

Score-Based Diffusion Models

LECTURE 16

CS236: Deep Generative Models

Stanford University

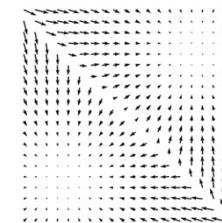
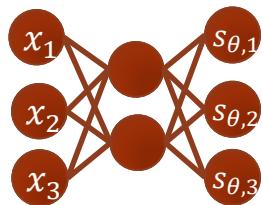
Plan for today

1. Recap on score-based models
2. Diffusion models as score-based models
3. Diffusion models as (hierarchical) VAEs
4. Diffusion models as normalizing flow models
5. Efficient sampling strategies
6. Controllable generation

Score-based models

- A model that represents the score function

$$s_\theta(\mathbf{x})$$

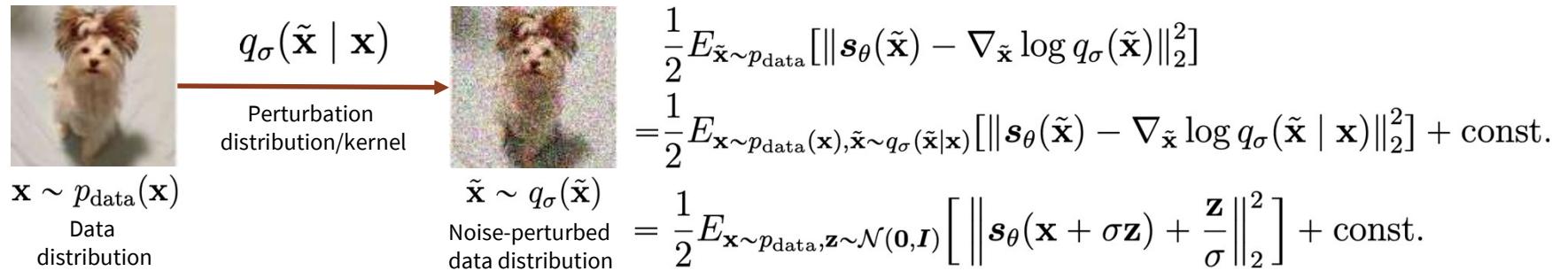


- **Score estimation:** training the score-based model from datapoints
- Score matching

$$\begin{aligned} & \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}} [\|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x})\|_2^2] \\ &= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}} \left[\frac{1}{2} \|\mathbf{s}_\theta(\mathbf{x})\|_2^2 + \text{tr}(\nabla_{\mathbf{x}} \mathbf{s}_\theta(\mathbf{x})) \right] + \text{const.} \end{aligned}$$

- Not scalable for deep score-based models and high dimensional data

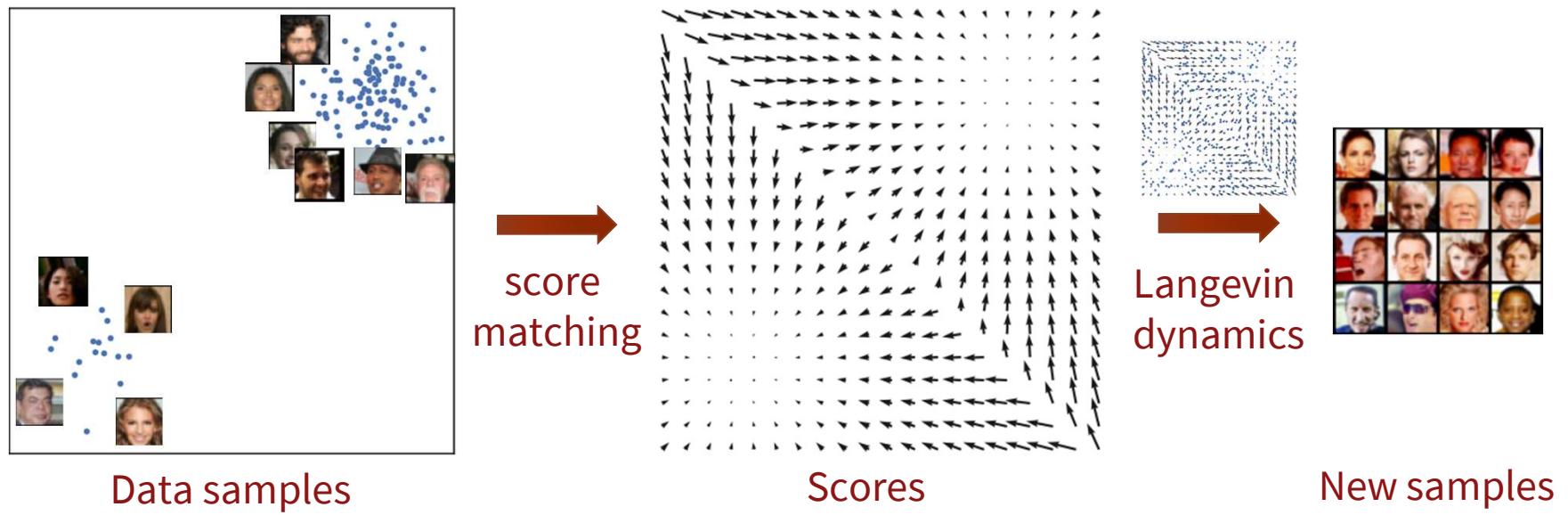
Denoising score matching



- **Pros:**
 - Much more scalable than score matching
 - Reduces score estimation to a denoising task
- **Con:** estimates score of noise-perturbed data

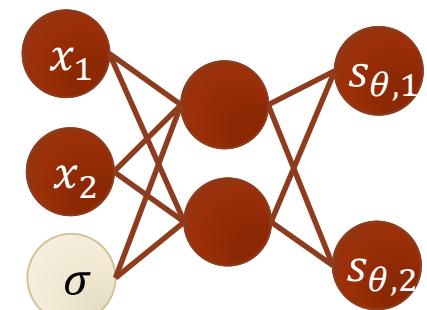
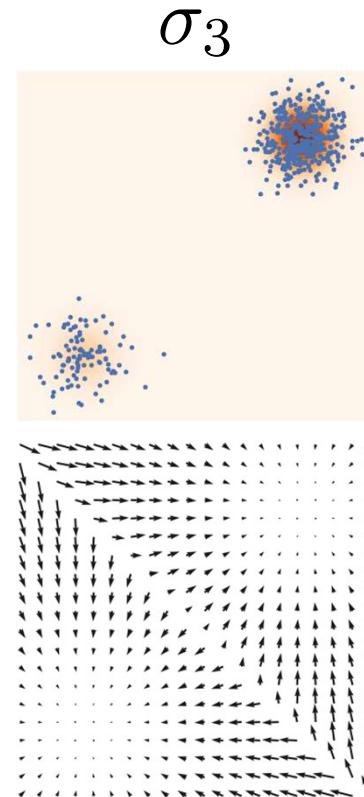
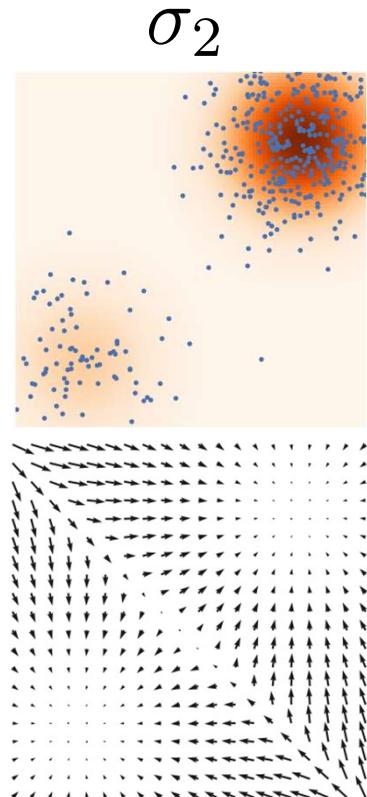
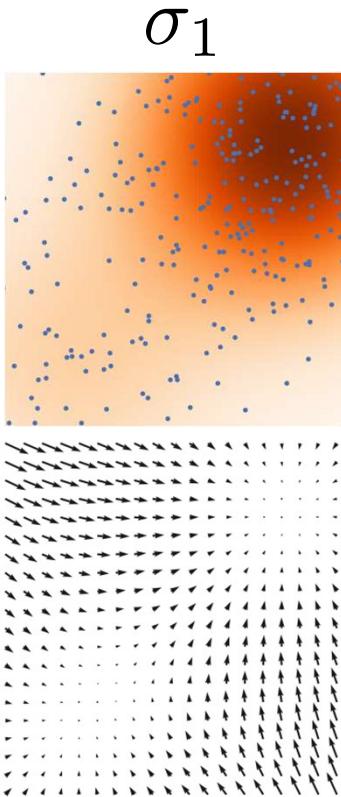
$$\mathbf{s}_\theta(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log q_\sigma(\mathbf{x}) \neq \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$$

Score-based generative modeling



Stanford University

Joint Score Estimation via Noise Conditional Score Networks



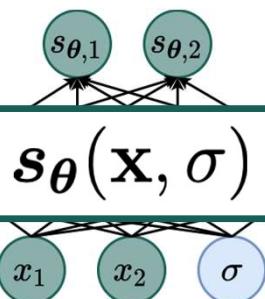
Noise Conditional
Score Network
(NCSN)

Stanford University

Using multiple noise levels

$$p_{\sigma_1}(\mathbf{x}) < p_{\sigma_2}(\mathbf{x}) < p_{\sigma_3}(\mathbf{x})$$

Data



Noise Conditional
Score Model

Annealed Langevin dynamics

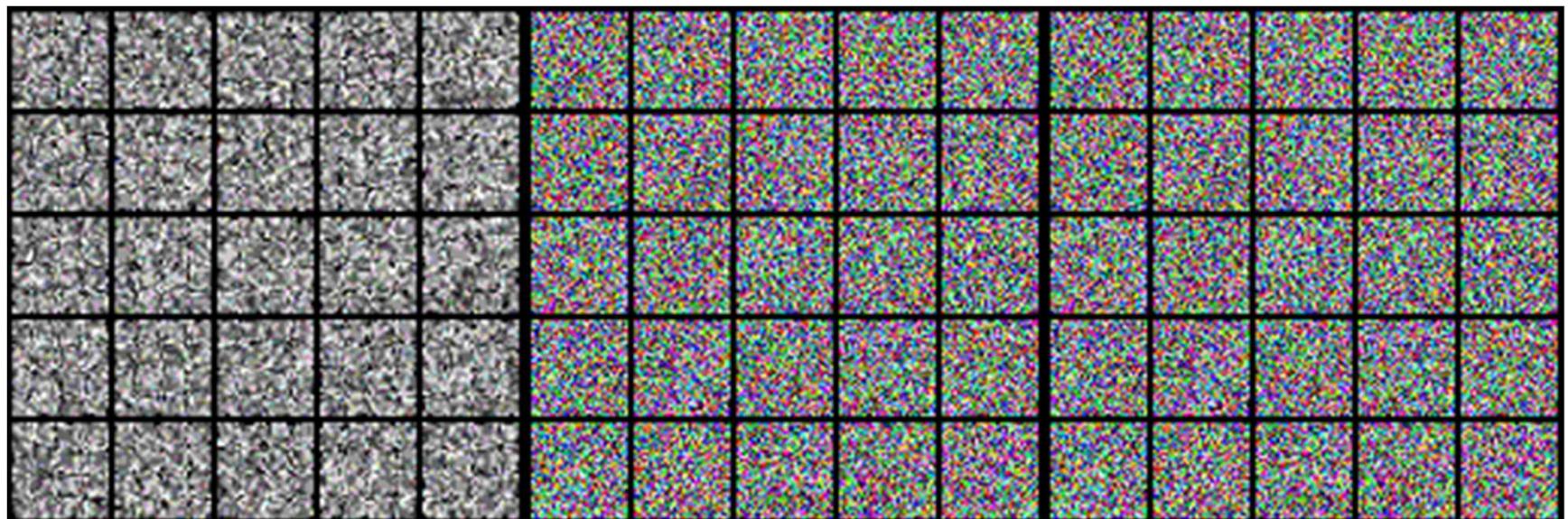
$$\frac{1}{N} \sum_{i=1}^N \lambda(\sigma_i) \mathbb{E}_{p_{\sigma_i}(\mathbf{x})} [\nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x})] - \frac{1}{2} \| \nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x}) \|_2^2$$

Positive weighting function

$s_{\theta}(\mathbf{x}, \sigma_1)$ $s_{\theta}(\mathbf{x}, \sigma_2)$ $s_{\theta}(\mathbf{x}, \sigma_3)$

Stanford University

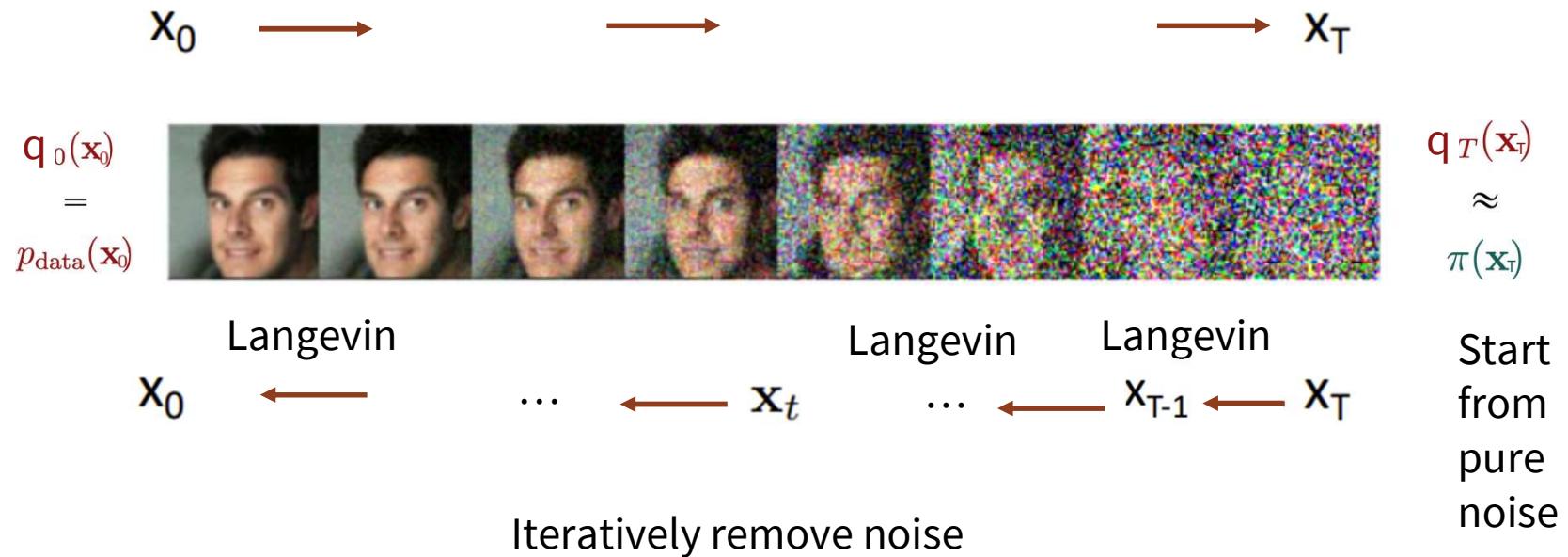
Experiments: Sampling



Stanford University

Sampling as Iterative Denoising

Inverse process: iteratively add Gaussian noise



Iterative noising process



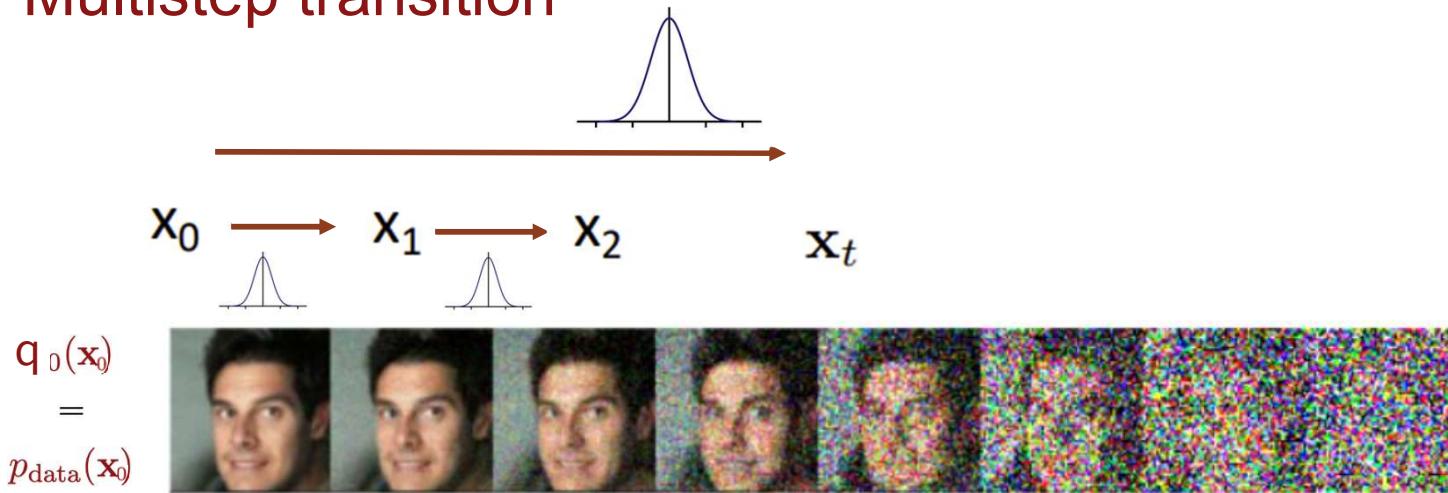
Noise perturbed densities are obtained by adding Gaussian noise

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

Defines a joint distribution $q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$

Stanford University

Multistep transition

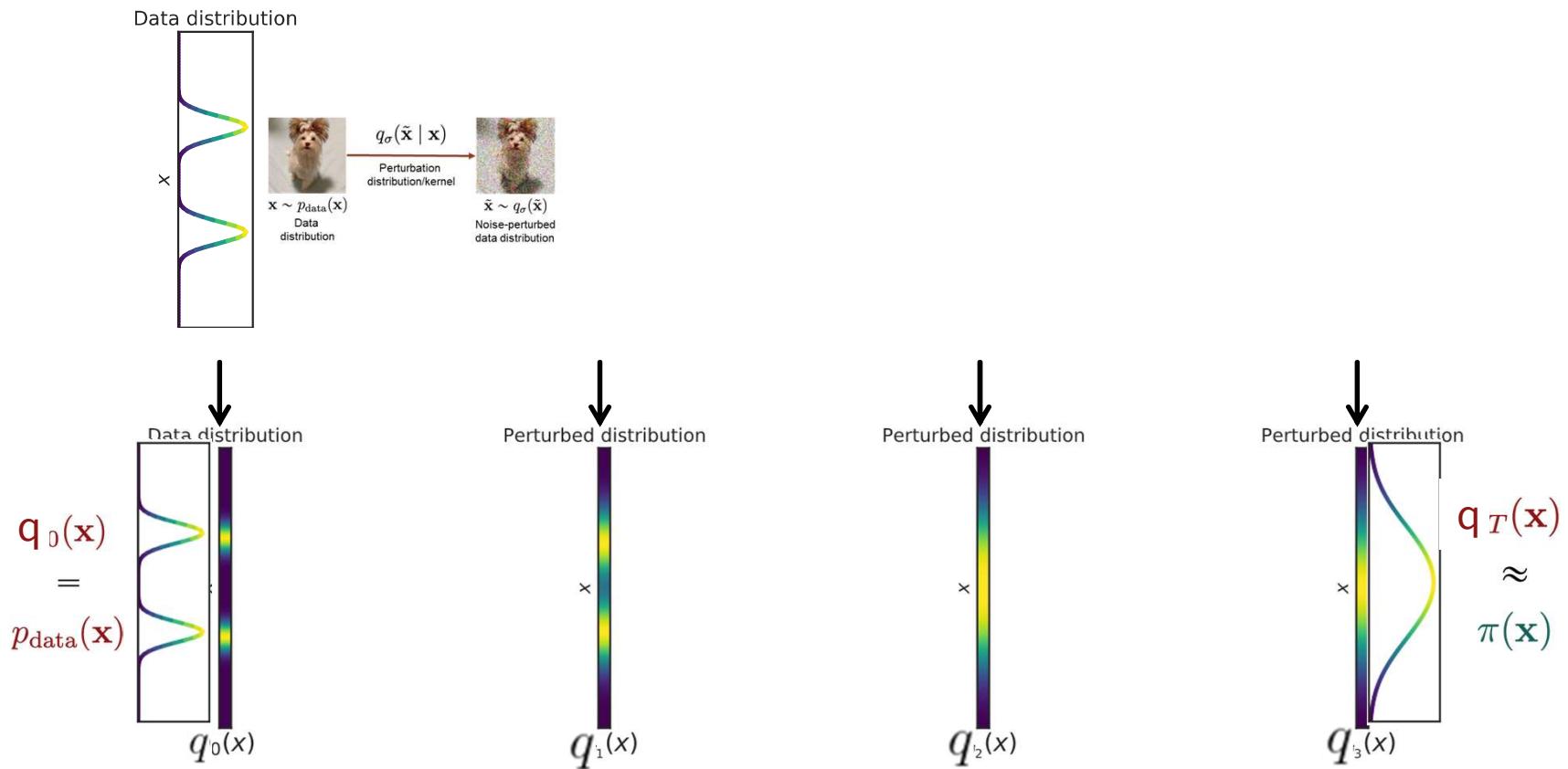


Multi-step transitions are also Gaussian and can be computed in closed form

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad \bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$$

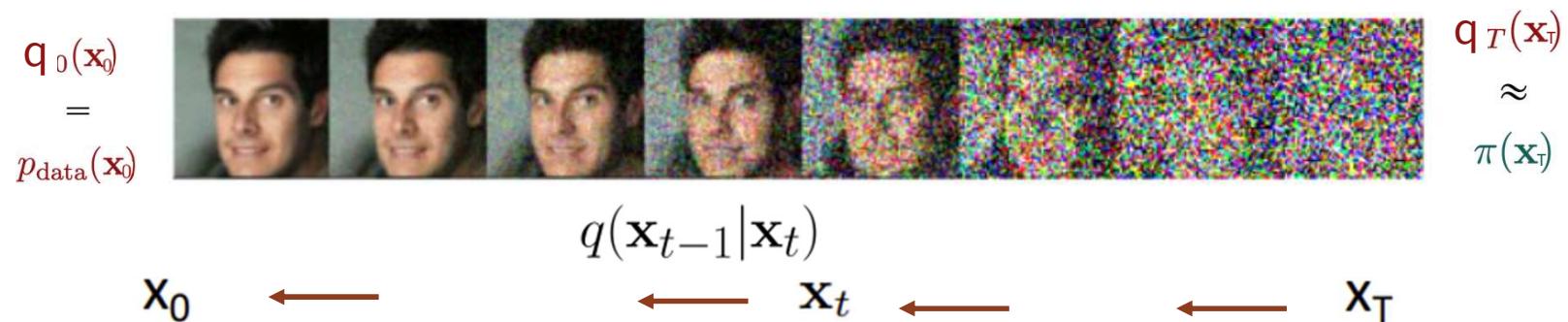
- 1) Same as noise-perturbed data distributions in score-based models
- 2) Efficient sampling at any t

Diffusion perspective



Stanford University

Iterative Denoising



Ideal sampling process:

1. Sample x_T from $\pi(x)$, i.e. from pure noise
2. Iteratively sample from the true denoising distribution $q(x_{t-1}|x_t)$

Issue: Exact denoising distribution is unknown!

Solution: Learn a variational approximation

Iterative Denoising

$$\begin{aligned} q_0(\mathbf{x}_0) &= p_{\text{data}}(\mathbf{x}_0) \\ &\quad \text{---} \qquad \text{---} \qquad \text{---} \qquad \text{---} \qquad \text{---} \qquad \text{---} \end{aligned}$$

$$\begin{aligned} q_T(\mathbf{x}_T) &\approx \\ \pi(\mathbf{x}_T) & \end{aligned}$$

$$\mathbf{x}_0 \quad \leftarrow \quad \approx \quad \mathbf{x}_t \quad \leftarrow \quad \quad \quad \leftarrow \quad \mathbf{x}_T$$
$$q(\mathbf{x}_{t-1}|\mathbf{x}_t)$$
$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

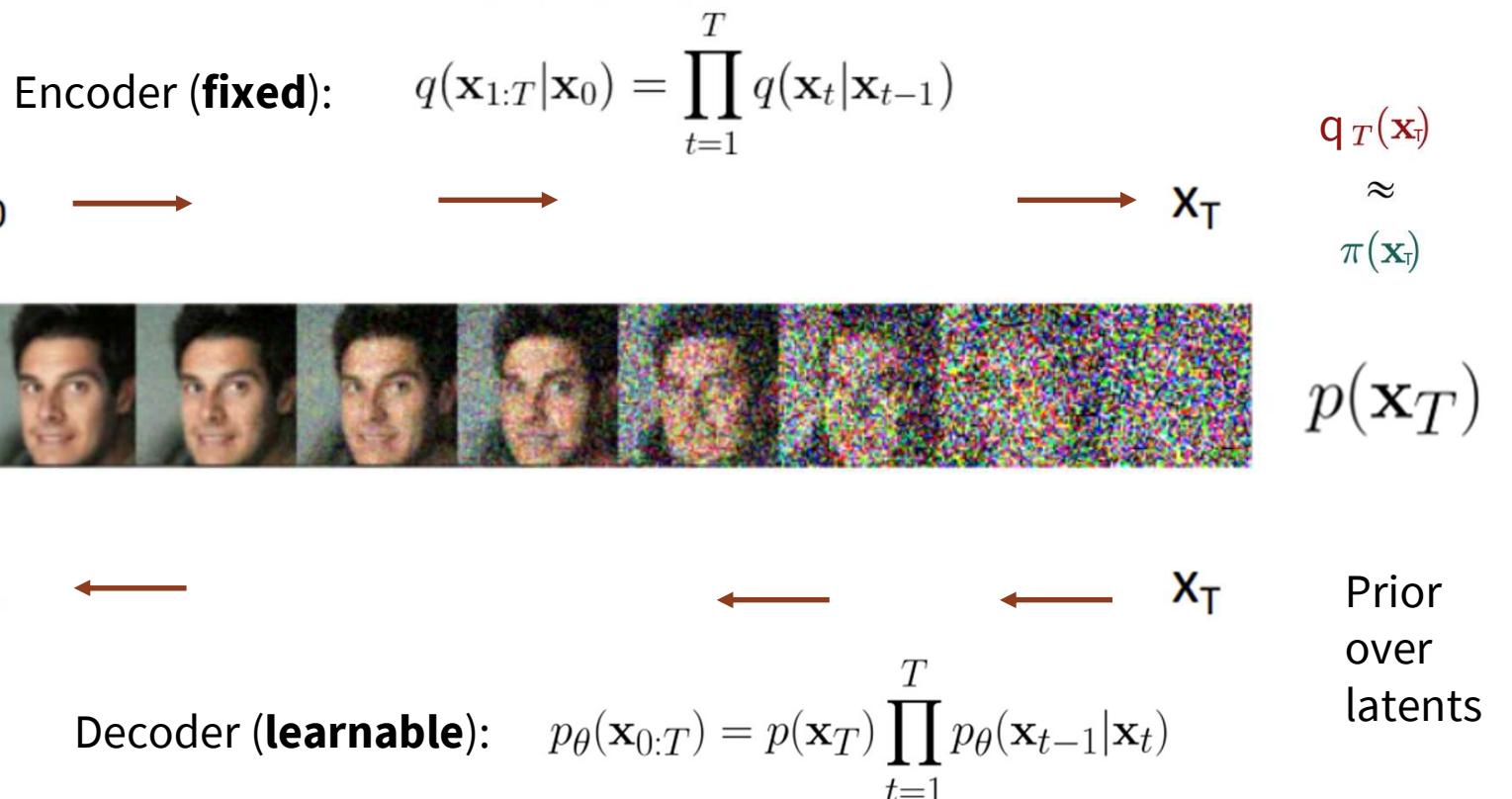
Sample \mathbf{x}_T from $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}) = \pi$

Iteratively sample from $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$

Joint distribution: $p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$

Stanford University

Diffusion model as a hierarchical VAE



From VAE to Hierarchical VAE

Basic VAE:

A mixture of an infinite number of Gaussians:

$$① \quad \mathbf{z} \sim \mathcal{N}(0, I)$$

$$② \quad p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mu_\theta(\mathbf{z}), \Sigma_\theta(\mathbf{z})) \text{ where } \mu_\theta, \Sigma_\theta \text{ are neural networks}$$

ELBO training:

$$E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{z}, \mathbf{x}; \theta) - \log q_\phi(\mathbf{z}|\mathbf{x})]$$



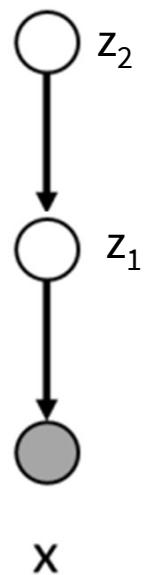
From VAE to Hierarchical VAE

Hierachical VAE (decoder): $p(x, z_1, z_2) = p(z_2) p(z_1 | z_2) p(x | z_1)$

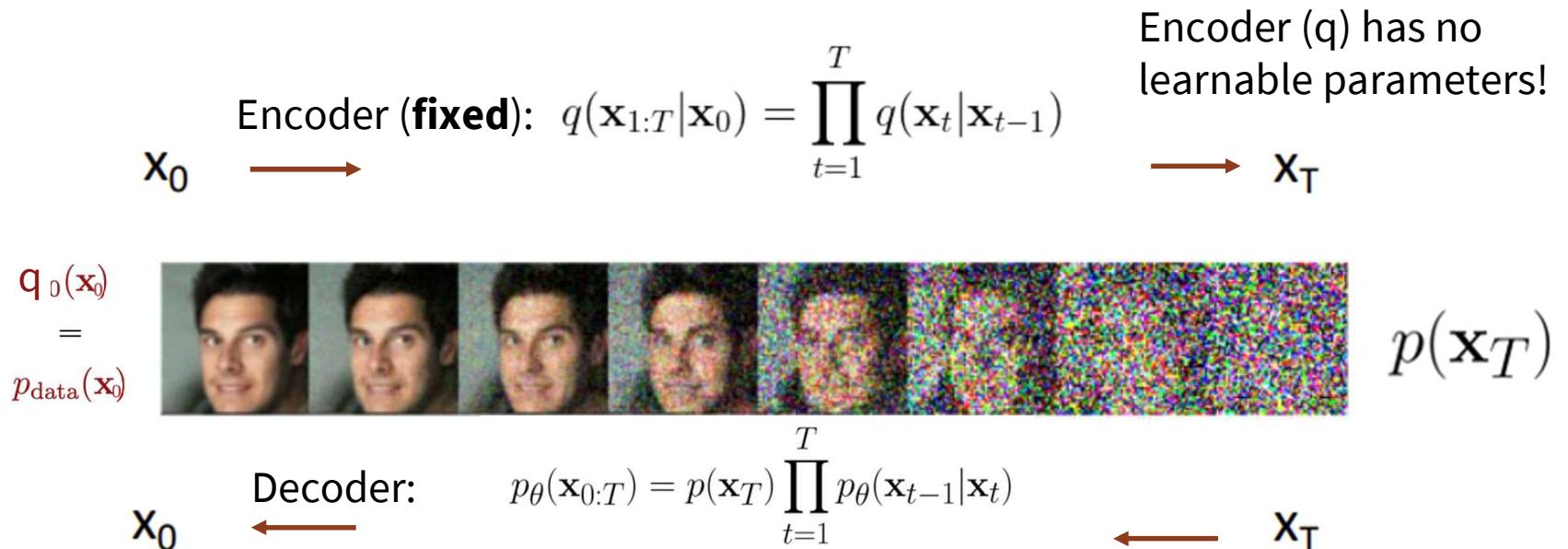
1. Sample z_2 from a simple prior $N(0, I)$
2. Sample z_1 from a decoder $p(z_1 | z_2)$
3. Sample x from a decoder $p(x | z_1)$

Encoder: $q(z_1, z_2 | x)$

ELBO training: $E_{q(z_1, z_2 | x)} \left[\log \left(\frac{p(x, z_1, z_2)}{q(z_1, z_2 | x)} \right) \right]$



Training a denoising diffusion probabilistic model



ELBO loss (negative ELBO averaged over data distribution):

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] =: L$$

Diffusion models as score-based models

The ELBO objective is

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] =: L$$

Decoder parameterization: $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right)$$

Up to scaling, predict noise that was added and subtract it

$$L = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), t \sim \mathcal{U}\{1, T\}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\lambda_t \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2 \right]$$

ELBO loss reduces to denoising score matching!

Stanford University

Training and inference

Algorithm 1 Training

```

1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
        $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$ 
6: until converged

```

Denoising score
matching training

$$\frac{1}{L} \sum_{i=1}^L E_{\mathbf{x} \sim p_{\text{data}}, \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\|\boldsymbol{\epsilon}_{\theta}(\mathbf{x} + \sigma_i \mathbf{z}, \sigma_i) + \mathbf{z}\|_2^2 \right]$$

Algorithm 2 Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

```

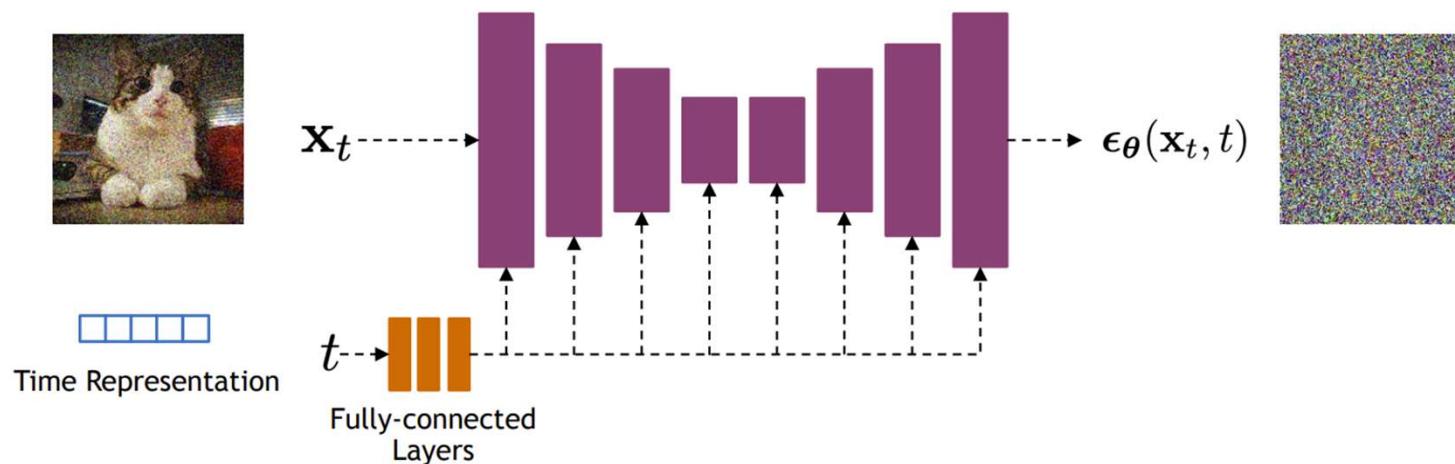
Iteratively Sample from decoders $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$
 (Annealed Langevin Dynamics)

$$[\boldsymbol{\epsilon}_{\theta}(\cdot, \sigma_i) := \sigma_i \mathbf{s}_{\theta}(\cdot, \sigma_i)]$$

Stanford University

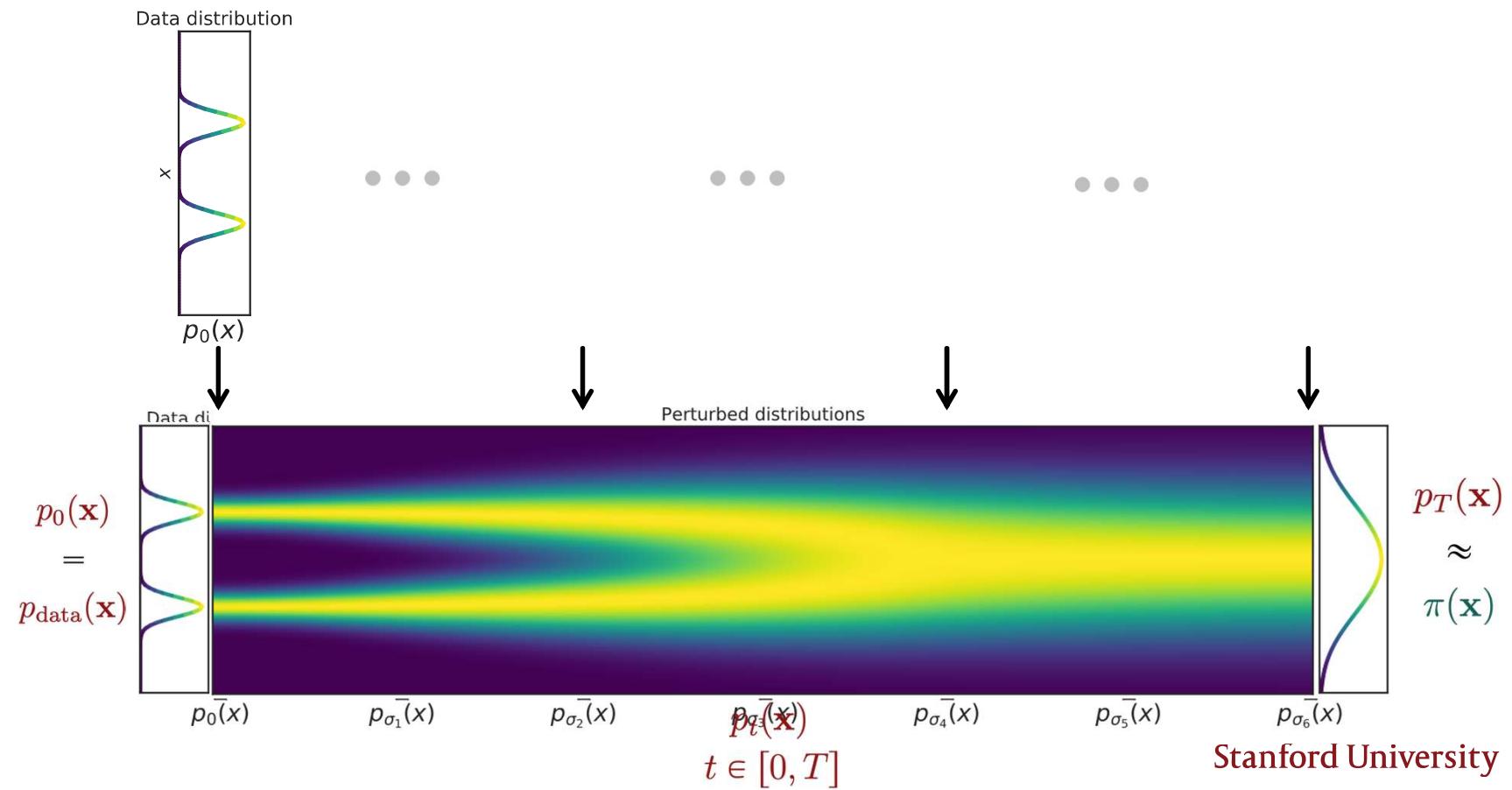
Architecture for the denoiser

Unet architecture used in practice

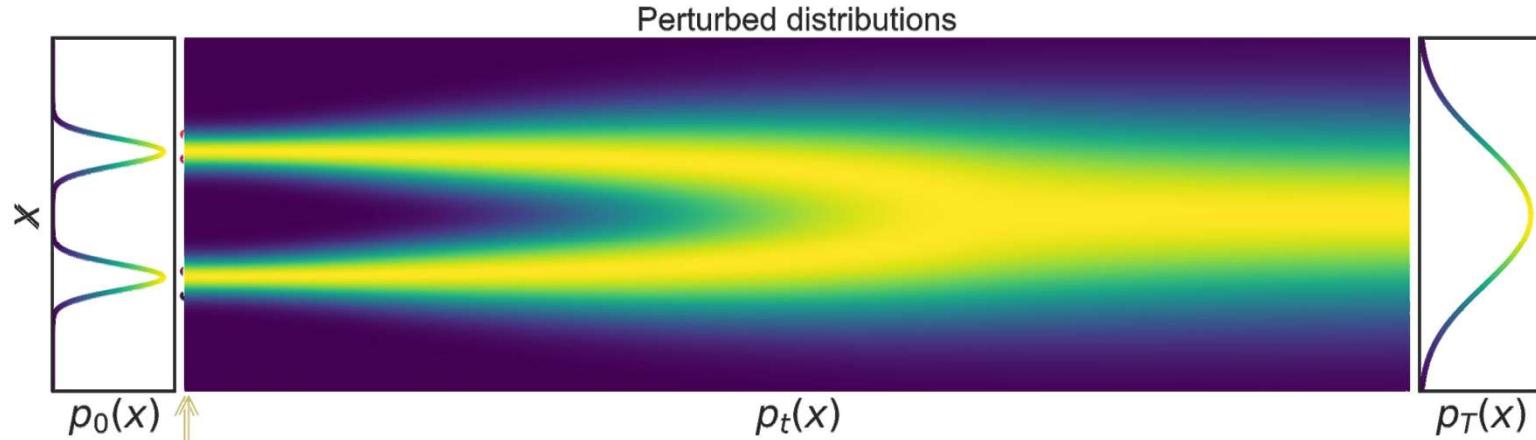


Same as the noise-conditional score network (t represents the time index or equivalently the noise level)

Infinite noise levels



Perturbing data with stochastic processes



Stochastic differential equation (SDE)

$$dx_t = [f(x_t, t)]dt + g(t)d\omega_t$$

Deterministic drift

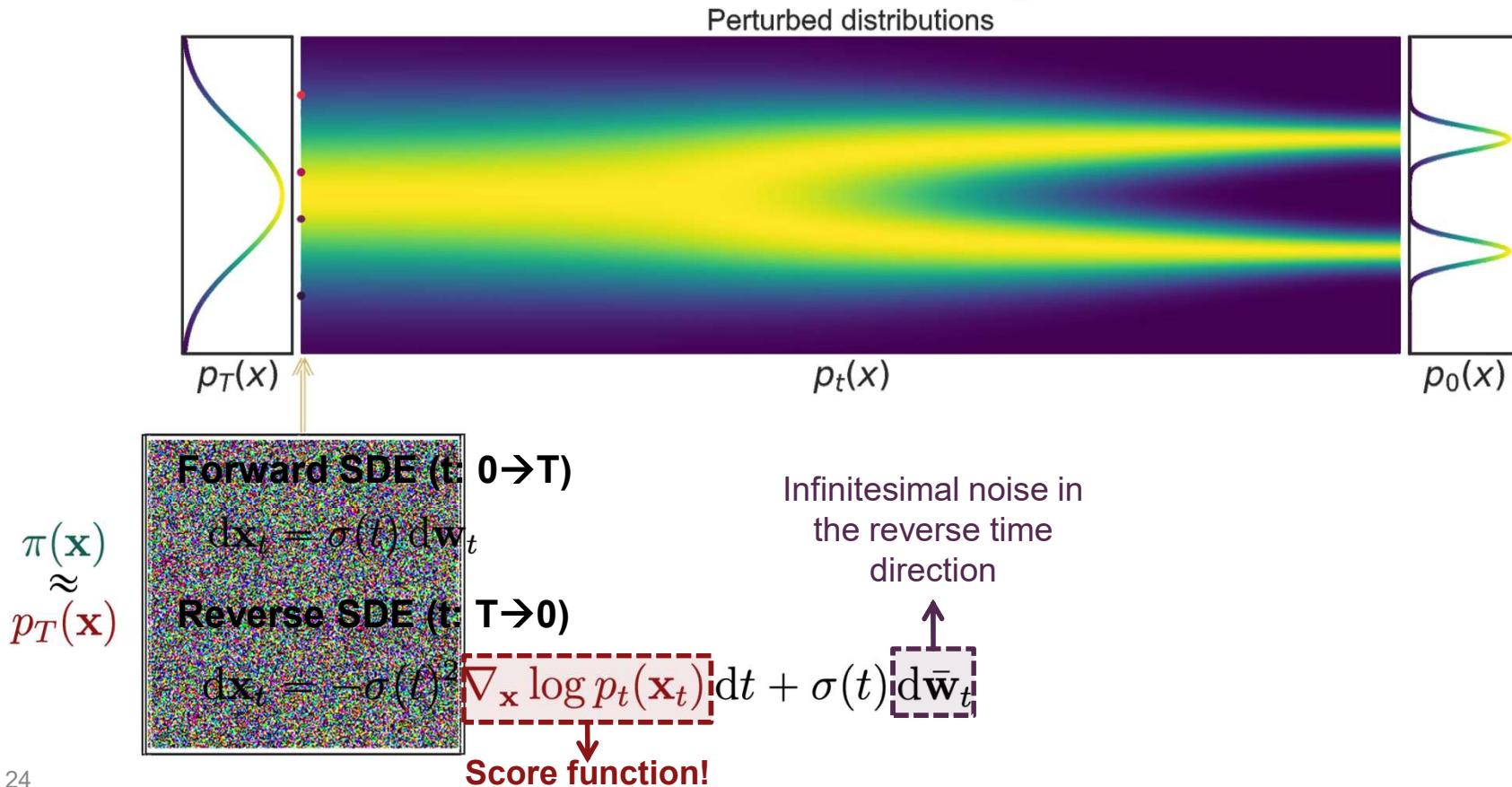
Infinitesimal noise

$$p_T(x) \approx \pi(x)$$

WLOG: Toy SDE

$$dx_t = \sigma(t) d\omega_t$$

Generation via reverse stochastic processes



Score-based generative modeling via SDEs

- Time-dependent score-based model

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

- Training:

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

- 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)

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

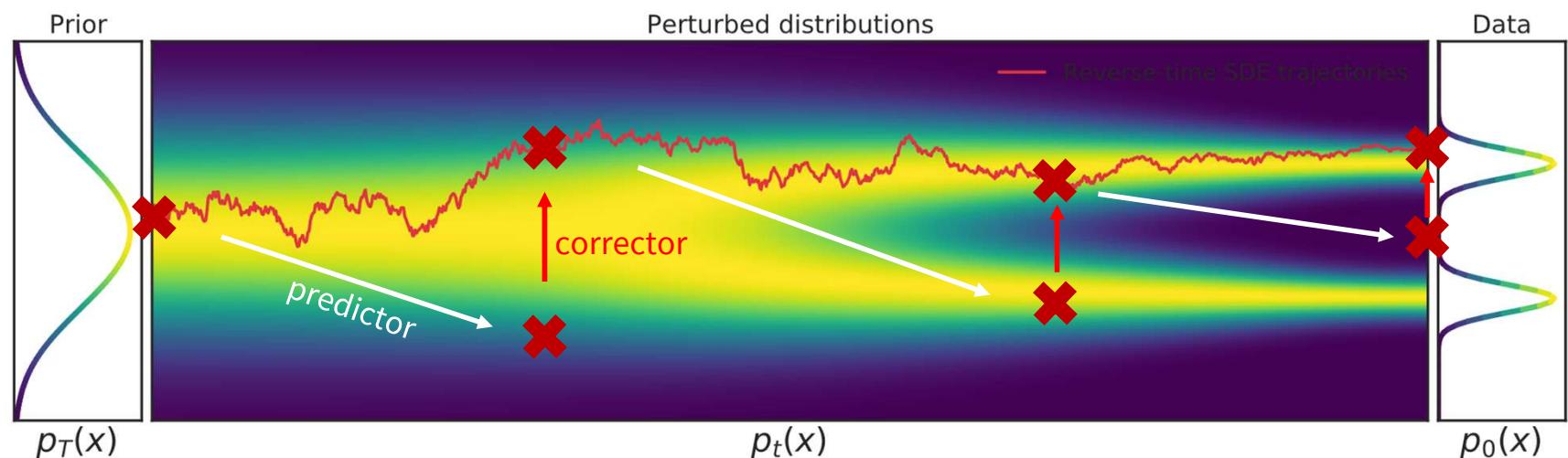
$$t \leftarrow t + \Delta t$$

Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole. "Score-Based Generative Modeling through Stochastic Differential Equations." ICLR 2021.

Stanford University

Predictor-Corrector sampling methods

- Predictor-Corrector sampling.
 - **Predictor:** Numerical SDE solver
 - **Corrector:** Score-based MCMC



Stanford University

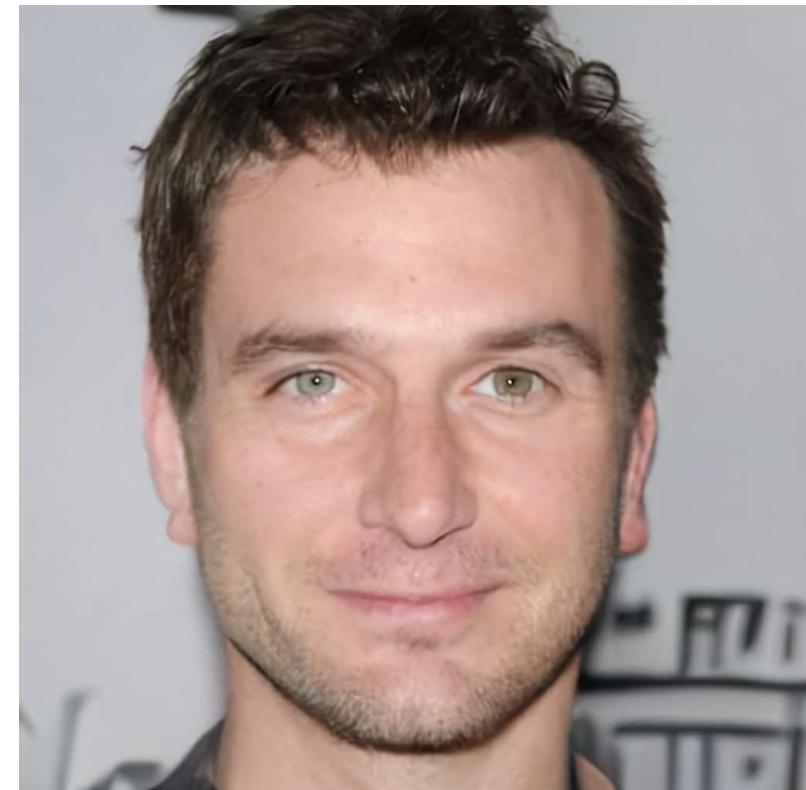
Results on predictor-corrector sampling

Model	FID↓	IS↑
Conditional		
BigGAN (Brock et al., 2018)	14.73	9.22
StyleGAN2-ADA (Karras et al., 2020a)	2.42	10.14
Unconditional		
StyleGAN2-ADA (Karras et al., 2020a)	2.92	9.83
NCSN (Song & Ermon, 2019)	25.32	$8.87 \pm .12$
NCSNv2 (Song & Ermon, 2020)	10.87	$8.40 \pm .07$
DDPM (Ho et al., 2020)	3.17	$9.46 \pm .11$
DDPM++	2.78	9.64
DDPM++ cont. (VP)	2.55	9.58
DDPM++ cont. (sub-VP)	2.61	9.56
DDPM++ cont. (deep, VP)	2.41	9.68
DDPM++ cont. (deep, sub-VP)	2.41	9.57
NCSN++	2.45	9.73
NCSN++ cont. (VE)	2.38	9.83
NCSN++ cont. (deep, VE)	2.20	9.89

Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole. "Score-Based Generative Modeling through Stochastic Differential Equations." ICLR 2021.

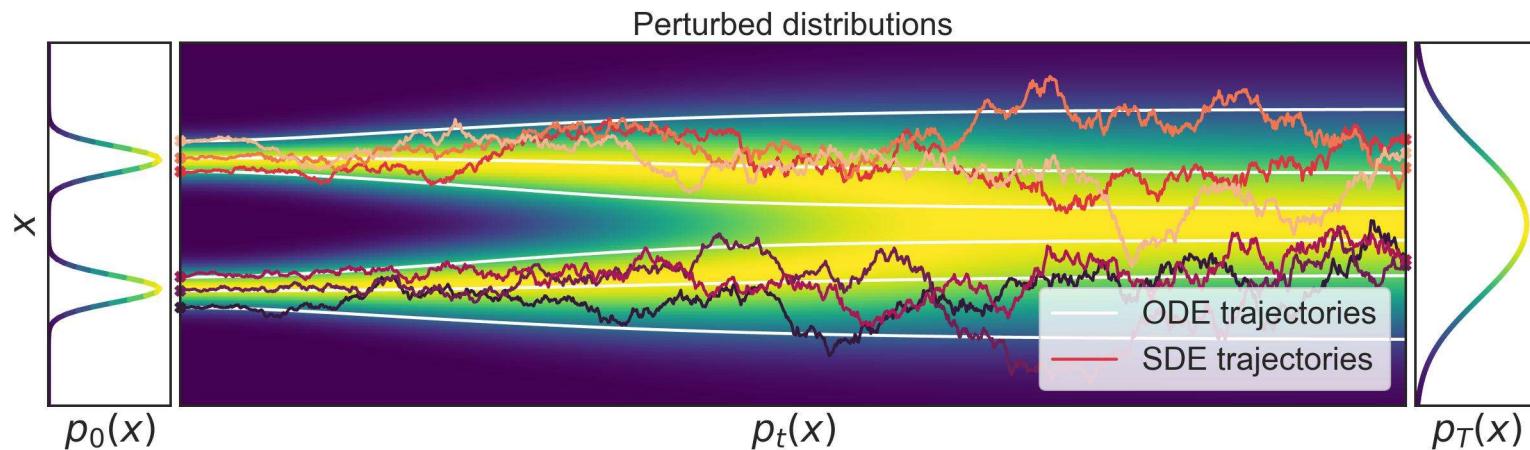
Stanford University

High-Fidelity Generation for 1024x1024 Images



Stanford University

Converting the SDE to an ODE



SDE

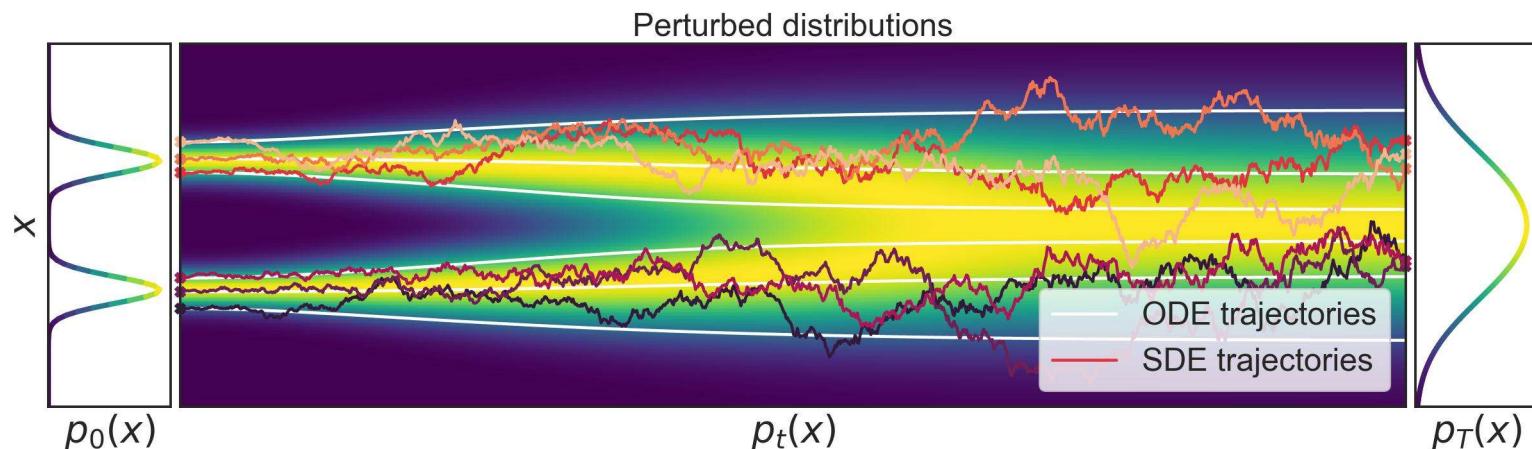
$$d\mathbf{x}_t = \sigma(t) d\mathbf{w}_t$$

Ordinary differential equation (ODE)

$$\frac{d\mathbf{x}_t}{dt} = -\frac{1}{2}\sigma(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)$$

Score function
 $\approx s_{\theta}(\mathbf{x}, t)$
Stanford University

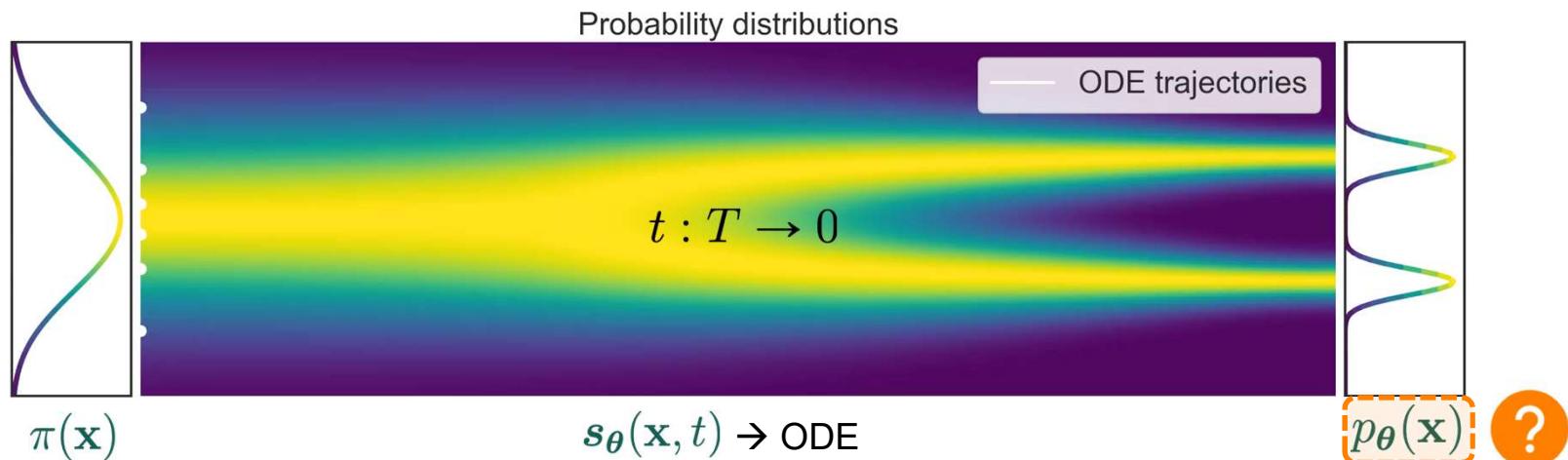
Converting the SDE to an ODE



We can think of this as a (continuous time, infinite depth) normalizing flow

1. Unique ODE solution \rightarrow Invertible mapping
2. To invert, solve ODE backwards from T to 0

Evaluating likelihoods with ODEs (flow model)



Computing the probability density function (change of variables formula)

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

ODE solver

Computable in polynomial time

Stanford University

Competitive likelihoods on test data

Negative log-probability ↓ (**bits/dim**)

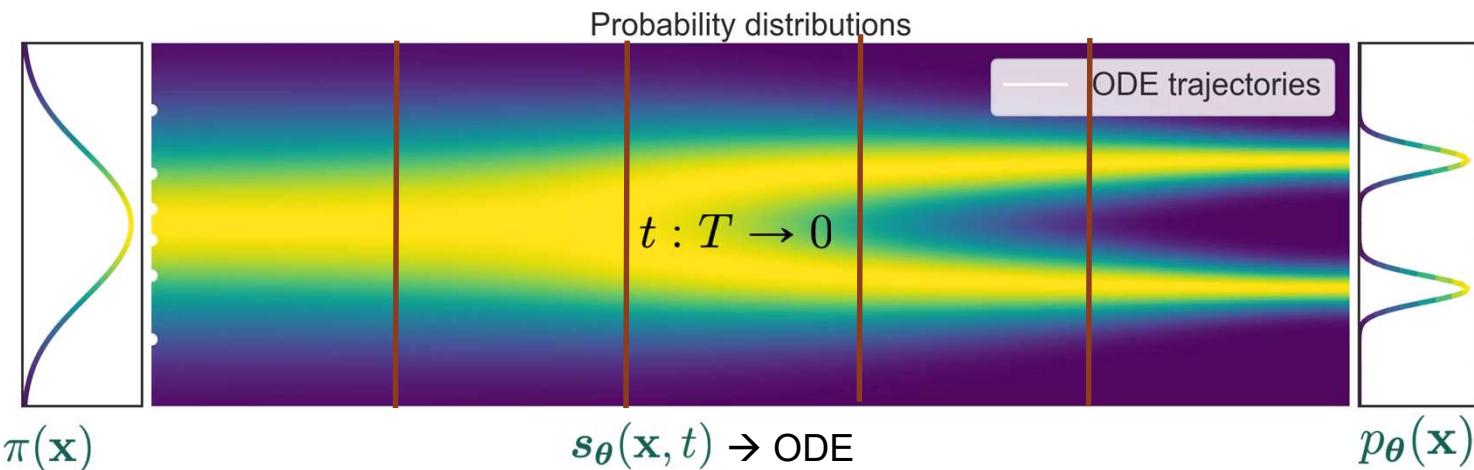
Method	CIFAR-10	ImageNet 32x32
PixelSNAIL [Chen et al. 2018]	2.85	3.80
Delta-VAE [Razavi et al. 2019]	2.83	3.77
Sparse Transformer [Child et al. 2019]	2.80	–



Challenges years of dominance of autoregressive models and VAEs

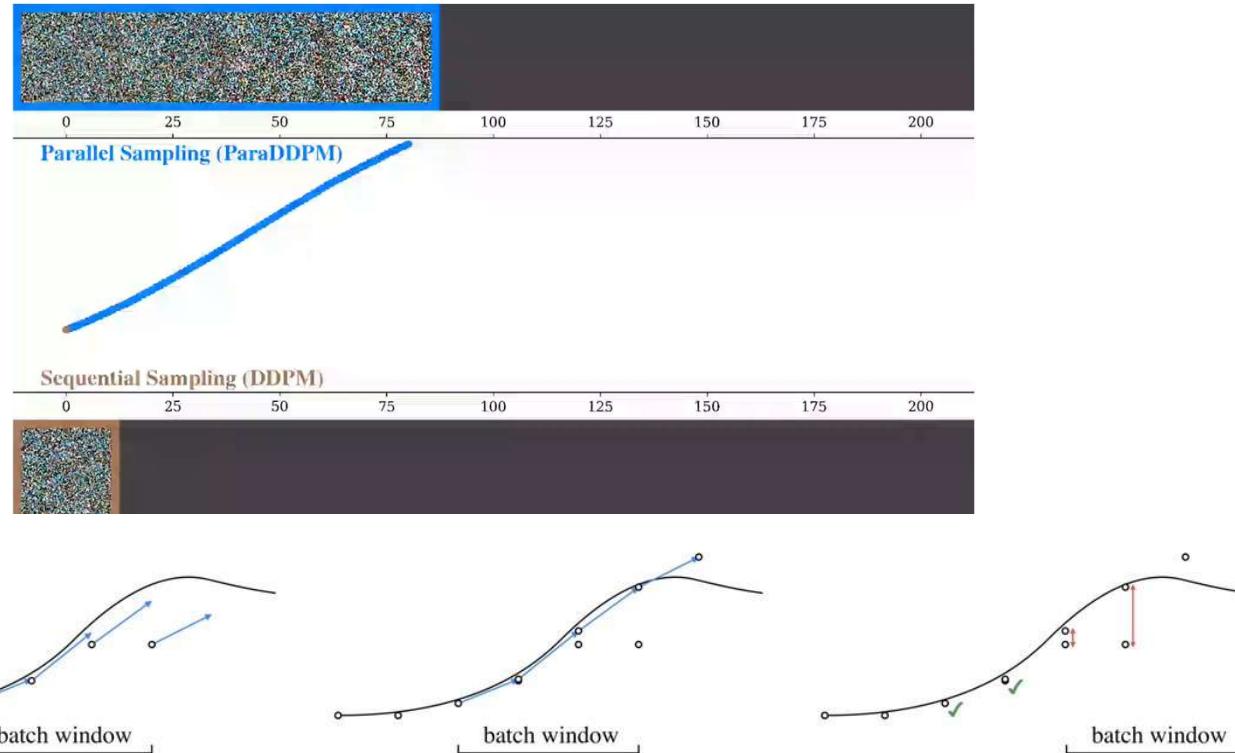
Stanford University

Accelerated sampling



- Numerical methods + ODE formulation to accelerate sampling
- DDIM [Song and Ermon, 2021]:
 - Coarsely discretize the time axis, take big steps
 - Corresponds to exponential integrator (semi-linear ODE) [Lu et al, 2022; Zhang and Chen, 2022]
 - 10x-50x speedups, comparable sample quality

Parallel ODE solving

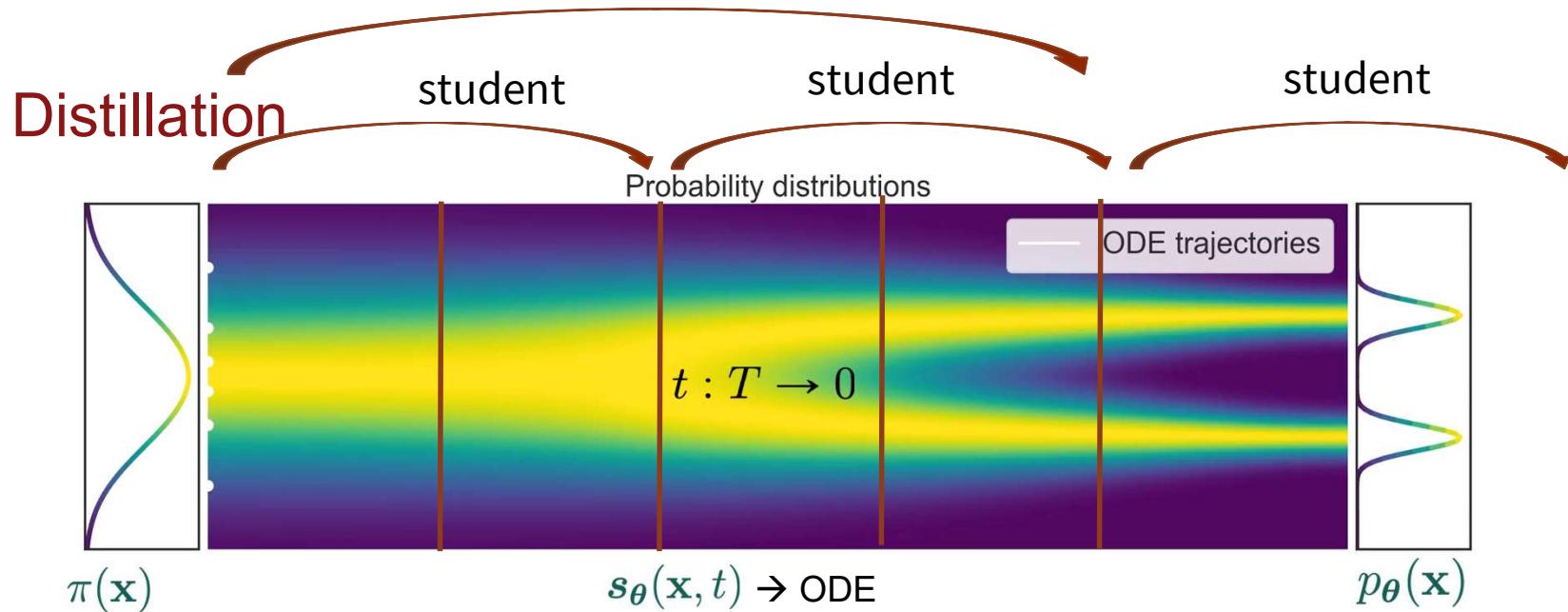


(a) Compute the drift of $x_{t:t+p}^k$ on a batch window of size $p = 4$, in parallel

(b) Update the values to $x_{t:t+p}^{k+1}$ using the cumulative drift of points in the window

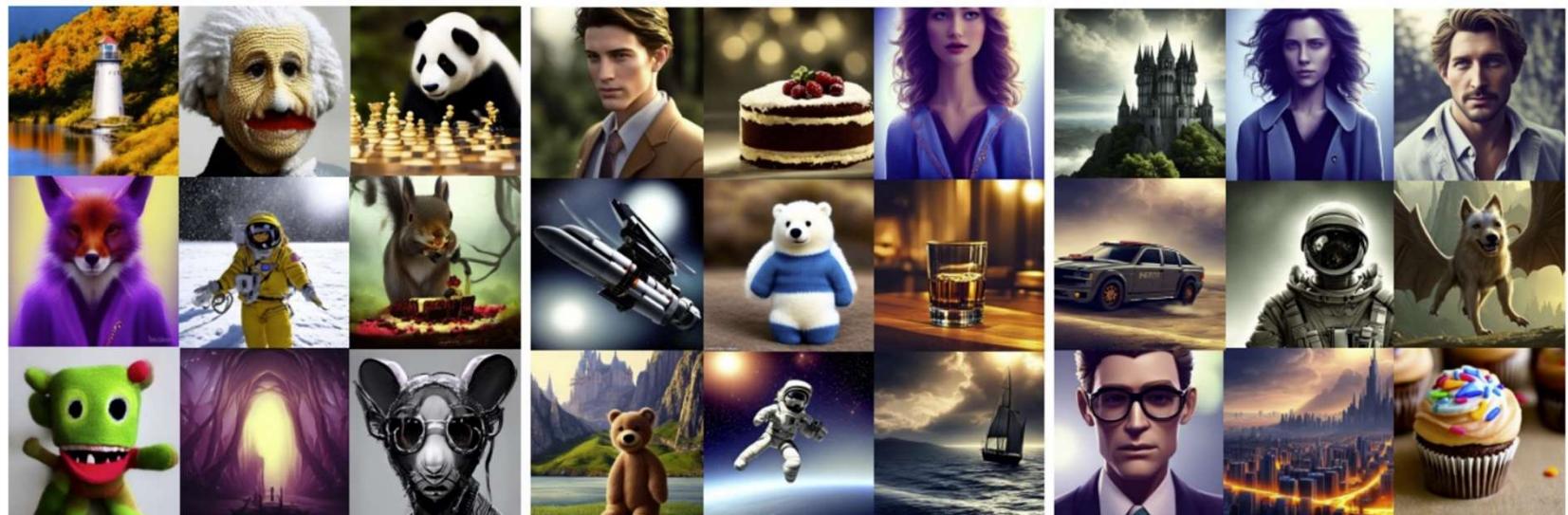
(c) Determine how far to slide the window forward, based on the error $\|x_i^{k+1} - x_i^k\|^2$.

versity



- Progressive distillation [Salimans, Ho 2022]
 - DDIM sampler as a teacher model
 - Student model trained to do in 1 step what DDIM achieves in 2 steps
 - Applied recursively to drastically reduce the number of steps required

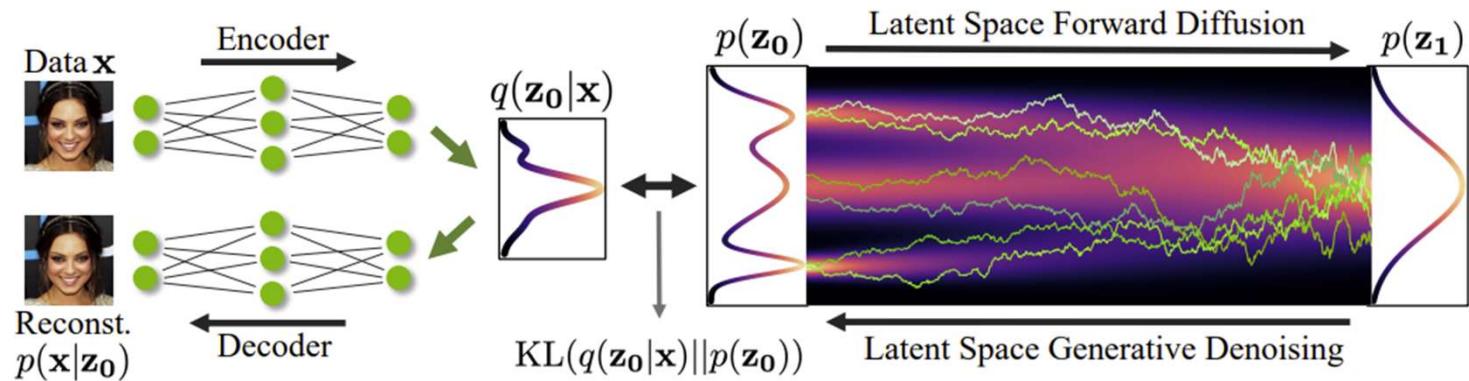
On distillation of guided diffusion models



Meng, Rombach, Gao, Kingma, Ermon, Ho, Salimans “On distillation of guided diffusion models.” 2023.

Stanford University

Latent diffusion model



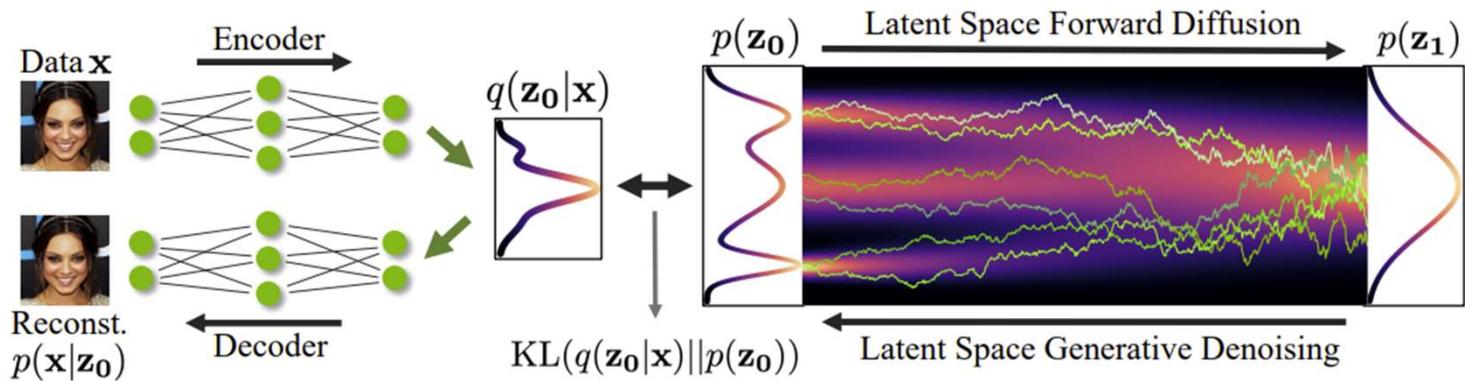
VAE mapping data to lower dimensional space

- 1. Faster**
- 2. Can be applied to any data type (e.g. discrete)**

Diffusion model prior over the latent space of the autoencoder

Stanford University

Stable diffusion text2image model



VAE mapping data to lower dimensional space

1. Pre-trained, focus on reconstruction (autoencoder)

Diffusion model prior over the latent space of the autoencoder

2. Trained in the second stage, keeping initial autoencoder fixed

Large scale, open source model, widely adopted

Stanford University

Conditional generation

User input:

An astronaut riding a horse



Stanford University

Conditional generation

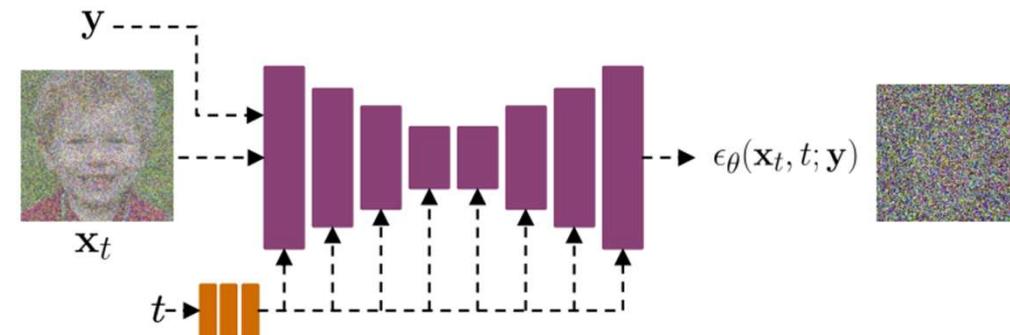
Let (x, y) denote (image,caption) pairs

Training a conditional generative model involves learning $p(x | y)$

Train score model for the image x conditional on caption y

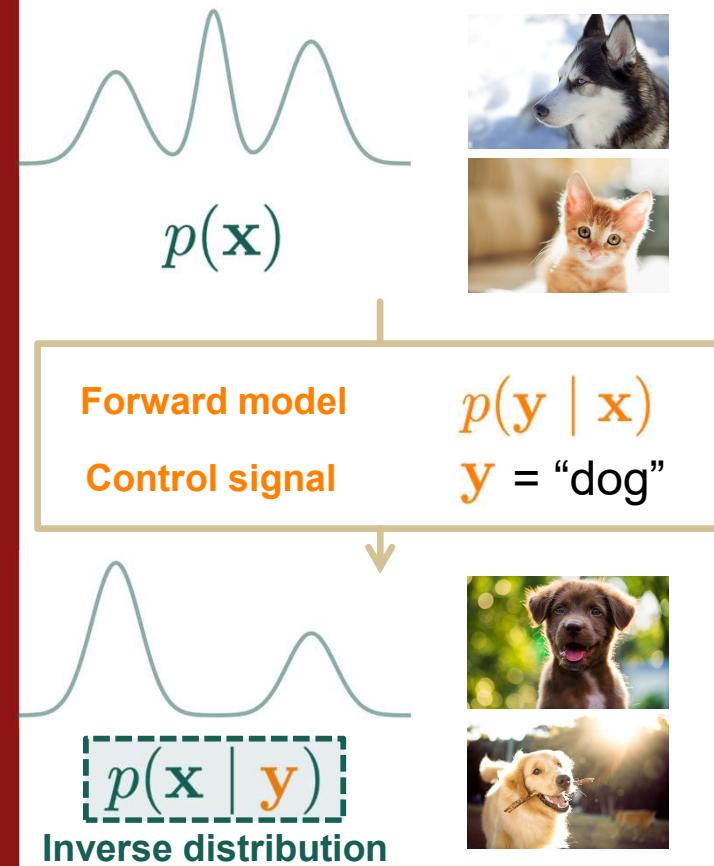
$$\mathbb{E}_{(x,y) \sim p_{\text{data}}(x,y)} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{t \sim \mathcal{U}[0,T]} \| \epsilon_\theta(x_t, t; y) - \epsilon \|_2^2$$

Need a suitable architecture



Stanford University

Control the generation process



Bayes' rule:

$$p(x | y) = \frac{p(x)p(y | x)}{p(y)}$$

Annotations: $p(x)$ and $p(y | x)$ are marked with green checkmarks; $p(y)$ is marked with a red X.

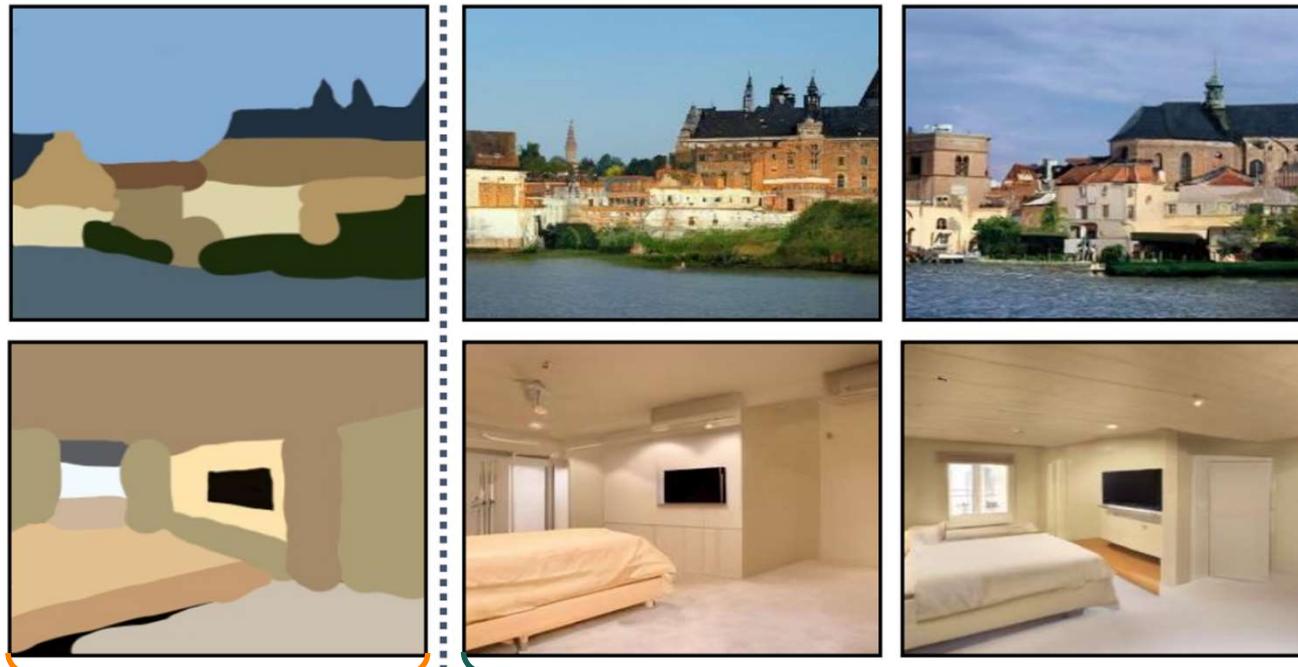
Bayes' rule for score functions:

$$\begin{aligned} \nabla_x \log p(x | y) &= \nabla_x \log p(x) \\ &\quad + \nabla_x \log p(y | x) \\ &\quad - \boxed{\nabla_x \log p(y)} 0 \\ &= \boxed{\nabla_x \log p(x)} + \boxed{\nabla_x \log p(y | x)} \end{aligned}$$

Plug in different forward models for the same score model

Stroke to image synthesis

Stroke Painting to Image



Stroke paintings
 y

Sampled images
 $x \mid y$

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

Forward model
 $p(y \mid x)$
can be specified.

Stanford University

Language-guided image generation

y

(Prompt)

Treehouse in the
style of Studio
Ghibli animation

x | y



Forward model

$p(y | x)$

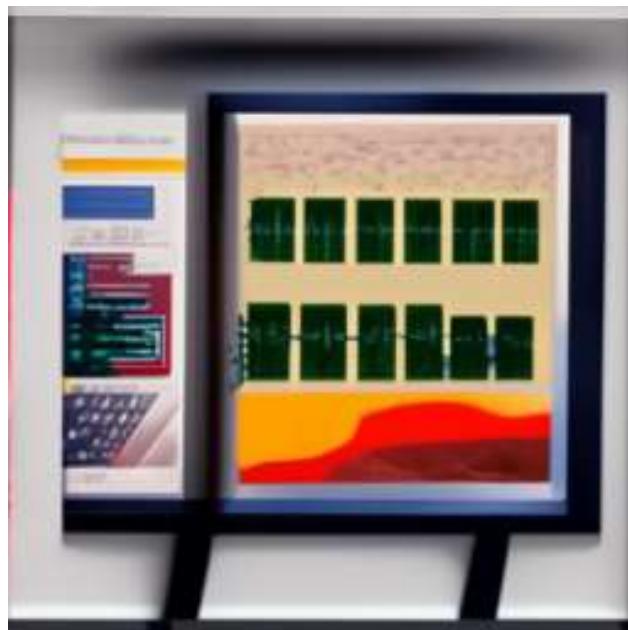
is an image
captioning neural
network.

[Work by @danielrussruss]

Stanford University

Controllable generation: Text-guided generation

An abstract painting of computer science:



A painting of the starry night by van Gogh



<https://colab.research.google.com/drive/1FuOobQOmDJuG7rGsMWfQa883A9r4HxEO?usp=sharing>

Stanford University

Classifier-free guidance

Bayes' rule:

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{x}) p(\mathbf{y} | \mathbf{x})}{p(\mathbf{y})}$$

Bayes' rule for score functions:

$$\begin{aligned}\nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{y}) &= \nabla_{\mathbf{x}} \log p(\mathbf{x}) \\&\quad + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) \\&\quad - \boxed{\nabla_{\mathbf{x}} \log p(\mathbf{y})} \ 0 \\&= \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) \\&\quad \qquad \qquad \qquad \text{Classifier obtained as the difference}\end{aligned}$$

Conditional score

Unconditional score

Classifier-free guidance

Train both a conditional and an unconditional score model (by randomly dropping the caption during training)

Combine the two models as follows

$$\begin{aligned}(1 + w) \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) &= (1 + w) (\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= (1 + w) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - w \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\end{aligned}$$

w is the classifier-guidance strength

Effect of classifier guidance



Increased classifier guidance strength (w)

Stanford University

Summary

Discrete time diffusion model as a hierarchical VAE

- Connections with score-based models (ELBO equivalence to score matching)

Continuous time diffusion models

- SDE perspective: VAE with an infinite number of latent variables
- ODE perspective: normalizing flow (exact likelihoods!)

Accelerated sampling

- Advanced numerical methods for solving ODE/SDEs
- Distillation

Controllable generation

- Classifier guidance
- Classifier-free guidance