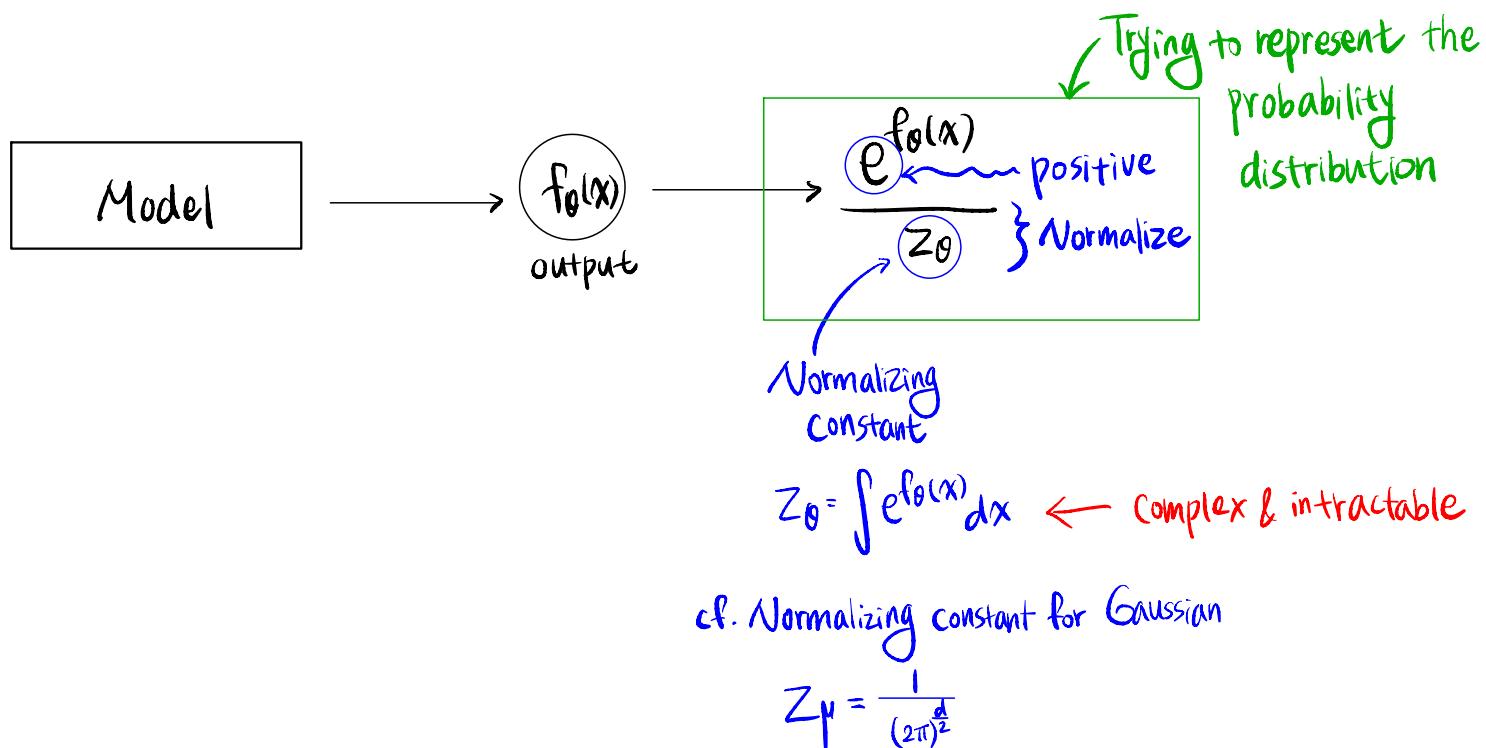


Diffusion and Score-based Generative Models



< Tackling the intractable normalizing constant >

■ Approximating the normalizing constant

- Energy-based models ← inaccurate probability evaluation

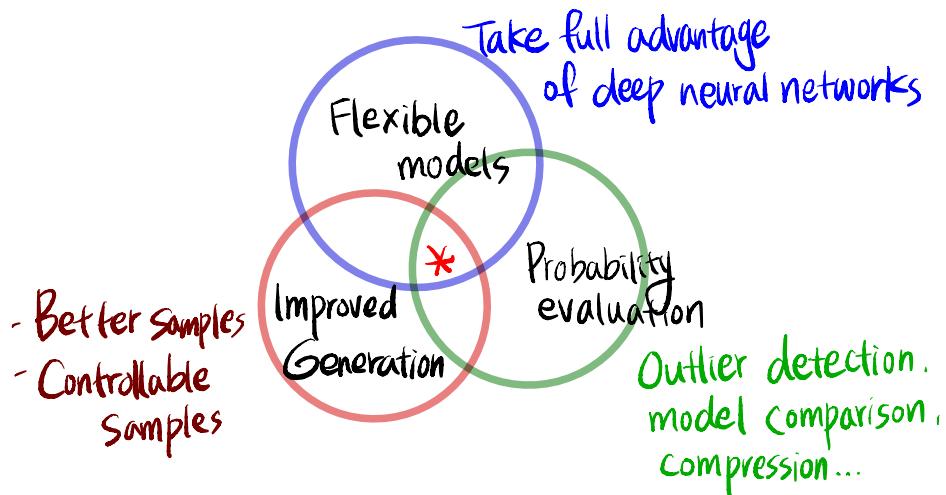
■ Using restricted neural network models

- Autoregressive models
 - Normalizing flow models
 - Variational auto-encoders
- } ← restricted model family
(=limit flexibility)

■ Modeling the Generation Process Only

- Generative Adversarial Networks (GANs) ← cannot evaluate probabilities

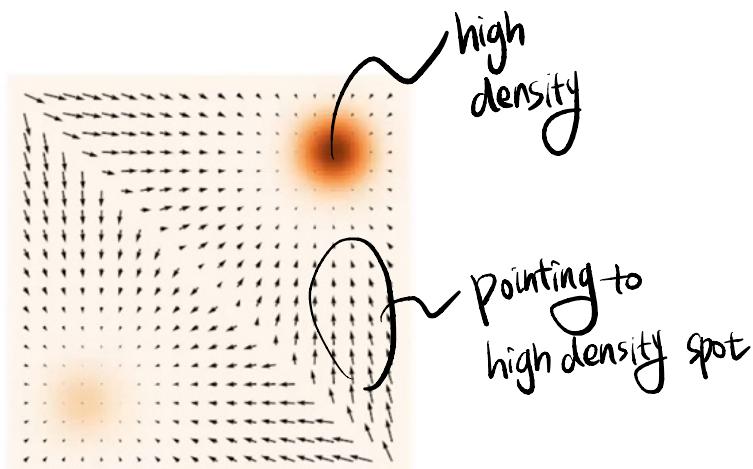
<Desiderata of better generative modeling>



<Our proposal: working with score functions>

$p(x)$
Probability density function

$\nabla_x \log p(x)$
(Stein) Score function



Score vs. density function (color coded)
gives direction to where the density function grows quickly

(\rightarrow Score functions preserve every information about density function)

[Score-based generative modeling: outline]

Flexible models



- Bypass the normalizing constant
- Principled statistical methods

Improved Generation



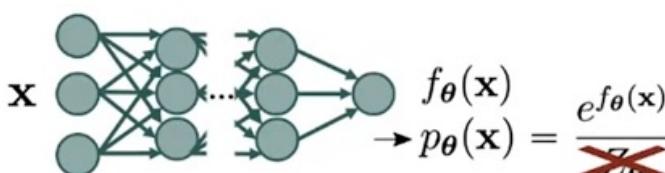
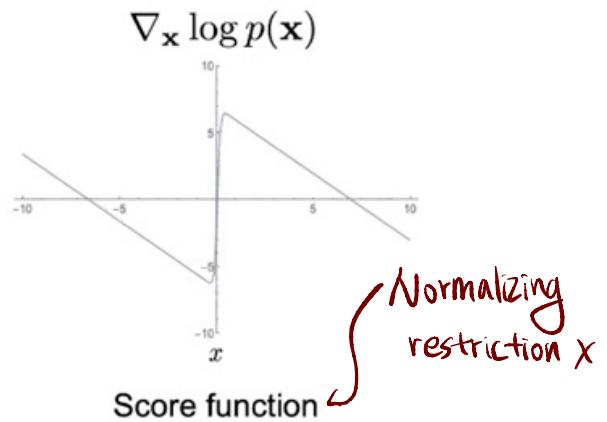
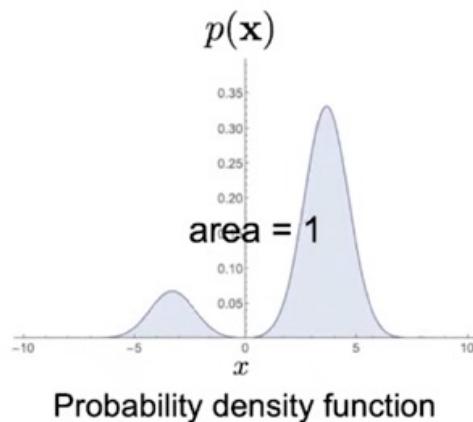
- Higher sample quality than GANs
- Controllable generation

Probability evaluation



- Accurate probability evaluation
- Better estimation of data probabilities

* Score functions bypass the normalizing constant



$$\begin{aligned} \nabla_x \log p_{\theta}(x) &= \nabla_x f_{\theta}(x) - [\nabla_x \log Z_{\theta}] \\ &= \nabla_x f_{\theta}(x) \\ &= [s_{\theta}(x)] \end{aligned}$$

Score model

gradient of constant

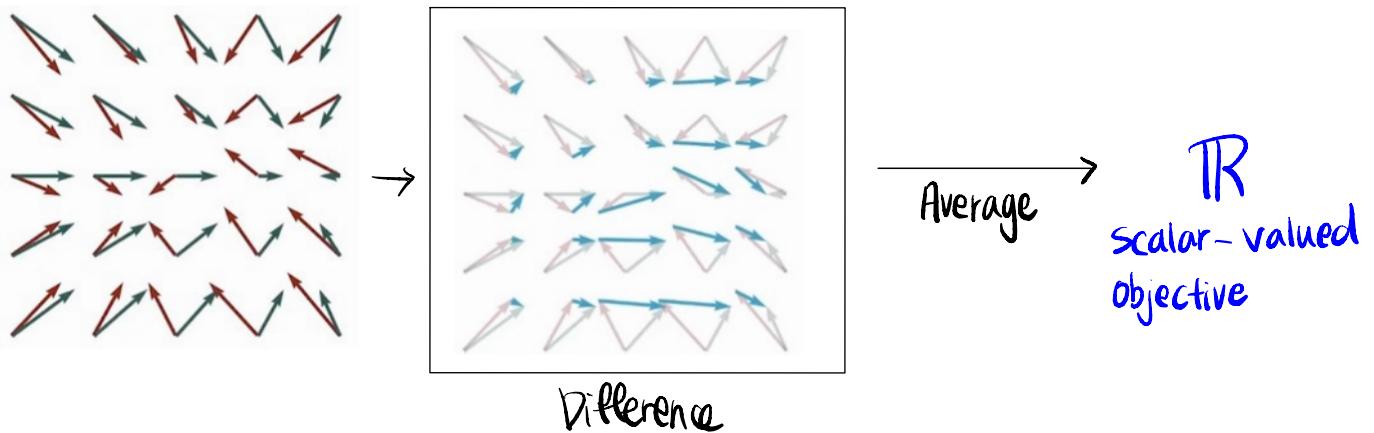
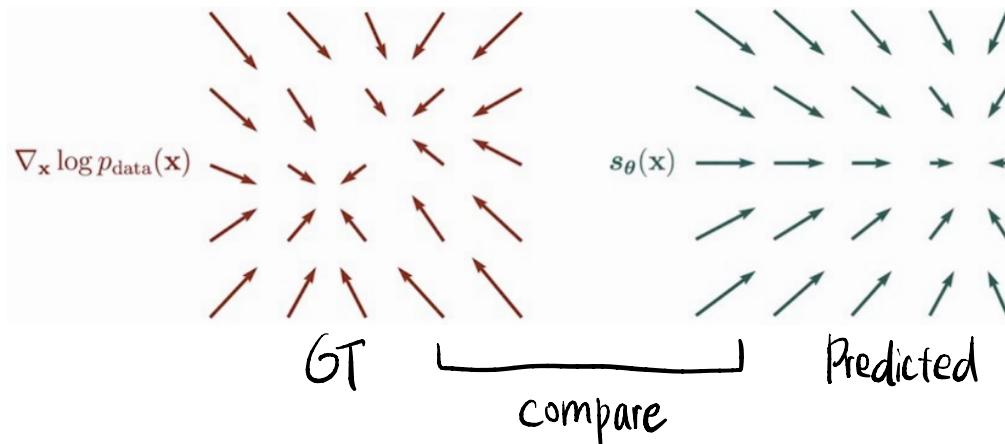
[Score models can be estimated from data]

Given: $\{x_1, x_2, \dots, x_N\}^{i.i.d} \sim p_{\text{data}}(x)$

Goal: $\nabla_x \log p_{\text{data}}(x)$

Score Model: $s_\theta(x) : \mathbb{R}^d \rightarrow \mathbb{R}^d \approx \nabla_x \log p_{\text{data}}(x)$

Objective: How to compare two vector fields of scores?



Objective:

$$\frac{1}{2} \mathbb{E}_{p_{\text{data}}(x)} \left[\|\nabla_x \log p_{\text{data}}(x) - s_\theta(x)\|_2^2 \right]$$

(Fisher divergence)

Cannot be computed because we don't know the ground truth value

Score Matching [Hyvärinen 2005]:

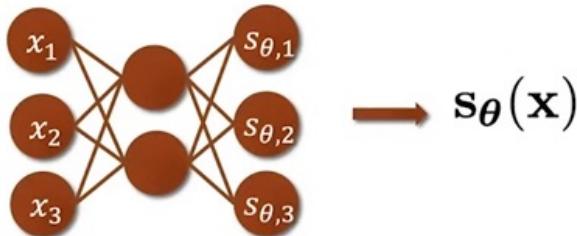
$$\mathbb{E}_{P_{\text{data}}(x)} \left[\frac{1}{2} \|S_\theta(x)\|_2^2 + \text{trace}(\underbrace{\nabla_x S_\theta(x)}_{\text{Jacobian of } S_\theta(x)}) \right]$$

) can be approximated using the empirical mean

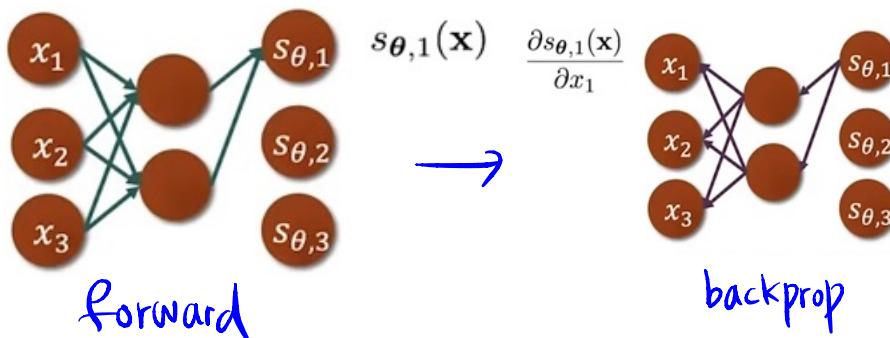
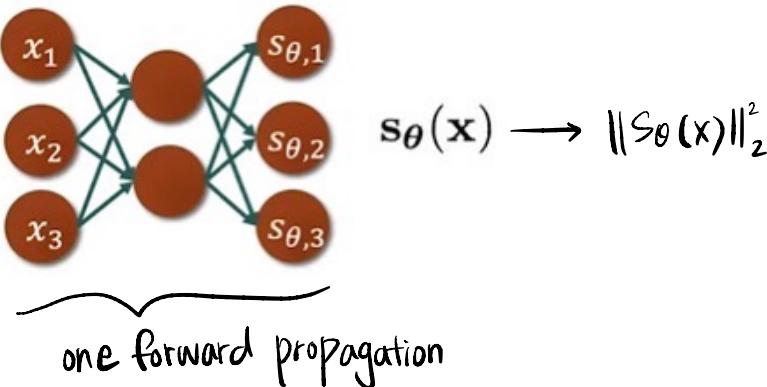
$$\approx \frac{1}{N} \sum_{i=1}^N \left[\frac{1}{2} \|S_\theta(x_i)\|_2^2 + \text{trace}(\nabla_x S_\theta(x_i)) \right]$$

Score Matching is not scalable

- Deep score models



- Need to compute: $\|S_\theta(x)\|_2^2$ and $\text{trace}(\nabla_x S_\theta(x))$



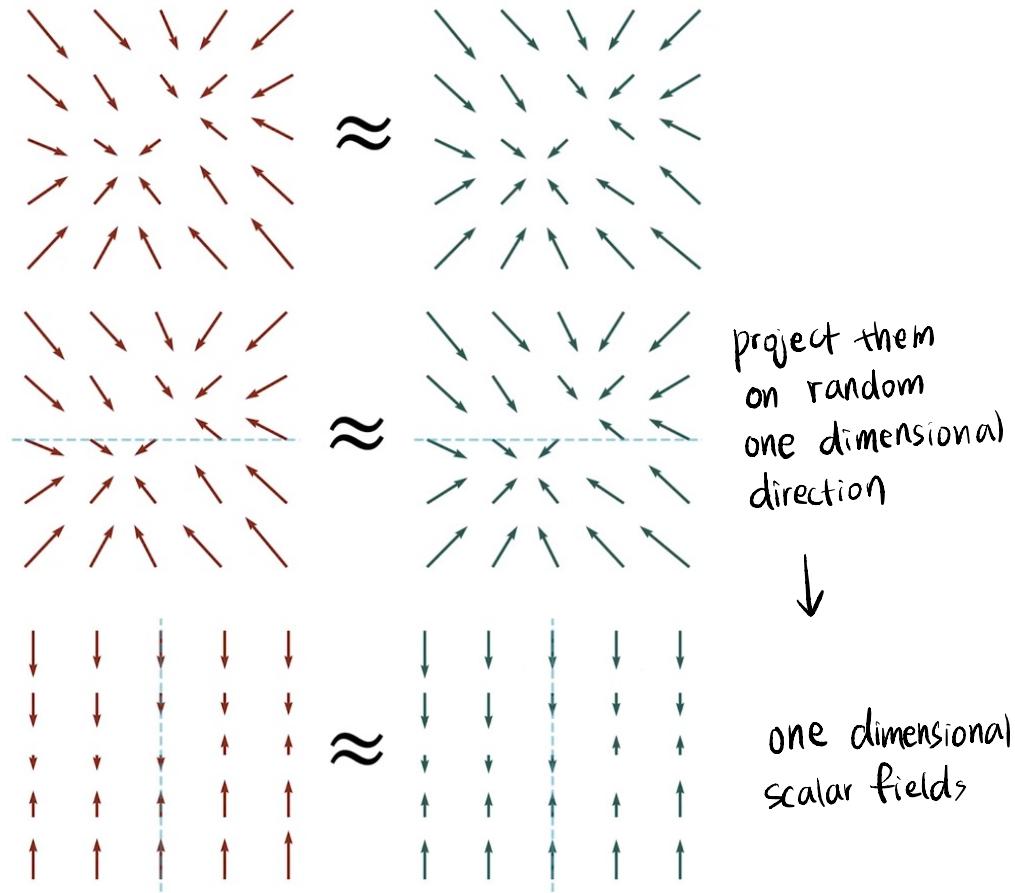
$O(\# \text{ dimensions of } x)$

$$\nabla_x S_\theta(x) = \begin{pmatrix} \frac{\partial s_{\theta,1}(x)}{\partial x_1} & \frac{\partial s_{\theta,1}(x)}{\partial x_2} & \frac{\partial s_{\theta,1}(x)}{\partial x_3} \\ \frac{\partial s_{\theta,2}(x)}{\partial x_1} & \frac{\partial s_{\theta,2}(x)}{\partial x_2} & \frac{\partial s_{\theta,2}(x)}{\partial x_3} \\ \frac{\partial s_{\theta,3}(x)}{\partial x_1} & \frac{\partial s_{\theta,3}(x)}{\partial x_2} & \frac{\partial s_{\theta,3}(x)}{\partial x_3} \end{pmatrix}$$

Σ : trace

Sliced score matching

- Intuition: one dimensional problem should be easier
- Idea: project onto random directions



Randomized Objective: Sliced Fisher Divergence

$$\frac{1}{2} \mathbb{E}_{P_V} \mathbb{E}_{P_{\text{data}}(x)} \left[(\nabla^T \nabla \log_{P_{\text{data}}} (x) - \nabla^T S_\theta(x))^2 \right]$$

V : projection direction $\in \mathbb{R}^d$

P_V : distribution of projection direction

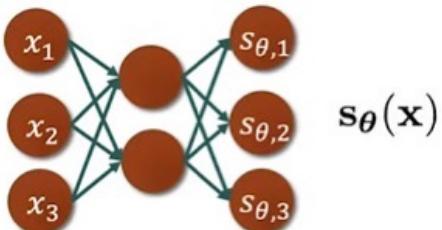
Integration by parts \rightarrow Sliced Score Matching

$$\mathbb{E}_{P_V} \mathbb{E}_{P_{\text{data}}(x)} \left[\nabla^T \nabla S_\theta(x) V + \frac{1}{2} (\nabla^T S_\theta(x))^2 \right]$$

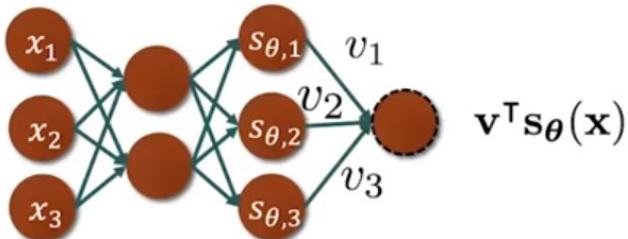
scalable!

Computing Jacobian-vector products is Scalable

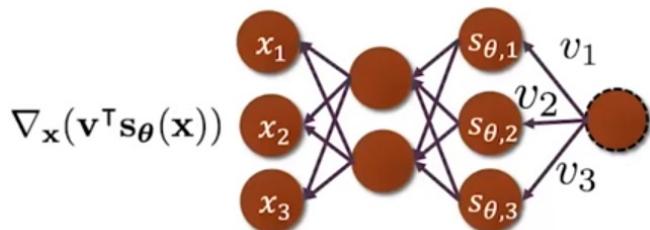
$$V^T \nabla_x S_\theta(x) V = V^T \nabla_x (V^T S_\theta(x))$$



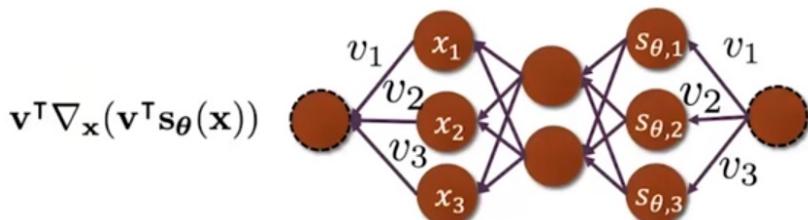
$$V^T \nabla_x (V^T S_\theta(x))$$



$$V^T \nabla_x (V^T S_\theta(x))$$



$$V^T \nabla_x (V^T S_\theta(x))$$



$$V^T \nabla_x (V^T S_\theta(x))$$

One
Backprop!
(efficient)

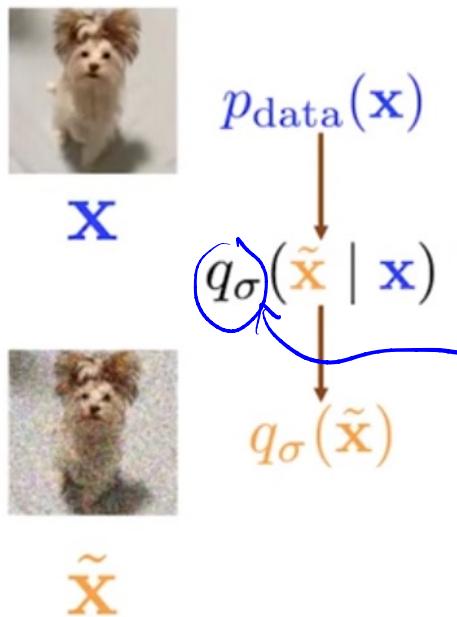
How sliced score matching works in practice

- Sample a minibatch of data points $\{x_1, x_2, \dots, x_n\} \sim P_{\text{data}}(x)$
- Sample a minibatch of projection directions $\{v_1, v_2, \dots, v_n\} \sim P_v$
- Estimate the sliced score matching loss with empirical means

$$\frac{1}{n} \sum_{i=1}^n [v_i^\top \nabla_x s_\theta(x_i) v_i + \frac{1}{2} (v_i^\top s_\theta(x_i))^2]$$

- The projection distribution is typically Gaussian or Rademacher
- Stochastic gradient descent
- Can use more projections per datapoint to boost performance

Denoising score matching



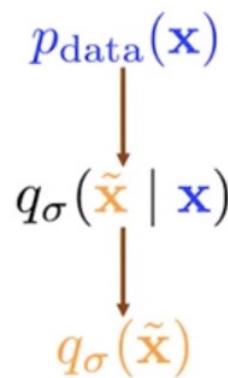
- Matching the score of a noise-perturbed distribution

$$\frac{1}{2} \mathbb{E}_{q_\theta(\tilde{x})} [\|\nabla_{\tilde{x}} \log q_\theta(\tilde{x}) - s_\theta(\tilde{x})\|_2^2]$$

- Denoising score matching

$$\frac{1}{2} \mathbb{E}_{\substack{\mathbb{E}_{q_\theta(\tilde{x}|x)} \\ p_{\text{data}}(x)}} [\|\nabla_{\tilde{x}} \log q_\theta(\tilde{x}|x) - s_\theta(\tilde{x})\|_2^2]$$

* Cannot estimate scores of noise-free distributions
scalable



- Denoising score matching:
matching the score of a noise-perturbed distribution

$\tilde{\mathbf{x}}$

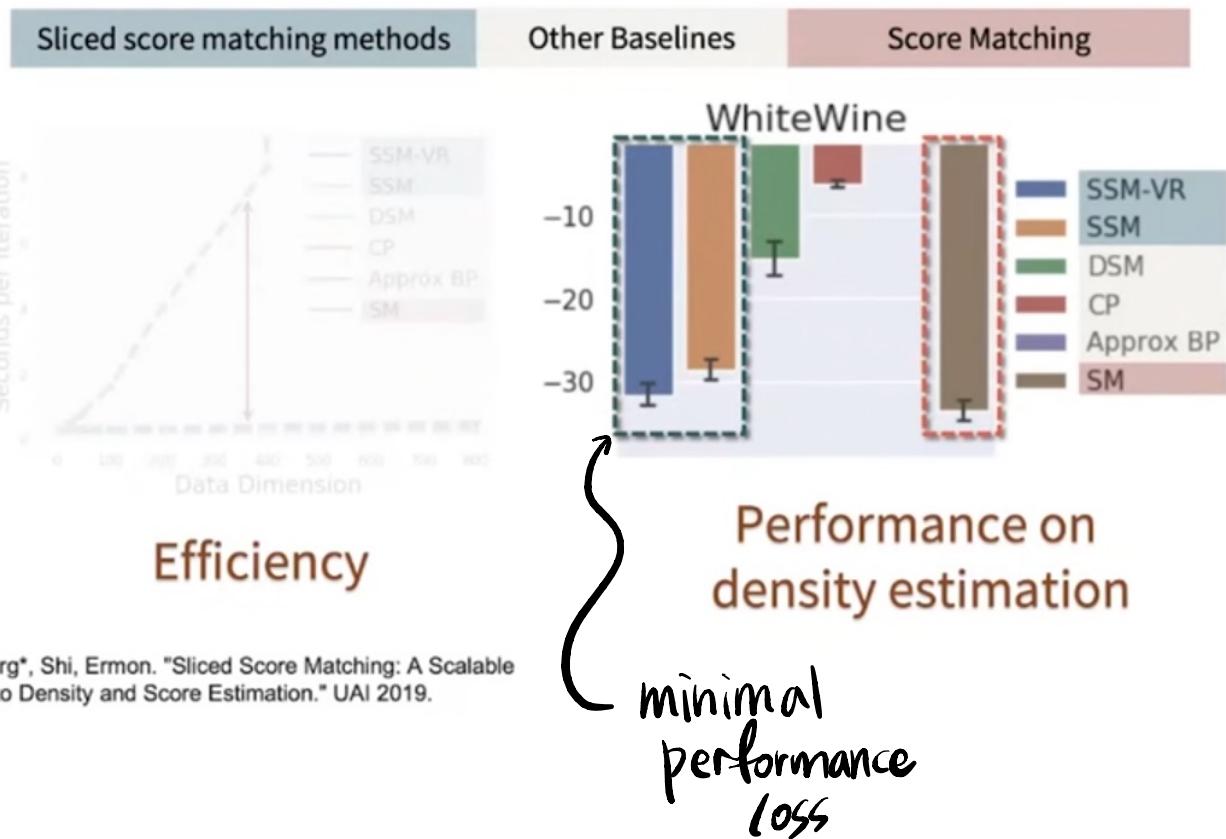
$\langle \text{Procedure} \rangle$

- Sample a minibatch of datapoints $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \sim p_{\text{data}}(\mathbf{x})$
- Sample a minibatch of perturbed datapoints $\{\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_n\} \sim q_\sigma(\tilde{\mathbf{x}})$
 $\tilde{\mathbf{x}}_i \sim q_\sigma(\tilde{\mathbf{x}}_i | \mathbf{x}_i)$
- Estimate the denoising score matching loss with empirical means

$$\frac{1}{2n} \sum_{i=1}^n \left[\| S_\sigma(\tilde{\mathbf{x}}_i) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}}_i | \mathbf{x}_i) \|_2^2 \right]$$
- If Gaussian perturbation

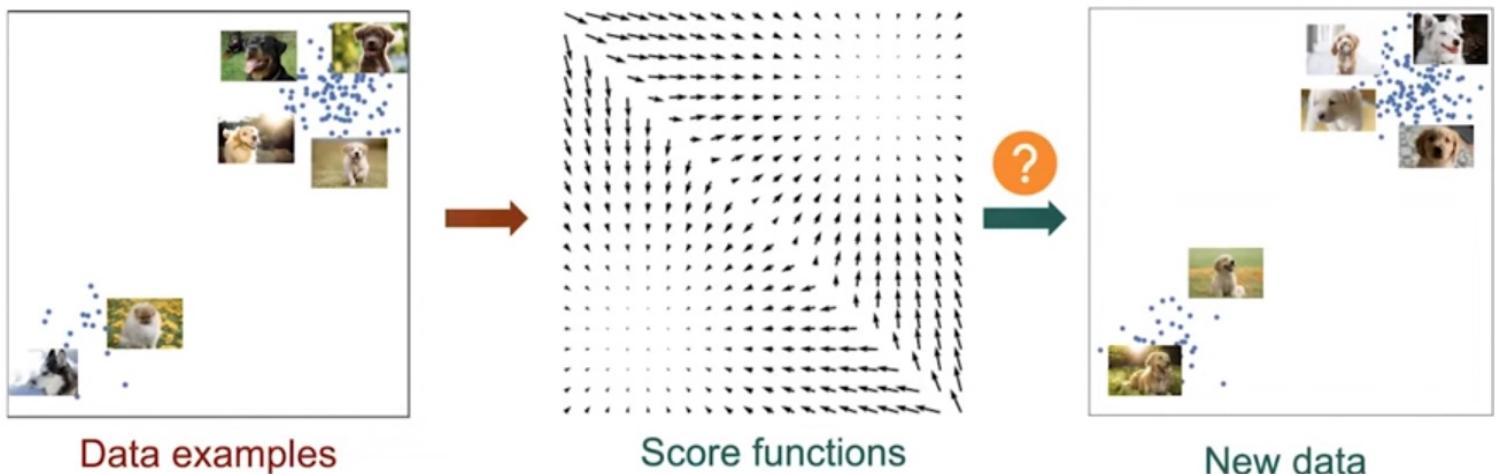
$$\frac{1}{2n} \sum_{i=1}^n \left[\| S_\sigma(\tilde{\mathbf{x}}_i) + \frac{\tilde{\mathbf{x}}_i - \mathbf{x}_i}{\sigma^2} \|_2^2 \right]$$
- Stochastic gradient descent
- Need to choose a very small σ !

Results: Sliced Score Matching for EBMs

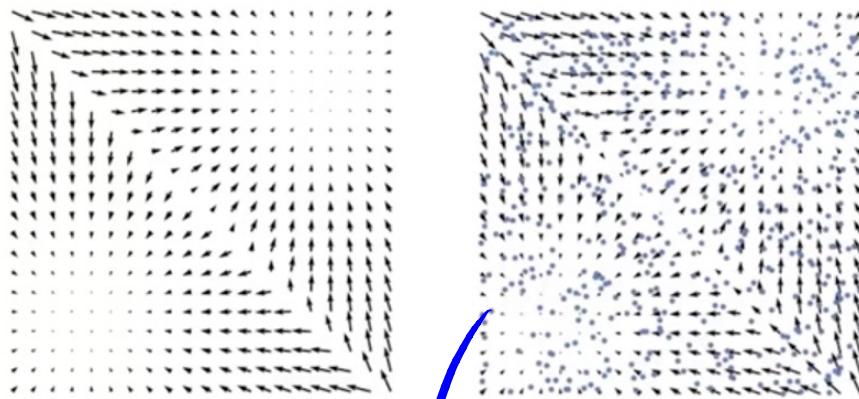


Song*, Garg*, Shi, Ermon. "Sliced Score Matching: A Scalable Approach to Density and Score Estimation." UAI 2019.

Sampling from score functions

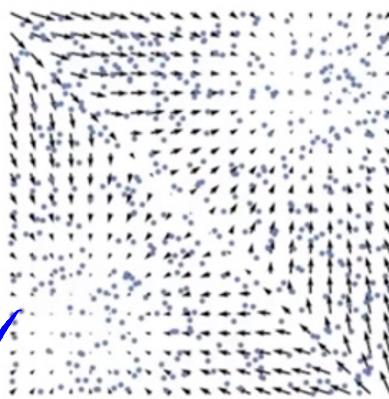


Sampling from score functions: Langevin dynamics



Score function

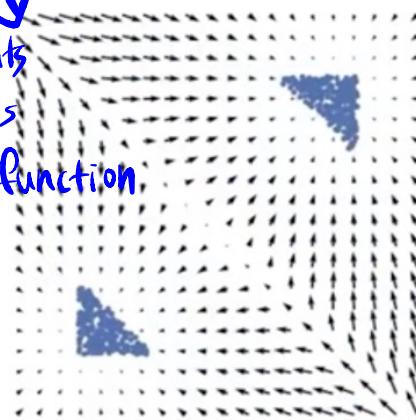
$$s_\theta(x)$$



Follow the scores

$$\tilde{x}_{t+1} \leftarrow \tilde{x}_t + \frac{\epsilon}{2} s_\theta(\tilde{x}_t)$$

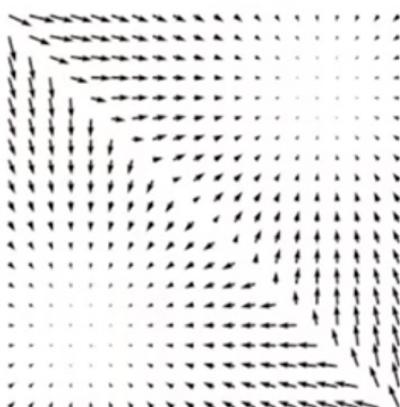
moving the points
following the directions
predicted by the score function
→ points will collapse
into each other



follow
noisy-version of
score functions

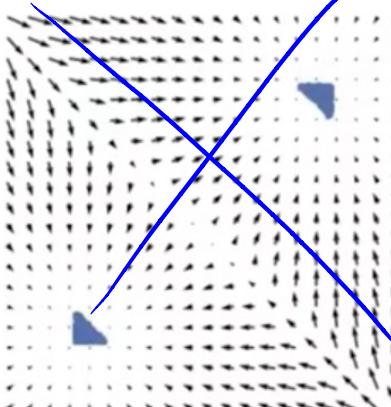
→ Correct samples guaranteed

Sampling from score functions: Langevin dynamics



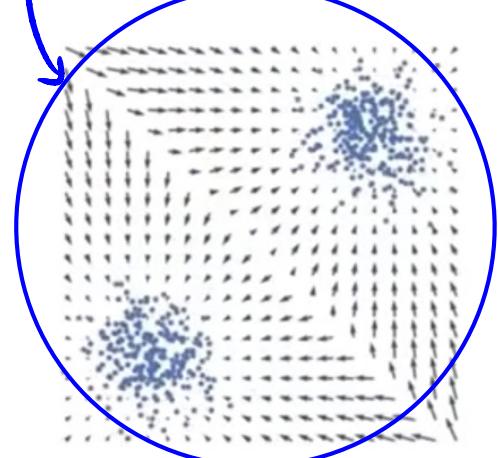
Score function

$$s_\theta(x)$$



Follow the scores

$$\tilde{x}_{t+1} \leftarrow \tilde{x}_t + \frac{\epsilon}{2} s_\theta(\tilde{x}_t)$$



Follow the noisy scores

$$z_t \sim \mathcal{N}(0, I)$$

$$\tilde{x}_{t+1} \leftarrow \tilde{x}_t + \frac{\epsilon}{2} s_\theta(\tilde{x}_t) + \sqrt{\epsilon} z_t$$

inject
Gaussian
Noise

Langevin dynamics sampling

- Sample from $p(x)$ using only the score $D_x \log p(x)$

- Initialize $x^0 \sim \pi(x)$
 \approx prior distribution

- Repeat for $t \leftarrow 1, 2, \dots, T$

$$z^t \sim N(0, 1)$$

$$x^t \leftarrow x^{t-1} + \frac{\epsilon}{2} D_x \log p(x^{t-1}) + \sqrt{\epsilon} z^t$$

- If $\epsilon \rightarrow 0$ and $T \rightarrow \infty$, we are guaranteed to have $x^T \sim p(x)$

- Langevin dynamics + score estimation

$$S_\theta(x) \approx D_x \log p(x)$$

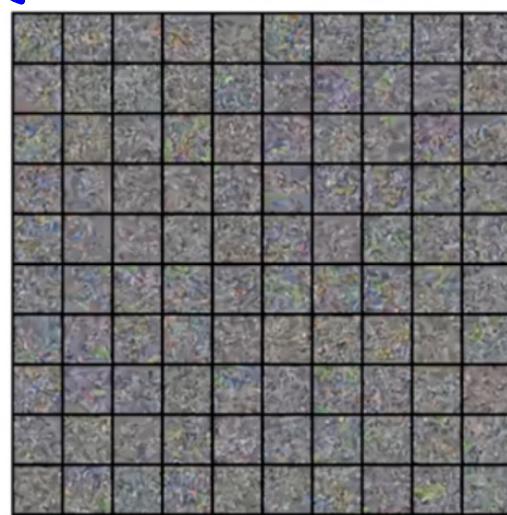
doesn't work well in practice
with naive approach

Score matching + Langevin dynamics

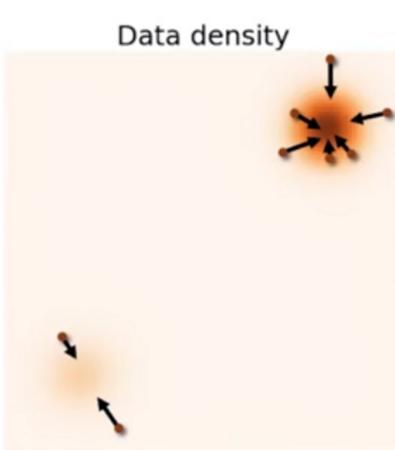
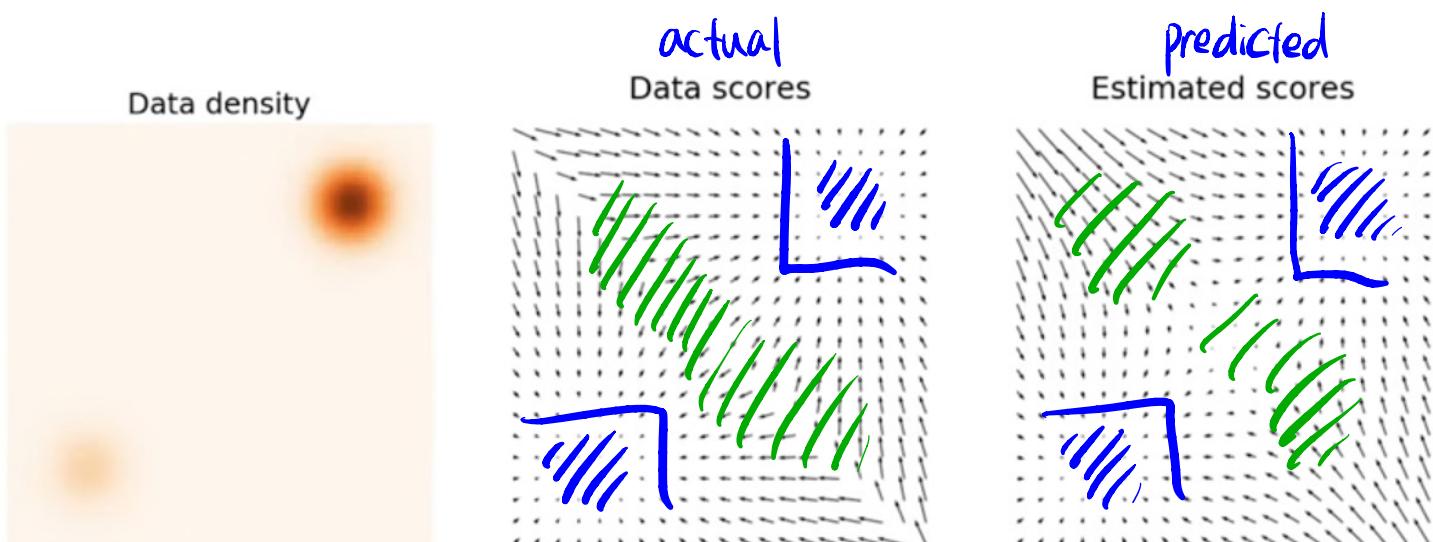
CIFAR-10 data



Model samples

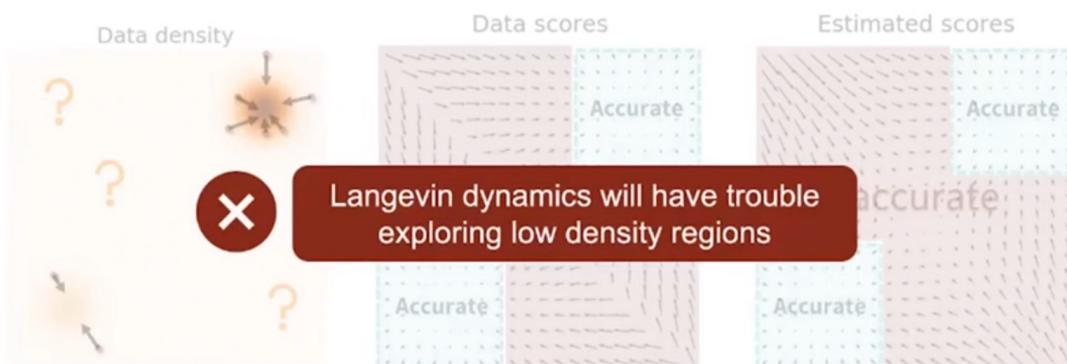


Challenge in low data density regions



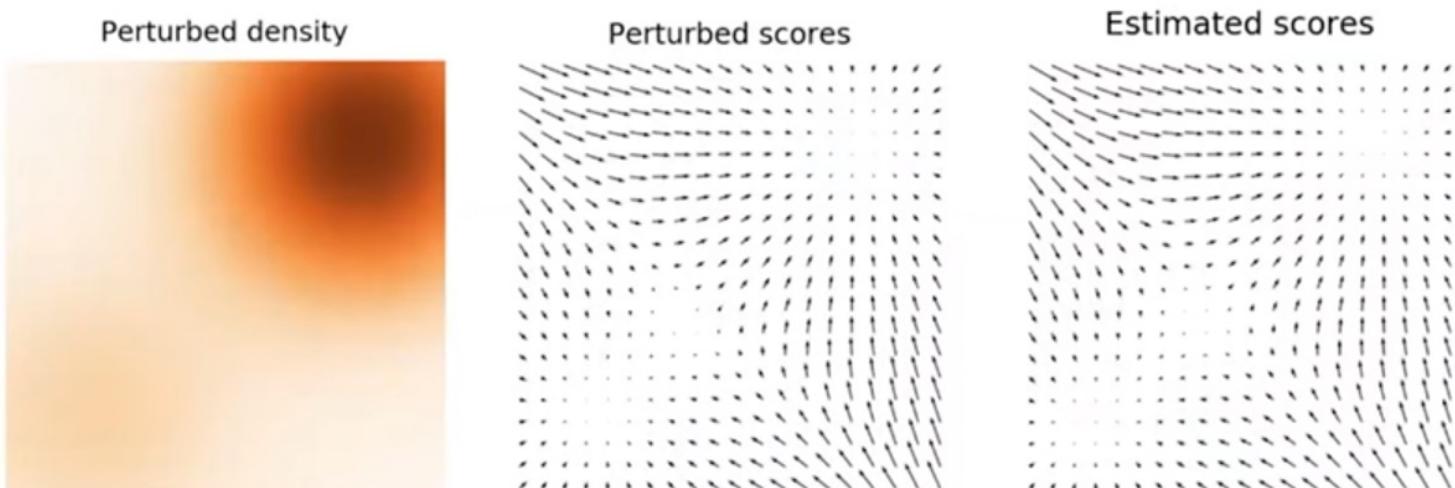
- : Score matching compares the difference between ground truth only at data samples
- In low density regions, we don't have enough samples
- \times enough information to infer the true score functions

Challenge in low data density regions



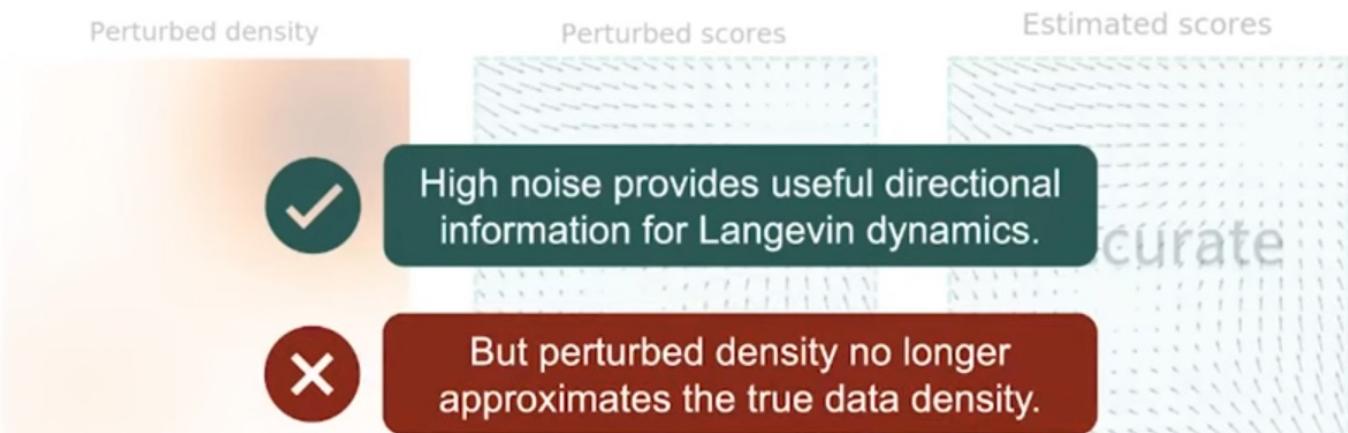
$$\frac{1}{2} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - s_{\theta}(\mathbf{x})\|_2^2]$$

Improving score estimation by adding noise



After adding enough Gaussian noise, perturb the datapoints to everywhere in space
→ the size of lower density regions become smaller

Improving score estimation by adding noise

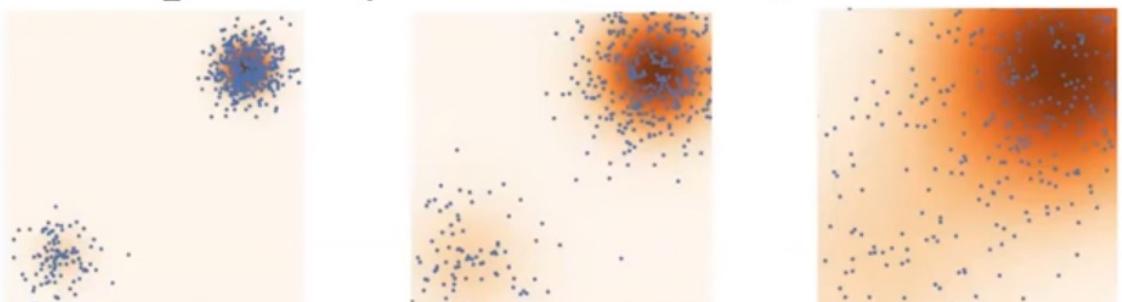


[Song and Ermon, NeurIPS 2019 (oral)]

Using multiple noise levels

Std: $\sigma_1 < \text{std: } \sigma_2 < \text{std: } \sigma_3$

Data
 $\sim N(0, \sigma)$



Corresponding
noisy-data
density
 $p_\sigma(x)$

$p_{\sigma_1}(x)$

$p_{\sigma_2}(x)$

$p_{\sigma_3}(x)$

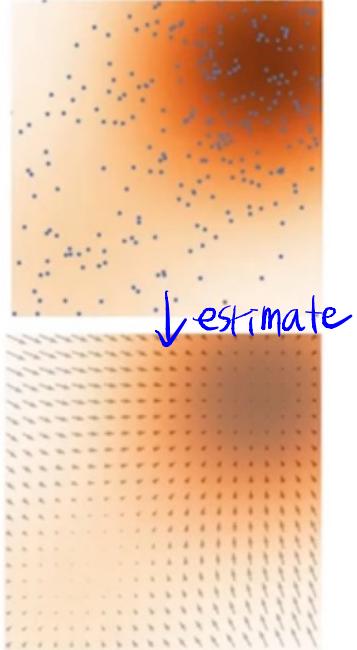
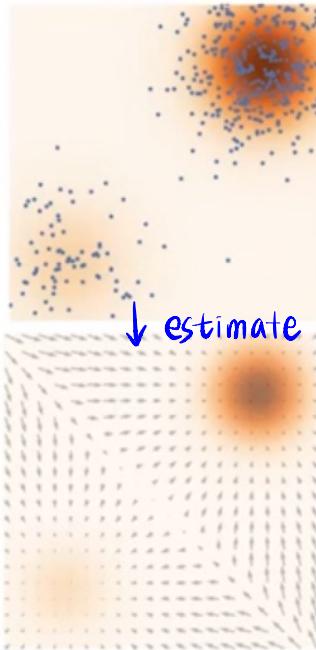
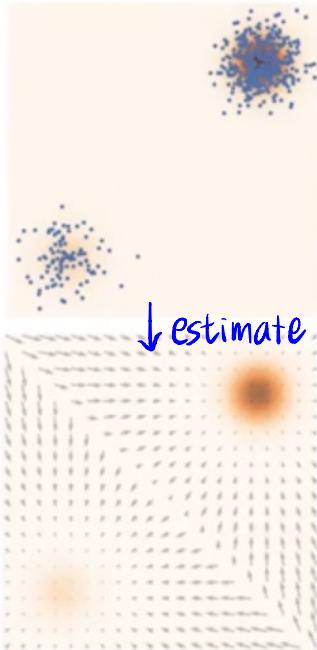
noisy
Data

Using multiple noise levels

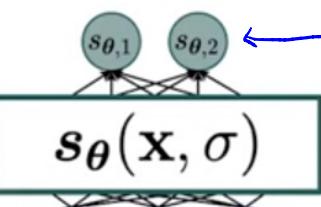
$p_{\sigma_1}(x)$

$p_{\sigma_2}(x)$

$p_{\sigma_3}(x)$



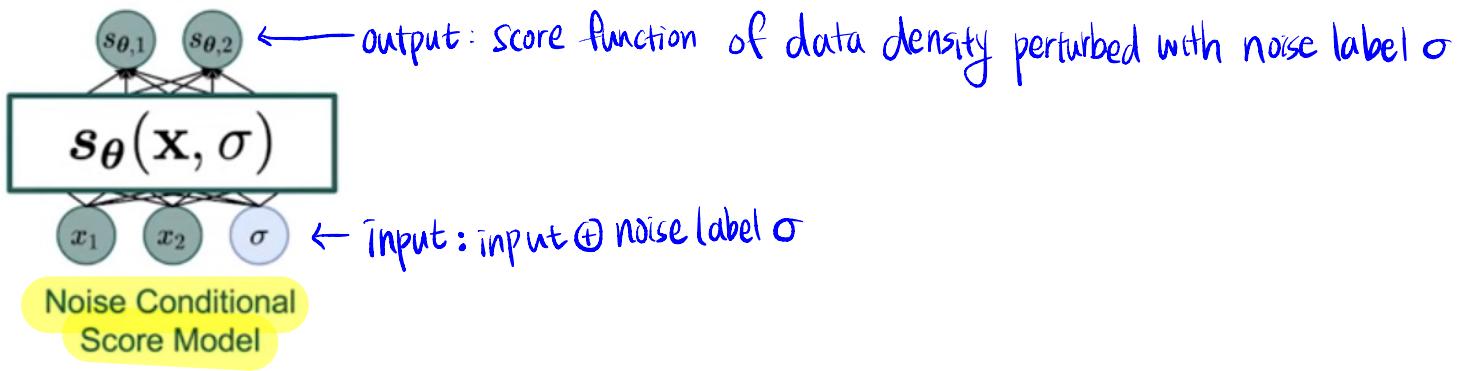
estimate
the underlying
score
function



← output: score function of data density perturbed with noise label σ

← Input: input \oplus noise label σ

Noise Conditional
Score Model



how
can we
train this
model
?

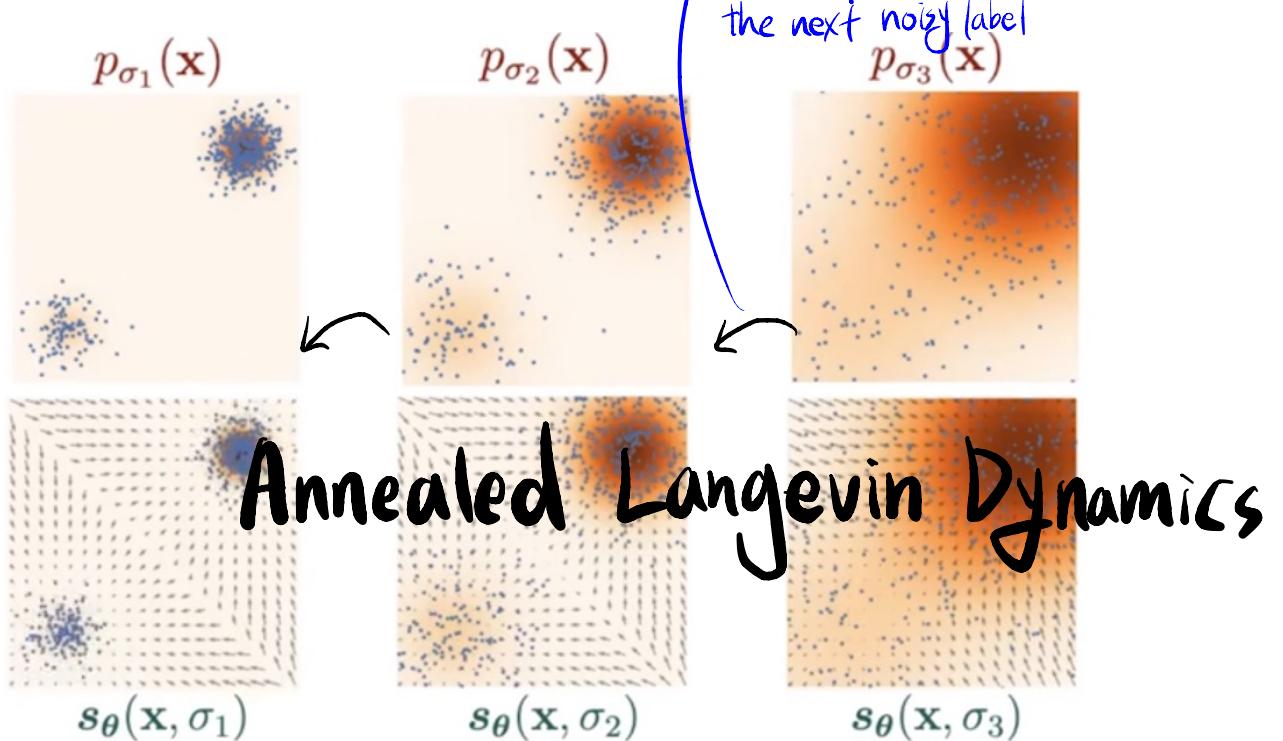
\Rightarrow
leverage
the idea of
Score-matching

$$\frac{1}{N} \sum_{i=1}^N \lambda(\sigma_i) \mathbb{E}_{p_{\sigma_i}(x)} [\| \nabla_x \log p_{\sigma_i}(x) - s_{\theta}(x, \sigma_i) \|_2^2]$$

positive
weighting
function

Score matching loss
for all noise labels σ_i

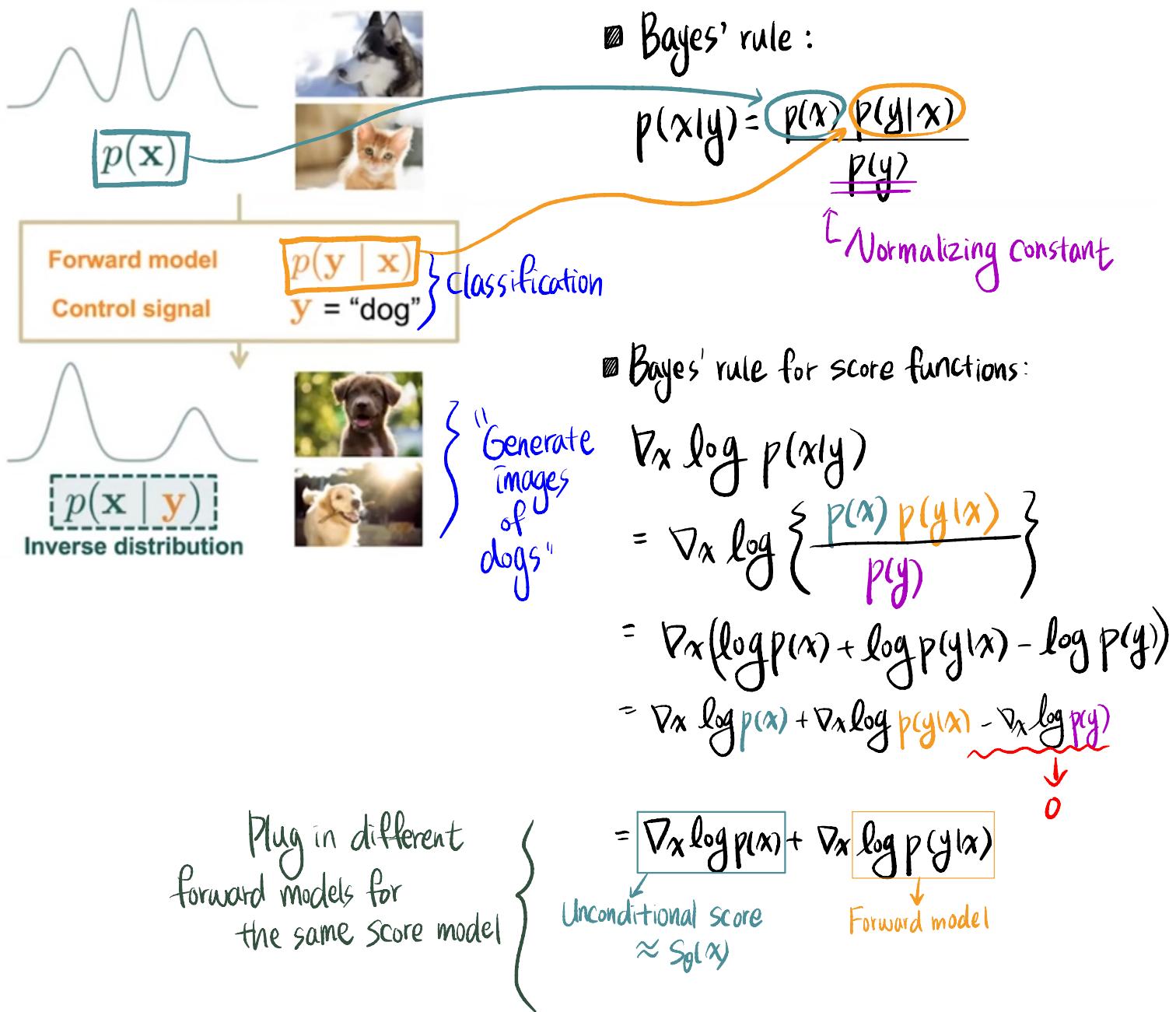
- A generalization to the training objective of diffusion probabilistic models
- First unveiled by DDPM



Quantitative results on CIFAR0-10

Model	Inception	FID
CIFAR-10 Unconditional		
PixelCNN [59]	4.60	65.93
PixelIQN [42]	5.29	49.46
EBM [12]	6.02	40.58
WGAN-GP [18]	$7.86 \pm .07$	36.4
MoLM [45]	$7.90 \pm .10$	18.9
SNGAN [36]	$8.22 \pm .05$	21.7
ProgressiveGAN [25]	$8.80 \pm .05$	-
NCSN (Ours)	<u>$8.87 \pm .12$</u>	25.32

Control the generation process



Controllable Generation: class-conditional generation

- y is the **class label**
 - $p_t(y | x)$ is a time-dependent classifier (classifier-guidance)



Image inpainting

Given

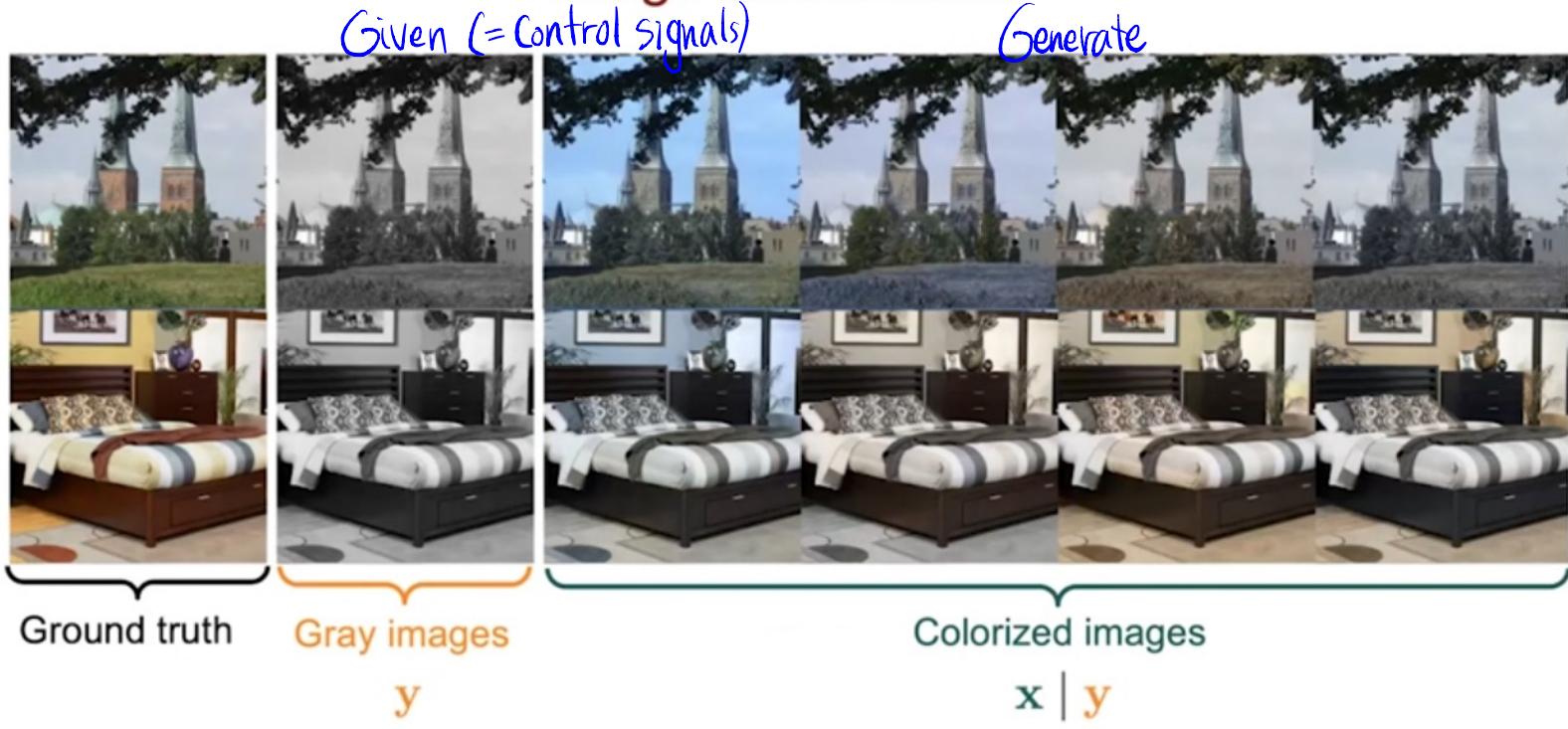
Generate



Forward model can be directly specified

[Song et al. ICLR 2021 (Outstanding Paper Award)]

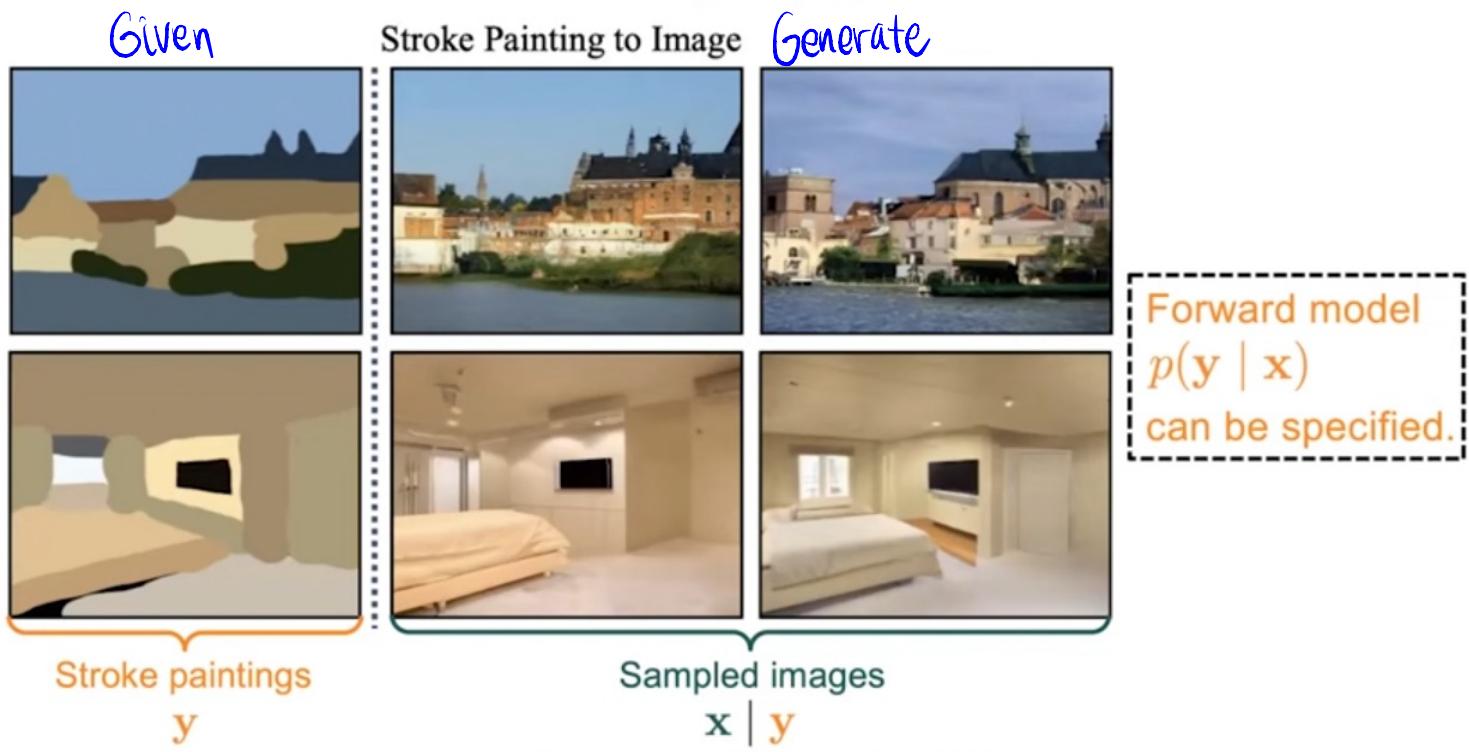
Image colorization



Forward model can be directly specified

→ One Score-based generative model for both image inpainting and colorization

Stroke to image synthesis



[Meng, He, Song, Song, Wu, Zhu, Ermon. ICLR 2022]

Language-guided image generation

y

(Prompt)

Treehouse in the
style of Studio
Ghibli animation

Given

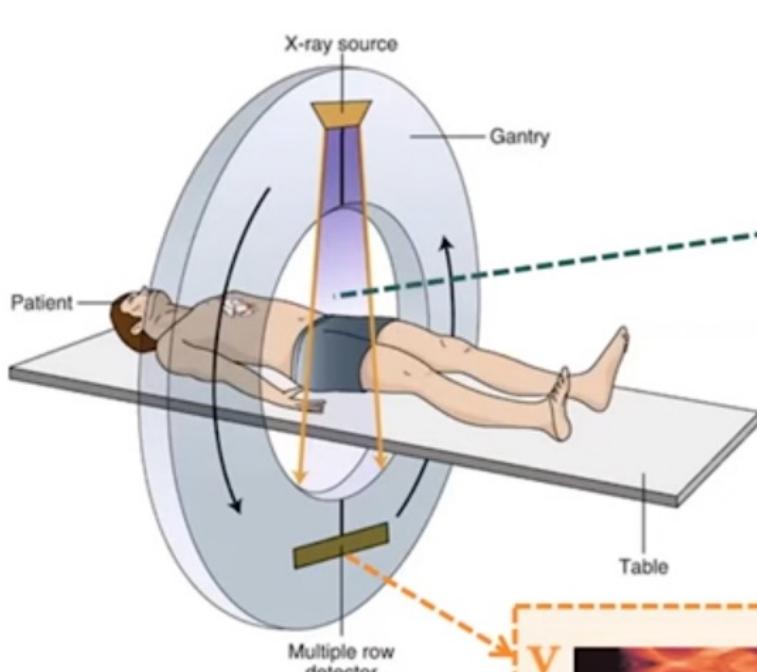
$x | y$



[Work by @danielrussruss]

Generate $x | y$

Medical image reconstruction



Cross-sectional image



$x | y$

Generate

Sparse-view computed
tomography (CT)



Given

Forward model $p(y|x)$ is given by physical simulation

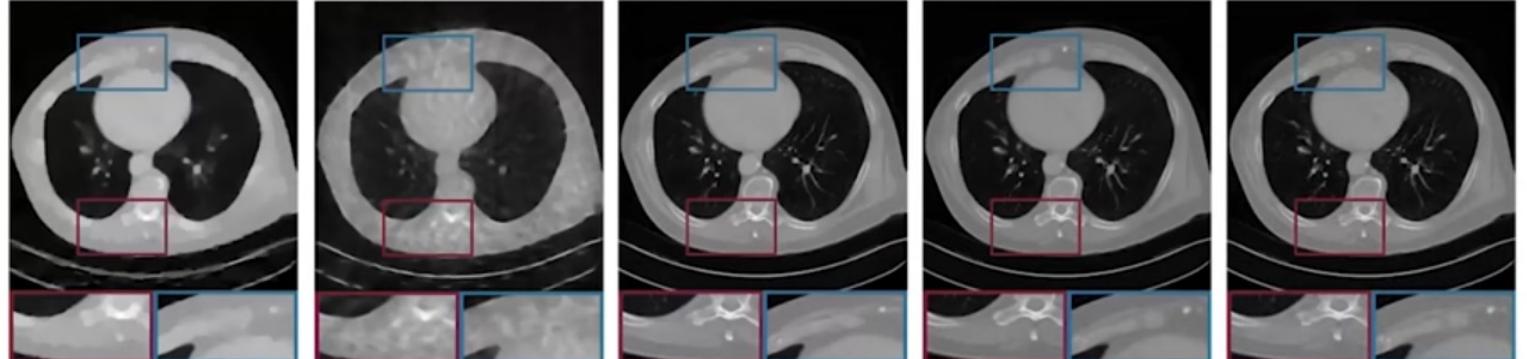
PSNR : \uparrow : better

SSIM : \downarrow : better

Medical image reconstruction

Sparse-view CT (just 23 projections)

PSNR: 20.30, SSIM: 0.778 PSNR: 22.78, SSIM: 0.603 PSNR: 31.76, SSIM: 0.882 PSNR: 35.23, SSIM: 0.912



FISTA

Neumann network

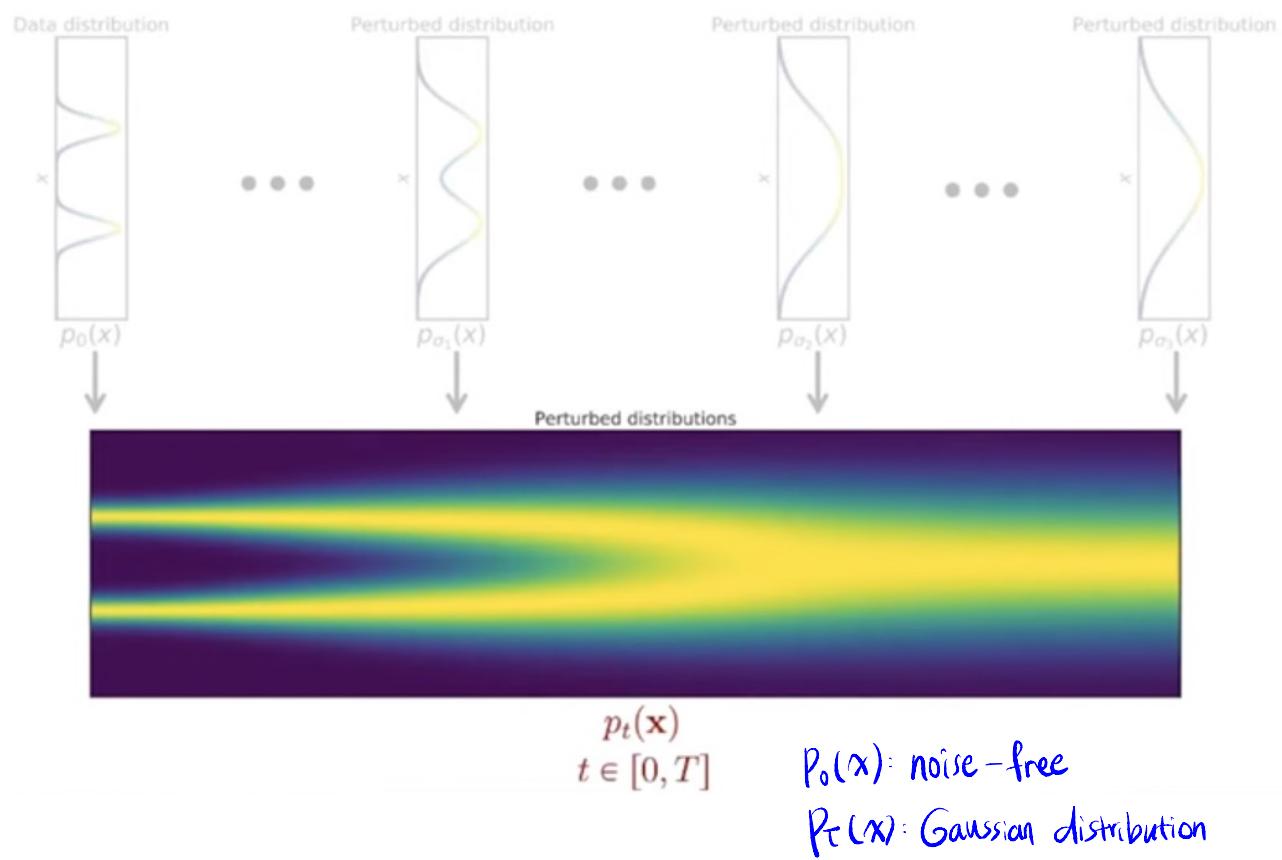
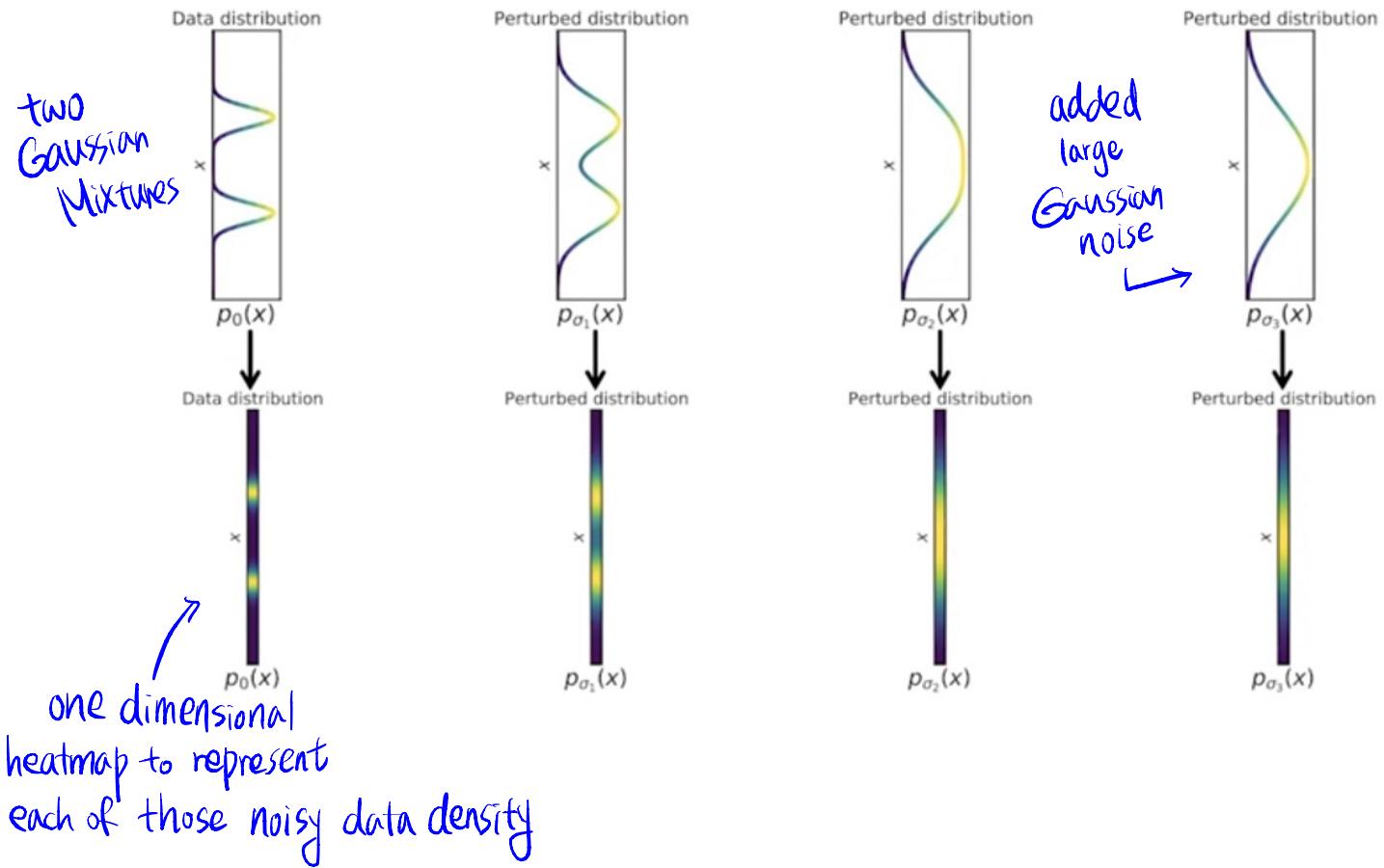
SIN-PRN

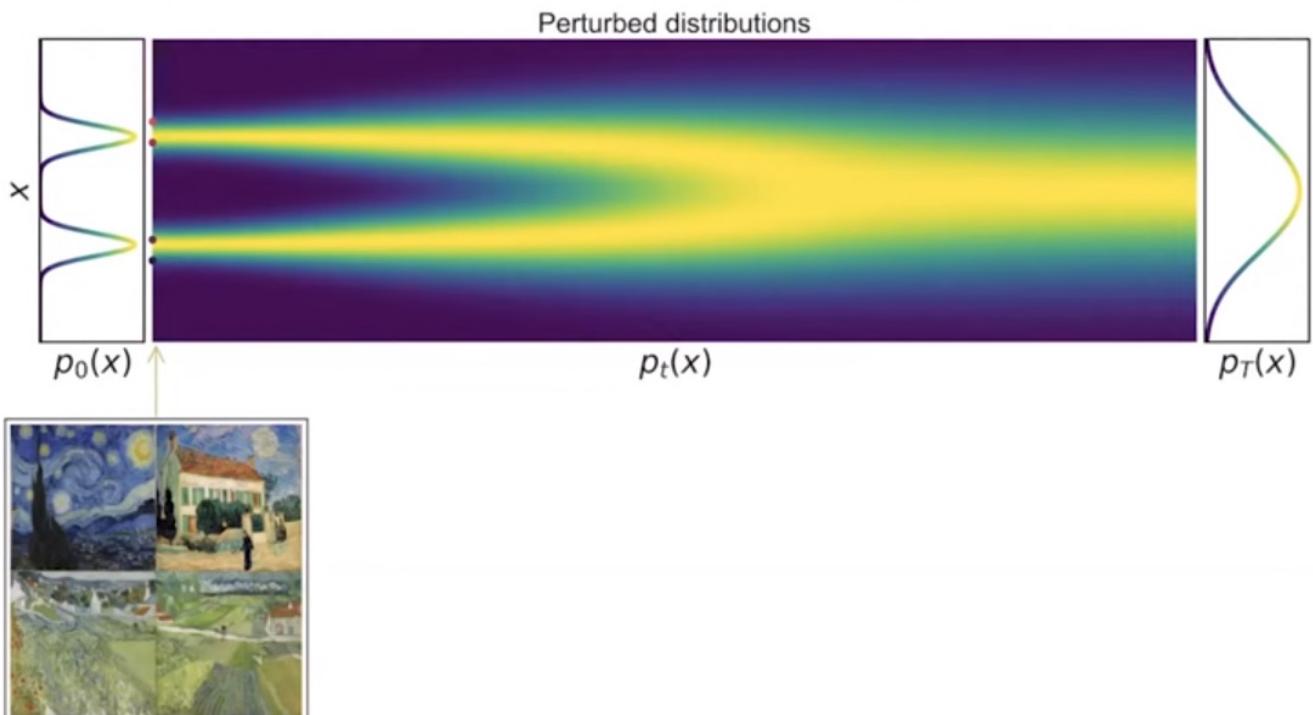
Ours

Ground truth

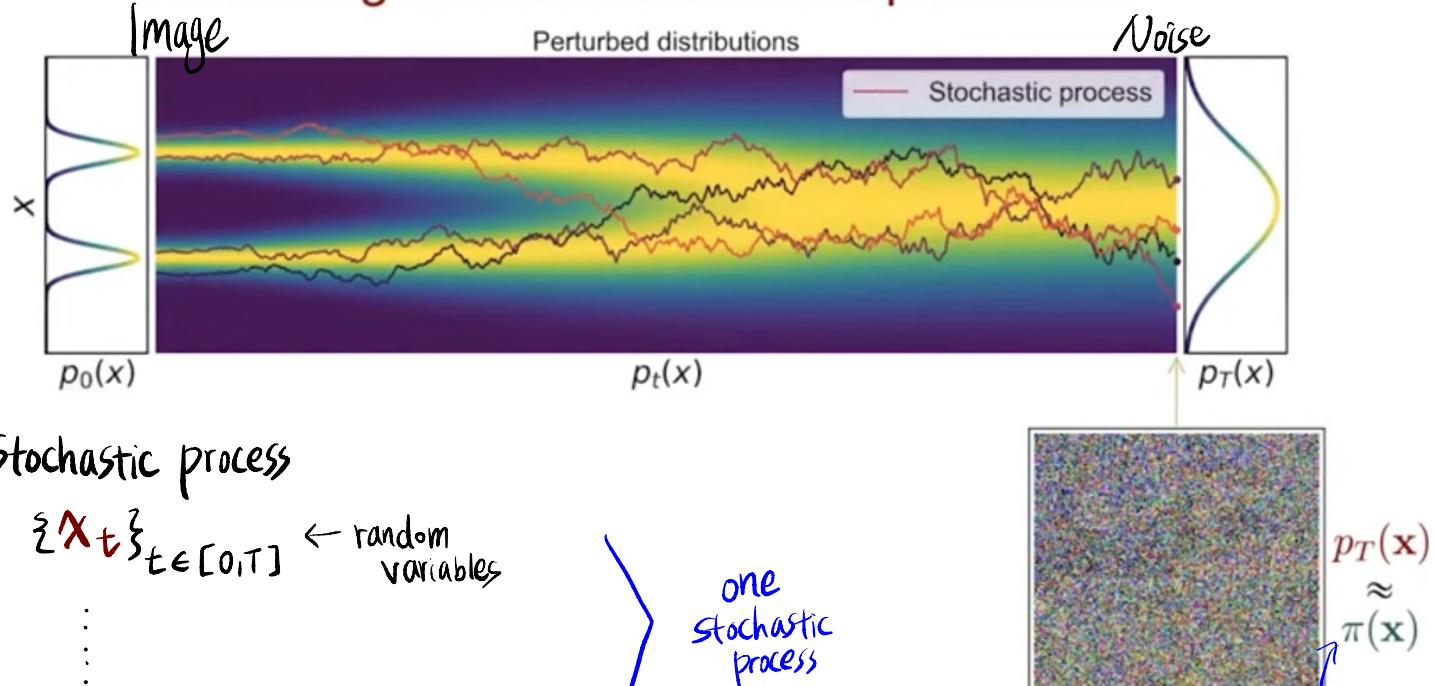
[Probability Evaluation]

Infinite noise levels





Perturbing data with stochastic processes



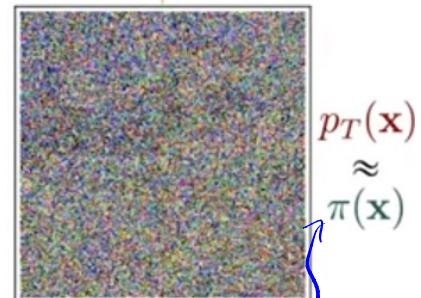
Stochastic process

$$\{X_t\}_{t \in [0, T]} \leftarrow \text{random variables}$$

Corresponding probability densities for random variables

$$\{p_t(x)\}_{t \in [0, T]}$$

one stochastic process
↔
infinite number of probability densities



Samples from simple Gaussian distribution

How to choose the right stochastic process such that
it represents an infinite number of noisy data densities

Stochastic process

$$\{X_t\}_{t \in [0, T]} \leftarrow \text{random variables}$$



Stochastic Differential Equation (SDE)

$$dx_t = f(x_t, t) dt + g(t) dW_t$$

Deterministic
drift

Infinitesimal
noise

} without
loss of generality

WLOG: Toy SDE

$$dx_t = \sigma(t) dW_t$$

stochastic
term

= a continuous

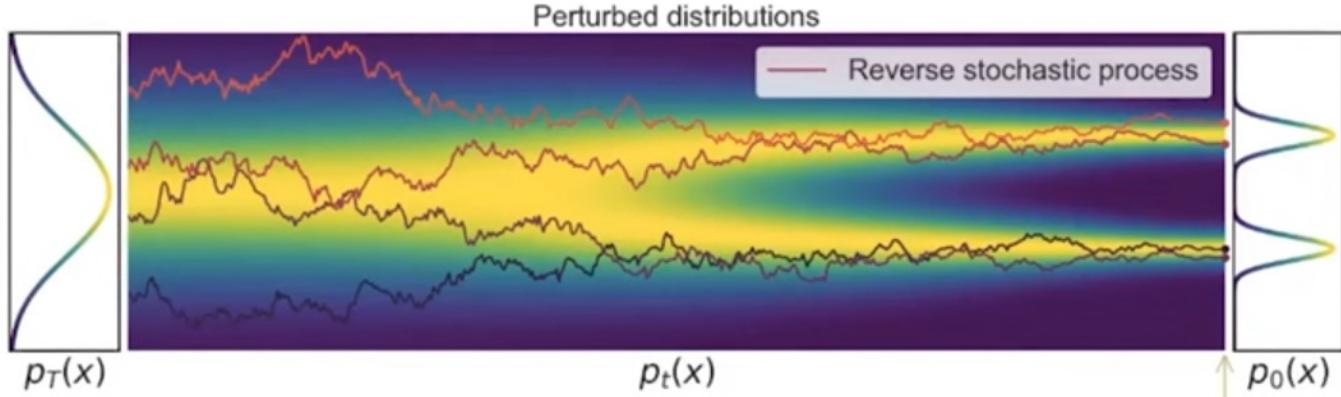
time generalization of noise label O_i

Corresponding
probability densities for random variables

$$\{p_t(x)\}_{t \in [0, T]}$$

Reverse Stochastic process

Generation via reverse stochastic processes



Forward SDE ($t: 0 \rightarrow T$)

$$dX_t = \sigma(t) dW_t$$

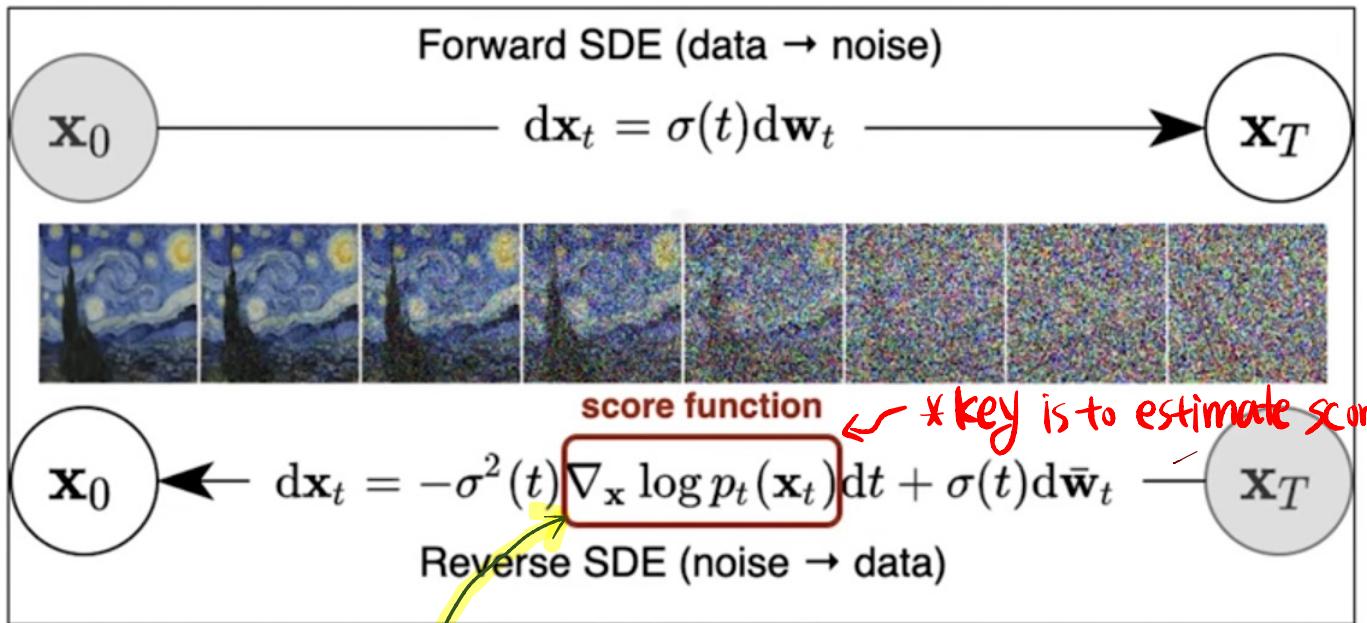
Reverse SDE ($t: T \rightarrow 0$)

$$dX_t = -\sigma(t)^2 \underbrace{\nabla_x \log p_t(X_t)}_{\text{Score function}} dt + \sigma(t) \boxed{d\bar{W}_t}$$

infinitesimal noise
in the reverse time
direction



Score-based generative modeling via SDEs



Time conditional Score model

$$S_\theta(\mathbf{x}, t) \approx \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)$$

train 'time conditional Score function' model
to approximate score function
of the data density and time instance t

Training objective:

$$\mathbb{E}_{t \sim \text{Uniform}[0, T]} [\lambda(t) \mathbb{E}_{p_t(x)} [\|\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) - S_\theta(\mathbf{x}, t)\|_2^2]]$$

↓
positive weighting
function

↓
Score matching loss for
any time instant t

Score-based generative modeling via SDEs

(After training, we achieve :)

- Time-dependent score-based model

$$\mathbf{s}_\theta(\mathbf{x}, t) \approx \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$$

- Training:

Minimize $\mathbb{E}_{t \in \mathcal{U}(0, T)} [\lambda(t) \mathbb{E}_{p_t(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, t)\|_2^2]]$
Objective :

Train efficiently by $\begin{cases} \textcircled{1} \text{ denoising score matching} \\ \textcircled{2} \text{ sliced score matching} \end{cases}$

- Reverse-time SDE

$$d\mathbf{x} = -\sigma^2(t) \mathbf{s}_\theta(\mathbf{x}, t) dt + \sigma(t) d\bar{\mathbf{w}}$$

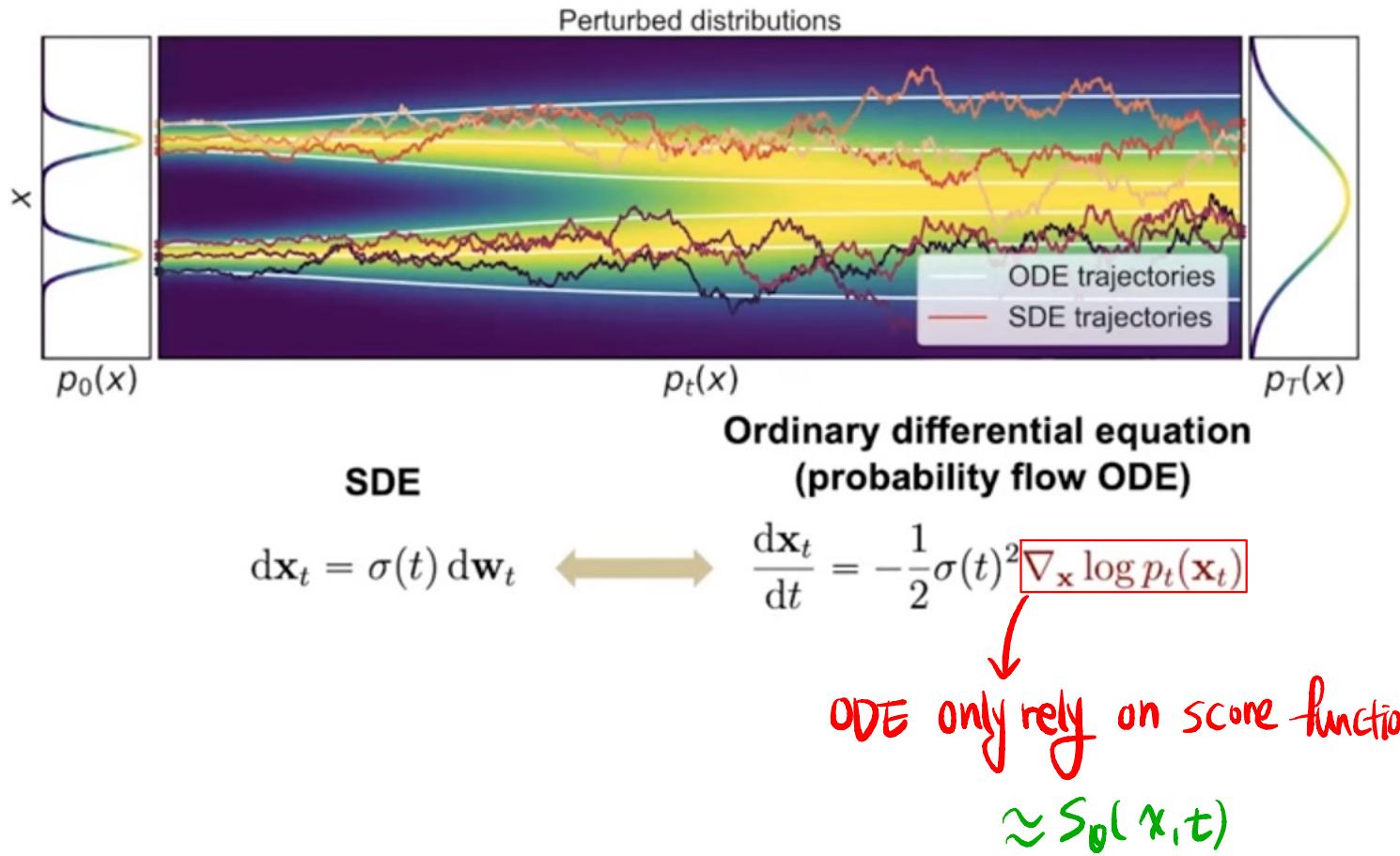
- Euler-Maruyama (analogous to Euler for ODEs)
 $\mathbf{x} \leftarrow \mathbf{x} - \sigma(t)^2 \mathbf{s}_\theta(\mathbf{x}, t) \Delta t + \sigma(t) \mathbf{z} \quad (\mathbf{z} \sim \mathcal{N}(\mathbf{0}, |\Delta t| \mathbf{I}))$
 $t \leftarrow t + \Delta t$

cf.

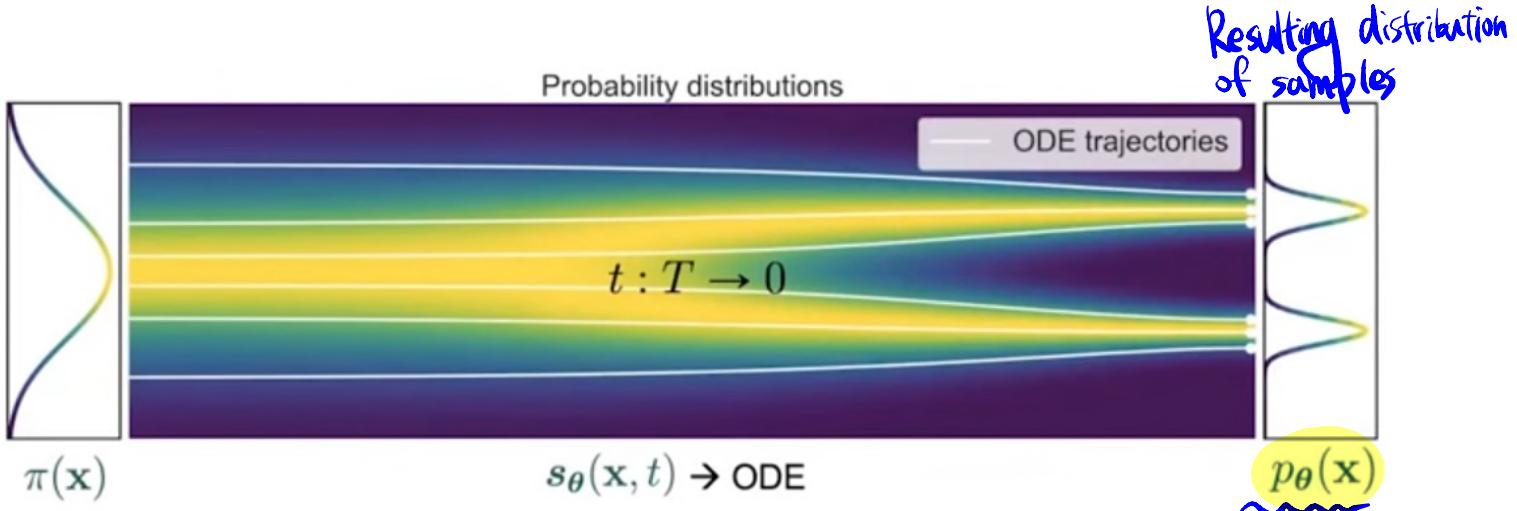
SDE vs. ODE

Markov non-Markov

Converting the SDE to an ODE



Evaluating the probabilities with ODEs



Computing the exact likelihood with ODEs

Likelihood:

$$\pi(\mathbf{x}) \xrightarrow[t : T \rightarrow 0]{s_{\theta}(\mathbf{x}, t)} p_{\theta}(\mathbf{x})$$

Applications of likelihood:

- Lossless compression (Witten et al. 1987, Townsend et al. 2019)
- Unsupervised anomaly detection (Pimental et al. 2014, Song et al. 2018)
- Generative classification (Ng & Jordan 2002, Zimmermann, Schott, Song, et al. 2021)
- Density estimation (Silverman, 1986)

Probability flow ODE allows exact likelihood computation:

$$\log p_{\theta}(\mathbf{x}_0) = \log \pi(\mathbf{x}_T) - \frac{1}{2} \int_0^T \sigma(t)^2 \text{trace}(\nabla_{\mathbf{x}} s_{\theta}(\mathbf{x}, t)) dt$$

ODE solver Unbiased estimator

(Change-of-variable formula for ODEs, Chen et al. 2018)
[Song et al. ICLR 2021 (Outstanding Paper Award)]

Efficient maximum likelihood training

Theorem (informal): Connection between the Kullback-Leibler (KL) divergence and score matching.

$$\begin{aligned} \text{KL}(p_{\text{data}} \| p_{\theta}) &\leq \frac{1}{2} \mathbb{E}_{t \sim \text{Uniform}[0, T]} [\sigma(t)^2 \mathbb{E}_{p_t(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) - s_{\theta}(\mathbf{x}, t)\|_2^2]] \\ &\quad + \boxed{\text{KL}(p_T \| \pi)} \approx 0 \end{aligned}$$

↑
“Likelihood weighting”

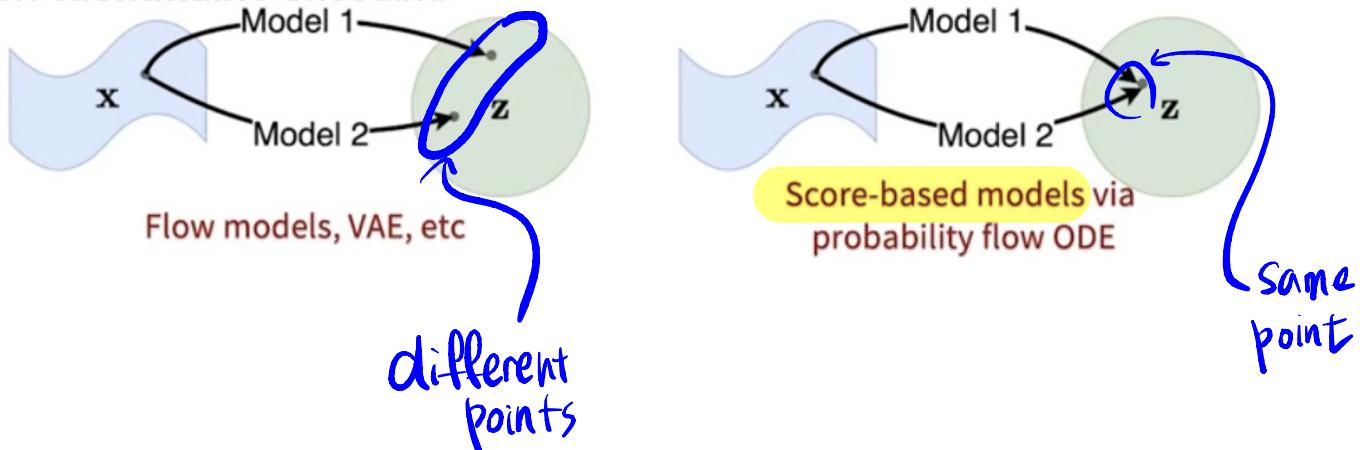


Efficient surrogate loss for
maximum likelihood training

[Song et al. NeurIPS 2021 (spotlight)]

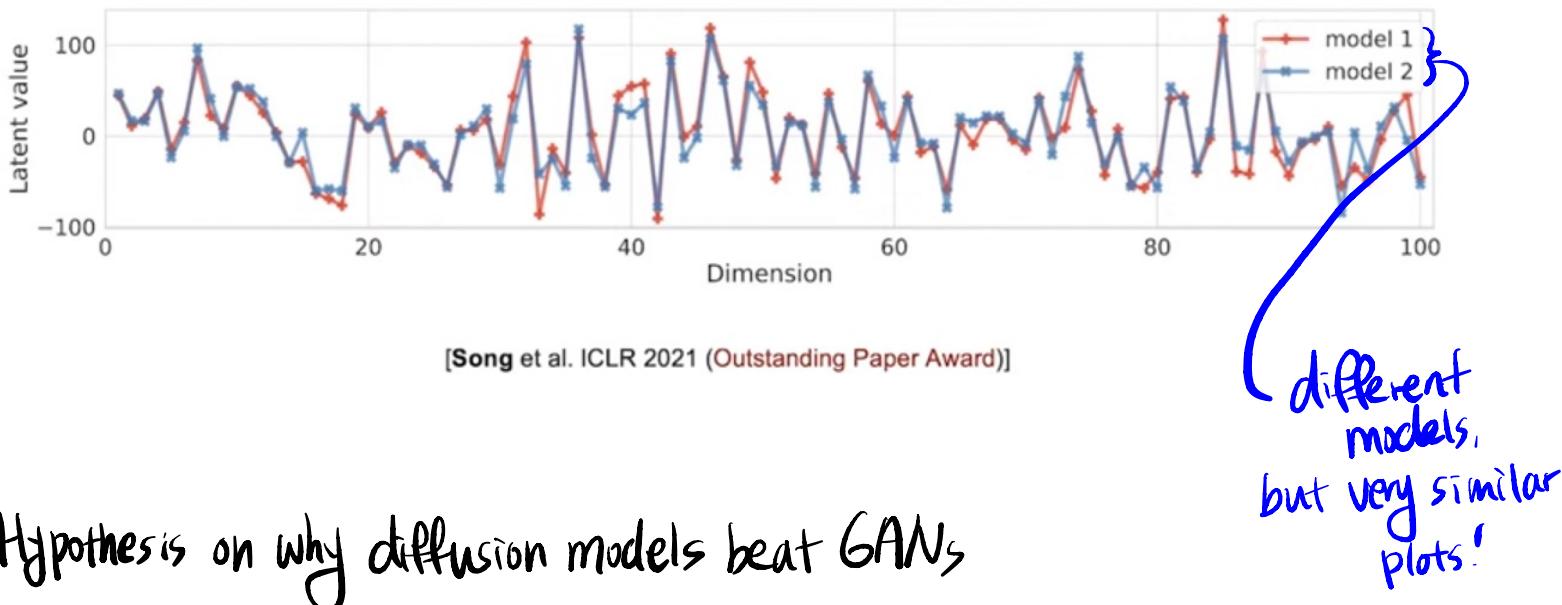
Probability flow ODE: uniquely identifiable encoding

- Uniquely **identifiable** encoding



- No trainable parameters in the probability flow ODE

$$dx = \frac{1}{2} \sigma(t)^2 \nabla_x \log p_t(x) dt$$



Hypothesis on why diffusion models beat GANs

- ◻ More flexible network architectures
- ◻ More diversity