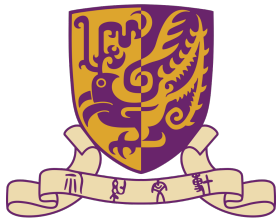
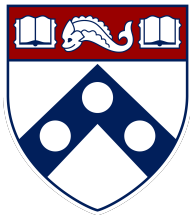


# Theoretical Foundations of Diffusion Models



Yuxin Chen, Gen Li, Yuting Wei

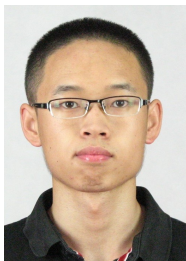
UPenn & CUHK

# Instructors

---



Yuxin Chen  
UPenn

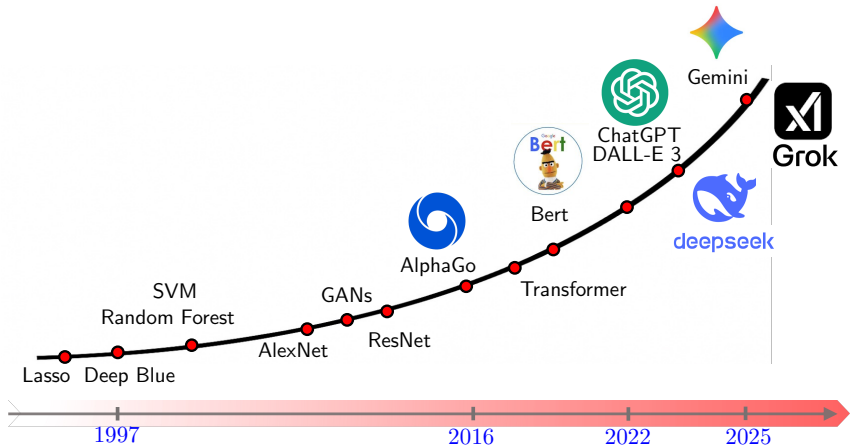


Gen Li  
CUHK

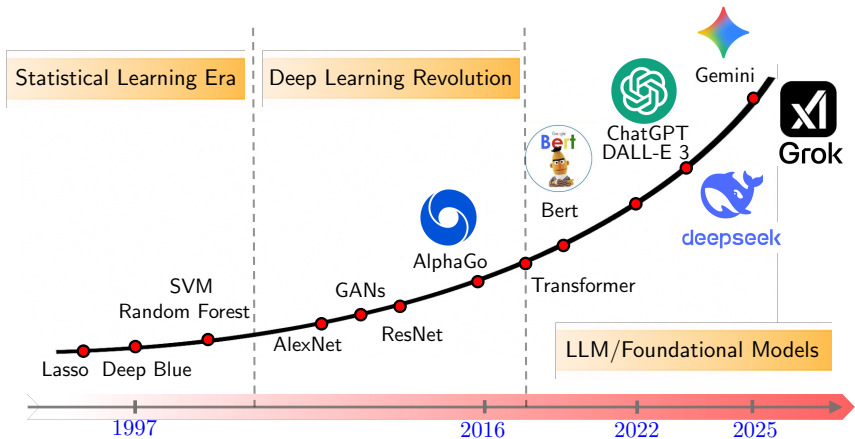


Yuting Wei  
UPenn

# Evolution of AI models



# Evolution of AI models





# Diffusion models hold great promises

---

---

## Diffusion Models Beat GANs on Image Synthesis

---

Prafulla Dhariwal\*  
OpenAI  
prafulla@openai.com

Alex Nichol\*  
OpenAI  
alex@openai.com

NeurIPS 2021, cited 10023

# Diffusion models hold great promises

---

---

## Diffusion Models Beat GANs on Image Synthesis

---

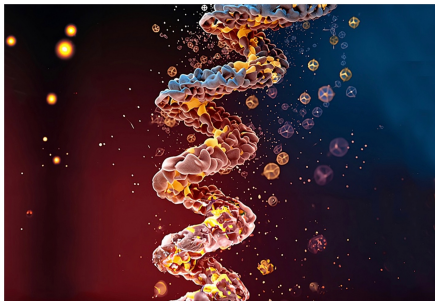
**Prafulla Dhariwal\***  
OpenAI  
prafulla@openai.com

**Alex Nichol\***  
OpenAI  
alex@openai.com

### Generative AI imagines new protein structures

“FrameDiff” is a computational tool that uses generative AI to craft new protein structures, with the aim of accelerating drug development and improving gene therapy.

Rachel Gordon | MIT CSAIL  
July 12, 2023



# Diffusion models hold great promises

## Diffusion Models Be

Diffusion models are now turbocharging reinforcement learning systems

By Ben Dickson - March 4, 2024

**Prafulla Dhariwal\***

OpenAI

prafulla@openai.com

Like 75



Facebook



Twitter



Reddit



LinkedIn

## Generative AI imagin

“FrameDiff” is a computational protein structures, with the aim of improving gene therapy.

Rachel Gordon | MIT CSAIL

July 12, 2023



# Diffusion models hold great promises

## Diffusion Models Be:

Diffusion models are now turbocharging reinforcement learning systems

By Ben Dickson · March 4, 2024

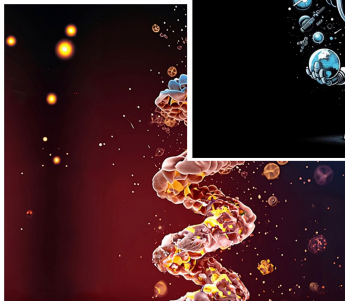
Prafulla Dhariwal\*  
OpenAI  
prafulla@openai.com



## Generative AI imagin

“FrameDiff” is a computational to protein structures, with the aim of improving gene therapy.

Rachel Gordon | MIT CSAIL  
July 12, 2023



## ScienceAdvances

HOME > SCIENCE ADVANCES > VOL. 10, NO. 13 > GENERATIVE EMULATION OF WEATHER FORECAST ENSEMBLES WITH DIFFUSION MODELS

RESEARCH ARTICLE ATMOSPHERIC SCIENCE



## Generative emulation of weather forecast ensembles with diffusion models

LIZAO LI · ROBERT CARVER · IGNACIO LOPEZ-GOMEZ · FEI SHA · AND JOHN ANDERSON [Authors Info & Affiliations](#)

SCIENCE ADVANCES · 29 Mar 2024 · Vol 10, Issue 13 · DOI:10.1126/sciadv.adf4483

## Stable Diffusion Based Solution for Medical Imaging

By NETSOL Technologies, on October 23, 2023

In this blog post, we'll deeply dive into the fascinating world of stable diffusion-based solutions and explore how they work alongside generative AI to take medical imaging to the next level.

INNOVATION > AI

## Microsoft Team Uses Diffusion Model For Materials Science

By John Werner, Contributor. © I am an MIT Senior Fellow & Lecturer, 5x...

Follow Author

Published Jan 21, 2025, 03:40am EST

# Generative modeling

---

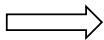
training data



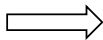
- Given training data  $\underbrace{X^{\text{train},i} \sim p_{\text{data}}}_{\text{from a general distribution}} (1 \leq i \leq N)$  in  $\mathbb{R}^d$

# Generative modeling

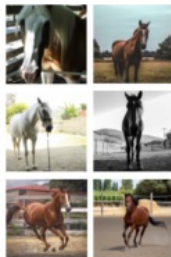
training data



Generative  
modeling



new samples



- Given training data  $\underbrace{X^{\text{train},i} \sim p_{\text{data}}}_{\text{from a general distribution}} (1 \leq i \leq N)$  in  $\mathbb{R}^d$
- Generate **new** samples  $Y \sim p_{\text{data}}$

# A natural approach: density estimation

---

- learn the distribution directly (parameterized by  $\theta$ ):

$$p(x | \theta) = \frac{e^{-f_{\theta}(x)}}{Z_{\theta}},$$

where  $Z_{\theta}$  is a normalizing constant depending on  $\theta$

# A natural approach: density estimation

---

- learn the distribution directly (parameterized by  $\theta$ ):

$$p(x | \theta) = \frac{e^{-f_{\theta}(x)}}{Z_{\theta}},$$

where  $Z_{\theta}$  is a normalizing constant depending on  $\theta$

- Use maximum likelihood (or posterior) to estimate  $\theta$ :

$$\max_{\theta} \sum_{i=1}^N \log p(X_i | \theta)$$

# A natural approach: density estimation

---

- learn the distribution directly (parameterized by  $\theta$ ):

$$p(x | \theta) = \frac{e^{-f_{\theta}(x)}}{Z_{\theta}},$$

where  $Z_{\theta}$  is a normalizing constant depending on  $\theta$

- Use maximum likelihood (or posterior) to estimate  $\theta$ :

$$\max_{\theta} \sum_{i=1}^N \log p(X_i | \theta) \longrightarrow \text{Intractable!}$$

## Another approach: score function

---

The **(Stein) score function** of a distribution  $p(x)$  is defined as

$$s(x) := \nabla_x \log p_X(x).$$

Note that

$$\begin{aligned} \nabla \log p(x | \theta) &= \nabla_x \log \frac{e^{-f_\theta(x)}}{Z_\theta} \\ &= -\nabla_x f_\theta(x) \end{aligned}$$

getting rid of the annoying  $Z_\theta$ !

## Another approach: score function

---

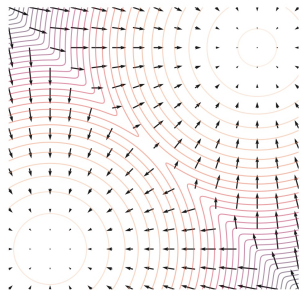
The **(Stein) score function** of a distribution  $p(x)$  is defined as

$$s(x) := \nabla_x \log p_X(x).$$

Note that

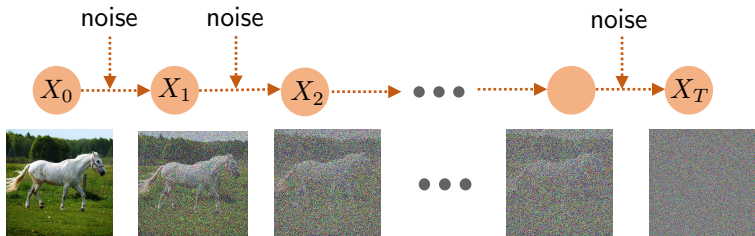
$$\begin{aligned} \nabla \log p(x | \theta) &= \nabla_x \log \frac{e^{-f_\theta(x)}}{Z_\theta} \\ &= -\nabla_x f_\theta(x) \end{aligned}$$

getting rid of the annoying  $Z_\theta$ !

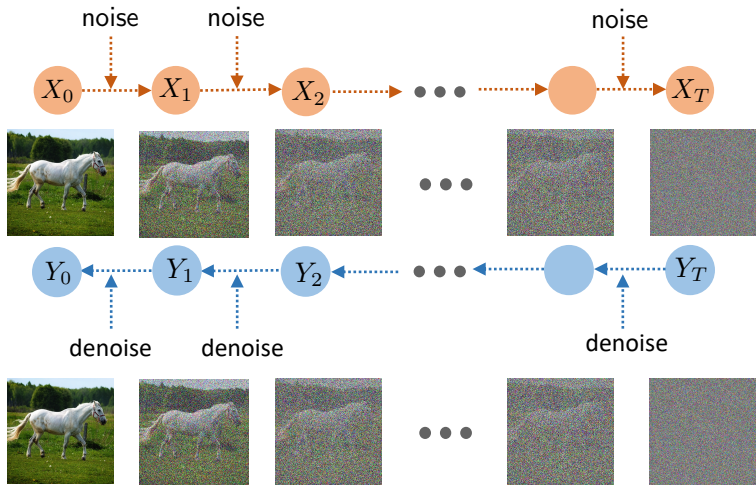


Score function of Gaussian mixtures

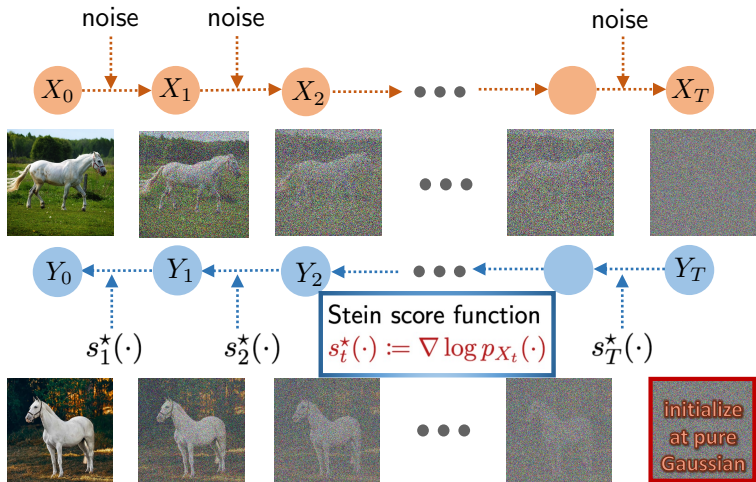
The score function points towards regions of higher probability.



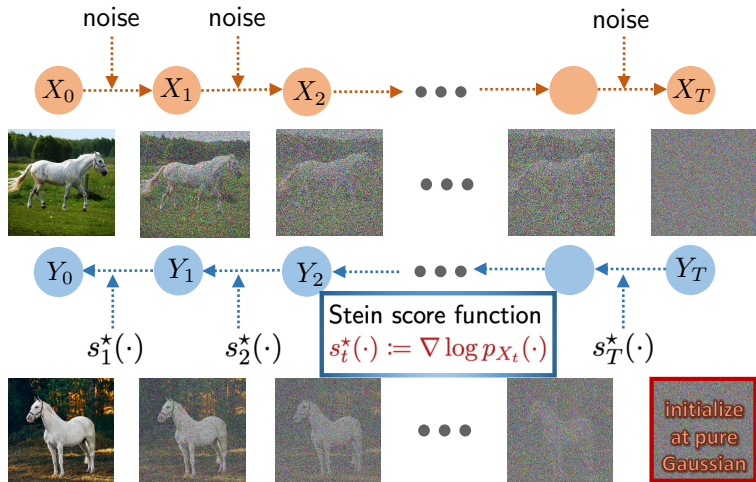
- **forward process:** diffuse data into noise



- **forward process:** diffuse data into noise
- **reverse process:** convert pure noise into data-like distributions



- **forward process:** diffuse data into noise
- **reverse process:** convert pure noise into data-like distributions



- **forward process:** diffuse data into noise
- **reverse process:** convert pure noise into data-like distributions

**Goal:**  $Y_t \stackrel{d}{\approx} X_t, \quad t = T, \dots, 1$

# Score is all you need

---

How to learn a reverse process s.t.  $Y_t \stackrel{d}{\approx} X_t$ , for  $t = T, \dots, 1$ ?

# Score is all you need

---

How to learn a reverse process s.t.  $Y_t \stackrel{d}{\approx} X_t$ , for  $t = T, \dots, 1$ ?

It is feasible as long as one knows the score function  $\underbrace{\nabla \log p_{X_t}(x)}_{\text{w.r.t. } X}$

# Score is all you need

How to learn a reverse process s.t.  $Y_t \stackrel{d}{\approx} X_t$ , for  $t = T, \dots, 1$ ?

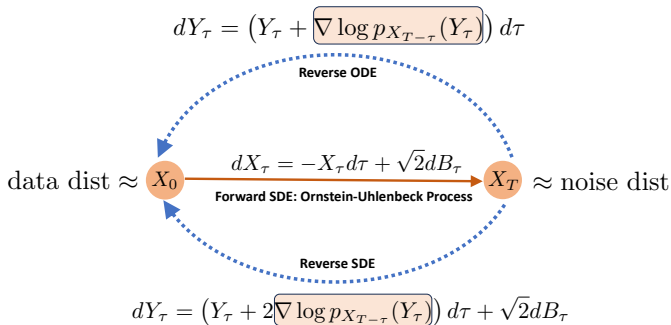
It is feasible as long as one knows the score function  $\underbrace{\nabla \log p_{X_t}(x)}_{\text{w.r.t. } X}$

$$\text{data dist} \approx X_0 \xrightarrow[\text{Forward SDE: Ornstein-Uhlenbeck Process}]{dX_\tau = -X_\tau d\tau + \sqrt{2}dB_\tau} X_T \approx \text{noise dist}$$

# Score is all you need

How to learn a reverse process s.t.  $Y_t \stackrel{d}{\approx} X_t$ , for  $t = T, \dots, 1$ ?

It is feasible as long as one knows the score function  $\underbrace{\nabla \log p_{X_t}(x)}_{\text{w.r.t. } X}$




# Towards mathematical theory for diffusion models

---

- hard to develop full-fledged **end-to-end** theory

# Towards mathematical theory for diffusion models

---

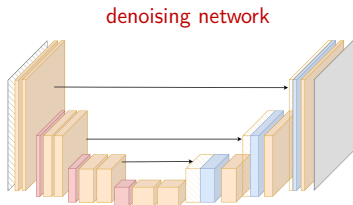
- hard to develop full-fledged **end-to-end** theory
- “divide-and-conquer”: score learning  decouple  $\leftarrow \color{red}{X} \rightarrow$  sampling

# A divide-and-conquer approach

— Li, Lu, Tan '22

— Chen, Chewi, Li, Li, Salim, Zhang '22

— Benton, De Bortoli, Doucet, Deligiannidis '23



learn  $s_t(\cdot) \approx \nabla \log p_{X_t}(\cdot)$

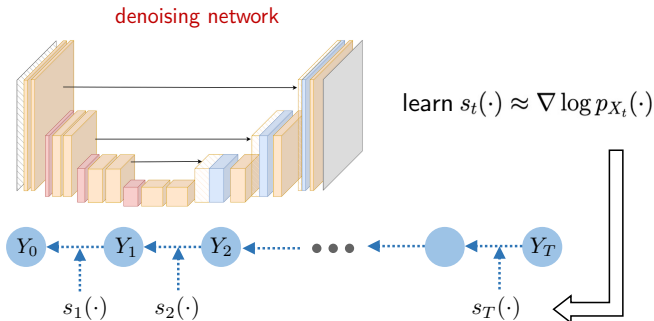
1. **score learning/matching:** learn estimates  $s_t(\cdot)$  for  $\nabla \log p_{X_t}(\cdot)$

# A divide-and-conquer approach

— Li, Lu, Tan '22

— Chen, Chewi, Li, Li, Salim, Zhang '22

— Benton, De Bortoli, Doucet, Deligiannidis '23



1. **score learning/matching:** learn estimates  $s_t(\cdot)$  for  $\nabla \log p_{X_t}(\cdot)$
2. **data generation:** sampling w/ aid of score estimates  $\{s_t(\cdot)\}_{11/124}$

## Score matching via denoising

---

$$X_0 \sim p_{\text{data}}, \quad X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \mathcal{N}(0, I_d), \quad 1 \leq t \leq T$$

## Score matching via denoising

---

$$X_0 \sim p_{\text{data}}, \quad X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \mathcal{N}(0, I_d), \quad 1 \leq t \leq T$$

Tweedie's formula ([Hyvarinen'05](#); [Vincent'11](#)):

$$s^*(x) = \frac{1}{\sqrt{1 - \bar{\alpha}_t}} \mathbb{E} \left[ W \mid \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} W = x \right],$$

where the expectation is taken over  $W \sim \mathcal{N}(0, I_d)$ ,  $X_0 \sim p_{\text{data}}$ .

# Score matching via denoising

---

$$X_0 \sim p_{\text{data}}, \quad X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \mathcal{N}(0, I_d), \quad 1 \leq t \leq T$$

Tweedie's formula ([Hyvarinen'05](#); [Vincent'11](#)):

$$s^*(x) = \frac{1}{\sqrt{1 - \bar{\alpha}_t}} \mathbb{E} \left[ W \mid \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} W = x \right],$$

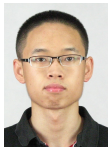
where the expectation is taken over  $W \sim \mathcal{N}(0, I_d)$ ,  $X_0 \sim p_{\text{data}}$ .

- nonparametric methods [Wibisono et al.'24](#); [Zhang et al.'24](#); [Dou et al.'24](#)
- AMP [Wu & Montanari'23](#)
- neural networks [Cole and Lu'24](#), [Mei and Wu'23](#), [Okon et al.'23](#)

## **Agenda:**

1. non-asymptotic convergence theory
2. adaptation to (unknown) low dimensionality
3. acceleration via higher-order approximation
4. provable benefits of diffusion guidance
5. diffusion models for inverse problems
6. discrete diffusion (or diffusion language models)

## *Part 1: nonasymptotic convergence theory*



Gen Li  
CUHK



Yuling Yan  
DE Shaw



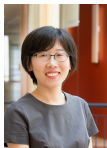
Yuchen Jiao  
CUHK



Yuchen Zhou  
UIUC



Yuxin Chen  
UPenn



Yuejie Chi  
Yale



Yuting Wei  
UPenn

# Two mainstream approaches

---

## Denoising Diffusion Probabilistic Models

**Jonathan Ho**  
UC Berkeley  
jonathanho@berkeley.edu

**Ajay Jain**  
UC Berkeley  
ajayj@berkeley.edu

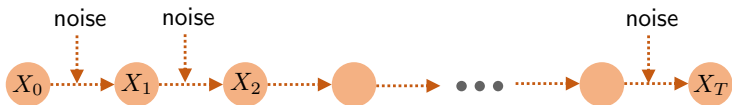
**Pieter Abbeel**  
UC Berkeley  
pabbeel@cs.berkeley.edu

## DENOISING DIFFUSION IMPLICIT MODELS

**Jiaming Song, Chenlin Meng & Stefano Ermon**  
Stanford University  
{tsong, chenlin, ermon}@cs.stanford.edu

# DDPM vs. DDIM

---



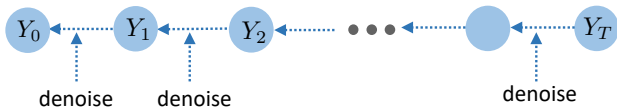
forward process:  $X_0 \sim p_{\text{data}}$  (target distribution)

$$X_t = \sqrt{\alpha_t} X_{t-1} + \sqrt{1 - \alpha_t} \mathcal{N}(0, I_d) \quad t \geq 1$$

- $\beta_t := 1 - \alpha_t$  controls variance of injected noise

# DDPM vs. DDIM

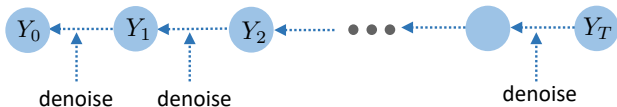
---



— *Ho, Jain, Abbeel '20*

1. A stochastic sampler: **denoising diffusion probabilistic models**  
DDPM

# DDPM vs. DDIM



— Ho, Jain, Abbeel '20

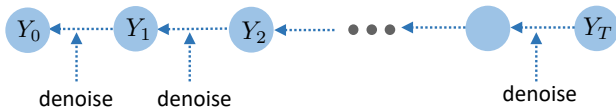
1. A stochastic sampler: **denoising diffusion probabilistic models**  
DDPM

$$Y_T \sim \mathcal{N}(0, I_d)$$

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \underbrace{Y_t + \eta_t^{\text{ddpm}} s_t(Y_t)}_{\text{deterministic}} + \underbrace{\sigma_t^{\text{ddpm}} \mathcal{N}(0, I_d)}_{\text{stochastic}} \right), \quad t = T, \dots, 1$$

# DDPM vs. DDIM

---

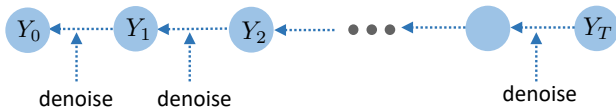


— Song, Meng, Ermon '20

— Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole '20

2. A deterministic sampler: **denoising diffusion implicit models**  
DDIM; or probability flow ODE

# DDPM vs. DDIM



— Song, Meng, Ermon '20

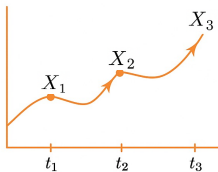
— Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole '20

2. A deterministic sampler: **denoising diffusion implicit models**  
DDIM; or probability flow ODE

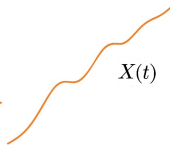
$$Y_T \sim \mathcal{N}(0, I_d)$$
$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \underbrace{Y_t + \eta_t^{\text{ddim}} s_t(Y_t)}_{\text{deterministic}} \right), \quad t = T, \dots, 1$$

# Interpretations from lens of SDE/ODE

---

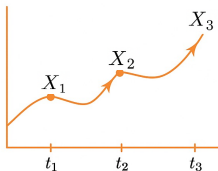


discrete-time

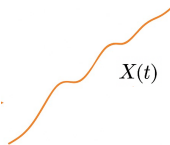


continuous-time

# Interpretations from lens of SDE/ODE



discrete-time

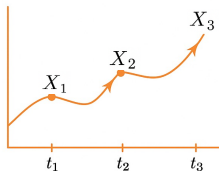


continuous-time

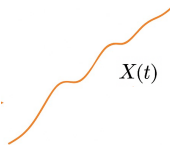
forward process

$$\begin{aligned} X_t &= \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d) \\ \Rightarrow dX_t &= -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE}) \end{aligned}$$

# Interpretations from lens of SDE/ODE



discrete-time



continuous-time

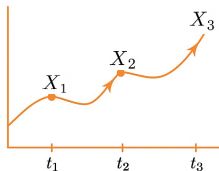
forward process

$$\begin{aligned} X_t &= \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d) \\ \Rightarrow dX_t &= -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE}) \end{aligned}$$

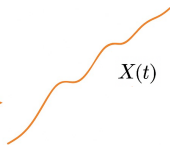
- $\exists$  reverse-time SDE w/ same path distribution ([Anderson '82](#))

$$dY_t = (Y_t + 2s_{T-t}^*(Y_t)) \beta(T-t) dt + \sqrt{2\beta(T-t)} dW_t$$

# Interpretations from lens of SDE/ODE



discrete-time



continuous-time

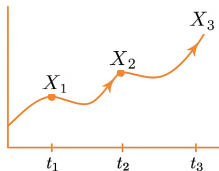
forward process

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d)$$
$$\Rightarrow dX_t = -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE})$$

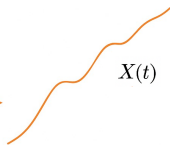
- $\exists$  reverse-time SDE w/ same path distribution ([Anderson '82](#))

time discretization  $\longrightarrow$  DDPM

# Interpretations from lens of SDE/ODE



discrete-time



continuous-time

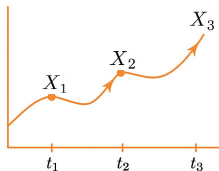
forward process

$$\begin{aligned} X_t &= \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d) \\ \Rightarrow dX_t &= -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE}) \end{aligned}$$

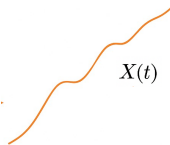
- $\exists$  reverse-time ODE w/ same *marginal* dist (Song et al. '20)

$$dY_t = (Y_t + s_{T-t}^*(Y_t)) \beta(T - t) dt$$

# Interpretations from lens of SDE/ODE



discrete-time



continuous-time

forward process

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d)$$
$$\Rightarrow dX_t = -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE})$$

- $\exists$  reverse-time ODE w/ same *marginal* dist (Song et al. '20)

time discretization  $\longrightarrow$  DDIM

**Key takeaway:** in continuous-time limits, sampling is feasible once exact score functions are available

— *almost no restriction on target data distributions*

**Key takeaway:** in continuous-time limits, sampling is feasible once exact score functions are available

— *almost no restriction on target data distributions*

**Questions:**

- what happens in discrete time? — effect of discretization error
- what if we only have imperfect scores? — effect of score error

# An incomplete list of prior theory

---

- Lee, Lu, Tan '22
- Chen, Chewi, Li, Li, Salim, Zhang '22
- Chen, Lee, Lu '22
- Lee, Lu, Tan '23
- Chen, Daras, Dimakis '23
- Chen, Chewi, Lee, Li, Lu, Salim '23
- Benton, De Bortoli, Doucet, Deligiannidis '23
- Li, Wei, Chen, Chi '23
- Benton, Deligiannidis, Doucet '23
- Cheng, Lu, Tan, Xie '23
- Tang '23
- Li, Wei, Chi, Chen '24
- Li, Yan '24a, '24b
- Azangulov, Deligiannidis, Rousseau '24
- Potaptchik, Azangulov, Deligiannidis '24
- Huang, Wei, Chen '24
- Gao, Zhu '24
- Huang, Huang, Lin '24
- Li, Jiao '24
- Li, Di, Gu '24
- Liang, Ju, Liang, Shroff '24
- Tang, Zhao '24
- Liang, Huang, Chen '25
- Li, Cai, Wei '25
- Tang, Yan '25
- Yu, Yu '25
- Gentiloni-Silveri, Ocello '25
- Li, Zhou, Wei, Chen '25
- Jain, Zhang '25
- Jiao, Zhou, Li '25
- ...

# Typical assumptions

---

- log-concavity of target distribution  $p_{\text{data}}$
- smoothness of score function  $s_t^* = \nabla \log p_t(x)$
- $\ell_\infty / \ell_2$  loss of the score estimate

## Our assumptions: target distribution $p_{\text{data}}$

---

$\mathbb{P}(\|X_0\|_2 \leq T^{c_R}) = 1$  for arbitrarily large const  $c_R > 0$

## Our assumptions: target distribution $p_{\text{data}}$

---

$$\mathbb{P}(\|X_0\|_2 \leq T^{c_R}) = 1 \text{ for arbitrarily large const } c_R > 0$$

- support size can be very large

## Our assumptions: target distribution $p_{\text{data}}$

---

$$\mathbb{P}(\|X_0\|_2 \leq T^{c_R}) = 1 \text{ for arbitrarily large const } c_R > 0$$

- support size can be very large
- very general: *no need of log-concavity, smoothness, etc*

## Our assumptions: target distribution $p_{\text{data}}$

---

$$\mathbb{P}(\|X_0\|_2 \leq T^{c_R}) = 1 \text{ for arbitrarily large const } c_R > 0$$

- support size can be very large
- very general: *no need of log-concavity, smoothness, etc*
- can also be replaced by  $\mathbb{E}[\|X_0\|_2] \leq T^{c_M}$  for large const  $c_M$

## Our assumptions: score estimates $\{s_t(\cdot)\}$

---

- $\ell_2$  score estimation error:  $s_t^*(X) := \nabla \log p_{X_t}(X)$ ,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[ \|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2$$

## Our assumptions: score estimates $\{s_t(\cdot)\}$

---

- $\ell_2$  score estimation error:  $s_t^*(X) := \nabla \log p_{X_t}(X)$ ,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[ \|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2$$

- much weaker than *pointwise* score error assumption

## Our assumptions: score estimates $\{s_t(\cdot)\}$

---

- $\ell_2$  score estimation error:  $s_t^*(X) := \nabla \log p_{X_t}(X)$ ,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[ \|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2$$

- much weaker than *pointwise* score error assumption
- suffices for DDPM **but not DDIM** (counterexample in [Li et al. '24](#))

## Our assumptions: score estimates $\{s_t(\cdot)\}$

---

- $\ell_2$  score estimation error:  $s_t^*(X) := \nabla \log p_{X_t}(X)$ ,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[ \|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2$$

- much weaker than *pointwise* score error assumption
  - suffices for DDPM **but not DDIM** (counterexample in [Li et al. '24](#))
- *Jacobian estimation error (for DDIM only)*:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[ \left\| \frac{\partial s_t}{\partial x}(X) - \frac{\partial s_t^*}{\partial x}(X) \right\| \right] \leq \varepsilon_{\text{Jacobi}}$$

# Assumptions: learning rates

---

$$X_0 \sim p_{\text{data}}, \quad X_t = \sqrt{\alpha_t} X_{t-1} + \sqrt{1 - \alpha_t} \mathcal{N}(0, I_d)$$

- **learning rates:** for some consts  $c_0, c_1 > 0$ ,

$$1 - \alpha_1 = \frac{1}{T^{c_0}}$$

$$1 - \alpha_t = \frac{c_1 \log T}{T} \min \left\{ (1 - \alpha_1) \left( 1 + \frac{c_1 \log T}{T} \right)^t, 1 \right\}$$

2 phases: exp growth  $\rightarrow$  flat

## A glimpse of convergence theory (up to log factor)

---

**DDPM:**

$$\begin{aligned} \text{KL}(p_{X_1} \parallel p_{Y_1}) &\lesssim d^2/T^2 + \varepsilon_{\text{score}}^2 && \text{(Jiao, Zhou, Li '25)} \\ \text{TV}(p_{X_1}, p_{Y_1}) &\lesssim d/T + \varepsilon_{\text{score}} && \text{(Li, Yan '24)} \end{aligned}$$

# A glimpse of convergence theory (up to log factor)

---

DDPM: 
$$\begin{aligned} \text{KL}(p_{X_1} \parallel p_{Y_1}) &\lesssim d^2/T^2 + \varepsilon_{\text{score}}^2 && \text{(Jiao, Zhou, Li '25)} \\ \text{TV}(p_{X_1}, p_{Y_1}) &\lesssim d/T + \varepsilon_{\text{score}} && \text{(Li, Yan '24)} \end{aligned}$$

- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{KL} \leq \varepsilon^2$
- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{TV} \leq \varepsilon$

# A glimpse of convergence theory (up to log factor)

---

DDPM: 
$$\begin{aligned} \text{KL}(p_{X_1} \parallel p_{Y_1}) &\lesssim d^2/T^2 + \varepsilon_{\text{score}}^2 && \text{(Jiao, Zhou, Li '25)} \\ \text{TV}(p_{X_1}, p_{Y_1}) &\lesssim d/T + \varepsilon_{\text{score}} && \text{(Li, Yan '24)} \end{aligned}$$

- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{KL} \leq \varepsilon^2$

- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{TV} \leq \varepsilon$

- Pinsker inequality ( $\text{TV} \leq \sqrt{\frac{1}{2}\text{KL}}$ ) is loose when bounding TV

# A glimpse of convergence theory (up to log factor)

---

DDPM: 
$$\begin{aligned} \text{KL}(p_{X_1} \parallel p_{Y_1}) &\lesssim d^2/T^2 + \varepsilon_{\text{score}}^2 && \text{(Jiao, Zhou, Li '25)} \\ \text{TV}(p_{X_1}, p_{Y_1}) &\lesssim d/T + \varepsilon_{\text{score}} && \text{(Li, Yan '24)} \end{aligned}$$

- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{KL} \leq \varepsilon^2$
- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{TV} \leq \varepsilon$
- **stability:** degrades gracefully as  $\varepsilon_{\text{score}} \uparrow$

# A glimpse of convergence theory (up to log factor)

---

DDPM: 
$$\begin{aligned} \text{KL}(p_{X_1} \parallel p_{Y_1}) &\lesssim d^2/T^2 + \varepsilon_{\text{score}}^2 && \text{(Jiao, Zhou, Li '25)} \\ \text{TV}(p_{X_1}, p_{Y_1}) &\lesssim d/T + \varepsilon_{\text{score}} && \text{(Li, Yan '24)} \end{aligned}$$

- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{KL} \leq \varepsilon^2$
- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{TV} \leq \varepsilon$
- **stability:** degrades gracefully as  $\varepsilon_{\text{score}} \uparrow$  and  $\varepsilon_{\text{Jacobi}} \uparrow$
- **general data distribution:** no need of smoothness, log-concavity

# A glimpse of convergence theory (up to log factor)

---

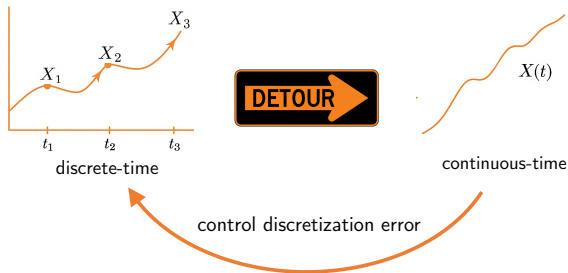
**DDPM:**  $\text{KL}(p_{X_1} \| p_{Y_1}) \lesssim d^2/T^2 + \varepsilon_{\text{score}}^2$  (Jiao, Zhou, Li '25)  
 $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}$  (Li, Yan '24)

**DDIM:**  $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \sqrt{d}\varepsilon_{\text{score}} + d\varepsilon_{\text{Jacobi}}$  (Li et al. '24)

- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{KL} \leq \varepsilon^2$
- **iteration complexity:**  $d/\varepsilon$  (assuming accurate scores)  
to yield  $\text{TV} \leq \varepsilon$
- **stability:** degrades gracefully as  $\varepsilon_{\text{score}} \uparrow$  and  $\varepsilon_{\text{Jacobi}} \uparrow$
- **general data distribution:** no need of smoothness, log-concavity

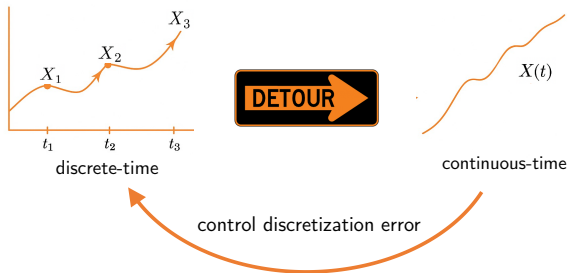
# Analysis strategy # 1 (for DDPM)

- *Chen, Chewi, Li, Li, Salim, Zhang '22, Chen, Lee, Lu '22, Tang, Zhao '24*
- *Benton, De Bortoli, Doucet, Deligiannidis '23, Huang, Wei, Chen '24*



# Analysis strategy # 1 (for DDPM)

- *Chen, Chewi, Li, Li, Salim, Zhang '22, Chen, Lee, Lu '22, Tang, Zhao '24*
- *Benton, De Bortoli, Doucet, Deligiannidis '23, Huang, Wei, Chen '24*



*Analogy: (stochastic) gradient descent vs. gradient flow, TD learning via ODE*

# Analysis strategy # 1 (for DDPM)

- *Chen, Chewi, Li, Li, Salim, Zhang '22, Chen, Lee, Lu '22, Tang, Zhao '24*  
— *Benton, De Bortoli, Doucet, Deligiannidis '23, Huang, Wei, Chen '24*



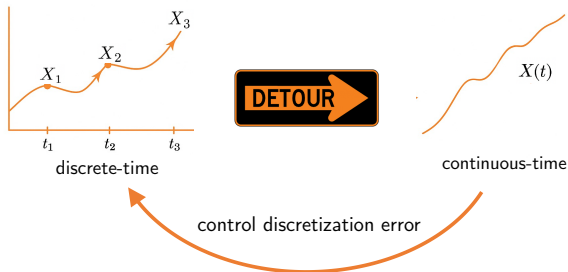
2 key steps:

- apply change of measure (e.g. Girsanov thm) to show

$$\text{KL}(P^{\text{true}} \parallel P^{\text{ddpm}}) \leq \int w(t) \underbrace{\mathbb{E} \left[ \left\| \text{drift}^{\text{true}}(t) - \text{drift}^{\text{ddpm}}(t) \right\|^2 \right]}_{\text{score error} + \text{discretization error}} dt + \text{small-term}$$

# Analysis strategy # 1 (for DDPM)

- *Chen, Chewi, Li, Li, Salim, Zhang '22, Chen, Lee, Lu '22, Tang, Zhao '24*  
— *Benton, De Bortoli, Doucet, Deligiannidis '23, Huang, Wei, Chen '24*



2 key steps:

- leverage stochastic localization to characterize

$$\text{discretization error} \xleftrightarrow{\text{link}} \mathbb{E}[\text{Cov}(X_0 | X_t)]$$

## Analysis strategy # 2 (for DDIM & DDPM)

---

— *Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24*

— *Li, Yan '24, Liang, Huang, Chen '25, Jiao, Zhou, Li '25*

Tackle discrete-time process directly & track changes of, e.g., TV distance

## Analysis strategy # 2 (for DDIM & DDPM)

---

— *Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24*

— *Li, Yan '24, Liang, Huang, Chen '25, Jiao, Zhou, Li '25*

Tackle discrete-time process directly & track changes of, e.g., TV distance

yields state-of-the-art theory for DDIM & DDPM!

## Analysis strategy # 2 (for DDIM & DDPM)

---

— Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24

— Li, Yan '24, Liang, Huang, Chen '25, Jiao, Zhou, Li '25

Tackle discrete-time process directly & track changes of, e.g., TV distance

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I), \quad Y_{t-1} = \underbrace{\frac{1}{\sqrt{1 - \beta_t}} Y_t + \frac{\beta_t}{2\sqrt{1 - \beta_t}} s_t(Y_t)}_{=: \Phi_t(Y_t)}$$

## Analysis strategy # 2 (for DDIM & DDPM)

---

— Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24

— Li, Yan '24, Liang, Huang, Chen '25, Jiao, Zhou, Li '25

Tackle discrete-time process directly & track changes of, e.g., TV distance

$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0 \quad \iff \quad \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)} \approx 1 \quad \forall y_t \in \mathcal{E}_t \text{ (some "typical" set)}$$

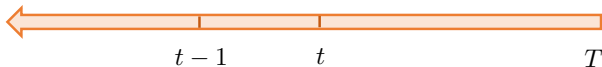
# Analysis strategy # 2 (for DDIM & DDPM)

— Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24

— Li, Yan '24, Liang, Huang, Chen '25, Jiao, Zhou, Li '25

Tackle discrete-time process directly & track changes of, e.g., TV distance

$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0 \iff \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)} \approx 1 \quad \forall y_t \in \mathcal{E}_t \text{ (some "typical" set)}$$



$$\frac{p_{Y_{t-1}}(\Phi_t(y_t))}{p_{X_{t-1}}(\Phi_t(y_t))} = \underbrace{\frac{p_{Y_{t-1}}(\Phi_t(y_t))}{p_{Y_t}(y_t)}}_{\text{relation btw } Y_t \text{ \& } Y_{t-1}} \left( \underbrace{\frac{p_{X_{t-1}}(\Phi_t(y_t))}{p_{X_t}(y_t)}}_{\text{relation btw } X_t \text{ \& } X_{t-1}} \right)^{-1} \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)}$$

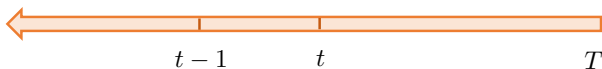
# Analysis strategy # 2 (for DDIM & DDPM)

— Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24

— Li, Yan '24, Liang, Huang, Chen '25, Jiao, Zhou, Li '25

Tackle discrete-time process directly & track changes of, e.g., TV distance

$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0 \iff \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)} \approx 1 \quad \forall y_t \in \mathcal{E}_t \text{ (some "typical" set)}$$



$$\frac{p_{Y_{t-1}}(\Phi_t(y_t))}{p_{X_{t-1}}(\Phi_t(y_t))} = \underbrace{\frac{p_{Y_{t-1}}(\Phi_t(y_t))}{p_{Y_t}(y_t)}}_{\text{relation btw } Y_t \text{ \& } Y_{t-1}} \left( \underbrace{\frac{p_{X_{t-1}}(\Phi_t(y_t))}{p_{X_t}(y_t)}}_{\text{relation btw } X_t \text{ \& } X_{t-1}} \right)^{-1} \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)}$$

$$\frac{p_{\Phi_t(Y_t)}(\Phi_t(y_t))}{p_{Y_t}(y_t)} = \det \left( \frac{\partial \Phi_t}{\partial y_t} \right)^{-1}$$

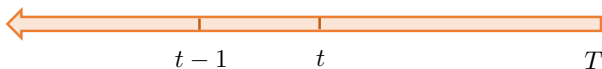
# Analysis strategy # 2 (for DDIM & DDPM)

— Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24

— Li, Yan '24, Liang, Huang, Chen '25, Jiao, Zhou, Li '25

Tackle discrete-time process directly & track changes of, e.g., TV distance

$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0 \iff \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)} \approx 1 \quad \forall y_t \in \mathcal{E}_t \text{ (some "typical" set)}$$

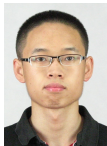


$$\frac{p_{Y_{t-1}}(\Phi_t(y_t))}{p_{X_{t-1}}(\Phi_t(y_t))} = \underbrace{\frac{p_{Y_{t-1}}(\Phi_t(y_t))}{p_{Y_t}(y_t)}}_{\text{relation btw } Y_t \text{ \& } Y_{t-1}} \left( \underbrace{\frac{p_{X_{t-1}}(\Phi_t(y_t))}{p_{X_t}(y_t)}}_{\text{relation btw } X_t \text{ \& } X_{t-1}} \right)^{-1} \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)}$$

$$\frac{p_{\Phi_t(Y_t)}(\Phi_t(y_t))}{p_{Y_t}(y_t)} = \det\left(\frac{\partial \Phi_t}{\partial y_t}\right)^{-1}$$

some concentration bounds

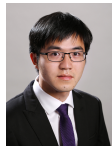
## *Part 2: adaptation to (unknown) low dimensionality*



Gen Li  
CUHK



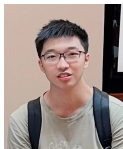
Yuling Yan  
DE Shaw



Changxiao Cai  
UMich



Jiadong Liang  
UPenn



Zhihan Huang  
UPenn



Yuxin Chen  
UPenn



Yuting Wei  
UPenn

# Recap: theory for mainstream diffusion models

## Denoising Diffusion Probabilistic Models

Jonathan Ho  
UC Berkeley  
jonathanho@berkeley.edu

Ajay Jain  
UC Berkeley  
ajayj@berkeley.edu

Pieter Abbeel  
UC Berkeley  
pabbeel@cs.berkeley.edu

## DENOISING DIFFUSION IMPLICIT MODELS

Jiaming Song, Chenlin Meng & Stefano Ermon  
Stanford University  
{tsong, chenlin, ermon}@cs.stanford.edu

## Theorem 1 (Li, Wei, Chi, Chen '24, Li, Yan '24)

With perfect scores, both DDIM & DDPM yield  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in  $\tilde{O}(d/\varepsilon)$  iterations

- $d$ : ambient dimension

$d/\varepsilon$  iterations are too slow ...

---



ImageNet:  $d = 150,528$  pixels per image

$d/\varepsilon$  iterations are too slow ...

---



ImageNet:  $d = 150,528$  pixels per image (so  $\frac{d}{\varepsilon} > 10^6$  for moderate  $\varepsilon$ )

$d/\varepsilon$  iterations are too slow ...

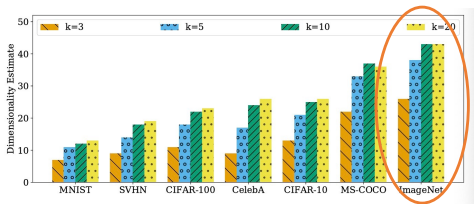
---



ImageNet:  $d = 150,528$  pixels per image (so  $\frac{d}{\varepsilon} > 10^6$  for moderate  $\varepsilon$ )

*In practice, DDIM/DDPM yield good samples in hundreds (or tens) of iterations ...*

# $d/\varepsilon$ iterations are too slow ...



ImageNet:  $d = 150,528$  pixels per image (so  $\frac{d}{\varepsilon} > 10^6$  for moderate  $\varepsilon$ )  
 $k = 43$  intrinsic dimension (Pope et al. '21)

*In practice, DDIM/DDPM yield good samples in hundreds (or tens) of iterations ...*

*Can diffusion models adapt to intrinsic low dimensionality?*

# Intrinsic dimension

---

The target distribution  $p_{\text{data}}$  is said to have **intrinsic dimension**  $k$  if

$$\log \underbrace{N^{\text{cover}}(\text{support}(p_{\text{data}}), \|\cdot\|_2, \varepsilon_0)}_{\text{covering number of support of } p_{\text{data}}} \lesssim k \log \frac{1}{\varepsilon_0}$$

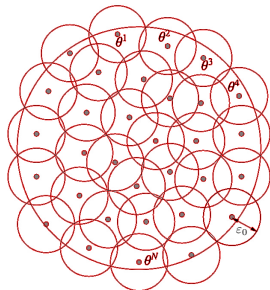
# Intrinsic dimension

---

The target distribution  $p_{\text{data}}$  is said to have **intrinsic dimension**  $k$  if

$$\log \underbrace{N^{\text{cover}}(\text{support}(p_{\text{data}}), \|\cdot\|_2, \varepsilon_0)}_{\text{covering number of support of } p_{\text{data}}} \lesssim k \log \frac{1}{\varepsilon_0}$$

- k-dimensional linear subspaces
- k-dimensional manifolds
- union of the above
- ...



— see [Li, Yan '24](#), [Huang, Wei, Chen '24](#)

# Assumptions

---

- **minimal data assumptions:**

$$\mathbb{P}(\|X_0\|_2 \leq \underbrace{T^{c_R}}_{\text{polynomially large diameter}}) = 1$$

for arbitrarily large constant  $c_R > 0$

- **perfect score estimates:**  $s_t(\cdot) = \nabla \log p_{X_t}(\cdot)$   
→ not needed; only to simplify presentation

# Convergence theory in KL divergence

---

## Theorem 2 (Li and Yan '24; Huang, Wei, Chen '24)

*DDPM sampler (its original form) yields  $\text{KL}(p_{X_1} \parallel p_{Y_1}) \leq \varepsilon$  in*

*$\tilde{O}(k/\varepsilon)$  iterations*

— related work [Azangulov et al. '24](#), [Potapchik et al.'24](#)

# Convergence theory in KL divergence

---

## Theorem 2 (Li and Yan '24; Huang, Wei, Chen '24)

DDPM sampler (its original form) yields  $\text{KL}(p_{X_1} \parallel p_{Y_1}) \leq \varepsilon$  in

$\tilde{O}(k/\varepsilon)$  iterations

— related work [Azangulov et al. '24](#), [Potapchik et al.'24](#)

- optimal scaling in  $k$

# Convergence theory in total variation

## Theorem 3 (Liang, Huang, Chen '25; Li, Yan '25; Tang, Yan '25)

Both DDPM & DDIM (their original form) yield  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in  $\tilde{O}(k/\varepsilon)$  iterations

$$\text{DDIM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( Y_t + \eta_t^{\text{ddim}} s_t(Y_t) \right)$$

$$\text{DDPM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( Y_t + \eta_t^{\text{ddpm}} s_t(Y_t) + \sigma_t^{\text{ddpm}} \mathcal{N}(0, I_d) \right)$$

# Convergence theory in total variation

## Theorem 3 (Liang, Huang, Chen '25; Li, Yan '25; Tang, Yan '25)

Both DDPM & DDIM (their original form) yield  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in  $\tilde{O}(k/\varepsilon)$  iterations

$$\text{DDIM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( Y_t + \frac{1 - \alpha_t}{1 + \sqrt{\frac{\alpha_t - \bar{\alpha}_t}{1 - \bar{\alpha}_t}}} s_t(Y_t) \right)$$

$$\text{DDPM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( Y_t + (1 - \alpha_t) s_t(Y_t) + \sqrt{\frac{(1 - \alpha_t)(\alpha_t - \bar{\alpha}_t)}{1 - \bar{\alpha}_t}} \mathcal{N}(0, I_d) \right)$$

where  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

# Convergence theory in total variation

## Theorem 3 (Liang, Huang, Chen '25; Li, Yan '25; Tang, Yan '25)

Both DDPM & DDIM (their original form) yield  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in  $\tilde{O}(k/\varepsilon)$  iterations

$$\text{DDIM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( Y_t + \frac{1 - \alpha_t}{1 + \sqrt{\frac{\alpha_t - \bar{\alpha}_t}{1 - \bar{\alpha}_t}}} s_t(Y_t) \right)$$

$$\text{DDPM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( Y_t + (1 - \alpha_t) s_t(Y_t) + \sqrt{\frac{(1 - \alpha_t)(\alpha_t - \bar{\alpha}_t)}{1 - \bar{\alpha}_t}} \mathcal{N}(0, I_d) \right)$$

where  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

- originally derived to optimize variational lower bounds!

# Interpretation from lens of SDE/ODE

---

reverse-time SDE (same distribution as  $X_t$ ):

$$dY_t = (Y_t + 2s_{T-t}^*(Y_t))dt + \sqrt{2}dB_t$$

# Interpretation from lens of SDE/ODE

---

reverse-time SDE (same distribution as  $X_t$ ):

$$dY_t = (Y_t + 2s_{T-t}^*(Y_t))dt + \sqrt{2} dB_t$$

Tweedie's formula  $\Downarrow$

$$dY_t = \left( \underbrace{c_{1,t}Y_t}_{\text{linear drift}} + c_{2,t} \underbrace{\mathbb{E}[X_0 \mid X_{T-t} = Y_t]}_{\text{cond. mean of } X_0} \right) dt + \sqrt{2} dB_t$$

- **key enabler:**  $\mathbb{E}[X_0 \mid X_t]$  is “projection” onto low-dimensional structure

# Interpretation from lens of SDE/ODE

---

reverse-time SDE (same distribution as  $X_t$ ):

$$dY_t = (Y_t + 2s_{T-t}^*(Y_t))dt + \sqrt{2} dB_t$$

Tweedie's formula  $\Downarrow$

$$dY_t = \left( \underbrace{c_{1,t}Y_t}_{\text{linear drift}} + c_{2,t} \underbrace{\mathbb{E}[X_0 \mid X_{T-t} = Y_t]}_{\text{cond. mean of } X_0} \right) dt + \sqrt{2} dB_t$$

time-discretize  $\Downarrow$

DDIM or DDPM

- **key enabler:**  $\mathbb{E}[X_0 \mid X_t]$  is “projection” onto low-dimensional structure
- **discretization scheme matters:** time-discretize carefully to retain low-dimensional adaptation

Can diffusion models adapt to other structures,  
e.g. *Gaussian mixture models*?



figure credit: Dall-E 3 from OpenAI

# An incomplete list of prior art

---

Gaussian mixture models ([Pearson'94](#))

$$X_0 \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma), \quad \sum_{i=1}^K \pi_i = 1.$$

# An incomplete list of prior art

---

Gaussian mixture models ([Pearson'94](#))

$$X_0 \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma), \quad \sum_{i=1}^K \pi_i = 1.$$

- Dasgupta'99
- Vempala & Wang'04
- Arora & Kannan'05
- Kalai et al.'10
- Moitra & Valiant'10
- Hsu & Kakade'13
- Diakonikolas et al.'18
- Hopkins & Li'18
- Ghosal & Van Der Vaart'01
- Dwivedi, Wainwright, et al.'20
- Chen'95
- Heinrich & Kahn'18
- Wu & Yang'20
- Doss et al.'23
- Saha & Guntuboyina'20
- Ashtiani et al.'18
- [Shah et al.'23](#)
- [Liang et al.'24](#)
- [Wu et al.'24](#)
- [Chidambaram et al.'24](#)
- [Chen et al.'24](#)
- [Gatmiry et al.'24](#)
- [Wang et al.'24](#)

# Dimension-free convergence of DDPM for GMMs

Gaussian mixture models (Pearson'94)

$$X_0 \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma), \quad \sum_{i=1}^K \pi_i = 1$$

## Theorem 4 (Li, Cai, Wei '25)

For spherical Gaussian mixture with  $\Sigma = \sigma^2 I_d$ , DDPM yields

$$\text{TV}(X_0, Y_0) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \varepsilon_{\text{score}} \sqrt{\log T}$$

# Dimension-free convergence of DDPM for GMMs

Gaussian mixture models (Pearson'94)

$$X_0 \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma), \quad \sum_{i=1}^K \pi_i = 1$$

## Theorem 4 (Li, Cai, Wei '25)

For spherical Gaussian mixture with  $\Sigma = \sigma^2 I_d$ , DDPM yields

$$\text{TV}(X_0, Y_0) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \varepsilon_{\text{score}} \sqrt{\log T}$$

- to yield  $\varepsilon$ -accuracy, it requires  $\underbrace{\tilde{O}(1/\varepsilon)}_{\text{dimension-free}}$  iterations

# Dimension-free convergence of DDPM for GMMs

Gaussian mixture models (Pearson'94)

$$X_0 \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma), \quad \sum_{i=1}^K \pi_i = 1$$

## Theorem 4 (Li, Cai, Wei '25)

For spherical Gaussian mixture with  $\Sigma = \sigma^2 I_d$ , DDPM yields

$$\text{TV}(X_0, Y_0) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \epsilon_{\text{score}} \sqrt{\log T}$$

- to yield  $\epsilon$ -accuracy, it requires  $\underbrace{\tilde{O}(1/\epsilon)}_{\text{dimension-free}}$  iterations
- stable vis-a-vis score errors

# Dimension-free convergence of DDPM for GMMs

Gaussian mixture models (Pearson'94)

$$X_0 \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma), \quad \sum_{i=1}^K \pi_i = 1$$

## Theorem 4 (Li, Cai, Wei '25)

For spherical Gaussian mixture with  $\Sigma = \sigma^2 I_d$ , DDPM yields

$$\text{TV}(X_0, Y_0) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \varepsilon_{\text{score}} \sqrt{\log T}$$

*Even in ultra-high-dimension, diffusion models are highly effective in sampling GMMs!*

## Comparison w/ prior theory of GMMs

---

$$\text{(our theory)} \quad \text{TV}(X_0, Y_0) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \varepsilon_{\text{score}} \sqrt{\log T}.$$

# Comparison w/ prior theory of GMMs

---

$$\text{(our theory)} \quad \text{TV}(X_0, Y_0) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \varepsilon_{\text{score}} \sqrt{\log T}.$$

- Liang, Shi, Song, Zhou'24:  $\underbrace{\tilde{O}(d/\varepsilon^2)}_{\text{show no adaptation phenomenon}}$  complexity bound

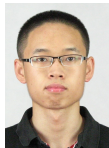
# Comparison w/ prior theory of GMMs

---

$$\text{(our theory)} \quad \text{TV}(X_0, Y_0) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \varepsilon_{\text{score}} \sqrt{\log T}.$$

- Liang, Shi, Song, Zhou'24:  $\underbrace{\tilde{O}(d/\varepsilon^2)}_{\text{show no adaptation phenomenon}}$  complexity bound
- Chen et al.'24; Gatmiry et al.'24: focus on score estimation using piecewise polynomial regression + existing DDPM convergence theory

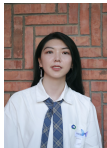
## *Part 3: acceleration via higher-order approximation*



Gen Li  
CUHK



Yuchen Zhou  
UIUC



Yu Huang  
UPenn



Timofey Efimov  
CMU



Yuchen Wu  
Cornell



Yuxin Chen  
UPenn



Yuejie Chi  
Yale



Yuting Wei  
UPenn

# DDPM and DDIM are still slow ...

---

Low sampling speed!

100s-1000s steps



initialize  
at pure  
Gaussian

— *Song, Meng, Ermon '20*

# DDPM and DDIM are still slow ...

---

Low sampling speed!

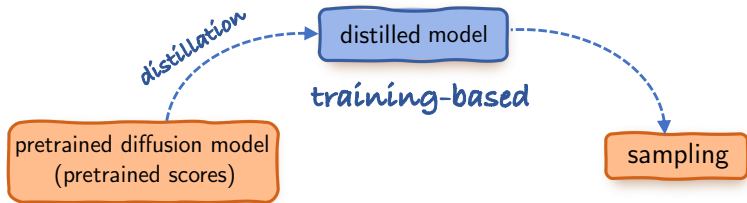
100s-1000s steps



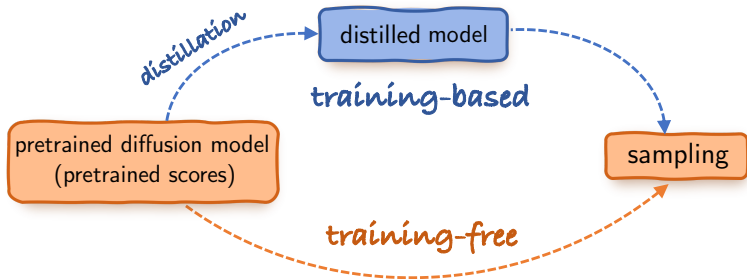
initialize  
at pure  
Gaussian

50K  $32 \times 32$  images: DDPM (20h) vs. single-step GANs (< 1min)

— *Song, Meng, Ermon '20*



- **training-based:** distill pre-trained diffusion model into another  
requires additional training  
model that can be executed rapidly
  - e.g., progressive distillation, consistency model



- **training-free:** directly invoke pre-trained score estimates for sampling w/o additional training
  - e.g., DPM-Solver/++ (Lu et al. '22), UniPC (Zhao et al. '23), ...

Can we design a *training-free* sampler that is provably faster than DDIM/DDPM?

# Discretization $\longleftrightarrow$ approximation

---

**A starting point:** equiv solution to probability flow ODE

$$\underbrace{Y_{\bar{\alpha}_{t-1}}^{\text{ode}}}_{\text{represent } Y_{t-1}} = \frac{1}{\sqrt{\alpha_t}} \underbrace{Y_{\bar{\alpha}_t}^{\text{ode}}}_{\text{represent } Y_t} + \underbrace{\int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) d\gamma}_{}$$

where  $s_{\gamma}^*(x) := \nabla \log p_{\sqrt{\gamma}X_0 + \sqrt{1-\gamma}\mathcal{N}(0, I_d)}(x)$

# Discretization $\longleftrightarrow$ approximation

---

**A starting point:** equiv solution to probability flow ODE

$$\underbrace{Y_{\bar{\alpha}_{t-1}}^{\text{ode}}}_{\text{represent } Y_{t-1}} = \frac{1}{\sqrt{\alpha_t}} \underbrace{Y_{\bar{\alpha}_t}^{\text{ode}}}_{\text{represent } Y_t} + \underbrace{\int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) d\gamma}_{\text{infinitely many score evaluations!}}$$

where  $s_{\gamma}^*(x) := \nabla \log p_{\sqrt{\gamma}X_0 + \sqrt{1-\gamma}\mathcal{N}(0, I_d)}(x)$

# Discretization $\longleftrightarrow$ approximation

---

**A starting point:** equiv solution to probability flow ODE

$$\underbrace{Y_{\bar{\alpha}_{t-1}}^{\text{ode}}}_{\text{represent } Y_{t-1}} = \frac{1}{\sqrt{\alpha_t}} \underbrace{Y_{\bar{\alpha}_t}^{\text{ode}}}_{\text{represent } Y_t} + \underbrace{\int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) d\gamma}_{\text{infinitely many score evaluations!}}$$

where  $s_{\gamma}^*(x) := \nabla \log p_{\sqrt{\gamma}X_0 + \sqrt{1-\gamma}\mathcal{N}(0, I_d)}(x)$

- can we approximate the integral by a few score evals?

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\alpha_t}^*(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\alpha_t}^*(Y_{\alpha_t}^{\text{ode}}) \approx s_t(Y_t)$

$$\Rightarrow Y_{t-1} \approx \frac{1}{\sqrt{\alpha_t}} \left( Y_t + \frac{1 - \alpha_t}{2} s_t(Y_t) \right) \quad \text{original DDIM}$$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\bar{\alpha}_t}^*(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\epsilon}$  iterations; 1 score eval per iteration (DDIM)

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\bar{\alpha}_t}^*(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\epsilon}$  iterations; 1 score eval per iteration (DDIM)

**refined approximation?**

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\bar{\alpha}_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\bar{\alpha}_t}^*(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\varepsilon}$  iterations; 1 score eval per iteration (DDIM)

**refined approximation?**

$$s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_t(Y_t) + \frac{\gamma - \bar{\alpha}_t}{\bar{\alpha}_t - \bar{\alpha}_{t+1}} \left( s_t(Y_t) - s_{t+1}(Y_{t+1}) \right)$$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\alpha_t}^*(Y_{\alpha_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\epsilon}$  iterations; 1 score eval per iteration (DDIM)

**2nd order approx:** (Li, Huang, Efimov, Wei, Chi, Chen '24)

$$\sqrt{\alpha_t} Y_{t-1} \approx Y_t + \frac{1 - \alpha_t}{2} s_t(Y_t) + \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} (s_t(Y_t) - \sqrt{\alpha_{t+1}} s_{t+1}(Y_{t+1}))$$

— similar in spirit to DPM-Solver-2 (Lu et al '22)

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\alpha_t}^*(Y_{\alpha_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\epsilon}$  iterations; 1 score eval per iteration (DDIM)

**2nd order approx:** (Li, Huang, Efimov, Wei, Chi, Chen '24)

$\frac{\text{poly}(d)}{\sqrt{\epsilon}}$  iterations; 2 score evals per iteration

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\alpha_t}^*(Y_{\alpha_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\epsilon}$  iterations; 1 score eval per iteration (DDIM)

**even higher-order approximation?**

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\bar{\alpha}_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\alpha_t}^*(Y_{\alpha_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\epsilon}$  iterations; 1 score eval per iteration (DDIM)

**even higher-order approximation?** for order  $K$ :

$$\frac{1}{\gamma^{3/2}} s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx \sum_{0 \leq i < K} \psi_i(\gamma) \frac{s_{\gamma_{t,i}}^*(Y_{\gamma_{t,i}}^{\text{ode}})}{(\gamma_{t,i})^{3/2}}$$

- $K$  anchor points:  $\gamma_{t,0}, \dots, \gamma_{t,K-1}$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$

**1st order approx:**  $s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx s_{\alpha_t}^*(Y_{\alpha_t}^{\text{ode}}) \approx s_t(Y_t)$

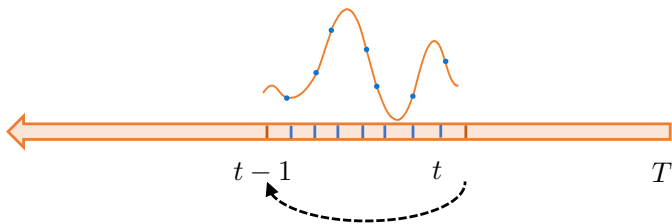
$\frac{d}{\epsilon}$  iterations; 1 score eval per iteration (DDIM)

**even higher-order approximation?** for order  $K$ :

$$\frac{1}{\gamma^{3/2}} s_{\gamma}^*(Y_{\gamma}^{\text{ode}}) \approx \sum_{0 \leq i < K} \psi_i(\gamma) \frac{s_{\gamma_{t,i}}^*(Y_{\gamma_{t,i}}^{\text{ode}})}{(\gamma_{t,i})^{3/2}}$$

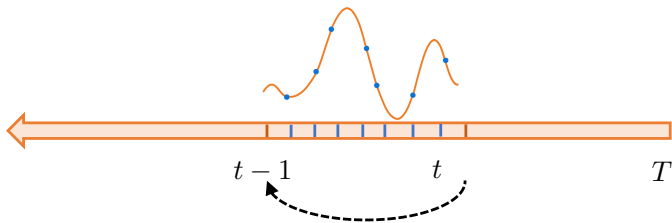
- $K$  anchor points:  $\gamma_{t,0}, \dots, \gamma_{t,K-1}$
- Lagrange basis polynomial:  $\psi_i(\gamma) := \frac{\prod_{i': i' \neq i} (\gamma - \gamma_{t,i'})}{\prod_{i': i' \neq i} (\gamma_{t,i} - \gamma_{t,i'})}$

# Proposed $K$ -th order sampler (Li et al. '25)



$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approx by deg-}(K-1) \text{ Lagrange polynomials}} d\gamma$$

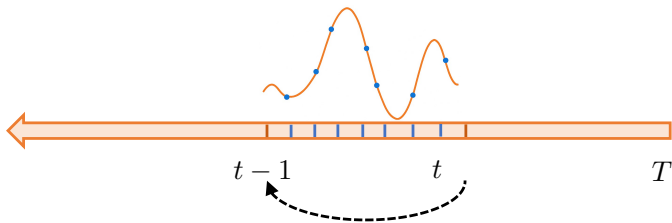
# Proposed $K$ -th order sampler (Li et al. '25)



$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approx by deg-}(K-1) \text{ Lagrange polynomials}} d\gamma$$

- successively, alternately refine  $Y_{\gamma_{t,i}}^{\text{ode}}$  and  $s_{\gamma_{t,i}}(Y_{\gamma_{t,i}}^{\text{ode}})$

## Proposed $K$ -th order sampler (Li et al. '25)



$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approx by deg-}(K-1) \text{ Lagrange polynomials}} d\gamma$$

- successively, alternately refine  $Y_{\gamma_{t,i}}^{\text{ode}}$  and  $s_{\gamma_{t,i}}(Y_{\gamma_{t,i}}^{\text{ode}})$

$K$  score evals per iteration;  $\tilde{O}(1)$  rounds of refinements

# Convergence theory for our accelerated sampler

---

## Theorem 5 (Li, Zhou, Wei, Chen '25)

Consider any  $K = O(1)$ . With perfect scores, our accelerated deterministic sampler yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in

$$\tilde{O}(d^{1+2/K} / \varepsilon^{1/K}) \text{ iterations}$$

# Convergence theory for our accelerated sampler

---

## Theorem 5 (Li, Zhou, Wei, Chen '25)

Consider any  $K = O(1)$ . With perfect scores, our accelerated deterministic sampler yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in

$$\tilde{O}(d^{1+2/K} / \varepsilon^{1/K}) \text{ iterations}$$

- # score function evaluations:  $\frac{d^{1+o(1)}}{\varepsilon^{1/K}}$

# Convergence theory for our accelerated sampler

## Theorem 5 (Li, Zhou, Wei, Chen '25)

Consider any  $K = O(1)$ . With perfect scores, our accelerated deterministic sampler yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in

$$\tilde{O}(d^{1+2/K} / \varepsilon^{1/K}) \text{ iterations}$$

- **# score function evaluations:**  $\frac{d^{1+o(1)}}{\varepsilon^{1/K}}$
- outperforms vanilla DDIM ( $d/\varepsilon$ )
  - substantially improved  $\varepsilon$ -dependency
  - almost no loss in  $d$ -dependency;

# Convergence theory for our accelerated sampler

## Theorem 5 (Li, Zhou, Wei, Chen '25)

Consider any  $K = O(1)$ . With perfect scores, our accelerated deterministic sampler yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in

$$\tilde{O}(d^{1+2/K} / \varepsilon^{1/K}) \text{ iterations}$$

- **# score function evaluations:**  $\frac{d^{1+o(1)}}{\varepsilon^{1/K}}$
- outperforms vanilla DDIM ( $d/\varepsilon$ )
  - substantially improved  $\varepsilon$ -dependency
  - almost no loss in  $d$ -dependency;
- **minimal assumptions** on data distributions
  - see also [Huang et al. '24, '25](#) (Runge-Kutta; stronger assumptions)

*Can we design algorithms to improve the dependency on dimension  $d$ ?*

## Non-uniform Lipschitz property

---

Non-uniform Lipschitz constant  $L \geq 1$ : for every  $\gamma \in (0, 1)$ ,

$$\mathbb{P} \left( (1 - \gamma) \|\nabla s_\gamma^*(X_\gamma)\|_2 \leq L \right) \geq 1 - \frac{1}{d^4}.$$

# Non-uniform Lipschitz property

---

Non-uniform Lipschitz constant  $L \geq 1$ : for every  $\gamma \in (0, 1)$ ,

$$\mathbb{P}((1 - \gamma)\|\nabla s_\gamma^*(X_\gamma)\|_2 \leq L) \geq 1 - \frac{1}{d^4}.$$

- much weaker than the global  $\|\nabla s_\gamma^*(x)\|_2 \leq \tilde{L}$  for all  $x \in \mathbb{R}^d$  and  $\tau$

# Non-uniform Lipschitz property

---

Non-uniform Lipschitz constant  $L \geq 1$ : for every  $\gamma \in (0, 1)$ ,

$$\mathbb{P}((1 - \gamma)\|\nabla s_\gamma^*(X_\gamma)\|_2 \leq L) \geq 1 - \frac{1}{d^4}.$$

- much weaker than the global  $\|\nabla s_\gamma^*(x)\|_2 \leq \tilde{L}$  for all  $x \in \mathbb{R}^d$  and  $\tau$
- Examples that  $L$  is small but  $\tilde{L}$  is large
  - $X_0 \sim N(\mu, \Sigma) \Rightarrow L = 1, \tilde{L} = \infty$  for singular  $\Sigma$
  - $X_0 \sim \sum_{h=1}^H \pi_h \mathcal{N}(\mu_h, \sigma_h^2 I_d) \Rightarrow L \lesssim \log(d) \log(H), \tilde{L}$  can be  $\Omega(\|\mu\|_2^2) \approx \mathbb{E}[\|X_0\|_2^2]$  even when  $\sigma_h^2 = 1$
  - $X_0$  has indep. entries with  $\mathbb{E}[|X_{0,i}|] \leq d^{c_R} \Rightarrow L \lesssim \log d, \tilde{L}$  can be infinite

# Convergence theory for DDPM

---

## Theorem 6 (Jiao, Zhou, Li '25)

DDPM yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  and  $\text{KL}(p_{X_1} \| p_{Y_1}) \leq \varepsilon^2$  in

$$\frac{\min\{L\sqrt{d}, d\}}{\varepsilon} \text{ iterations}$$

# Convergence theory for DDPM

---

## Theorem 6 (Jiao, Zhou, Li '25)

DDPM yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  and  $\text{KL}(p_{X_1} \| p_{Y_1}) \leq \varepsilon^2$  in

$$\frac{\min\{L\sqrt{d}, d\}}{\varepsilon} \text{ iterations}$$

- When  $L = \tilde{O}(1)$ , iteration complexity:  $\sqrt{d}/\varepsilon$

# Convergence theory for DDPM

---

## Theorem 6 (Jiao, Zhou, Li '25)

DDPM yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  and  $\text{KL}(p_{X_1} \| p_{Y_1}) \leq \varepsilon^2$  in

$$\frac{\min\{L\sqrt{d}, d\}}{\varepsilon} \text{ iterations}$$

- When  $L = \tilde{O}(1)$ , iteration complexity:  $\sqrt{d}/\varepsilon$
- When  $L = O(\sqrt{d})$ , sublinear dependency on  $d$

# Convergence theory for DDPM

---

## Theorem 6 (Jiao, Zhou, Li '25)

DDPM yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  and  $\text{KL}(p_{X_1} \| p_{Y_1}) \leq \varepsilon^2$  in

$$\frac{\min\{L\sqrt{d}, d\}}{\varepsilon} \text{ iterations}$$

- When  $L = \tilde{O}(1)$ , iteration complexity:  $\sqrt{d}/\varepsilon$
- When  $L = O(\sqrt{d})$ , sublinear dependency on  $d$
- When  $L = \tilde{O}(1)$ , the lower bound is  $\tilde{\Omega}(\sqrt{d}/\varepsilon)$

# Sampler based on randomized midpoint

---

For integral  $\int_0^1 f(x)dx$ , interval width  $\delta$ :

- Fixed endpoints: error  $\propto \|f'\|\delta$
- Randomized midpoints: error  $\propto \|f'\|\delta^{3/2}$

— *Gupta et al '24, Shen and Lee '19*

# Sampler based on randomized midpoint

---

For integral  $\int_0^1 f(x)dx$ , interval width  $\delta$ :

- Fixed endpoints: error  $\propto \|f'\|\delta$
- Randomized midpoints: error  $\propto \|f'\|\delta^{3/2}$

— Gupta et al '24, Shen and Lee '19

For ODE in diffusion models:

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_{\gamma}^*(Y_{\gamma}^{\text{ode}})}_{\text{approximate by?}} d\gamma$$
$$\approx \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \frac{\bar{\alpha}_{t-1} - \bar{\alpha}_t}{\sqrt{\gamma_t^3}} s_{\gamma_t}^*(Y_{\gamma_t}^{\text{ode}}), \quad \gamma_t \sim \text{Unif}([\bar{\alpha}_t, \bar{\alpha}_{t-1}])$$

# Convergence theory for our accelerated sampler

---

## Theorem 7 (Jiao & Li '25)

*Our accelerated sampler yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in*

$$\tilde{O}\left(\frac{\min\{d, d^{2/3}L^{1/3}, d^{1/3}L\}}{\varepsilon^{2/3}}\right) \text{ iterations}$$

# Convergence theory for our accelerated sampler

---

## Theorem 7 (Jiao & Li '25)

Our accelerated sampler yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in

$$\tilde{O}\left(\frac{\min\{d, d^{2/3}L^{1/3}, d^{1/3}L\}}{\varepsilon^{2/3}}\right) \text{ iterations}$$

- To achieve sublinear dependency on  $d$ , our accelerated sampler requires  $L = O(d)$  vs. DDPM requires  $L = O(\sqrt{d})$

# Convergence theory for our accelerated sampler

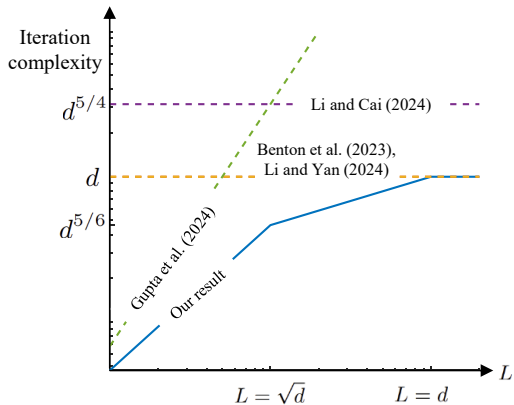
---

## Theorem 7 (Jiao & Li '25)

Our accelerated sampler yields  $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$  in

$$\tilde{O}\left(\frac{\min\{d, d^{2/3}L^{1/3}, d^{1/3}L\}}{\varepsilon^{2/3}}\right) \text{ iterations}$$

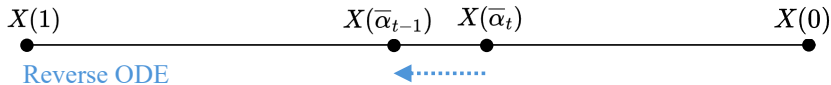
- When  $L = \tilde{O}(1)$ ,  $O(d^{1/3})$  dependency vs.  $O(d^{1/2})$  dependency of DDPM



achieve improvements over a full range of  $L$

# Forward-value discretization of diffusion ODE

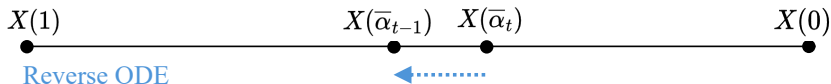
---



$$X(\bar{\alpha}_{t-1}) = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} X(\bar{\alpha}_t) + \sqrt{1 - \bar{\alpha}_{t-1}} \int_{\lambda(\bar{\alpha}_{t-1})}^{\lambda(\bar{\alpha}_t)} e^{\lambda} \underbrace{\mu(X(\bar{\alpha}(\lambda)), \bar{\alpha}(\lambda))}_{\text{data predictor}} d\lambda$$

$$\lambda = \frac{1}{2} \log \frac{\bar{\alpha}}{1 - \bar{\alpha}}, \quad \mu(x, \bar{\alpha}) = \frac{x + (1 - \bar{\alpha})s_{\bar{\alpha}}(x)}{\sqrt{\bar{\alpha}}}$$

# Forward-value discretization of diffusion ODE



$$X(\bar{\alpha}_{t-1}) = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} X(\bar{\alpha}_t) + \sqrt{1 - \bar{\alpha}_{t-1}} \int_{\lambda(\bar{\alpha}_{t-1})}^{\lambda(\bar{\alpha}_t)} e^\lambda \underbrace{\mu(X(\bar{\alpha}(\lambda)), \bar{\alpha}(\lambda))}_{\text{data predictor}} d\lambda$$

Reverse-value discretization of diffusion ODE:

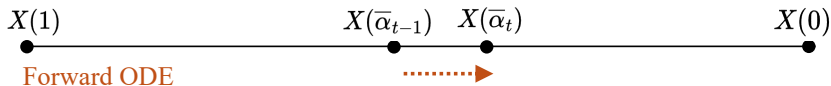
Approx. integrand  $\mu(X(\bar{\alpha}(\lambda)), \bar{\alpha}(\lambda))$  with **reverse value at  $\bar{\alpha}_t$**

$$Y_{t-1} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} Y_t + \left( \sqrt{\bar{\alpha}_{t-1}} - \sqrt{\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{1 - \bar{\alpha}_t}} \right) \mu(Y_t, \bar{\alpha}_t),$$

$t = T, \dots, 2.$

# Forward-value discretization of diffusion ODE

---



Forward-value discretization of diffusion ODE:

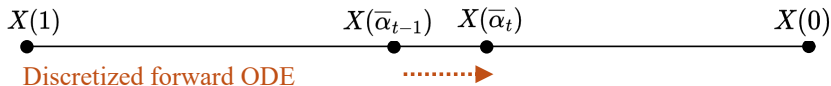
Approx. integrand  $\mu(X(\bar{\alpha}(\lambda)), \bar{\alpha}(\lambda))$  with **forward value at  $\bar{\alpha}_{t-1}$**

$$Y_t^{\text{fw}} = \sqrt{\frac{1 - \bar{\alpha}_t}{1 - \bar{\alpha}_{t-1}}} Y_{t-1}^{\text{fw}} - \left( \sqrt{\bar{\alpha}_t} - \sqrt{\frac{(1 - \bar{\alpha}_t)\bar{\alpha}_{t-1}}{1 - \bar{\alpha}_{t-1}}} \right) \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}),$$

$t = T, \dots, 2.$

# Motivation: diffusion ODE

---

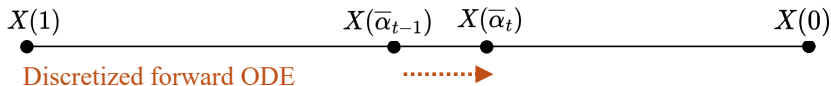


$$Y_{t-1}^{\text{fw}} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} Y_t^{\text{fw}} + \left( \sqrt{\bar{\alpha}_{t-1}} - \sqrt{\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{1 - \bar{\alpha}_t}} \right) \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}),$$

$t = T, \dots, 2.$

## Motivation: diffusion ODE

---



$$Y_{t-1}^{\text{fw}} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} Y_t^{\text{fw}} + \left( \sqrt{\bar{\alpha}_{t-1}} - \sqrt{\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{1 - \bar{\alpha}_t}} \right) \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}),$$

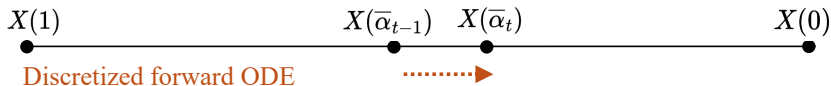
$t = T, \dots, 2.$

**Valid in two extreme cases:**

- $T = 2$ :  $Y_1^{\text{fw}} \approx \mu(Y_1^{\text{fw}}, \bar{\alpha}_1) \approx X_1$  ( $\bar{\alpha}_1 \approx 1$ ,  $\bar{\alpha}_2 \approx 0$ )

# Motivation: diffusion ODE

---



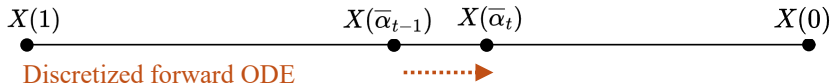
$$Y_{t-1}^{\text{fw}} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} Y_t^{\text{fw}} + \left( \sqrt{\bar{\alpha}_{t-1}} - \sqrt{\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{1 - \bar{\alpha}_t}} \right) \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}),$$

$t = T, \dots, 2.$

**Valid in two extreme cases:**

- $T = 2$ :  $Y_1^{\text{fw}} \approx \mu(Y_1^{\text{fw}}, \bar{\alpha}_1) \approx X_1$  ( $\bar{\alpha}_1 \approx 1$ ,  $\bar{\alpha}_2 \approx 0$ )
- $T \rightarrow \infty$ :  $\mu(Y^{\text{fw}}, \bar{\alpha}) \approx \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}) \Rightarrow$  solve the ODE

## Motivation: diffusion ODE



$$Y_{t-1}^{\text{fw}} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} Y_t^{\text{fw}} + \left( \sqrt{\bar{\alpha}_{t-1}} - \sqrt{\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{1 - \bar{\alpha}_t}} \right) \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}),$$

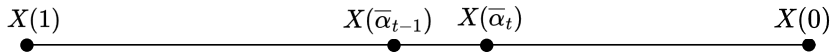
$t = T, \dots, 2.$

**Valid in two extreme cases:**

- $T = 2$ :  $Y_1^{\text{fw}} \approx \mu(Y_1^{\text{fw}}, \bar{\alpha}_1) \approx X_1$  ( $\bar{\alpha}_1 \approx 1$ ,  $\bar{\alpha}_2 \approx 0$ )
- $T \rightarrow \infty$ :  $\mu(Y^{\text{fw}}, \bar{\alpha}) \approx \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}) \Rightarrow$  solve the ODE

**Conjecture: forward-value discretized ODE  $Y_1^{\text{fw}} \stackrel{\text{d}}{\approx} X(\bar{\alpha}_1)$  for any  $T \geq 2$**

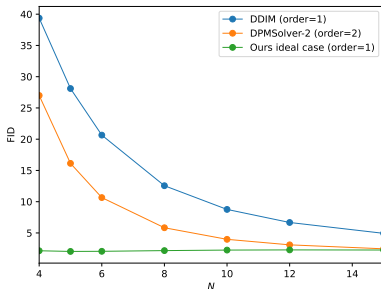
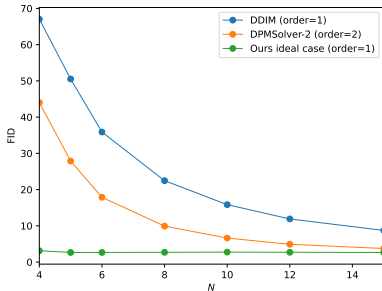
# Motivation: diffusion ODE



Discretized forward ODE

$$Y_{t-1}^{\text{fw}} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} Y_t^{\text{fw}} + \left( \sqrt{\bar{\alpha}_{t-1}} - \sqrt{\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{1 - \bar{\alpha}_t}} \right) \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}),$$

$t = T, \dots, 2.$



## Order of our sampler

---

Recall forward-value discretization of diffusion ODE: for  $t = T, \dots, 2$ ,

$$Y_{t-1}^{\text{fw}} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} Y_t^{\text{fw}} + \left( \sqrt{\bar{\alpha}_{t-1}} - \sqrt{\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{1 - \bar{\alpha}_t}} \right) \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}).$$

### Theorem 8 (Jiao, Li, Cai, Li. '25)

*The forward-value discretization of diffusion ODE satisfies*

$$\|Y_1^{\text{fw}} - X(\bar{\alpha}_1) + Y_1^{\text{ddim}} - X(\bar{\alpha}_1)\|_2 = \tilde{O}(1/T^2).$$

## Order of our sampler

---

Recall forward-value discretization of diffusion ODE: for  $t = T, \dots, 2$ ,

$$Y_{t-1}^{\text{fw}} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} Y_t^{\text{fw}} + \left( \sqrt{\bar{\alpha}_{t-1}} - \sqrt{\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{1 - \bar{\alpha}_t}} \right) \mu(Y_{t-1}^{\text{fw}}, \bar{\alpha}_{t-1}).$$

### Theorem 8 (Jiao, Li, Cai, Li. '25)

*The forward-value discretization of diffusion ODE satisfies*

$$\|Y_1^{\text{fw}} - X(\bar{\alpha}_1) + Y_1^{\text{ddim}} - X(\bar{\alpha}_1)\|_2 = \tilde{O}(1/T^2).$$

**Our proposed sampler is first-order**

( DDIM error  $\|Y_1^{\text{ddim}} - X(\bar{\alpha}_1)\|_2 = \tilde{O}(1/T)$  )

## Experiments: quantitative comparisons

---

FIDs↓ on ImageNet64 dataset (EDM2-S):

NFE	First-order Algorithms		High-order Algorithms		
	Ours	DDIM	DPMSolver-2	DPMSolver-3	UniPC-3
4	<b>22.35</b>	43.86	29.91	23.66	50.00
5	<b>11.98</b>	31.41	18.16	12.57	26.93
6	<b>7.20</b>	23.22	11.89	7.63	15.26
8	<b>3.64</b>	14.20	6.40	3.92	5.78
10	<b>2.51</b>	9.85	4.26	2.70	2.71

## Experiments: quantitative comparisons

---

FIDs $\downarrow$  on ImageNet64 dataset (EDM2-S):

NFE	First-order Algorithms		High-order Algorithms		
	Ours	DDIM	DPMSolver-2	DPMSolver-3	UniPC-3
4	<b>22.35</b>	43.86	29.91	23.66	50.00
5	<b>11.98</b>	31.41	18.16	12.57	26.93
6	<b>7.20</b>	23.22	11.89	7.63	15.26
8	<b>3.64</b>	14.20	6.40	3.92	5.78
10	<b>2.51</b>	9.85	4.26	2.70	2.71

- Better than first- and second-order algorithms
- Comparable with third-order algorithms

## Experiments: quantitative comparisons

---

FIDs $\downarrow$  on ImageNet512 dataset (EDM2-XXL):

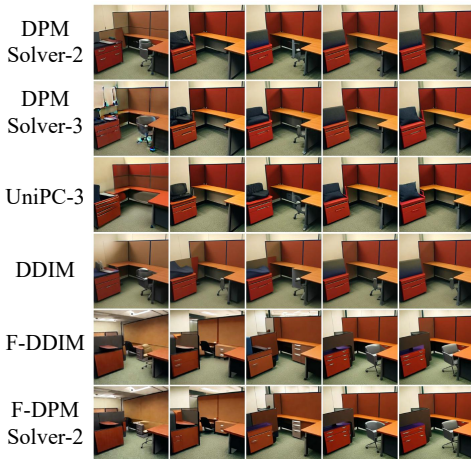
NFE	First-order Algorithms		High-order Algorithms		
	Ours	DDIM	DPMSolver-2	DPMSolver-3	UniPC-3
4	<b>50.96</b>	87.78	66.21	59.34	217.22
5	<b>22.69</b>	61.62	38.02	28.05	94.36
6	<b>14.33</b>	53.88	28.23	16.11	30.65
8	<b>6.18</b>	31.16	12.18	7.01	8.51
10	<b>4.57</b>	20.97	7.44	4.93	4.90

- Better than first- and second-order algorithms
- Comparable with third-order algorithms

# Experiments: qualitative comparisons

---

Stable diffusion: a desk and **chair** in an office cubicle



Only our sampler can generate **chair**

# Experiments: qualitative comparisons

Stable diffusion: **four** tennis players with rackets on a court

DPM  
Solver-2



DPM  
Solver-3



UniPC-3



DDIM



F-DDIM



F-DPM  
Solver-2



Only our sampler can generate **four** players

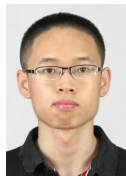
*Part 4: provable benefits of diffusion guidance*



Yuchen Jiao  
CUHK



Yuxin Chen  
UPenn



Gen Li  
CUHK

## Guided/controlled data generation

---

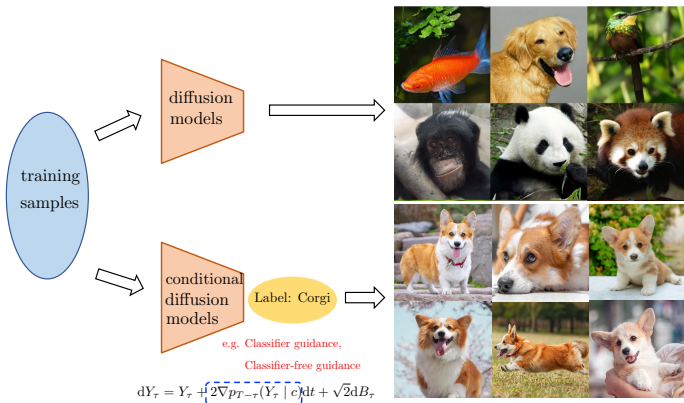
external conditions / specific instructions (via, e.g., prompts)

→ steer generation toward specific prompts

# Guided/controlled data generation

external conditions / specific instructions (via, e.g., prompts)

→ steer generation toward specific prompts



**class-conditional sampling:** generate samples for a specified class  $c$

sample from  $p_{\text{data}|c}$

# Conditional diffusion models?

---

**class-conditional sampling:** generate samples for a specified class  $c$   
sample from  $p_{\text{data}|c}$

**a natural conditioning approach:** replace scores w/ cond. scores

# Conditional diffusion models?

---

**class-conditional sampling:** generate samples for a specified class  $c$   
sample from  $p_{\text{data}|c}$

**a natural conditioning approach:** replace scores w/ cond. scores

$$\text{(unguided)} \quad dY_t = \left( \frac{1}{2} Y_t + \nabla \log p_{X_{1-t}}(Y_t) \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

# Conditional diffusion models?

---

**class-conditional sampling:** generate samples for a specified class  $c$   
sample from  $p_{\text{data}|c}$

**a natural conditioning approach:** replace scores w/ cond. scores

$$\text{(guided)} \quad dY_t = \left( \frac{1}{2} Y_t + \nabla \log p_{X_{1-t}|c}(Y_t | c) \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

# Conditional diffusion models?

---

**class-conditional sampling:** generate samples for a specified class  $c$   
sample from  $p_{\text{data}|c}$

**a natural conditioning approach:** replace scores w/ cond. scores

$$\text{(guided)} \quad dY_t = \left( \frac{1}{2} Y_t + \nabla \log p_{X_{1-t}|c}(Y_t | c) \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

- mathematically sound

# Conditional diffusion models?

---

**class-conditional sampling:** generate samples for a specified class  $c$   
sample from  $p_{\text{data}|c}$

**a natural conditioning approach:** replace scores w/ cond. scores

$$\text{(guided)} \quad dY_t = \left( \frac{1}{2} Y_t + \nabla \log p_{X_{1-t}|c}(Y_t | c) \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

- mathematically sound
- *unsatisfactory perceptual quality . . .*

sampling from this model via a standard diffusion sampler (e.g. DDPM). Interestingly, this standard way of conditioning usually does not perform well for diffusion models, for reasons that are unclear. In the text-to-image case for example, the generated samples tend to be visually incoherent and not faithful to the prompt, even for large-scale models (Ho and Salimans, 2022; Rombach et al., 2022).

## Practically more appealing: diffusion guidance

---

Add classifier probability to the drift to amplify guidance

$$dY_t^w = \left( \frac{1}{2} Y_t^w + \nabla \log p_{X_{1-t} | c}(Y_t^w | c) \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

## Practically more appealing: diffusion guidance

---

Add classifier probability to the drift to amplify guidance

$$dY_t^w = \left( \frac{1}{2} Y_t^w + \nabla \log p_{X_{1-t}|c}(Y_t^w | c) + \underbrace{w \nabla \log p_c | X_{1-t}(c | Y_t^w)}_{\text{guidance}} \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

# Practically more appealing: diffusion guidance

---

Add classifier probability to the drift to amplify guidance

$$dY_t^w = \left( \frac{1}{2} Y_t^w + \nabla \log p_{X_{1-t} | c}(Y_t^w | c) + \underbrace{w \nabla \log p_c | X_{1-t}(c | Y_t^w)}_{\text{guidance}} \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

---

## Diffusion Models Beat GANs on Image Synthesis

---

Prafulla Dhariwal\*  
OpenAI  
prafula@openai.com

Alex Nichol\*  
OpenAI  
alex@openai.com

classifier guidance



## CLASSIFIER-FREE DIFFUSION GUIDANCE

Jonathan Ho & Tim Salimans  
Google Research, Brain team  
{jonathanho,salimans}@google.com

# Mystery of diffusion guidance

---

$$dY_t^w = \left( \frac{1}{2} Y_t^w + \nabla \log p_{X_{1-t}|c}(Y_t^w | c) + \underbrace{w \nabla \log p_{c|X_{1-t}}(c | Y_t^w)}_{\text{guidance}} \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

- Bayes' rule  $\nabla \log p_{X_{1-t}|c} + w \nabla \log p_{c|X_{1-t}} = \nabla \log (p_{X_{1-t}} \cdot p_c^{1+w})$

# Mystery of diffusion guidance

---

$$dY_t^w = \left( \frac{1}{2} Y_t^w + \nabla \log p_{X_{1-t}|c}(Y_t^w | c) + \underbrace{w \nabla \log p_{c|X_{1-t}}(c | Y_t^w)}_{\text{guidance}} \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

- Bayes' rule  $\nabla \log p_{X_{1-t}|c} + w \nabla \log p_{c|X_{1-t}} = \nabla \log (p_{X_{1-t}} \cdot p_{c|X_{1-t}}^{1+w})$
- *appears to sample from*

$$p_{X_0} p_{c|X_0}^{1+w}$$

# Mystery of diffusion guidance

---

- appears to sample from

$$p_{X_0} p_c^{1+w} | X_0$$

- contrasts w/ Bayes

$$p_{X_0 | c} \propto p_{X_0} p_c | X_0$$

# Mystery of diffusion guidance

- appears to sample from

$$p_{X_0} p_c^{1+w} | X_0$$

- contrasts w/ Bayes

$$p_{X_0 | c} \propto p_{X_0} p_c | X_0$$

Mallat's Journey  
The ML Explosion  
Two Realities  
Ingredients of EML  
Clashing Mindsets  
Conclusion

Clashing 'Effectiveness' Narratives ...  
Clash of Elites

## Pivotal Moment, 1

- ▶ Norvig and Co-authors point out/claim (*in talks, not paper*):  
'Bayes Theorem is Empirically Wrong'.
- ▶ They mean: when used *in a certain dataset, in a certain way*, it is **empirically outperformed** by a different rule! i.e. **not**

$$P(B|A) \approx \frac{P(B)}{P(A)} \cdot P(A|B) \quad (\text{BAYES})$$

instead

$$P(B|A) \approx \frac{P(B)}{P(A)} \cdot P(A|B)^{1.5} \quad (\text{NOT-BAYES})$$

David Donoho - The Bridge from Mathematical to Digital, and Back

# Mystery of diffusion guidance

- appears to sample from

$$p_{X_0} p_c^{1+w} | X_0$$

- contrasts w/ Bayes


$$p_{X_0 | c} \propto p_{X_0} p_c | X_0$$

Mallat's Journey  
The ML Explosion  
Two Realities  
Ingredients of EML  
**Clashing Mindsets**  
Conclusion

Clashing 'Effectiveness' Narratives ...  
Clash of Elites

Pivotal Moment, 2

- ▶ Norvig et al.  
$$P(B|A) \approx \frac{P(B)}{P(A)} \cdot P(A|B)^{1.5} \quad (\text{NOT-BAYES})$$
- ▶ Empirical Machine Learning mindset:
  - ▶ *Mathematics is merely a source of gadgets.*
  - ▶ **Anyone** is free to grab **any** 'Math gadget' w/o motivation, amputate, mutilate and mashup *ad libitum*
  - ▶ *Measured Task Performance* on a URL-available dataset is objective arbiter of **what works**.
  - ▶ Poor Leaderboard ranking *deprecates the gadget*.



David Donoho - The Bridge from Mathematical to Digital, and Back

# Mystery of diffusion guidance

---

- appears to sample from

$$p_{X_0} p_c^{1+w} | X_0$$

- contrasts w/ Bayes

$$p_{X_0 | c} \propto p_{X_0} p_c | X_0$$

- *actually, even more complicated than*  $p_{X_0} p_c^{1+w} | X_0$  (Bradley et al. '24)

# Recent progress

---

Some recent progress for special distributions

- boosts classification confidence, diminishes diversity in Gaussian mixtures (Wu, Chen, Li, Wang, Wei '24)

# Recent progress

---

Some recent progress for special distributions

- boosts classification confidence, diminishes diversity in Gaussian mixtures (Wu, Chen, Li, Wang, Wei '24)
- samples more heavily from boundary of support of cond. dist (Chidambaram, Gutmiry, Chen, Lee, Lu '24)

# Recent progress

---

Some recent progress for special distributions

- boosts classification confidence, diminishes diversity in Gaussian mixtures (Wu, Chen, Li, Wang, Wei '24)
- samples more heavily from boundary of support of cond. dist (Chidambaram, Gutmiry, Chen, Lee, Lu '24)

Can we clarify benefits of guidance for general distributions?

# Our results: effectiveness of diffusion guidance

## Theorem 9 (Jiao, Chen, Li '25)

Guided diffusion w/ strength  $w$  yields

$$\underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_0}(c|Y_{output})}\right]}_{\text{unguided SDE}} \geq \underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_0}(c|Y_{output}^w)}\right]}_{\text{guided SDE w/ strength } w}$$

# Our results: effectiveness of diffusion guidance

## Theorem 9 (Jiao, Chen, Li '25)

Guided diffusion w/ strength  $w$  yields

$$\underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_0}(c|Y_{output})}\right]}_{\text{unguided SDE}} \geq \underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_0}(c|Y_{output}^w)}\right]}_{\text{guided SDE w/ strength } w}$$

- guidance improves  $\underbrace{\text{avg. reciprocal of classifier probability}}_{\text{related to Inception Score } \mathbb{E}[\log p_{c|X_0}(c|Y_{output})]}$

# Our results: effectiveness of diffusion guidance

## Theorem 9 (Jiao, Chen, Li '25)

For any  $\delta \in (0, 1)$ , *guided diffusion w/ strength w yields*  
*early stopping*

$$\underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_\delta}(c|Y_{1-\delta})}\right]}_{\text{unguided SDE}} \geq \underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_\delta}(c|Y_{1-\delta}^w)}\right]}_{\text{guided SDE w/ strength w}}$$

- guidance improves avg. reciprocal of classifier probability  
related to Inception Score  $\mathbb{E}[\log p_{c|X_0}(c|Y_{\text{output}})]$

# Our results: effectiveness of diffusion guidance

## Theorem 9 (Jiao, Chen, Li '25)

For any  $\delta \in (0, 1)$  and init  $y$ , guided diffusion  $w$ / strength  $w$  yields   
 *early stopping*

$$\underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_\delta}(c|Y_{1-\delta})} \mid Y_0 = y\right]}_{\text{unguided SDE}} \geq \underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_\delta}(c|Y_{1-\delta}^w)} \mid Y_0^w = y\right]}_{\text{guided SDE } w/\text{ strength } w}$$

- guidance improves avg. reciprocal of classifier probability  
related to Inception Score  $\mathbb{E}[\log p_{c|X_0}(c|Y_{\text{output}})]$

# Our results: effectiveness of diffusion guidance

## Theorem 9 (Jiao, Chen, Li '25)

For any  $\delta \in (0, 1)$  and init  $y$ , guided diffusion  $w$ / strength  $w$  yields   
 early stopping

$$\underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_\delta}(c|Y_{1-\delta})} \mid Y_0 = y\right]}_{\text{unguided SDE}} \geq \underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_\delta}(c|Y_{1-\delta}^w)} \mid Y_0^w = y\right]}_{\text{guided SDE } w/\text{ strength } w}$$

- guidance improves avg. reciprocal of classifier probability  
related to Inception Score  $\mathbb{E}[\log p_{c|X_0}(c|Y_{\text{output}})]$
- clarifies (rigorously) which metric guidance can improve

# Our results: effectiveness of diffusion guidance

## Theorem 9 (Jiao, Chen, Li '25)

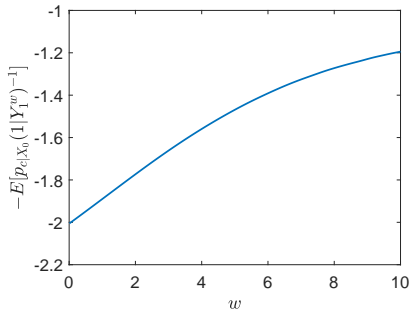
For any  $\delta \in (0, 1)$  and init  $y$ , guided diffusion  $w$ / strength  $w$  yields   
*early stopping*

$$\underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_\delta}(c|Y_{1-\delta})} \mid Y_0 = y\right]}_{\text{unguided SDE}} \geq \underbrace{\mathbb{E}\left[\frac{1}{p_{c|X_\delta}(c|Y_{1-\delta}^w)} \mid Y_0^w = y\right]}_{\text{guided SDE } w/\text{ strength } w}$$

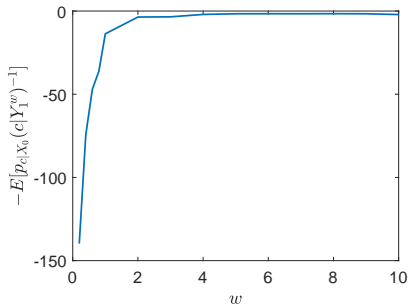
- guidance improves avg. reciprocal of classifier probability  
related to Inception Score  $\mathbb{E}[\log p_{c|X_0}(c|Y_{\text{output}})]$
- clarifies (rigorously) which metric guidance can improve
- holds for most distributions

# Experiments

---



Gaussian Mixture Models



ImageNet

## More generally: reward-guided diffusion

---

improve reward  $\mathbb{E}[r(Y^{\text{sample}})]$  by fine-tuning diffusion model  $\mathcal{D}$

## More generally: reward-guided diffusion

---

improve reward  $\mathbb{E}[r(Y^{\text{sample}})]$  by fine-tuning diffusion model  $\mathcal{D}$

$$dY_t^w = \left( \frac{1}{2} Y_t^w + \nabla \log p_{X_{1-t}}(Y_t^w) + w \cdot \text{guidance} \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

- guidance term:  $\underbrace{\nabla \log p_{X_{1-t}^{r\text{-wt}}}(Y_t^w)}_{\text{reward-reweighted score}} - \nabla \log p_{X_{1-t}}(Y_t^w)$

## More generally: reward-guided diffusion

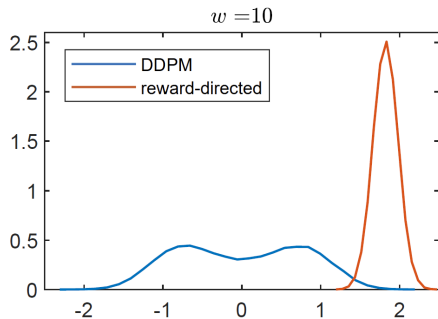
---

improve reward  $\mathbb{E}[r(Y^{\text{sample}})]$  by fine-tuning diffusion model  $\mathcal{D}$

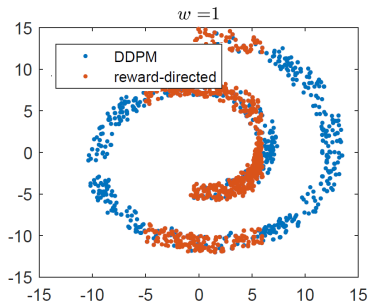
$$dY_t^w = \left( \frac{1}{2} Y_t^w + \nabla \log p_{X_{1-t}}(Y_t^w) + w \cdot \text{guidance} \right) \frac{dt}{t} + \frac{1}{\sqrt{t}} dB_t$$

- guidance term:  $\underbrace{\nabla \log p_{X_{1-t}^{r-wt}}(Y_t^w)}_{\text{reward-reweighted score}} - \nabla \log p_{X_{1-t}}(Y_t^w)$
- supported by theory; easy to train

# Experiments



$$r(x) = -(x - 2)^2$$



$$r(x) = 10\mathbb{1}(x_1 \in [-5, 6])$$

## *Part 5: diffusion models for inverse problems*



Yuchen Jiao  
CUHK



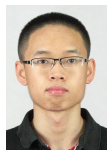
Na Li  
ZJU



Changxiao Cai  
UMich



Yuxin Chen  
UPenn

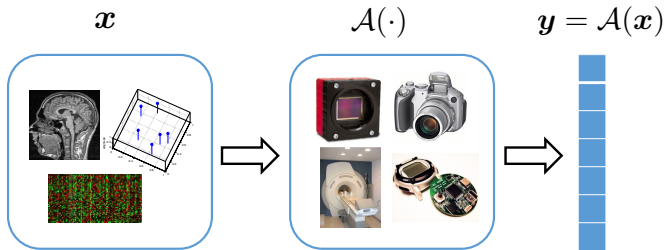


Gen Li  
CUHK

# Inverse problems

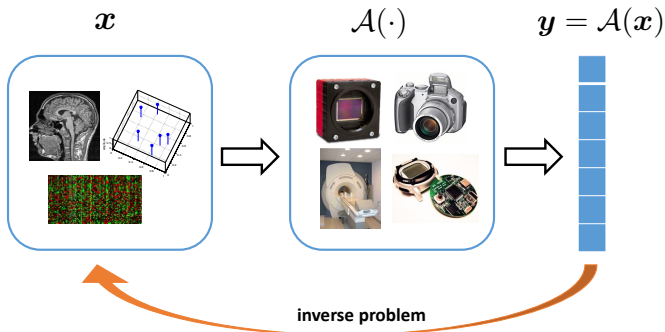
---

**Forward model:** we interrogate the signal of interest  $x$  through forward model  $\mathcal{A}$  and make measurements  $y$ .



# Inverse problems

**Forward model:** we interrogate the signal of interest  $x$  through forward model  $\mathcal{A}$  and make measurements  $y$ .

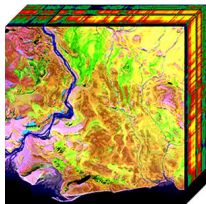


**Inverse problem:** recover the signal of interest  $x$  from  $y$ .

# Ubiquitous, but often ill-posed



healthcare



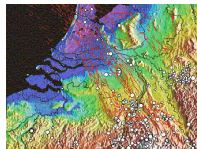
hyperspectral



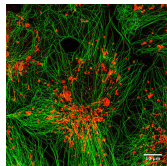
Internet traffic



Radio astronomy



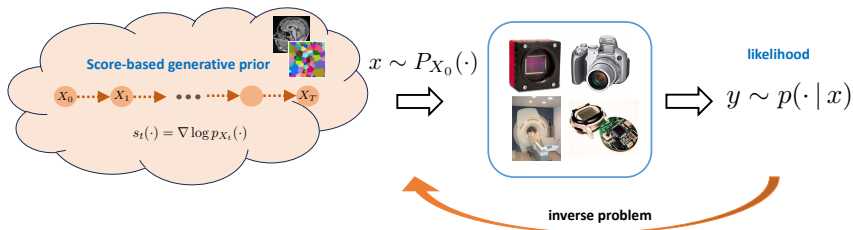
seismic imaging



microscopy

Can we exploit flexible / expressive data priors prescribed by diffusion models?

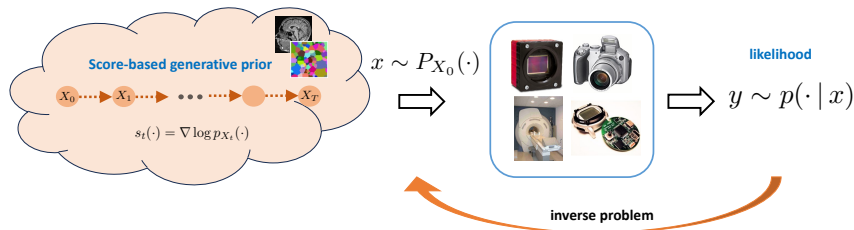
# Diffusion model for inverse problems



**Posterior sampling:** sample from

$$p(\cdot | y) \propto p(\cdot) p(y | x) = \underbrace{p(\cdot)}_{\text{prior}} \exp \underbrace{(\mathcal{L}(\cdot; y))}_{\text{log-likelihood}}$$

# Diffusion model for inverse problems



**Posterior sampling:** sample from

$$p(\cdot | y) \propto p(\cdot) p(y | x) = \underbrace{p(\cdot)}_{\text{prior}} \exp \left( \underbrace{\mathcal{L}(\cdot; y)}_{\text{log-likelihood}} \right)$$

**Score-based implicit prior:** the data prior  $p(\cdot)$  is accessed through its score functions  $s_t(\cdot) = \nabla \log p_{X_t}(\cdot)$ .

## Recall: Tweedie's formula

---

$$X_0 \sim p_{\text{data}}, \quad X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \mathcal{N}(0, I_d), \quad 1 \leq t \leq T$$

Tweedie's formula ([Hyvarinen'05](#); [Vincent'11](#)):

$$s_t^*(x) = \frac{1}{\sqrt{1 - \bar{\alpha}_t}} \mathbb{E} \left[ W \mid \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} W = x \right],$$

where the expectation is taken over  $W \sim \mathcal{N}(0, I_d)$ ,  $X_0 \sim p_{\text{data}}$ .

## Recall: Tweedie's formula

---

$$X_0 \sim p_{\text{data}}, \quad X_t = \sqrt{\bar{\alpha}_t}X_0 + \sqrt{1 - \bar{\alpha}_t}\mathcal{N}(0, I_d), \quad 1 \leq t \leq T$$

Tweedie's formula ([Hyvarinen'05](#); [Vincent'11](#)):

$$s_t^*(x) = \frac{1}{\sqrt{1 - \bar{\alpha}_t}} \mathbb{E} \left[ W \mid \sqrt{\bar{\alpha}_t}X_0 + \sqrt{1 - \bar{\alpha}_t}W = x \right],$$

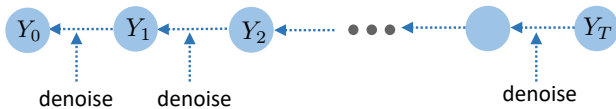
where the expectation is taken over  $W \sim \mathcal{N}(0, I_d)$ ,  $X_0 \sim p_{\text{data}}$ .

Data predictor (considered in inverse problems):

$$\begin{aligned} \mu_t(x) &:= \mathbb{E} \left[ X_0 \mid \sqrt{\bar{\alpha}_t}X_0 + \sqrt{1 - \bar{\alpha}_t}W = x \right] \\ &= \frac{1}{\sqrt{\bar{\alpha}_t}} (x - (1 - \bar{\alpha}_t)s_t^*(x)) \end{aligned}$$

## Recall: DDIM-type Sampler

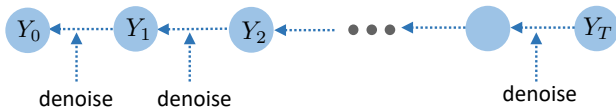
---



— *Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole '20*

## Recall: DDIM-type Sampler

---

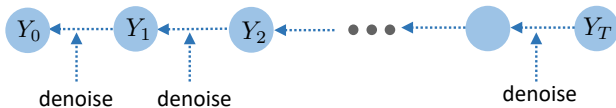


— Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole '20

$$\frac{\hat{X}_{t-1}}{\sigma_{t-1}} = \frac{\hat{X}_t}{\sigma_t} \sqrt{1 - \eta(1 - e^{-2\delta_{t-1}})} + e^{\lambda_{t-1}} \left( 1 - e^{-\delta_{t-1}} \sqrt{1 - \eta(1 - e^{-2\delta_{t-1}})} \right) \\ \cdot \mu_t(\hat{X}_t) + \sqrt{\eta(1 - e^{-2\delta_{t-1}})} \mathcal{N}(0, I_d), \quad 0 \leq \eta \leq 1$$

- $\delta_{t-1} := \lambda_{t-1} - \lambda_t$  where  $\lambda_t := \log(\alpha_t/\sigma_t)$  denotes signal-to-noise ratio

# Recall: DDIM-type Sampler



— *Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole '20*

$$\frac{\widehat{X}_{t-1}}{\sigma_{t-1}} = \frac{\widehat{X}_t}{\sigma_t} \sqrt{1 - \eta(1 - e^{-2\delta_{t-1}})} + e^{\lambda_{t-1}} \left( 1 - e^{-\delta_{t-1}} \sqrt{1 - \eta(1 - e^{-2\delta_{t-1}})} \right) \\ \cdot \mu_t(\widehat{X}_t) + \sqrt{\eta(1 - e^{-2\delta_{t-1}})} \mathcal{N}(0, I_d), \quad 0 \leq \eta \leq 1$$

- $\delta_{t-1} := \lambda_{t-1} - \lambda_t$  where  $\lambda_t := \log(\alpha_t/\sigma_t)$  denotes signal-to-noise ratio
- Common choices with  $\eta \in [0, 1]$ 
  - $\eta = 1$ : stochastic sampler (DDPM, [Ho, Jain, Abbeel '20](#))
  - $\eta = 0$ : deterministic sampler (probability flow ODE, [Song, Meng, Ermon '20](#))

# Denoising with linear observations

---

**Setting:** linear inverse problem:

$$y = Ax + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

**Goal:** sample from posterior  $p_{X_0|Y=y}$

# Denoising with linear observations

---

**Setting:** linear inverse problem:

$$y = Ax + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

**Goal:** sample from posterior  $p_{X_0|Y=y}$

**Basic idea:** replace  $\mu_t$  with  $\mu_{t,y}$  in a DDIM-type sampler

$$\mu_{t,y}(x) := \mathbb{E}[X_0 \mid X_t = x, AX_0 + \varepsilon = y]$$

immediate trial:  $\min_{\mu} \|\mu - \mu_t(x)\|_2^2 + \gamma \|y - A\mu\|_2^2$ .

# Prior algorithms

---

- Diffusion Posterior Sampling (DPS, [Chung et al'22](#))
- Denoising Diffusion Restoration Models (DDRM, [Kawar et al'22](#))
- Regularization by Denoising Diffusion Process (RED-diff, [Mardani et al'23](#))
- Denoising Diffusion Null-space Model (DDNM/DDNM+, [Wang et al'22](#))
- Projected Diffusion (ProjDiff, [Zhang et al'24](#))

# Prior algorithms

---

- Diffusion Posterior Sampling (DPS, [Chung et al'22](#))
- Denoising Diffusion Restoration Models (DDRM, [Kawar et al'22](#))
- Regularization by Denoising Diffusion Process (RED-diff, [Mardani et al'23](#))
- Denoising Diffusion Null-space Model (DDNM/DDNM+, [Wang et al'22](#))
- Projected Diffusion (ProjDiff, [Zhang et al'24](#))

## Limitations

- implementation-complicated
- lacks theoretical guarantees
- unsatisfactory performance

## Observation: two sources of uncertainty

---

1. Measurement uncertainty: Apply SVD to  $A = U\Sigma V^\top \in \mathbb{R}^{k \times d}$ :

$$y = Ax + \varepsilon \iff \Sigma_{ss}^{-1} u_s^\top y = v_s^\top x_0 + \mathcal{N}\left(0, \frac{\sigma^2}{\Sigma_{ss}^2}\right), \quad \Sigma_{ss} > 0$$

2. Diffusion uncertainty

$$v_s^\top X_t = \alpha_t v_s^\top X_0 + \sigma_t \mathcal{N}(0, 1)$$

## Observation: two sources of uncertainty

---

1. Measurement uncertainty: Apply SVD to  $A = U\Sigma V^\top \in \mathbb{R}^{k \times d}$ :

$$y = Ax + \varepsilon \iff \Sigma_{ss}^{-1} u_s^\top y = v_s^\top x_0 + \mathcal{N}\left(0, \frac{\sigma^2}{\Sigma_{ss}^2}\right), \quad \Sigma_{ss} > 0$$

2. Diffusion uncertainty

$$v_s^\top X_t = \alpha_t v_s^\top X_0 + \sigma_t \mathcal{N}(0, 1)$$

**Key insight:** divide based on SNR

**Direction sets:**

$\mathcal{S}_t := \{s : \alpha_t \sigma < \Sigma_{ss} \sigma_t\}$  (observation-dominated at time  $t$ )

$\mathcal{S}_t^c := \mathcal{S} \setminus \mathcal{S}_t$  (prior-dominated at time  $t$ )

$\mathcal{S}^c := [k] \setminus \mathcal{S}$  (no observations)

Define critical time  $\tau_s$  s.t.  $\alpha_{\tau_s} / \sigma_{\tau_s} = \Sigma_{ss} / \sigma$ .

# Our sampler

---

$$\frac{v_s^\top \widehat{X}_{t-1}}{\sigma_{t-1}} = \frac{v_s^\top \widehat{X}_t}{\sigma_t} \sqrt{1 - \eta(1 - e^{-2\delta_{t-1}})} + e^{\lambda_{t-1}} \left( 1 - e^{-\delta_{t-1}} \sqrt{1 - \eta(1 - e^{-2\delta_{t-1}})} \right) \cdot v_s^\top \mu_t(\widehat{X}_t) + \sqrt{\eta(1 - e^{-2\delta_{t-1}})} v_s^\top \mathcal{N}(0, I_d)$$

- **Observation-dominated** ( $s \in \mathcal{S}_t$ ): apply DDIM-type update with  $\eta = 0$  to the modified forward process

$$\xi_{t,s} := \alpha_t \xi_{0,s} + \sqrt{\sigma_t^2 - \alpha_t^2 \sigma^2 \Sigma_{ss}^{-2}} \mathcal{N}(0, 1), \quad \xi_{0,s} := \Sigma_{ss}^{-1} u_s^\top y$$

- **Prior-dominated directions** ( $s \in \mathcal{S}_t^c$ ): apply DDIM-type update with  $\eta = 1$  (corresponding to DDPM)
- **Unobserved directions** ( $s \in \mathcal{S}^c$ ): apply DDIM-type update with sufficiently large  $\eta = \eta_{\mathcal{S}^c}$ .

# Theoretical guarantee

---

## Theorem 10 (Jiao, Na, Cai, Chen, Li '26)

Assume  $X_0$  has bounded support. Let  $\delta = \max \delta_t$ . If  $\eta_{S^c}^3 \delta^2 \rightarrow 0$  and  $\eta_{S^c}^2 \delta / \log \frac{1}{\delta} \rightarrow \infty$ , then

$$\text{law}(\hat{X}_t) \rightarrow \text{law}(X_0 | Y = y) \quad \text{as } \lambda_t \rightarrow \infty.$$

# Theoretical guarantee

---

## Theorem 10 (Jiao, Na, Cai, Chen, Li '26)

Assume  $X_0$  has bounded support. Let  $\delta = \max \delta_t$ . If  $\eta_{S^c}^3 \delta^2 \rightarrow 0$  and  $\eta_{S^c}^2 \delta / \log \frac{1}{\delta} \rightarrow \infty$ , then

$$\text{law}(\hat{X}_t) \rightarrow \text{law}(X_0 | Y = y) \quad \text{as } \lambda_t \rightarrow \infty.$$

- provably solving inverse problems

# Theoretical guarantee

---

## Theorem 10 (Jiao, Na, Cai, Chen, Li '26)

Assume  $X_0$  has bounded support. Let  $\delta = \max \delta_t$ . If  $\eta_{S^c}^3 \delta^2 \rightarrow 0$  and  $\eta_{S^c}^2 \delta / \log \frac{1}{\delta} \rightarrow \infty$ , then

$$\text{law}(\hat{X}_t) \rightarrow \text{law}(X_0 | Y = y) \quad \text{as} \quad \lambda_t \rightarrow \infty.$$

- provably solving inverse problems
- providing guidance for parameter choices  $\eta$

# Experiments: Choices of $\eta$

Inpainting task

$\eta_{S_t^c}$	$\eta^{S^c}$	CelebA				ImageNet			
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
0	0	27.22	0.79	27.11	61.46	21.24	0.54	46.72	89.56
	1	31.32	0.88	16.97	30.24	26.13	0.78	26.17	37.76
	2	32.55	0.90	15.44	26.62	28.76	0.86	16.74	18.60
	4	33.16	<b>0.91</b>	15.07	25.42	30.72	<b>0.89</b>	14.45	<b>15.74</b>
	8	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.86
	16	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
	32	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
0	16	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
0.5		33.03	0.90	15.66	26.60	30.86	0.88	15.33	17.89
1		32.84	0.90	16.36	28.42	30.67	0.88	16.61	20.85
2		32.51	0.89	17.89	33.21	30.26	0.86	19.52	28.71

# Experiments: Choices of $\eta$

Inpainting task

$\eta_{S_t^c}$	$\eta_{S^c}$	CelebA				ImageNet			
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
0	0	27.22	0.79	27.11	61.46	21.24	0.54	46.72	89.56
	1	31.32	0.88	16.97	30.24	26.13	0.78	26.17	37.76
	2	32.55	0.90	15.44	26.62	28.76	0.86	16.74	18.60
	4	33.16	<b>0.91</b>	15.07	25.42	30.72	<b>0.89</b>	14.45	<b>15.74</b>
	8	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.86
	16	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
	32	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
0	16	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
0.5		33.03	0.90	15.66	26.60	30.86	0.88	15.33	17.89
1		32.84	0.90	16.36	28.42	30.67	0.88	16.61	20.85
2		32.51	0.89	17.89	33.21	30.26	0.86	19.52	28.71

- Larger  $\eta_{S^c}$  improves performance, aligning with our theory

# Experiments: Choices of $\eta$

## Inpainting task

$\eta_{S_t^c}$	$\eta_{S^c}$	CelebA				ImageNet			
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
0	0	27.22	0.79	27.11	61.46	21.24	0.54	46.72	89.56
	1	31.32	0.88	16.97	30.24	26.13	0.78	26.17	37.76
	2	32.55	0.90	15.44	26.62	28.76	0.86	16.74	18.60
	4	33.16	<b>0.91</b>	15.07	25.42	30.72	<b>0.89</b>	14.45	<b>15.74</b>
	8	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.86
	16	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
	32	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
0	16	<b>33.23</b>	<b>0.91</b>	<b>15.06</b>	<b>25.33</b>	<b>31.06</b>	<b>0.89</b>	<b>14.31</b>	15.85
0.5		33.03	0.90	15.66	26.60	30.86	0.88	15.33	17.89
1		32.84	0.90	16.36	28.42	30.67	0.88	16.61	20.85
2		32.51	0.89	17.89	33.21	30.26	0.86	19.52	28.71

- Larger  $\eta_{S^c}$  improves performance, aligning with our theory
- $\eta_{S_t^c} = 0$  achieves the best performance

# Experiments: Comparison with references

Super-resolution task

	Method	CelebA				ImageNet			
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
Single-sample	$A^\dagger y$	23.64	0.49	68.72	147.89	21.85	0.45	65.34	183.32
	DPS	27.98	0.78	23.10	39.91	24.44	0.67	31.81	<b>36.17</b>
	DDRM	29.20	0.82	21.92	40.14	25.66	0.72	34.88	55.71
	RED-diff	24.98	0.55	50.59	73.89	22.74	0.49	53.24	96.26
	DDNM+	29.20	0.82	21.91	39.96	25.62	0.72	34.39	53.78
	ProjDiff	29.49	0.83	<b>20.89</b>	36.61	25.73	0.72	33.03	49.70
	<b>Ours</b>	<u>29.84</u>	<u>0.84</u>	22.02	<b>34.16</b>	<u>25.90</u>	<u>0.73</u>	<b>32.84</b>	49.99
Posterior mean	DDNM+	30.22	<b>0.85</b>	21.97	43.54	26.11	0.73	34.38	54.27
	ProjDiff	30.31	<b>0.85</b>	22.55	38.34	26.10	<b>0.74</b>	33.13	50.27
	<b>Ours</b>	<b>30.35</b>	<b>0.85</b>	22.55	40.62	<b>26.13</b>	<b>0.74</b>	33.35	52.91

# Experiments: Comparison with references

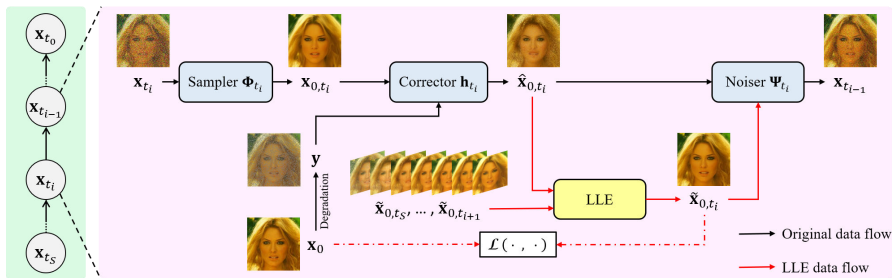
## Super-resolution task

	Method	CelebA				ImageNet			
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
Single-sample	$A^{\dagger}_y$	23.64	0.49	68.72	147.89	21.85	0.45	65.34	183.32
	DPS	27.98	0.78	23.10	39.91	24.44	0.67	31.81	<b>36.17</b>
	DDRM	29.20	0.82	21.92	40.14	25.66	0.72	34.88	55.71
	RED-diff	24.98	0.55	50.59	73.89	22.74	0.49	53.24	96.26
	DDNM+	29.20	0.82	21.91	39.96	25.62	0.72	34.39	53.78
	ProjDiff	29.49	0.83	<b>20.89</b>	36.61	25.73	0.72	33.03	49.70
	<b>Ours</b>	<u>29.84</u>	<u>0.84</u>	22.02	<b>34.16</b>	<u>25.90</u>	<u>0.73</u>	<b>32.84</b>	49.99
Posterior mean	DDNM+	30.22	<b>0.85</b>	21.97	43.54	26.11	0.73	34.38	54.27
	ProjDiff	30.31	<b>0.85</b>	22.55	38.34	26.10	<b>0.74</b>	33.13	50.27
	<b>Ours</b>	<b>30.35</b>	<b>0.85</b>	22.55	40.62	<b>26.13</b>	<b>0.74</b>	33.35	52.91

- Ours: best on **three** metrics
- Others: best on **at most one or two** metrics

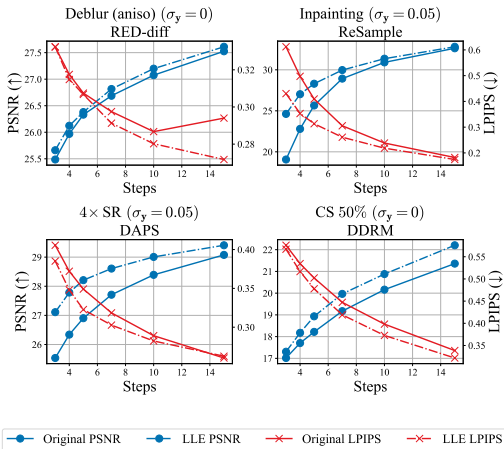
# Accelerated diffusion-based inverse algorithms

*Inspired by accelerated diffusion sampler*



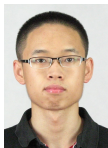
The proposed canonical form of diffusion-based inverse algorithms and the workflow of our Learnable Linear Extrapolation (LLE) method

# Numerical experiments

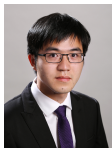


LLE achieves improvements across multiple tasks consistently

*Part 6: discrete diffusion (diffusion language models)*



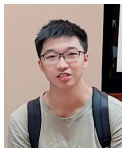
Gen Li  
CUHK



Changxiao Cai  
UMich



Daniil Dmitriev  
UPenn



Zhihan Huang  
UPenn



Yuting Wei  
UPenn

# Language generative models

training data

## How to Handle the “Terrible Twos”



Dr Becky Kennedy, Clinical Psychologist

The key to managing and understanding the two-year-old stage out-of-control behavior - and then, based on that understanding change.

The New York Times

GENERATION GRANDPARENT

### When a Darling Grandbaby Becomes a Devilish Terrible Two

The time comes in any relationship when the initial infatuation dampens a bit. But we're not breaking up.



Generative modeling



new samples

How to deal with terrible twos?

- Drink a coffee before engaging
- Lower your expectations
- Always carry snacks
- Pick your battles
- Master distraction
- Accept that logic is useless

# Language generative models

training data

## How to Handle the “Terrible Twos”



Dr. Becky Kennedy, Clinical Psychologist

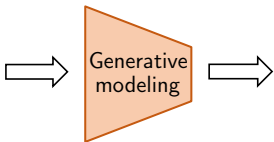
The key to managing and understanding the two-year-old stage out-of-control behavior - and then, based on that understanding change.

The New York Times

GENERATION GRANDPARENT

### When a Darling Grandbaby Becomes a Devilish Terrible Two

The time comes in any relationship when the initial infatuation dampens a bit. But we're not breaking up.



new samples

How to deal with terrible twos?

- Drink a coffee before engaging
- Lower your expectations
- Always carry snacks
- Pick your battles
- Master distraction
- Accept that logic is useless

- Given training data  $\underbrace{X^{\text{train},i} \sim p_{\text{data}}}_{\text{from a discrete distribution}} (1 \leq i \leq N)$  in  $[S]^d$

# Language generative models

training data

## How to Handle the “Terrible Twos”



Dr. Becky Kennedy, Clinical Psychologist

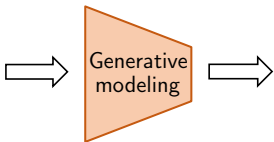
The key to managing and understanding the two-year-old stage out-of-control behavior - and then, based on that understanding change.

The New York Times

GENERATION GRANDPARENT

### When a Darling Grandbaby Becomes a Devilish Terrible Two

The time comes in any relationship when the initial infatuation dampens a bit. But we're not breaking up.



new samples

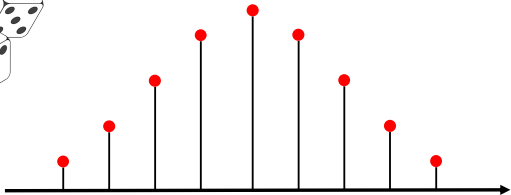
How to deal with terrible twos?

- Drink a coffee before engaging
- Lower your expectations
- Always carry snacks
- Pick your battles
- Master distraction
- Accept that logic is useless

- Given training data  $\underbrace{X^{\text{train},i} \sim p_{\text{data}}}_{\text{from a discrete distribution}} (1 \leq i \leq N)$  in  $[S]^d$
- Generate **new** samples  $Y \sim p_{\text{data}}$

# Challenges in discrete probabilistic modeling

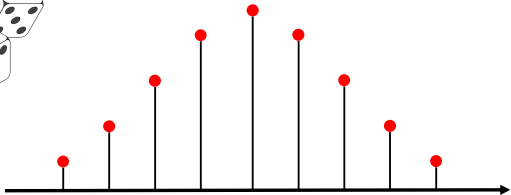
---



Discrete distribution:  $p(x) \geq 0$ , and  $\sum_{x \in \mathcal{X}} p(x) = 1$ .

# Challenges in discrete probabilistic modeling

---

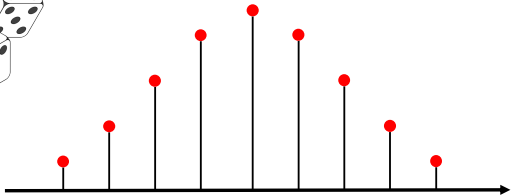


Discrete distribution:  $p(x) \geq 0$ , and  $\sum_{x \in \mathcal{X}} p(x) = 1$ .

- large sample space  $S^d$ , **exponentially** grow with  $d$

# Challenges in discrete probabilistic modeling

---



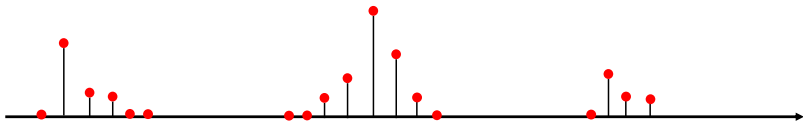
Discrete distribution:  $p(x) \geq 0$ , and  $\sum_{x \in \mathcal{X}} p(x) = 1$ .

- large sample space  $S^d$ , **exponentially** grow with  $d$
- no gradient  $\nabla_x p(x)$

# Challenges in discrete probabilistic modeling

---

Proper text: We thank the reviewer for the insightful comments



Discrete distribution:  $p(x) \geq 0$ , and  $\sum_{x \in \mathcal{X}} p(x) = 1$ .

- large sample space  $S^d$ , **exponentially** grow with  $d$
- no gradient  $\nabla_x p(x)$
- real texts are highly structured

# A successful paradigm: Autoregressive modeling

---

$\mathcal{X} = \{1, \dots, S\}^d$ ,  $x =$  All animals are equal, but some are more equal

$$\begin{aligned} p_{\text{data}}(x) &= p_{\text{data}}(x_1, \dots, x_d) \\ &= p_{\text{data}}(x_1) \underbrace{p_{\text{data}}(x_2 \mid x_1)}_{\substack{\downarrow \\ p_{\theta}(x_i \mid \text{previous words})}} \dots p_{\text{data}}(x_d \mid x_1, x_2, \dots, x_{d-1}) \end{aligned}$$

# A successful paradigm: Autoregressive modeling

---

$\mathcal{X} = \{1, \dots, S\}^d$ ,  $x =$  All animals are equal, but some are more equal

$$\begin{aligned} p_{\text{data}}(x) &= p_{\text{data}}(x_1, \dots, x_d) \\ &= p_{\text{data}}(x_1) \underbrace{p_{\text{data}}(x_2 \mid x_1)}_{\substack{\downarrow \\ p_{\theta}(x_i \mid \text{previous words})}} \dots p_{\text{data}}(x_d \mid x_1, x_2, \dots, x_{d-1}) \end{aligned}$$

advantages:

- only requires modeling  $p_{\theta}(x_i \mid \text{previous words})$  on  $[S]$
- approximate any probability distribution

# A successful paradigm: Autoregressive modeling

---

$\mathcal{X} = \{1, \dots, S\}^d$ ,  $x =$  All animals are equal, but some are more equal

$$\begin{aligned} p_{\text{data}}(x) &= p_{\text{data}}(x_1, \dots, x_d) \\ &= p_{\text{data}}(x_1) \underbrace{p_{\text{data}}(x_2 \mid x_1)} \dots p_{\text{data}}(x_d \mid x_1, x_2, \dots, x_{d-1}) \\ &\quad \downarrow \\ &\quad p_{\theta}(x_i \mid \text{previous words}) \end{aligned}$$

disadvantages:

- a rigid left-to-right order  $\rightarrow$  hard to control

# A successful paradigm: Autoregressive modeling

---

$\mathcal{X} = \{1, \dots, S\}^d$ ,  $x =$  All animals are equal, but some are more equal

$$\begin{aligned} p_{\text{data}}(x) &= p_{\text{data}}(x_1, \dots, x_d) \\ &= p_{\text{data}}(x_1) \underbrace{p_{\text{data}}(x_2 \mid x_1)}_{\substack{\downarrow \\ p_{\theta}(x_i \mid \text{previous words})}} \dots p_{\text{data}}(x_d \mid x_1, x_2, \dots, x_{d-1}) \end{aligned}$$

disadvantages:

- a rigid left-to-right order  $\rightarrow$  hard to control
- slow sampling speed

# A successful paradigm: Autoregressive modeling

---

$\mathcal{X} = \{1, \dots, S\}^d$ ,  $x =$  All animals are equal, but some are more equal

$$\begin{aligned} p_{\text{data}}(x) &= p_{\text{data}}(x_1, \dots, x_d) \\ &= p_{\text{data}}(x_1) \underbrace{p_{\text{data}}(x_2 \mid x_1)}_{\substack{\downarrow \\ p_{\theta}(x_i \mid \text{previous words})}} \dots p_{\text{data}}(x_d \mid x_1, x_2, \dots, x_{d-1}) \end{aligned}$$

disadvantages:

- a rigid left-to-right order  $\rightarrow$  hard to control
- slow sampling speed
- errors accumulate

*"autoregressive transformers are "doomed", as generation "drifts" from the data distribution and diverges during sampling." — Yann LeCun*

## An alternative: modeling score function

---

- write discrete distribution (parameterized by  $\theta$ ):

$$p_{\theta}(x) = \frac{e^{f_{\theta}(x)}}{Z_{\theta}}$$

where  $Z_{\theta} = \sum_{x \in \mathcal{X}} p_{\theta}(x)$  is a normalizing const. depending on  $\theta$

# An alternative: modeling score function

---

- write discrete distribution (parameterized by  $\theta$ ):

$$p_{\theta}(x) = \frac{e^{f_{\theta}(x)}}{Z_{\theta}}$$

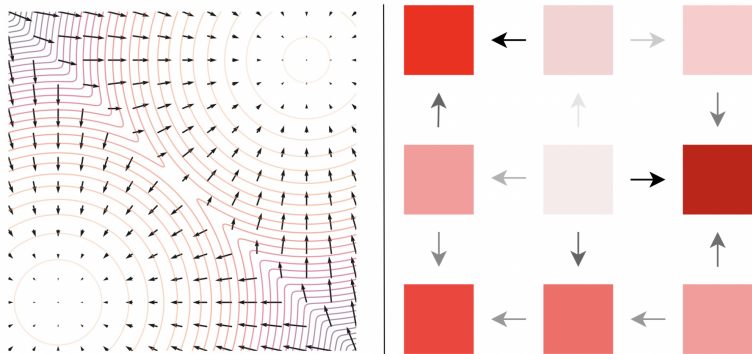
where  $Z_{\theta} = \sum_{x \in \mathcal{X}} p_{\theta}(x)$  is a normalizing const. depending on  $\theta$

- concrete score function:

$$s_{\theta}(y, x) = \frac{p_{\theta}(y)}{p_{\theta}(x)} = \frac{e^{f_{\theta}(y)}}{e^{f_{\theta}(x)}}$$

— analogous to  $\nabla \log p_{\theta}(x) = \nabla_x f_{\theta}(x)$  in cont. case

# An alternative: modeling score function



Left (continuous space): score function  $\nabla_x \log p(x)$  points to higher density regime

Right (discrete space): concrete score  $\frac{p(y)}{p(x)}$  generalizes for discrete spaces

— fig credit: Aaron Lou

## Towards a solid foundation for discrete diffusion models

*Austin, Johnson, Ho, Tarlow, van den Berg '21*

*Lou, Meng, Ermon '24*

*Sahoo, Arriola, Schiff, Gokaslan, Marroquin, Chiu, Rush, Kuleshov '24*

*Ou, Nie, Xue, Zhu, Sun, Li, Li '24*

*Campbell, Benton, De Bortoli, Rainforth, Deligiannidis, Doucet '22*

*Chen, Ying '25*

*Liang, Liang, Lai, Shroff '25*

*Li, Cai '25*

*Bach, Saremi '25*

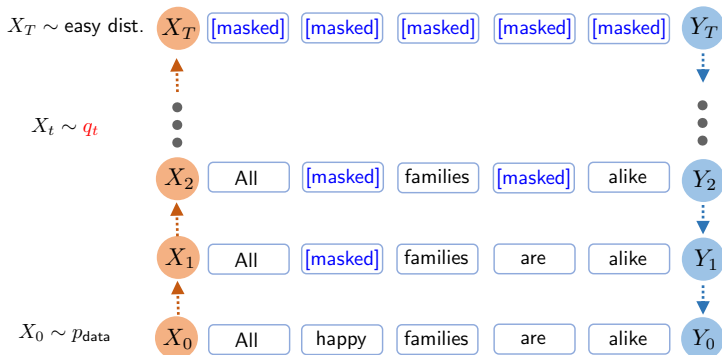
...



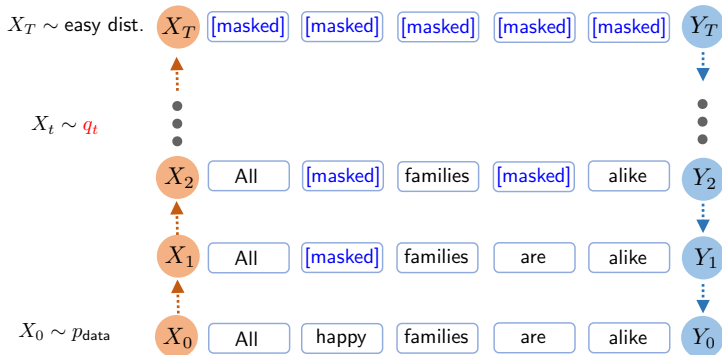
*A brief introduction to discrete diffusion under  
continuous-time Markov chain formulation*

*CTMC*

# A CTMC formulation for discrete case



# A CTMC formulation for discrete case

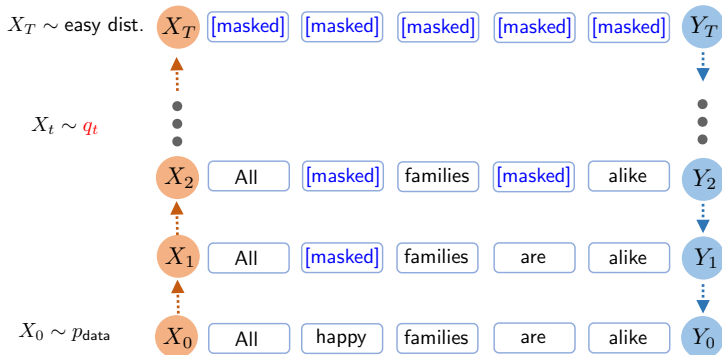


The transition probabilities satisfy, as  $\Delta t \rightarrow 0^+$ :

$$\Pr(x_{t+\Delta t} = y \mid x_t = x) = \mathbf{I}\{x = y\} + Q_t(x, y)\Delta t + o(\Delta t)$$

$$\text{rate matrix: } Q_t(x, y) \geq 0, y \neq x, Q_t(x, x) = -\sum_{y \neq x} Q_t(x, y)$$

# A CTMC formulation for discrete case

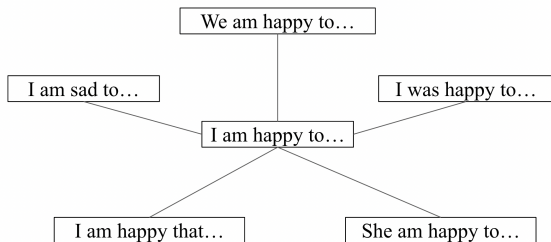


The transition probabilities satisfy, as  $\Delta t \rightarrow 0^+$ :

**Kolmogorov equation:** 
$$\frac{dq_t}{dt} = Q_t^\top q_t, \quad \text{for } 0 \leq t \leq T$$

## Two prevalent examples

---



- Uniform noising process  $\rightarrow \text{Unif}(\mathcal{X})$

$$Q_t(x, y) = 1/S, \quad x \sim y: \text{ differ at one coordinate}$$

- Masking noising process  $\rightarrow \delta_{\text{MASK}}$

$$Q_t(x, y) = 1, \quad \text{for some } i, x_i \neq \text{MASK}, y_i = \text{MASK}, x^{-i} = y^{-i}$$

## Score is all you need

---

How to learn a reverse process s.t.  $Y_t \stackrel{d}{\approx} X_t$ , for  $t = T, \dots, 1$ ?

# Score is all you need

---

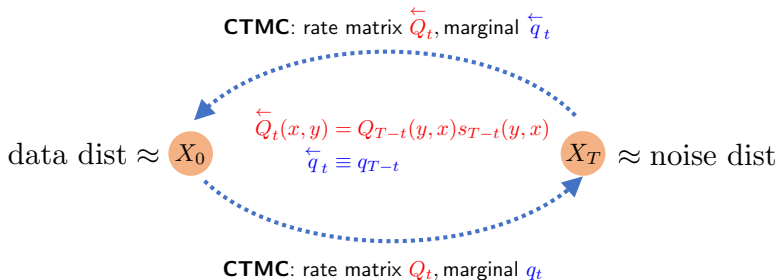
How to learn a reverse process s.t.  $Y_t \stackrel{d}{\approx} X_t$ , for  $t = T, \dots, 1$ ?

It is feasible as long as one knows score  $s_t(y, x) := p_t(y)/p_t(x)$

# Score is all you need

How to learn a reverse process s.t.  $Y_t \stackrel{d}{\approx} X_t$ , for  $t = T, \dots, 1$ ?

It is feasible as long as one knows score  $s_t(y, x) := p_t(y)/p_t(x)$



# Score estimation

---

Score entropy loss [Lou et al.'24](#):

$$\mathcal{L}_{\text{SE}}(t, \hat{s}, s) := \mathbb{E}_{x \sim q_t} \left[ \sum_{y \neq x} Q_t(y, x) s(y, x) \underbrace{\left( \frac{\hat{s}(y, x)}{s(y, x)} - 1 - \log \left( \frac{\hat{s}(y, x)}{s(y, x)} \right) \right)}_{\text{Bregman divergence } D(\hat{s}, s) \text{ with } \phi(x) = -\log x} \right]$$

# Score estimation

---

Score entropy loss [Lou et al.'24](#):

$$\mathcal{L}_{\text{SE}}(t, \hat{s}, s) := \mathbb{E}_{x \sim q_t} \left[ \sum_{y \neq x} Q_t(y, x) s(y, x) \underbrace{\left( \frac{\hat{s}(y, x)}{s(y, x)} - 1 - \log \left( \frac{\hat{s}(y, x)}{s(y, x)} \right) \right)}_{\text{Bregman divergence } D(\hat{s}, s) \text{ with } \phi(x) = -\log x} \right]$$

Key observations:

- if  $p(x) = \sum_{x_0} p(x | x_0) p_0(x_0)$ ,  $s(y, x) := \frac{p(y)}{p(x)}$  can be replaced by  $\frac{p(y|x_0)}{p(x|x_0)}$

# Score estimation

Score entropy loss [Lou et al.'24](#):

$$\mathcal{L}_{\text{SE}}(t, \hat{s}, s) := \mathbb{E}_{x \sim q_t} \left[ \sum_{y \neq x} Q_t(y, x) s(y, x) \underbrace{\left( \frac{\hat{s}(y, x)}{s(y, x)} - 1 - \log \left( \frac{\hat{s}(y, x)}{s(y, x)} \right) \right)}_{\text{Bregman divergence } D(\hat{s}, s) \text{ with } \phi(x) = -\log x} \right]$$

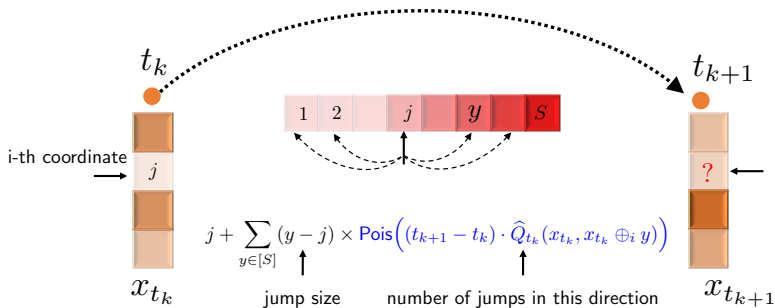
Key observations:

- if  $p(x) = \sum_{x_0} p(x | x_0) p_0(x_0)$ ,  $s(y, x) := \frac{p(y)}{p(x)}$  can be replaced by  $\frac{p(y|x_0)}{p(x|x_0)}$
- instead of evaluating  $s(y, x)$  for every  $y \neq x$ , only requires  $x \sim y$

— [SEE, Lou et al.'24](#), [Benton et al.'24](#), [Ou et al.'25](#)

# Score-based sampling: $\tau$ -leaping

— Campbell, et al.'22



Given  $\hat{s}_{t_i}(y, x)$  at points  $0 \leq t_0 < t_1 < \dots < t_N \leq T$ ,  $\tau$ -leaping is equivalent to a CTMC with

$$\hat{Q}_t^i(a, b) = \hat{Q}_{t_k}(x_{t_k}, x_{t_k} \oplus_i (b - a)), \quad \text{for } i \in [d],$$

with  $\hat{Q}_t(x, y) = Q_{T-t}(y, x) \hat{s}_{T-t}(y, x)$ .

# Prior theory

---

To obtain samples that are  $\varepsilon$ -close in KL to  $p_{\text{data}}$ , it takes

$$\tilde{O}\left(\frac{d^2 S}{\varepsilon}\right) \quad \text{and} \quad \tilde{O}\left(\frac{dS}{\varepsilon}\right)$$

for uniform and mask noising processes, respectively ([Liang et al.'25](#))

## Estimating the number of unseen species: How many words did Shakespeare know?

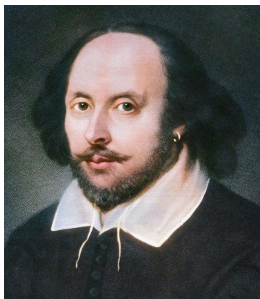
BRADLEY EFRON, RONALD THISTED

*Biometrika*, Volume 63, Issue 3, December 1976, Pages 435–447,

### Abstract

#### SUMMARY

Shakespeare wrote 31534 different words, of which 14376 appear only once, 4343 twice, etc. The question considered is how many words he knew but did not use. A parametric empirical Bayes model due to Fisher and a nonparametric model due to Good & Toulmin are examined. The latter theory is augmented using linear programming methods. We conclude that the models are equivalent to supposing that Shakespeare knew at least 35000 more words.



*Can we develop sharp dependences on  $d$  and  $S$ ?*

## Assumptions: score estimates $\{s_t(\cdot)\}$

---

- For discretization points:  $0 \leq t_0 < t_1 < \dots < t_N \leq T$ , assume

$$\sum_{k=0}^{N-1} (t_{k+1} - t_k) \mathcal{L}_{\text{SE}}(T - t_k, \widehat{s}_{T-t_k}, s_{T-t_k}) \leq \varepsilon_{\text{score}}$$

— *cont. case*:  $\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim q_t} [\|s_t(X) - s_t^*(X)\|_2^2] \leq \varepsilon_{\text{score}}^2$

**X** boundedness assump, e.g.  $\widehat{s}_{t_k}(x, y) \in [1/M, M]$

**X** regularity assump, e.g. *continuity of the score function*

# Uniform discrete diffusion

---

## Theorem 11 (Upper bound for uniform noising process)

For  $0 = t_0 < t_1 < \dots < t_N = T$ , let  $\Delta := \max_k \{t_{k+1} - t_k\} = O(1)$ .  
The  $\tau$ -leaping algorithm achieves

$$\text{KL}(p_{\text{data}} \| p_{\text{output}}) \lesssim \varepsilon_{\text{score}} + e^{-T} d \log S + \Delta d \log \left( \frac{S}{\Delta} \right).$$

# Uniform discrete diffusion

## Theorem 11 (Upper bound for uniform noising process)

For  $0 = t_0 < t_1 < \dots < t_N = T$ , let  $\Delta := \max_k \{t_{k+1} - t_k\} = O(1)$ .  
The  $\tau$ -leaping algorithm achieves

$$\text{KL}(p_{\text{data}} \| p_{\text{output}}) \lesssim \varepsilon_{\text{score}} + e^{-T} d \log S + \Delta d \log \left( \frac{S}{\Delta} \right).$$

Under constant step size schedule,  $t_{k+1} - t_k = T/N$ ,

$$N = \tilde{O} \left( \frac{d \log S}{\varepsilon} \right) \quad \text{discretization steps}$$

are enough to guarantee  $\text{KL}(p_{\text{data}} \| p_{\text{output}}) \leq \varepsilon_{\text{score}} + \varepsilon$ .

# Uniform discrete diffusion

## Theorem 11 (Upper bound for uniform noising process)

For  $0 = t_0 < t_1 < \dots < t_N = T$ , let  $\Delta := \max_k \{t_{k+1} - t_k\} = O(1)$ .  
The  $\tau$ -leaping algorithm achieves

$$\text{KL}(p_{\text{data}} \| p_{\text{output}}) \lesssim \varepsilon_{\text{score}} + e^{-T} d \log S + \Delta d \log \left( \frac{S}{\Delta} \right).$$

Under constant step size schedule,  $t_{k+1} - t_k = T/N$ ,

$$N = \tilde{O} \left( \frac{d \log S}{\varepsilon} \right) \quad \text{discretization steps}$$

are enough to guarantee  $\text{KL}(p_{\text{data}} \| p_{\text{output}}) \leq \varepsilon_{\text{score}} + \varepsilon$ .

— *matching lower bound for dist. far from uniform.*

# Comparisons with prior work

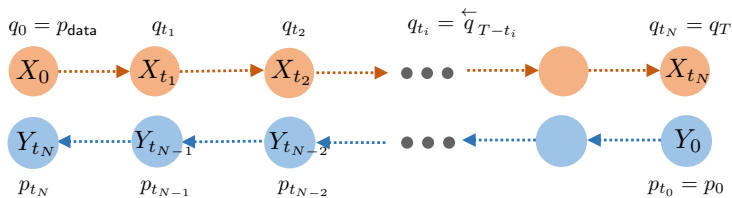
---

Paper	Score Est. Assump.	No Early Stopping	Iteration Complexity
Ren et al., '24	Bounded	$\times$	$d^2 S^2 / \epsilon$
Liang et al., '25a	Bounded	$\times$	$d^2 S / \epsilon$
This work	None	$\checkmark$	$d \log S / \epsilon$

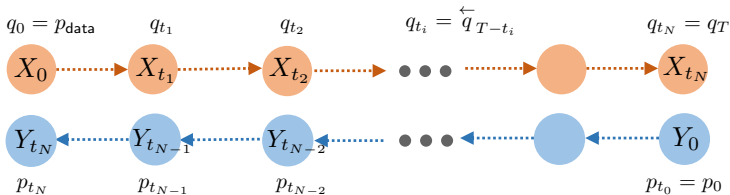
**No** early stopping and extra constraints are needed on the score estimator!

# Proof sketch: Step 1

---



# Proof sketch: Step 1



$$\begin{aligned}
 \text{KL}(q_0 \parallel p_T) &\leq \text{KL}(q_{T-t_0, \dots, T-t_N} \parallel p_{t_0, \dots, t_N}) \\
 &= \text{KL}(q_T \parallel p_0) + \sum_{k=0}^{N-1} \mathbb{E}_{x_{t_k} \sim \overleftarrow{q}_{t_k}} \text{KL}(\overleftarrow{q}_{t_{k+1}|t_k}(\cdot | x_{t_k}) \parallel p_{t_{k+1}|t_k}(\cdot | x_{t_k})) \\
 &\leq \text{KL}(q_T \parallel p_0) + \sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} \mathbb{E}_{x_t \sim \overleftarrow{q}_t} \sum_{y \neq x_t} \overleftarrow{Q}_t(x_t, y) \underbrace{D\left(\overleftarrow{Q}_t(x_t, y), \overleftarrow{Q}_t(x_t, y)\right)}_{D(a,b) := \frac{a}{b} - 1 - \log\left(\frac{a}{b}\right)} dt
 \end{aligned}$$

Girsanov's change-of-measure Theorem (e.g. [Liang et al., 2025](#)):

## Proof sketch: Step 2

---

For  $\tau$ -leaping sampler, consider  $t \in [t_k, t_{k+1})$  and write  $\ell = t_k$ .

$$\begin{aligned} & \sum_{y \neq x_t} \overleftarrow{Q}_t(x_t, y) D\left(\widehat{Q}_t(x_t, y), \overleftarrow{Q}_t(x_t, y)\right) \\ &= \frac{1}{S} \sum_{i \in [d]} \sum_{c \in [S]} s_{T-t}(x_t \oplus_i c, x_t) D\left(\widehat{s}_{T-\ell}(x_\ell \oplus_i c, x_\ell), s_{T-t}(x_t \oplus_i c, x_t)\right) \\ &= \frac{1}{S} \underbrace{\sum_{y_\ell \sim x_\ell} s_{T-\ell}(y_\ell, x_\ell) D\left(\widehat{s}_{T-\ell}(y_\ell, x_\ell), s_{T-\ell}(y_\ell, x_\ell)\right)}_{\text{controlled by small score entropy loss}} \\ &+ \frac{1}{S} \underbrace{\sum_{i \in [d]} \sum_{c \in [S]} \left(s_{T-\ell}(x_\ell \oplus_i c, x_\ell) - s_{T-t}(x_t \oplus_i c, x_t)\right) \log \widehat{s}_{T-\ell}(x_\ell \oplus_i c, x_\ell)}_{\text{expectation controlled by martingale property} = 0} \\ &+ \frac{1}{S} \sum_{y_t \sim x_t} \left(-\log s_{T-t}(y_t, x_t)\right) - \frac{1}{S} \sum_{y_\ell \sim x_\ell} \left(-\log s_{T-\ell}(y_\ell, x_\ell)\right). \end{aligned}$$

## Proof sketch: Step 3

---

It remains to control  $\varphi(T - t) - \varphi(T - t_k)$  with

$$\varphi(t) := \frac{1}{S} \mathbb{E}_{x_t \sim q_t} \left[ \sum_{i \in [d]} \sum_{c \in [S]} -\log s_t(x_t \oplus_i c, x_t) \right]$$

## Proof sketch: Step 3

---

It remains to control  $\varphi(T - t) - \varphi(T - t_k)$  with

$$\varphi(t) := \frac{1}{S} \mathbb{E}_{x_t \sim q_t} \left[ \sum_{i \in [d]} \sum_{c \in [S]} -\log s_t(x_t \oplus_i c, x_t) \right]$$

**Key observation:**  $\varphi(t)$  is (i) non-increasing and (ii) bounded  $0 \leq \varphi(t) \leq d(\log S + \max\{\log t^{-1}, 0\})$ .

## Proof sketch: Step 3

---

It remains to control  $\varphi(T - t) - \varphi(T - t_k)$  with

$$\varphi(t) := \frac{1}{S} \mathbb{E}_{x_t \sim q_t} \left[ \sum_{i \in [d]} \sum_{c \in [S]} -\log s_t(x_t \oplus_i c, x_t) \right]$$

**Key observation:**  $\varphi(t)$  is (i) non-increasing and (ii) bounded  $0 \leq \varphi(t) \leq d(\log S + \max\{\log t^{-1}, 0\})$ .

As a result, using (i), (ii), and the telescoping summation, one arrives at

$$\underbrace{\sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} (\varphi(T - t) - \varphi(T - t_k)) dt}_{\text{key quantity for discretization error}} \lesssim \Delta d \log \frac{S}{\Delta}.$$

*Sharp dependence but not yet adaptive...*  
*What about masking diffusion?*

— e.g., for  $p_{\text{data}} = \frac{1}{2}\delta_{\mathbf{0}} + \frac{1}{2}\delta_{\mathbf{1}}$  or  $p_{\text{data}} = \text{Unif}(\mathcal{X})$

# Fundamental information-theoretic quantities

Let  $P$  be a distribution over  $[S]^d$  and  $X := (X_1, \dots, X_d) \sim P$ . Define

**total correlation**  $\mathcal{C}(P) := \sum_{i=1}^d \underbrace{\mathcal{H}(X_i)}_{\text{entropy}} - \mathcal{H}(X_1, \dots, X_d)$

**dual total correlation**  $\mathcal{B}(P) := \mathcal{H}(X_1, \dots, X_d) - \sum_{i=1}^d \mathcal{H}(X_i | X_{-i})$

# Fundamental information-theoretic quantities

Let  $P$  be a distribution over  $[S]^d$  and  $X := (X_1, \dots, X_d) \sim P$ . Define

**total correlation**  $\mathcal{C}(P) := \sum_{i=1}^d \underbrace{\mathcal{H}(X_i)}_{\text{entropy}} - \mathcal{H}(X_1, \dots, X_d)$

**dual total correlation**  $\mathcal{B}(P) := \mathcal{H}(X_1, \dots, X_d) - \sum_{i=1}^d \mathcal{H}(X_i | X_{-i})$

- $\mathcal{C}(P) = \int_0^\infty (e^t - 1) \mathcal{I}(t) dt$  and  $\mathcal{B}(P) := \int_0^\infty \mathcal{I}(t) dt$  for

$$\mathcal{I}(t) := \sum_{i \neq j \in [d]} \underbrace{\mathcal{I}(x_t^i; x_t^j | x_t^{-(i,j)})}_{\text{conditional mutual information}}, \quad \text{for } P = q_0 = p_{\text{data}} \text{ and } x_t \sim q_t.$$

# Fundamental information-theoretic quantities

Let  $P$  be a distribution over  $[S]^d$  and  $X := (X_1, \dots, X_d) \sim P$ . Define

$$\text{total correlation} \quad \mathcal{C}(P) := \sum_{i=1}^d \underbrace{\mathcal{H}(X_i)}_{\text{entropy}} - \mathcal{H}(X_1, \dots, X_d)$$

$$\text{dual total correlation} \quad \mathcal{B}(P) := \mathcal{H}(X_1, \dots, X_d) - \sum_{i=1}^d \mathcal{H}(X_i | X_{-i})$$

- $\mathcal{C}(P) = \int_0^\infty (e^t - 1) \mathcal{I}(t) dt$  and  $\mathcal{B}(P) := \int_0^\infty \mathcal{I}(t) dt$  for

$$\mathcal{I}(t) := \sum_{i \neq j \in [d]} \underbrace{\mathcal{I}(x_t^i; x_t^j | x_t^{-(i,j)})}_{\text{conditional mutual information}}, \quad \text{for } P = q_0 = p_{\text{data}} \text{ and } x_t \sim q_t.$$

- define **effective total correlation**

$$\mathcal{D}(P) := \int_0^\infty \min(1, t) \cdot \mathcal{I}(t) dt \leq \min\{\mathcal{C}(P), \mathcal{B}(P)\}.$$

## Modified $\tau$ -leaping

Consider the following modification of  $\tau$ -leaping (for  $t \in [t_k, t_{k+1})$ ):

$$\widehat{Q}_t^i(a, b) = \frac{e^{T-t_k} - 1}{e^{T-t} - 1} \widehat{s}_{T-t_k}(x_{t_k} \odot_i b, x_{t_k}) \cdot \underbrace{\mathbf{I}\{a = x_{t_k}^i = \text{MASK}\}}_{\substack{\text{at most one transition at } x_t^i \\ \text{at the interval } [t_k, t_{k+1})}}$$

**Key observation:** for any  $t > 0$ ,  $b \in [S]$ ,  $x \in ([S] \cup \{\text{MASK}\})^d$ , and  $i$  such that  $x^i = \text{MASK}$ ,

$$s_t(x \odot_i b, x) = \frac{1}{e^t - 1} \cdot s_0(x \odot_i b, x)$$

## Masking discrete diffusion

### Theorem 12 (Upper bound for uniform noising process)

For  $0 = t_0 < t_1 < \dots < t_N = T$ , modified  $\tau$ -leaping satisfies

$$\text{KL}(p_{\text{data}} \| p_{\text{output}}) \lesssim \varepsilon_{\text{score}} + e^{-T} d \log S + \sum_{k=0}^{N-1} (t_{k+1} - t_k) \int_{T-t_{k+1}}^{T-t_k} \mathcal{I}(t) dt.$$

Under exponential-then-constant schedule,  $t_{k+1} - t_k \leq \kappa \min(1, T - t_{k+1})$ ,

$$N = \tilde{O} \left( \frac{\mathcal{D}(p_{\text{data}})}{\varepsilon} \right) \quad \text{discretization steps}$$

are enough to guarantee  $\text{KL}(p_{\text{data}} \| p_{\text{output}}) \leq \varepsilon_{\text{score}} + \varepsilon$ .

— it satisfies  $\mathcal{D}(P) \leq \min\{\mathcal{C}(P), \mathcal{B}(P)\} \leq d \log S$

# Comparisons with prior work

Paper	Score Est. Assump.	No Early Stopping	Sampler	Iteration Complexity	Adaptation
<a href="#">Liang et al., 25b</a>	Bounded	✗	$\tau$ -leaping	$dS/\varepsilon$	✗
<a href="#">Conforti et al., 25</a>	$\hat{s}_t \approx s_t$	✗	DMPM	$dS/\varepsilon$	✗
This work	None	✓	Modified $\tau$ -leaping	$\mathcal{D}/\varepsilon$	✓

**No** early stopping and **no** extra constraints on the score estimator!

# Distributions of small $\mathcal{D}(p_{\text{data}})$

---

## Completely independent of $d$

- product distributions on  $[S]^d$ :  $\mathcal{D}(p_{\text{data}}) = 0$  ( $N = 2$  suffices)
- mixture of two Dirac measures,  $\frac{1}{2}\delta_{k_1} + \frac{1}{2}\delta_{k_2}$ :  $\mathcal{D}(p_{\text{data}}) = \log 2$

## Sublinear dependence in $d$

- hidden Markov models
- cont. dist. with intrinsic dimension  $k$ , plus quantization
- sparse random regular graphs & stochastic block models

— *sublinear convergence rates for masking discrete diffusion!*

## Example: Latent parity model

---

Let  $d = 2k$  and consider the following distribution:

## Example: Latent parity model

---

Let  $d = 2k$  and consider the following distribution:

- sample  $b \sim \text{Bern}(1/2)$  & set  $x_1, \dots, x_k = b$

## Example: Latent parity model

---

Let  $d = 2k$  and consider the following distribution:

- sample  $b \sim \text{Bern}(1/2)$  & set  $x_1, \dots, x_k = b$
- sample  $x_{k+1}, \dots, x_{2k-1} \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(1/2)$

## Example: Latent parity model

---

Let  $d = 2k$  and consider the following distribution:

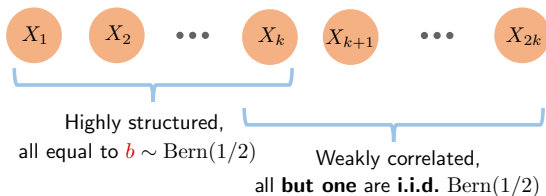
- sample  $b \sim \text{Bern}(1/2)$  & set  $x_1, \dots, x_k = b$
- sample  $x_{k+1}, \dots, x_{2k-1} \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(1/2)$
- set  $x_{2k} = (b + \sum_{i=k+1}^{2k-1} x_i) \bmod 2$

## Example: Latent parity model

---

Let  $d = 2k$  and consider the following distribution:

- sample  $b \sim \text{Bern}(1/2)$  & set  $x_1, \dots, x_k = b$
- sample  $x_{k+1}, \dots, x_{2k-1} \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(1/2)$
- set  $x_{2k} = (b + \sum_{i=k+1}^{2k-1} x_i) \bmod 2$



In this example,  $\mathcal{D}(q_0) = \Theta(1) \ll \Theta(d) = \min(\mathcal{B}(q_0), \mathcal{C}(q_0))$

# Proof sketch: Step 1

---

Similarly to the uniform case, writing  $x_t := (x_t \odot_i c, x_t)$ , and  $x_{t_k} := (x_{t_k} \odot_i c, x_{t_k})$ , we decompose

$$\begin{aligned} & \sum_{y \neq x_t} \overleftarrow{Q}_t(x_t, y) D\left(\widehat{Q}_t(x_t, y), \overleftarrow{Q}_t(x_t, y)\right) \\ &= \sum_{i \in m(x_t), c \in [S]} s_{T-t}(x_t) \cdot D\left(\frac{e^{T-t_k} - 1}{e^{T-t} - 1} \widehat{s}_{T-t_k}(x_{t_k}), s_{T-t}(x_t)\right) \\ &= \sum_{i \in m(x_t), c \in [S]} \underbrace{s_{T-t}(x_{t_k}) D(\widehat{s}_{T-t_k}(x_{t_k}), s_{T-t_k}(x_{t_k}))}_{\text{controlled by small score entropy loss}} \\ & \quad + \underbrace{(s_{T-t}(x_{t_k}) - s_{T-t}(x_t)) \log \frac{\widehat{s}_{T-t_k}(x_{t_k})}{s_{T-t_k}(x_{t_k})}}_{\text{expectation controlled by martingale property} = 0} + s_{T-t}(x_t) D(s_{T-t}(x_{t_k}), s_{T-t}(x_t)). \end{aligned}$$

## Proof sketch: Step 2

---

Using martingale property and changing measure under expectation:

$$\begin{aligned} & \mathbb{E}_{x_t, x_{t_k} \sim \bar{q}_{t, t_k}^{\leftarrow}} s_{T-t}(x_t \odot_i c, x_t) D(s_{T-t}(x_{t_k} \odot_i c, x_{t_k}), s_{T-t}(x_t \odot_i c, x_t)) \\ &= \mathbb{E}_{y_t, y_{t_k} \sim \bar{q}_{t, t_k}^{\leftarrow}} \log \frac{q_0(y_t) q_0(y_{t_k} \odot_i \text{MASK})}{q_0(y_t \odot_i \text{MASK}) q_0(y_{t_k} \odot_i y_t^i)} \\ &= \mathbb{E}_{y_t, y_{t_k} \sim \bar{q}_{t, t_k}^{\leftarrow}} (f_i(y_{t_k}) - f_i(y_t)), \end{aligned}$$

where we define

$$f_i(y) := \log \frac{q_0(y \odot_i y_t^i)}{q_0(y \odot_i \text{MASK})}.$$

Connection to the conditional mutual information (via *Dynkin's formula*)

$$\mathbb{E}_{y_t, y_{t_k} \sim \bar{q}_{t, t_k}^{\leftarrow}} (f_i(y_{t_k}) - f_i(y_t)) \lesssim \sum_{j \neq i} \int_{t_k}^t \mathbb{I}(y_v^i; y_v^j \mid y_v^{-(i,j)}) dv.$$

# Summary

---

- nonasymptotic convergence theory for diffusion models
- adaptation to unknown low dimensionality
- demystifying diffusion guidance
- provable training-free acceleration
- diffusion language models

# Summary

---

- nonasymptotic convergence theory for diffusion models
- adaptation to unknown low dimensionality
- demystifying diffusion guidance
- provable training-free acceleration
- diffusion language models

## **Future directions:**

- end-to-end theory that accounts for score learning + sampling?
- adaptive improvement under stylized statistical models
- design of high-order stochastic samplers

# Papers I

---

“A sharp convergence theory for the probability flow ODEs of diffusion models,” G. Li, Y. Wei, Y. Chi, Y. Chen, arXiv:2408.02320, 2024

“Towards non-asymptotic convergence for diffusion-based generative models,” G. Li, Y. Wei, Y. Chen, Y. Chi, ICLR 2024

“ $O(d/T)$  convergence theory for diffusion probabilistic models under minimal assumptions,” G. Li\*, Y. Yan\*, ICLR 2025

“Optimal convergence analysis of DDPM for general distributions,” Y. Jiao, Y. Zhou, G. Li, arXiv:2510.27562, 2025

“Minimax optimality of the probability flow ode for diffusion models,” C. Cai, G. Li, arXiv:2503.09583, 2025

“Adapting to unknown low-dimensional structures in score-based diffusion models,” G. Li\*, Y. Yan\*, NeurIPS 2024

“Low-dimensional adaptation of diffusion models: convergence in total variation,” J. Liang, Z. Huang, Y. Chen, COLT 2025

“Denoising diffusion probabilistic models are optimally adaptive to unknown low dimensionality,” Z. Huang, Y. Wei, Y. Chen, *Mathematics of Operations Research*, 2026+

“Dimension-free convergence of diffusion models for approximate Gaussian mixtures,” G. Li\*, C. Cai\*, Y. Wei, ICML 2026

# Papers II

---

“Accelerating convergence of score-based diffusion models, provably,” G. Li\*, Y. Huang\*, T. Efimov, Y. Wei, Y. Chi, Y. Chen, ICML 2024

“Stochastic Runge-Kutta methods: Provable acceleration of diffusion models,” Y. Wu, Y. Chen, Y. Wei, arXiv:2410.04760, 2024

“Faster diffusion models via higher-order approximation,” G. Li\*, Y. Zhou\*, Y. Wei, Y. Chen, arXiv:2506.24042, 2025

“Improved convergence rate for diffusion probabilistic models,” G. Li, Y. Jiao, ICLR 2025

“Provable acceleration for diffusion models under minimal assumptions,” G. Li\*, C. Cai\*, arXiv:2410.23285

“Theoretical insights for diffusion guidance: A case study for Gaussian mixture models,” Y. Wu, M. Chen, Z. Li, M. Wang, Y. Wei, ICML 2024

“Provable efficiency of guidance in diffusion models for general data distribution,” Y. Jiao, G. Li, ICML 2025

“Towards a unified framework for guided diffusion models,” Y. Jiao, Y. Chen, G. Li, arXiv:2512.04985, 2025

“Provably solving inverse problems with diffusion prior through DDIM-type sampler,” Y. Jiao, N. Li, C. Cai, Y. Chen, G. Li, 2026

# Papers III

---

“Breaking AR’s sampling bottleneck: provable acceleration via diffusion language models,” G. Li\*, C. Cai\*, NeurIPS 2025

“Efficient sampling with discrete diffusion models: sharp and adaptive guarantees,” Daniil Dmitriev\*, Zhihan Huang\*, Yuting Wei, COLT 2026

(\* = equal contributions)