar{\sigma}}$  [Song et al., 2021]
- ▶ Performance: EDM  $>$  other models  
⇒ “bell-shape” is preferable

# Total weighting: theory vs practice

In practice:

- ▶ EDM [Karras et al., 2022]



- ▶ Other total weighting functions used in practice (mostly monotone): e.g.,  $\beta_{\bar{\sigma}} = \frac{1}{\bar{\sigma}}$  [Song et al., 2021]
- ▶ Performance: EDM  $>$  other models  
 $\Rightarrow$  “bell-shape” is preferable
- ▶ Roughly, **Theory**  $\approx$  **Practice**

# Proof of convergence

- ▶ Recall: input data  $X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$ ,  
output data  $-\frac{\xi_{ij}}{\bar{\sigma}_{t_j}}$ ,  
where  $x_i \sim P_0$ ,  $\xi_{ij} \sim \mathcal{N}(0, I)$



# Proof of convergence

- ▶ Recall: input data  $X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$ ,  
output data  $-\frac{\xi_{ij}}{\bar{\sigma}_{t_j}}$ ,  
where  $x_i \sim P_0$ ,  $\xi_{ij} \sim \mathcal{N}(0, I)$
- ▶ Framework [Allen-Zhu et al., 2019]:

# Proof of convergence

- ▶ Recall: input data  $X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$ ,

output data  $-\frac{\xi_{ij}}{\bar{\sigma}_{t_j}}$ ,

where  $x_i \sim P_0$ ,  $\xi_{ij} \sim \mathcal{N}(0, I)$

- ▶ Framework [Allen-Zhu et al., 2019]:

semi-smoothness + local strongly convex  
key

# Proof of convergence

- ▶ Recall: input data  $X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$ ,

output data  $-\frac{\xi_{ij}}{\bar{\sigma}_{t_j}}$ ,

where  $x_i \sim P_0$ ,  $\xi_{ij} \sim \mathcal{N}(0, I)$

- ▶ Framework [Allen-Zhu et al., 2019]:

semi-smoothness + local strongly convex  
key

- ▶ No longer works in denoising diffusion models

# Proof of convergence

- ▶ Recall: input data  $X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$ ,  
output data  $-\frac{\xi_{ij}}{\bar{\sigma}_{t_j}}$ ,  
where  $x_i \sim P_0$ ,  $\xi_{ij} \sim \mathcal{N}(0, I)$
- ▶ Framework [Allen-Zhu et al., 2019]:  
semi-smoothness + local strongly convex  
key
- ▶ No longer works in denoising diffusion models
- ▶ Reason: (1) scaling ;  
bad  
~~small output data~~

# Proof of convergence

- ▶ Recall: input data  $X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$ ,

$$\text{output data } -\frac{\xi_{ij}}{\bar{\sigma}_{t_j}},$$

where  $x_i \sim P_0$ ,  $\xi_{ij} \sim \mathcal{N}(0, I)$

- ▶ Framework [Allen-Zhu et al., 2019]:

semi-smoothness + local strongly convex  
key

- ▶ No longer works in denoising diffusion models

- ▶ Reason: (1) scaling ; (2) correlation  
~~bad~~                      ~~“good”~~  
~~small output data~~      ~~data separability~~

# Proof of convergence

- ▶ Recall: input data  $X_{ij} = x_i + \bar{\sigma}_{t_j} \xi_{ij}$ ,

output data  $-\frac{\xi_{ij}}{\bar{\sigma}_{t_j}}$ ,

where  $x_i \sim P_0$ ,  $\xi_{ij} \sim \mathcal{N}(0, I)$

- ▶ Framework [Allen-Zhu et al., 2019]:

semi-smoothness + local strongly convex  
key

- ▶ No longer works in denoising diffusion models

- ▶ Reason: (1) scaling ; (2) correlation  
bad "good"  
~~small output data~~ ~~data separability~~

- ▶ Our proof: high probability bound using some high-dimensional geometry facts

# Theory: more about training

Recall:

## Theorem

For any  $\epsilon_{\text{train}} > 0$ , consider  $m \geq M(\epsilon_{\text{train}})$ . With high probability,

$$\begin{aligned} & \bar{\mathcal{L}}_{em}(\theta^{(k)}) \\ & \leq \prod_{s=0}^{k-1} \left( 1 - C_5 h w(t_{j^*(s)}) (t_{j^*(s)} - t_{j^*(s)-1}) \bar{\sigma}_{t_{j^*(s)}} \frac{m d^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \bar{\mathcal{L}}_{em}(\theta^{(0)}) \end{aligned}$$

Moreover, when  $K = \Theta\left(d^{\frac{1-a_0}{2}} n^2 N \log\left(\frac{d}{\epsilon_{\text{train}}}\right)\right)$ ,

$$\bar{\mathcal{L}}_{em}(\theta^{(K)}) \leq \epsilon_{\text{train}}.$$

# Theory: more about training

Properties of denoising score matching objective:



## Theory: more about training

Properties of denoising score matching objective:

- ▶  $\bar{\mathcal{L}}_{em}(\theta) \rightarrow \bar{\mathcal{L}}$  as  $n \rightarrow \infty$

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2} \sum_{j=1}^N w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \frac{1}{n} \sum_{i=1}^n \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

$$\bar{\mathcal{L}}(\theta) = \frac{1}{2} \sum_{j=1}^N w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{t_j}) + \xi\|^2$$

## Theory: more about training

Properties of denoising score matching objective:

- ▶  $\bar{\mathcal{L}}_{em}(\theta) \rightarrow \bar{\mathcal{L}}$  as  $n \rightarrow \infty$

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2} \sum_{j=1}^N w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \frac{1}{n} \sum_{i=1}^n \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

$$\bar{\mathcal{L}}(\theta) = \frac{1}{2} \sum_{j=1}^N w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{t_j}) + \xi\|^2$$

- ▶ Score matching  $\rightarrow$  denoising score matching:

## Theory: more about training

Properties of denoising score matching objective:

- ▶  $\bar{\mathcal{L}}_{em}(\theta) \rightarrow \bar{\mathcal{L}}$  as  $n \rightarrow \infty$

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2} \sum_{j=1}^N w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \frac{1}{n} \sum_{i=1}^n \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

$$\bar{\mathcal{L}}(\theta) = \frac{1}{2} \sum_{j=1}^N w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{t_j}) + \xi\|^2$$

- ▶ Score matching  $\rightarrow$  denoising score matching:

$$\begin{aligned} 0 &\leq \mathbb{E}_{X_{t_j}} \|S(\theta; t_j, X_{t_j}) - \nabla_x \log p_{t_j}(X_{t_j})\|^2 \\ &= \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{t_j}) + \xi\|^2 + C_{t_j} \end{aligned}$$

## Theory: more about training

Properties of denoising score matching objective:

- ▶  $\bar{\mathcal{L}}_{em}(\theta) \rightarrow \bar{\mathcal{L}}$  as  $n \rightarrow \infty$

$$\bar{\mathcal{L}}_{em}(\theta) = \frac{1}{2} \sum_{j=1}^N w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \frac{1}{n} \sum_{i=1}^n \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{ij}) + \xi_{ij}\|^2$$

$$\bar{\mathcal{L}}(\theta) = \frac{1}{2} \sum_{j=1}^N w(t_j)(t_j - t_{j-1}) \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{t_j}) + \xi\|^2$$

- ▶ Score matching  $\rightarrow$  denoising score matching:

$$\begin{aligned} 0 &\leq \mathbb{E}_{X_{t_j}} \|S(\theta; t_j, X_{t_j}) - \nabla_x \log p_{t_j}(X_{t_j})\|^2 \\ &= \frac{1}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \|\bar{\sigma}_{t_j} S(\theta; t_j, X_{t_j}) + \xi\|^2 + C_{t_j} \end{aligned}$$

- ▶  $\Rightarrow \bar{\mathcal{L}}(\theta) + \bar{C} \geq 0$

- ▶  $\Rightarrow \bar{\mathcal{L}}(\theta) \geq -\bar{C} > 0$ , i.e., positive lower bound of  $\bar{\mathcal{L}}(\theta)$

## Theory: more training

- ▶ Theorem: choose arbitrarily large width  $m$ ,  $\epsilon_{\text{train}}$  can be as small as possible

## Theory: more training

- ▶ Theorem: choose arbitrarily large width  $m$ ,  $\epsilon_{\text{train}}$  can be as small as possible

$$0 \overset{?}{\leftarrow} \epsilon_{\text{train}} \approx \bar{\mathcal{L}}_{em}(\theta) \rightarrow \bar{\mathcal{L}}(\theta) \geq -\bar{C} > 0$$

## Theory: more training

- ▶ Theorem: choose arbitrarily large width  $m$ ,  $\epsilon_{\text{train}}$  can be as small as possible

$$0 \stackrel{?}{\leftarrow} \epsilon_{\text{train}} \approx \bar{\mathcal{L}}_{em}(\theta) \rightarrow \bar{\mathcal{L}}(\theta) \geq -\bar{C} > 0$$

- ▶ Contradiction?

# Interesting generalization setting #1

No contradiction:



# Interesting generalization setting #1

No contradiction:

can have both  $\bar{\mathcal{L}}_{em}(\theta) \approx 0$  and  $\bar{\mathcal{L}} \approx -\bar{C} > 0$

# Interesting generalization setting #1

No contradiction:

can have both  $\bar{\mathcal{L}}_{em}(\theta) \approx 0$  and  $\bar{\mathcal{L}} \approx -\bar{C} > 0$

Reason:

# Interesting generalization setting #1

No contradiction:

can have both  $\bar{\mathcal{L}}_{em}(\theta) \approx 0$  and  $\bar{\mathcal{L}} \approx -\bar{C} > 0$

Reason:

- ▶ Overparameterized setting: fix  $m$ ,

# Interesting generalization setting #1

No contradiction:

can have both  $\bar{\mathcal{L}}_{em}(\theta) \approx 0$  and  $\bar{\mathcal{L}} \approx -\bar{C} > 0$

Reason:

- ▶ Overparameterized setting: fix  $m$ ,
- ▶  $\bar{\mathcal{L}}_{em}(\theta) \leq \epsilon_{\text{train}}$ : sample size  $n \ll m \Rightarrow n$  is bounded

# Interesting generalization setting #1

No contradiction:

can have both  $\bar{\mathcal{L}}_{em}(\theta) \approx 0$  and  $\bar{\mathcal{L}} \approx -\bar{C} > 0$

Reason:

- ▶ Overparameterized setting: fix  $m$ ,
  - ▶  $\bar{\mathcal{L}}_{em}(\theta) \leq \epsilon_{\text{train}}$ : sample size  $n \ll m \Rightarrow n$  is bounded
  - ▶  $\bar{\mathcal{L}}_{em} \rightarrow \bar{\mathcal{L}}$  as  $n \rightarrow \infty$

# Interesting generalization setting #1

No contradiction:

can have both  $\bar{\mathcal{L}}_{em}(\theta) \approx 0$  and  $\bar{\mathcal{L}} \approx -\bar{C} > 0$

Reason:

- ▶ Overparameterized setting: fix  $m$ ,
  - ▶  $\bar{\mathcal{L}}_{em}(\theta) \leq \epsilon_{\text{train}}$ : sample size  $n \ll m \Rightarrow n$  is bounded
  - ▶  $\bar{\mathcal{L}}_{em} \rightarrow \bar{\mathcal{L}}$  as  $n \rightarrow \infty$

Consequence:

# Interesting generalization setting #1

No contradiction:

can have both  $\bar{\mathcal{L}}_{em}(\theta) \approx 0$  and  $\bar{\mathcal{L}} \approx -\bar{C} > 0$

Reason:

- ▶ Overparameterized setting: fix  $m$ ,
  - ▶  $\bar{\mathcal{L}}_{em}(\theta) \leq \epsilon_{\text{train}}$ : sample size  $n \ll m \Rightarrow n$  is bounded
  - ▶  $\bar{\mathcal{L}}_{em} \rightarrow \bar{\mathcal{L}}$  as  $n \rightarrow \infty$

Consequence:

- ▶  $\epsilon_n = |\bar{\mathcal{L}}_{em}(\theta) - \bar{\mathcal{L}}(\theta)|$

# Interesting generalization setting #1

No contradiction:

can have both  $\bar{\mathcal{L}}_{em}(\theta) \approx 0$  and  $\bar{\mathcal{L}} \approx -\bar{C} > 0$

Reason:

- ▶ Overparameterized setting: fix  $m$ ,
  - ▶  $\bar{\mathcal{L}}_{em}(\theta) \leq \epsilon_{\text{train}}$ : sample size  $n \ll m \Rightarrow n$  is bounded
  - ▶  $\bar{\mathcal{L}}_{em} \rightarrow \bar{\mathcal{L}}$  as  $n \rightarrow \infty$

Consequence:

- ▶  $\epsilon_n = |\bar{\mathcal{L}}_{em}(\theta) - \bar{\mathcal{L}}(\theta)|$
- ▶ If  $\epsilon_{\text{train}}$  is small  $\Rightarrow \epsilon_n$  is large



## Interesting generalization setting #

- ▶ Recall:

$$\text{KL}(p_\delta | q_{T-\delta}) \lesssim \underbrace{E_D + E_I}_{\text{sampling}} + E_S$$

where  $p_\delta$  is the true density at time  $\delta$ , and  $q_{T-\delta}$  is the approximated density of  $p_\delta$ .

## Interesting generalization setting #

- ▶ Recall:

$$\text{KL}(p_\delta | q_{T-\delta}) \lesssim \underbrace{E_D + E_I}_{\text{sampling}} + E_S$$

where  $p_\delta$  is the true density at time  $\delta$ , and  $q_{T-\delta}$  is the approximated density of  $p_\delta$ .

- ▶  $E_S$ : score error

# Interesting generalization setting #1

$$E_S \leq \max_{1 \leq j \leq N} \frac{\sigma_{t_{N-j}}^2}{w(t_{N-j})} \left( \epsilon_{\text{train}} + \underbrace{\epsilon_n}_{\bar{\mathcal{L}}_{em}(\theta) - \bar{\mathcal{L}}(\theta)} + \underbrace{\epsilon_{\text{est}}}_{\text{estimation error}} + \underbrace{\epsilon_{\text{approx}}}_{\text{approximation error}} \right)$$

# Interesting generalization setting #1

$$E_S \leq \max_{1 \leq j \leq N} \frac{\sigma_{t_{N-j}}^2}{w(t_{N-j})} \left( \epsilon_{\text{train}} + \underbrace{\epsilon_n}_{\bar{\mathcal{L}}_{em}(\theta) - \bar{\mathcal{L}}(\theta)} + \underbrace{\epsilon_{\text{est}}}_{\text{estimation error}} + \underbrace{\epsilon_{\text{approx}}}_{\text{approximation error}} \right)$$

- ▶  $\epsilon_n + \epsilon_{\text{est}} + \epsilon_{\text{approx}}$   
[Chen et al., 2023, Oko et al., 2023, Han et al., 2024]

# Interesting generalization setting #1

$$E_S \leq \max_{1 \leq j \leq N} \frac{\sigma_{t_{N-j}}^2}{w(t_{N-j})} \left( \epsilon_{\text{train}} + \underbrace{\epsilon_n}_{\bar{\mathcal{L}}_{em}(\theta) - \bar{\mathcal{L}}(\theta)} + \underbrace{\epsilon_{\text{est}}}_{\text{estimation error}} + \underbrace{\epsilon_{\text{approx}}}_{\text{approximation error}} \right)$$

- ▶  $\epsilon_n + \epsilon_{\text{est}} + \epsilon_{\text{approx}}$   
[Chen et al., 2023, Oko et al., 2023, Han et al., 2024]
- ▶ Regression generalization  $\rightarrow$  Diffusion models

# Interesting generalization setting #1

$$E_S \leq \max_{1 \leq j \leq N} \frac{\sigma_{t_{N-j}}^2}{w(t_{N-j})} \left( \epsilon_{\text{train}} + \underbrace{\epsilon_n}_{\bar{\mathcal{L}}_{em}(\theta) - \bar{\mathcal{L}}(\theta)} + \underbrace{\epsilon_{\text{est}}}_{\text{estimation error}} + \underbrace{\epsilon_{\text{approx}}}_{\text{approximation error}} \right)$$

- ▶  $\epsilon_n + \epsilon_{\text{est}} + \epsilon_{\text{approx}}$   
[Chen et al., 2023, Oko et al., 2023, Han et al., 2024]
- ▶ Regression generalization  $\rightarrow$  Diffusion models
- ▶ Open problem:
  - ▶  $\epsilon_{\text{train}} + \epsilon_n \geq -2\bar{C} > 0 \Rightarrow$  error bound  $\not\rightarrow 0$

# Interesting generalization setting #1

$$E_S \leq \max_{1 \leq j \leq N} \frac{\sigma_{t_{N-j}}^2}{w(t_{N-j})} \left( \epsilon_{\text{train}} + \underbrace{\epsilon_n}_{\bar{\mathcal{L}}_{em}(\theta) - \bar{\mathcal{L}}(\theta)} + \underbrace{\epsilon_{\text{est}}}_{\text{estimation error}} + \underbrace{\epsilon_{\text{approx}}}_{\text{approximation error}} \right)$$

- ▶  $\epsilon_n + \epsilon_{\text{est}} + \epsilon_{\text{approx}}$   
[Chen et al., 2023, Oko et al., 2023, Han et al., 2024]
- ▶ Regression generalization  $\rightarrow$  Diffusion models
- ▶ Open problem:
  - ▶  $\epsilon_{\text{train}} + \epsilon_n \geq -2\bar{C} > 0 \Rightarrow$  error bound  $\nrightarrow 0$   
change decomposition  $\Rightarrow$  tighter analysis?

# Interesting generalization setting #1

$$E_S \leq \max_{1 \leq j \leq N} \frac{\sigma_{t_{N-j}}^2}{w(t_{N-j})} \left( \epsilon_{\text{train}} + \underbrace{\epsilon_n}_{\bar{\mathcal{L}}_{em}(\theta) - \bar{\mathcal{L}}(\theta)} + \underbrace{\epsilon_{\text{est}}}_{\text{estimation error}} + \underbrace{\epsilon_{\text{approx}}}_{\text{approximation error}} \right)$$

- ▶  $\epsilon_n + \epsilon_{\text{est}} + \epsilon_{\text{approx}}$   
[Chen et al., 2023, Oko et al., 2023, Han et al., 2024]
- ▶ Regression generalization  $\rightarrow$  Diffusion models
- ▶ Open problem:
  - ▶  $\epsilon_{\text{train}} + \epsilon_n \geq -2\bar{C} > 0 \Rightarrow$  error bound  $\not\rightarrow 0$   
change decomposition  $\Rightarrow$  tighter analysis?
  - ▶  $S(\theta; t, X_t) \stackrel{?}{\rightarrow} \nabla_x \log p_t(X_t)$



## Example full error analysis

Theorem (EDM polynomial schedule [Karras et al., 2022])

$$KL(p_\delta | q_{T-\delta}) \lesssim \underbrace{\frac{m_2^2}{T^2}}_{E_I} + \underbrace{\frac{da^2 T^{\frac{1}{a}}}{\delta^{\frac{1}{a}} N} + (m_2^2 + d) \left( \frac{a^2 T^{\frac{1}{a}}}{\delta^{\frac{1}{a}} N} + \frac{a^3 T^{\frac{2}{a}}}{\delta^{\frac{2}{a}} N^2} \right)}_{E_D} + \underbrace{\frac{1}{N} \left( C_2 + \left( 1 - C_1 h \left( \frac{md^{\frac{a_0-1}{2}}}{n^3 N^2} \right) \right)^K \right)}_{E_S},$$

where  $\delta = t_0$ ,  $a = 7$ ,  $a_0 \in (1/2, 1)$ .

►  $C_2 = \epsilon_n + \epsilon_{\text{est}} + \epsilon_{\text{approx}}$

# Theory: full error analysis

$$E_S := \sum_{j=0}^{N-1} \alpha_j \left( \underbrace{t_j}_{\text{time schedule}}, \underbrace{\bar{\sigma}_{t_j}}_{\text{variance schedule}} \right) \mathbb{E}_{X_{t_j}} \|S(\theta; t_j, X_{t_j}) - \nabla \log p_{t_j}(X_{t_j})\|^2$$

sampling + optimization

# Theory: full error analysis

$$E_S := \sum_{j=0}^{N-1} \alpha_j \left( \underbrace{t_j}_{\text{time schedule}}, \underbrace{\bar{\sigma}_{t_j}}_{\text{variance schedule}} \right) \mathbb{E}_{X_{t_j}} \|S(\theta; t_j, X_{t_j}) - \nabla \log p_{t_j}(X_{t_j})\|^2$$

sampling + optimization

- ▶  $\alpha_j(t_j, \bar{\sigma}_{t_j}) \neq \beta_j$ , the weighting for training objective

# Theory: full error analysis

$$E_S := \sum_{j=0}^{N-1} \alpha_j \left( \underbrace{t_j}_{\text{time schedule}}, \underbrace{\bar{\sigma}_{t_j}}_{\text{variance schedule}} \right) \mathbb{E}_{X_{t_j}} \|S(\theta; t_j, X_{t_j}) - \nabla \log p_{t_j}(X_{t_j})\|^2$$

sampling + optimization

- ▶  $\alpha_j(t_j, \bar{\sigma}_{t_j}) \neq \beta_j$ , the weighting for training objective
- ▶ First choose total weighting  $\beta_j$ ; then apply the schedules  $t_j, \bar{\sigma}_{t_j}$

# Theory: full error analysis

$$E_S := \sum_{j=0}^{N-1} \alpha_j \left( \underbrace{t_j}_{\text{time schedule}}, \underbrace{\bar{\sigma}_{t_j}}_{\text{variance schedule}} \right) \mathbb{E}_{X_{t_j}} \|S(\theta; t_j, X_{t_j}) - \nabla \log p_{t_j}(X_{t_j})\|^2$$

sampling + optimization

- ▶  $\alpha_j(t_j, \bar{\sigma}_{t_j}) \neq \beta_j$ , the weighting for training objective
- ▶ First choose total weighting  $\beta_j$ ; then apply the schedules  $t_j, \bar{\sigma}_{t_j}$
- ▶ Next: focus on two concrete schedules used in practice

# Theory: full error analysis

$$E_S := \sum_{j=0}^{N-1} \alpha_j \left( \underbrace{t_j}_{\text{time schedule}}, \underbrace{\bar{\sigma}_{t_j}}_{\text{variance schedule}} \right) \mathbb{E}_{X_{t_j}} \|S(\theta; t_j, X_{t_j}) - \nabla \log p_{t_j}(X_{t_j})\|^2$$

sampling + optimization

- ▶  $\alpha_j(t_j, \bar{\sigma}_{t_j}) \neq \beta_j$ , the weighting for training objective
- ▶ First choose total weighting  $\beta_j$ ; then apply the schedules  $t_j, \bar{\sigma}_{t_j}$
- ▶ Next: focus on two concrete schedules used in practice
  - ▶ Theoretical implication: how to choose between two schedules

# Theory: full error analysis

Two most famous choices:

# Theory: full error analysis

Two most famous choices:

- ▶ Exponential schedules [Song et al., 2021]: first work



# Theory: full error analysis

Two most famous choices:

- ▶ Exponential schedules [Song et al., 2021]: first work
- ▶ Polynomial schedules [Karras et al., 2022]: improved design

# Theory: full error analysis

Two most famous choices:

- ▶ Exponential schedules [Song et al., 2021]: first work
- ▶ Polynomial schedules [Karras et al., 2022]: improved design

	Variance $\bar{\sigma}_t$	Time $t_k$
[Karras et al., 2022]	$t$	$\left( \bar{\sigma}_{\max}^{1/\rho} - \left( \bar{\sigma}_{\max}^{1/\rho} - \bar{\sigma}_{\min}^{1/\rho} \right) \frac{N-k}{N} \right)^\rho$
[Song et al., 2021]	$\sqrt{t}$	$\bar{\sigma}_{\max}^2 \left( \frac{\bar{\sigma}_{\min}^2}{\bar{\sigma}_{\max}^2} \right)^{\frac{N-k}{N}}$

# Theory: full error analysis

Two most famous choices:

- ▶ Exponential schedules [Song et al., 2021]: first work
- ▶ Polynomial schedules [Karras et al., 2022]: improved design

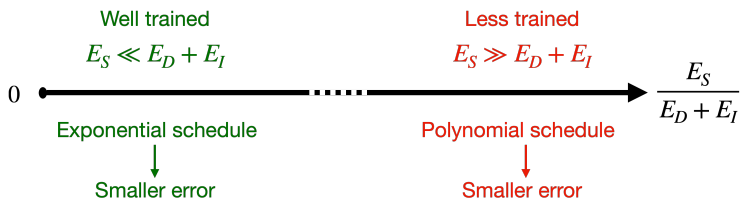
---

	Variance $\bar{\sigma}_t$	Time $t_k$
[Karras et al., 2022]	$t$	$\left( \bar{\sigma}_{\max}^{1/\rho} - \left( \bar{\sigma}_{\max}^{1/\rho} - \bar{\sigma}_{\min}^{1/\rho} \right) \frac{N-k}{N} \right)^\rho$
[Song et al., 2021]	$\sqrt{t}$	$\bar{\sigma}_{\max}^2 \left( \frac{\bar{\sigma}_{\min}^2}{\bar{\sigma}_{\max}^2} \right)^{\frac{N-k}{N}}$

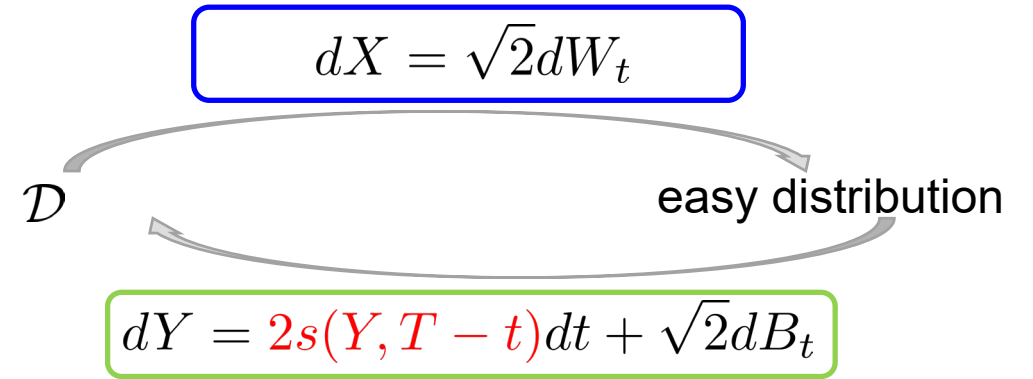
---

- ▶ Time schedule  $t_k$ : a function of  $k$
- ▶  $\rho = 7$

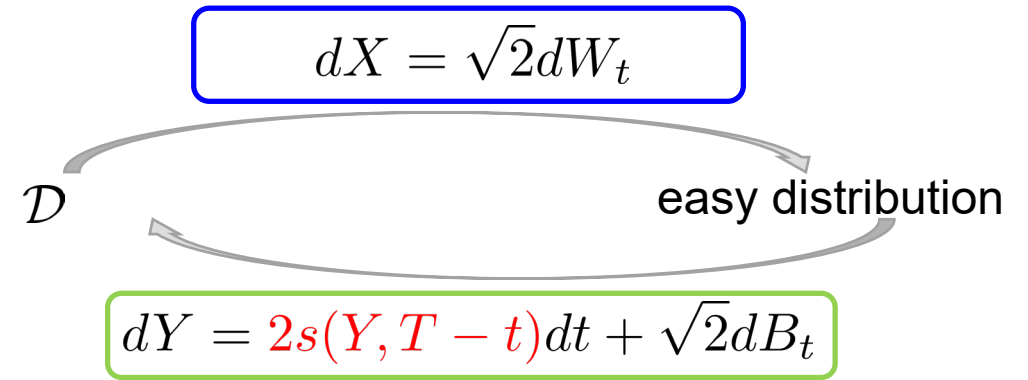
# Implication from theory: How to choose schedules?



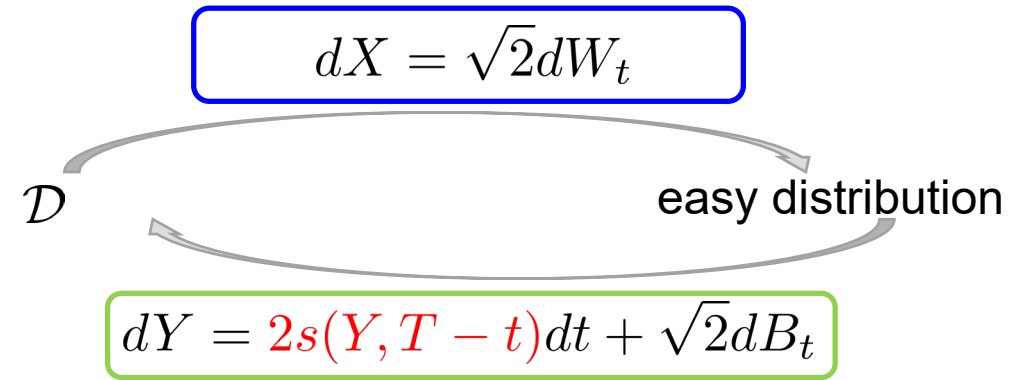
- 1st quantitative result that analyzes  
forward training + backward sampling



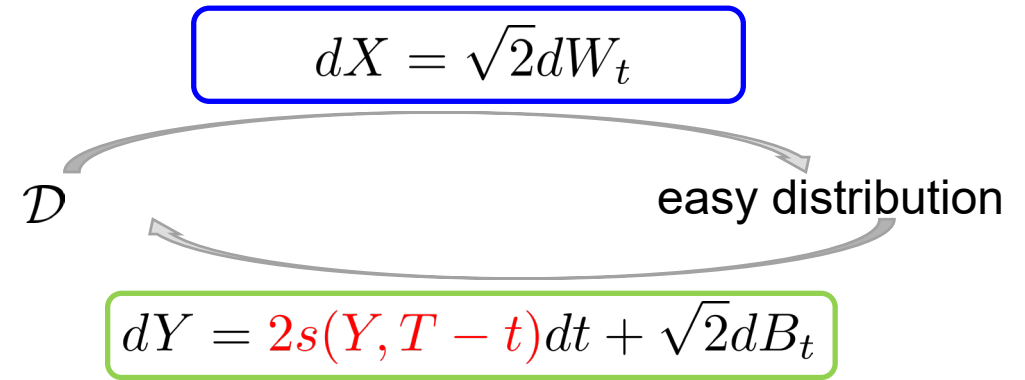
- 1st quantitative result that analyzes  
    **forward** training + **backward** sampling
- nontrivial analysis of training dynamics
  - overparameterized ReLU MLP



- 1st quantitative result that analyzes  
forward training + backward sampling
- nontrivial analysis of training dynamics
  - overparameterized ReLU MLP
  - weaker data assumptions (no separability, unbdd output, ...)

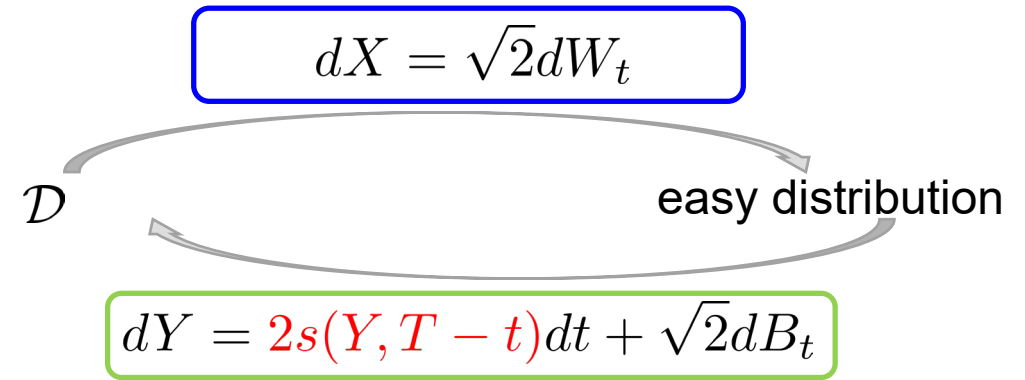


- 1st quantitative result that analyzes  
forward training + backward sampling
- nontrivial analysis of training dynamics
  - overparameterized ReLU MLP
  - weaker data assumptions (no separability, unbdd output, ...)
  - interesting generalization setting #1: SM - DSM gap

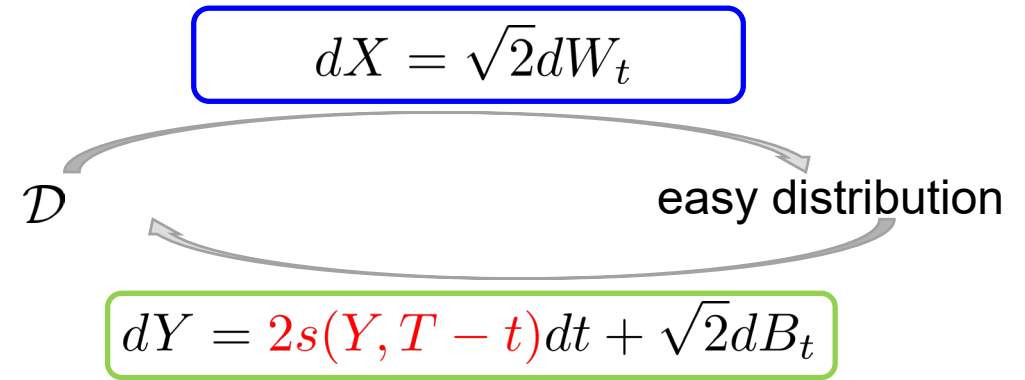




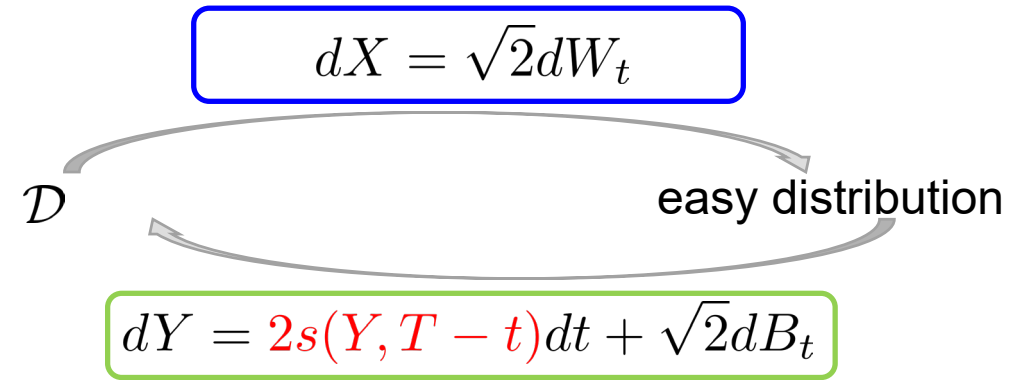
- 1st quantitative result that analyzes  
     **forward** training + **backward** sampling
- nontrivial analysis of training dynamics
  - overparameterized ReLU MLP
  - weaker data assumptions (no separability, unbdd output, ...)
  - interesting generalization setting #1: SM - DSM gap
- besides understanding, practical implication



- 1st quantitative result that analyzes  
     **forward** training + **backward** sampling
- nontrivial analysis of training dynamics
  - overparameterized ReLU MLP
  - weaker data assumptions (no separability, unbdd output, ...)
  - interesting generalization setting #1: SM - DSM gap
- besides understanding, practical implication
  - total weighting



- 1st quantitative result that analyzes  
     **forward** training + **backward** sampling
- nontrivial analysis of training dynamics
  - overparameterized ReLU MLP
  - weaker data assumptions (no separability, unbdd output, ...)
  - interesting generalization setting #1: SM - DSM gap
- besides understanding, practical implication
  - total weighting
  - variance and time schedules



Generalization Setting #2 (open)

## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

**accurate**, but **not innovative**

## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

Key: what exactly is this?

## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

Key: what exactly is this?

$$\min_{\theta} \int_0^T w(t) \underbrace{\mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2}_{\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2} dt$$



## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

Key: what exactly is this?

$$\min_{\theta} \int_0^T w(t) \underbrace{\mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2}_{\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2} dt$$

natural: empirical distribution, i.e. sum of Dirac's

## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

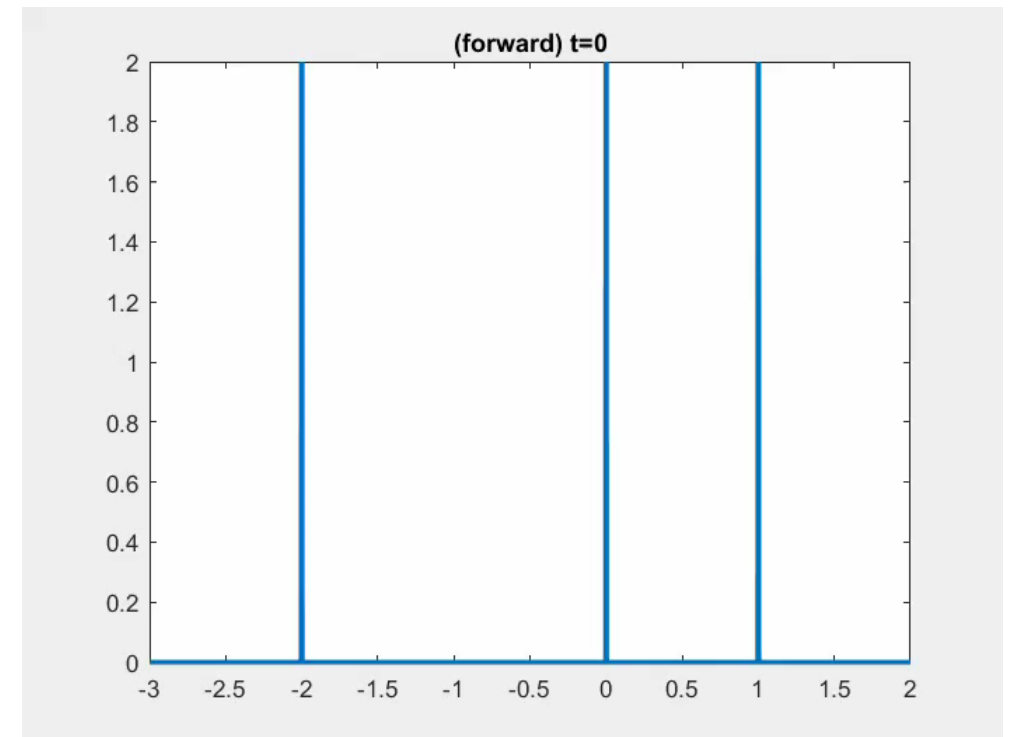
accurate, but not innovative

Key: what exactly is this?

$$\min_{\theta} \int_0^T w(t) \underbrace{\mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2}_{\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2} dt$$

natural: empirical distribution, i.e. sum of Dirac's

exact density evolution



## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

Key: what exactly is this?

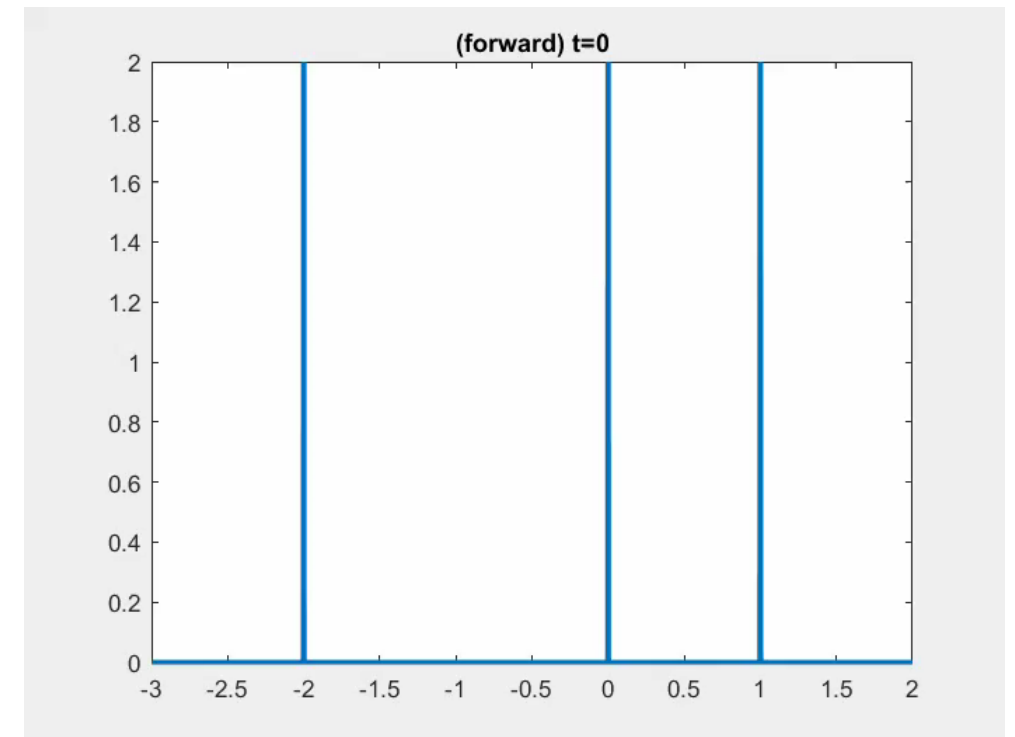
$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

$$\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2$$

natural: empirical distribution, i.e. sum of Dirac's

perfect score

exact density evolution



## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

Key: what exactly is this?

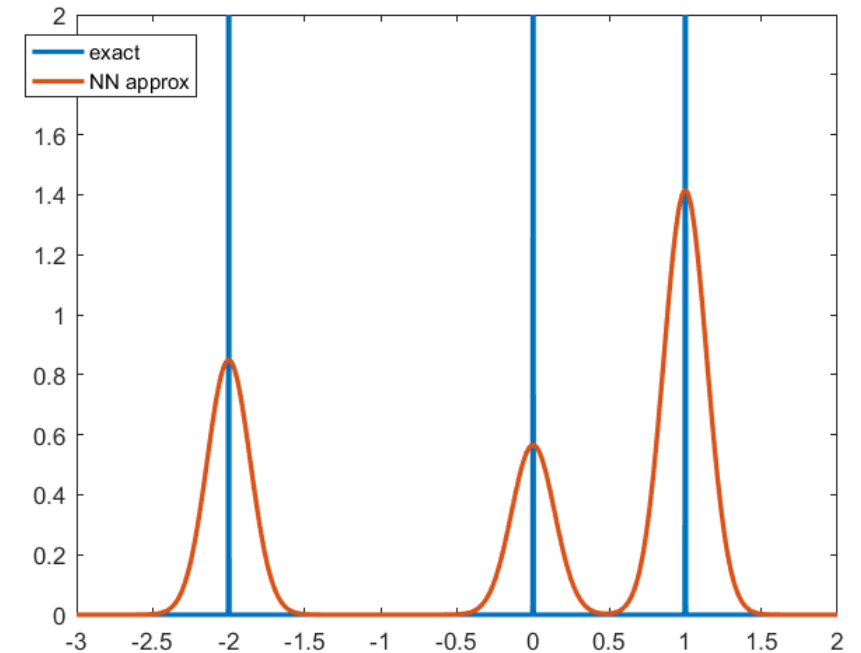
$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

$$\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2$$

natural: empirical distribution, i.e. sum of Dirac's

perfect score

approximation error



## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

- score error
- integration error
- initialization error
- early stopping

Key: what exactly is this?

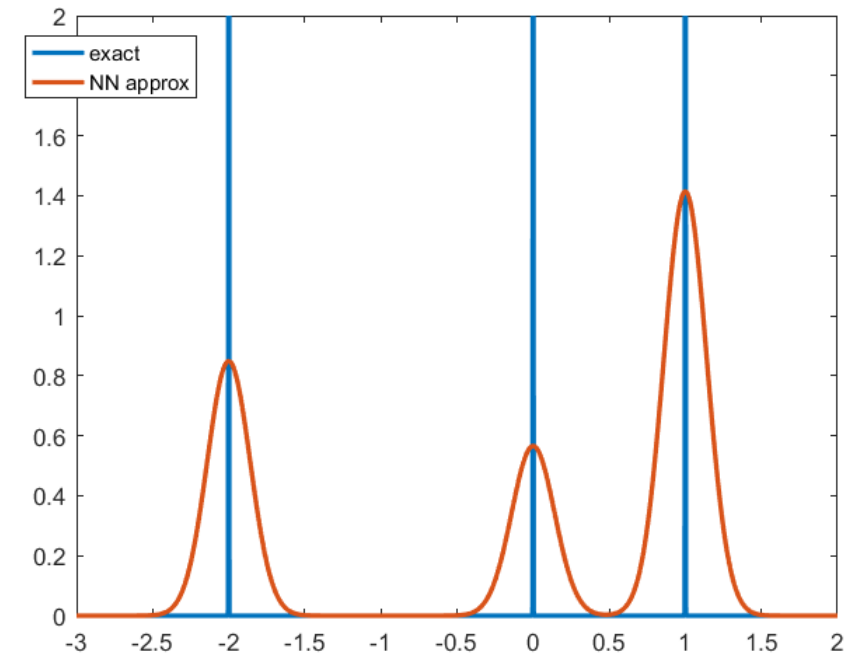
$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

$$\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2$$

natural: empirical distribution, i.e. sum of Dirac's

perfect score

4 types of errors



## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

Key: what exactly is this?

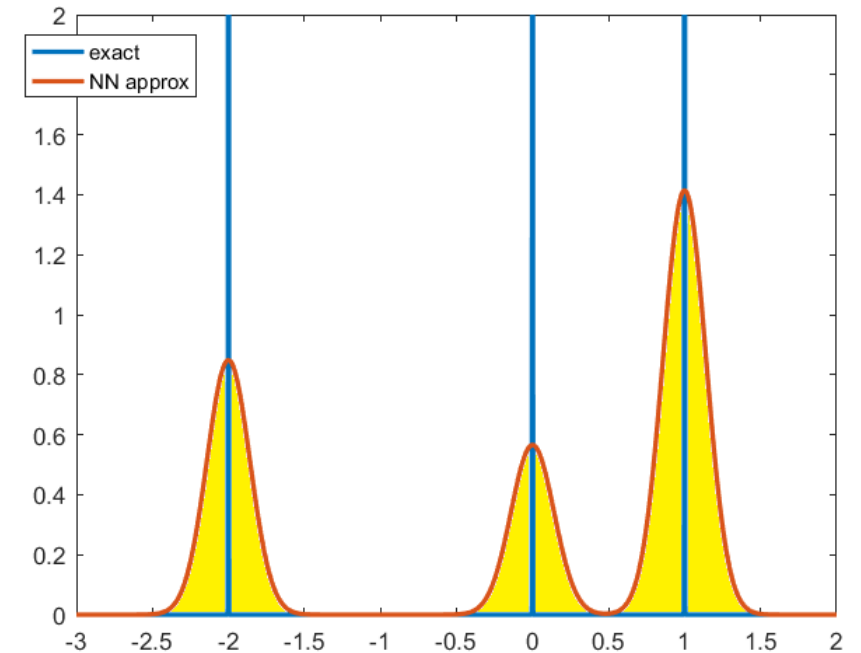
$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

$$\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2$$

natural: empirical distribution, i.e. sum of Dirac's

perfect score

4 types of errors → innovation?



## Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

Key: what exactly is this?

$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(X_t, t) - s_{\theta}(X_t, t)\|^2 dt$$

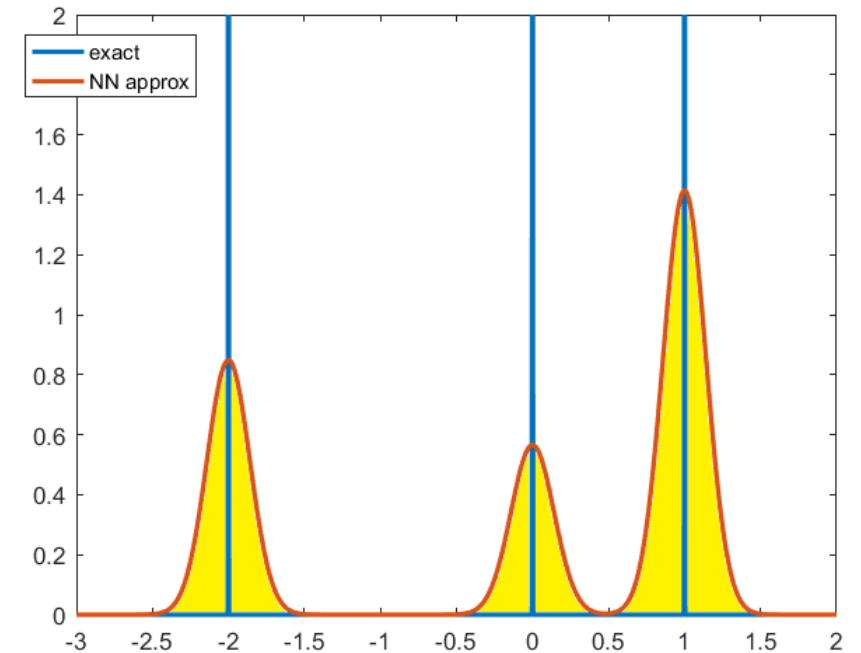
$$\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2$$

natural: empirical distribution, i.e. sum of Dirac's

perfect score

Kadkhodaie+ 23, Scarvelis+ 23, ...

4 types of errors → innovation?



Generalization Setting #2 (open)

Quantifications:  $d(p_{\text{training data}} | p_{\text{generated data}}) \leq \dots$

What if: generated data = uniformly drawn from training data ?

accurate, but not innovative

perfect score

Kadkhodaie+ 23, Scarvelis+ 23, ...

4 types of errors → innovation?

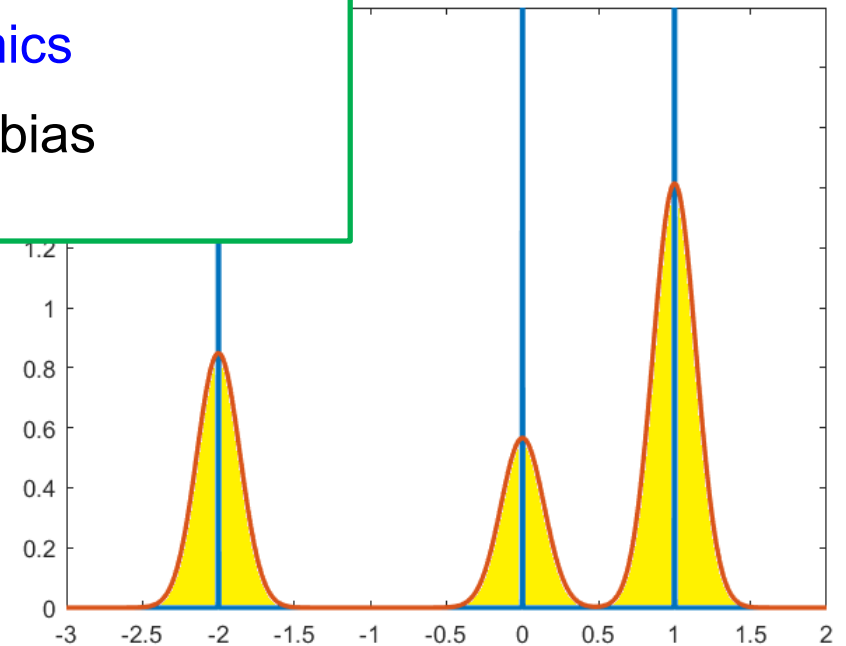
Key:

hope: understanding the training dynamics  
& tight analysis of sampling → inductive bias

$$\min_{\theta} \int_0^T w(t) \mathbb{E}_{X_t} \|\nabla_x \log p(\Lambda_t, t) - s_{\theta}(\Lambda_t, t)\|^2 dt$$

$$\approx \frac{1}{n} \sum_{i=1}^n \|\nabla_x \log p(X_t^i, t) - s_{\theta}(X_t^i, t)\|^2$$

natural: empirical distribution, i.e. sum of Dirac's







Can **discrete diffusion model** add to the success of **LLM**?

Can **discrete diffusion model** add to the success of **LLM**?

Quantitative error analysis (possible)

Can **discrete diffusion model** add to the success of **LLM**?

Quantitative error analysis (possible) →

what is it good at?

Can **discrete diffusion model** add to the success of **LLM**?

Quantitative error analysis (possible) →

what is it good at?

difference and similarity to autoregressive model

Can **discrete diffusion model** add to the success of **LLM**?

Quantitative error analysis (possible) →

what is it good at?

difference and similarity to autoregressive model

how to best deploy it?

...

Thank *you* for your attention and feedback!

Support:

**NSF** DMS-1847802, ECCS-1936776

**Cullen-Peck** Scholarship

**Emory-GT** AI.Humanity Award

**Simons Institute** Research Fellowship



MoleiTaoMath



**SCAN ME**



itsdynamical



**SCAN ME**