

# **CS 4644-DL / 7643-A**

## **ZSOLT KIRA**

Generative Models:

Denoising Diffusion Probabilistic Models (DDPMs)

Slides adapted from those by Danfei Xu

# Taxonomy of Generative Models

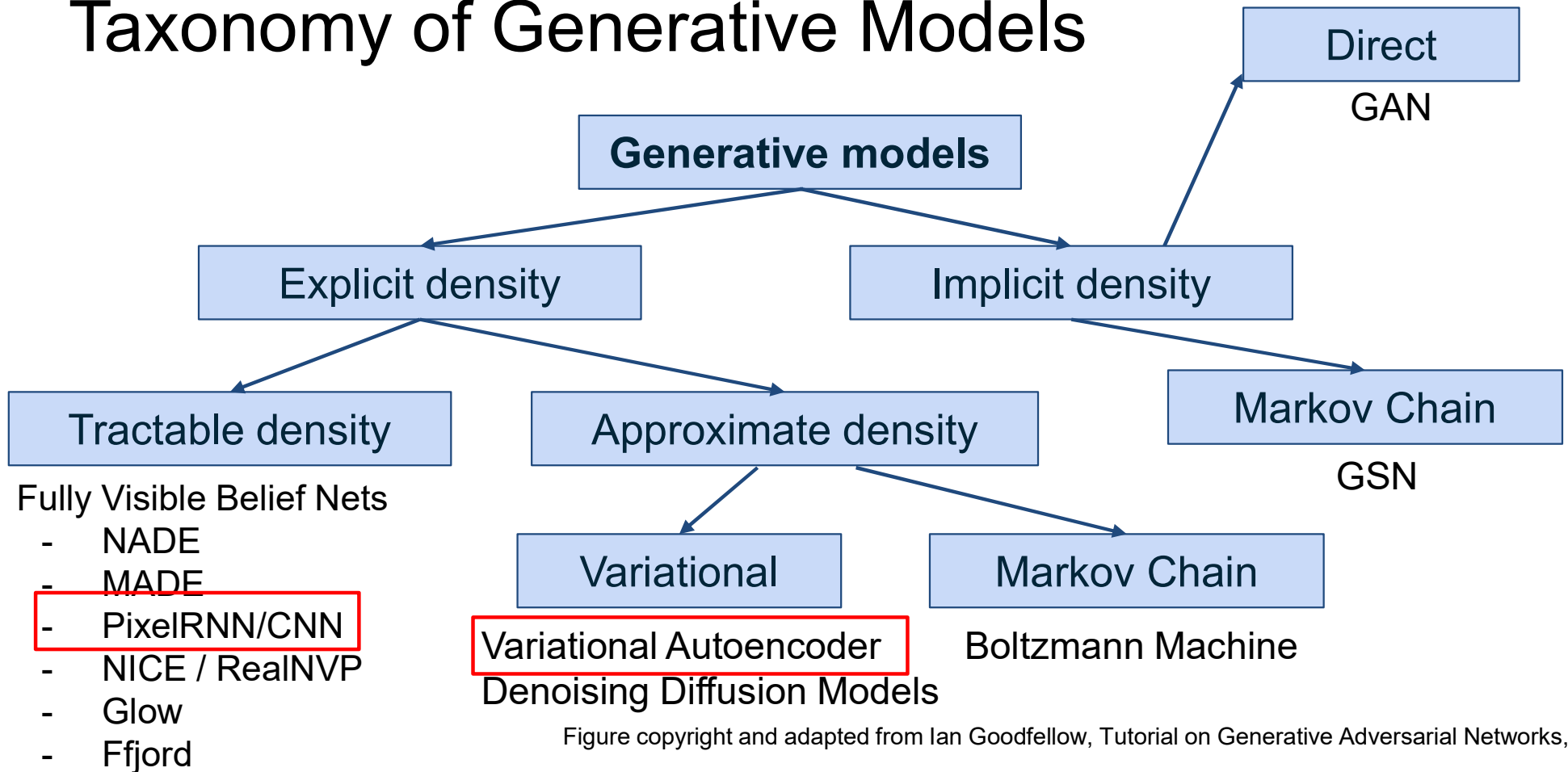


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

# PixelCNN *[van der Oord et al. 2016]*

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region (**masked convolution**)

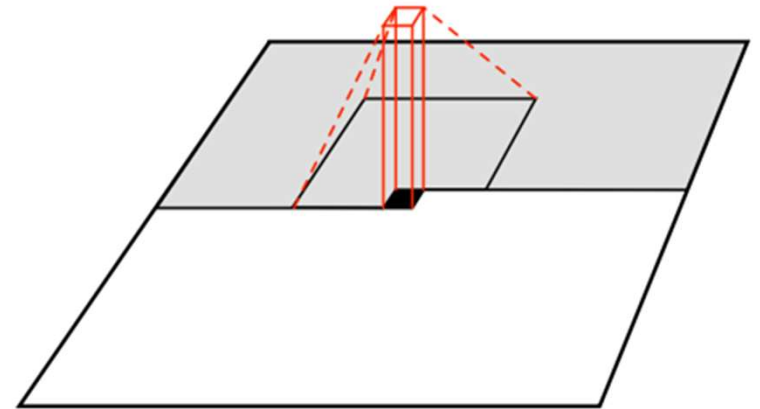


Figure copyright van der Oord et al., 2016. Reproduced with permission.

$$\begin{aligned}
\log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[ \log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\
&= \mathbf{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\
&= \mathbf{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z) q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)}) q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\
&= \mathbf{E}_z \left[ \log p_{\theta}(x^{(i)} | z) \right] - \mathbf{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\
&= \mathbf{E}_z \left[ \log p_{\theta}(x^{(i)} | z) \right] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)}))
\end{aligned}$$

↑  
Decoder network gives  $p_{\theta}(x|z)$ , can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick. see paper.)

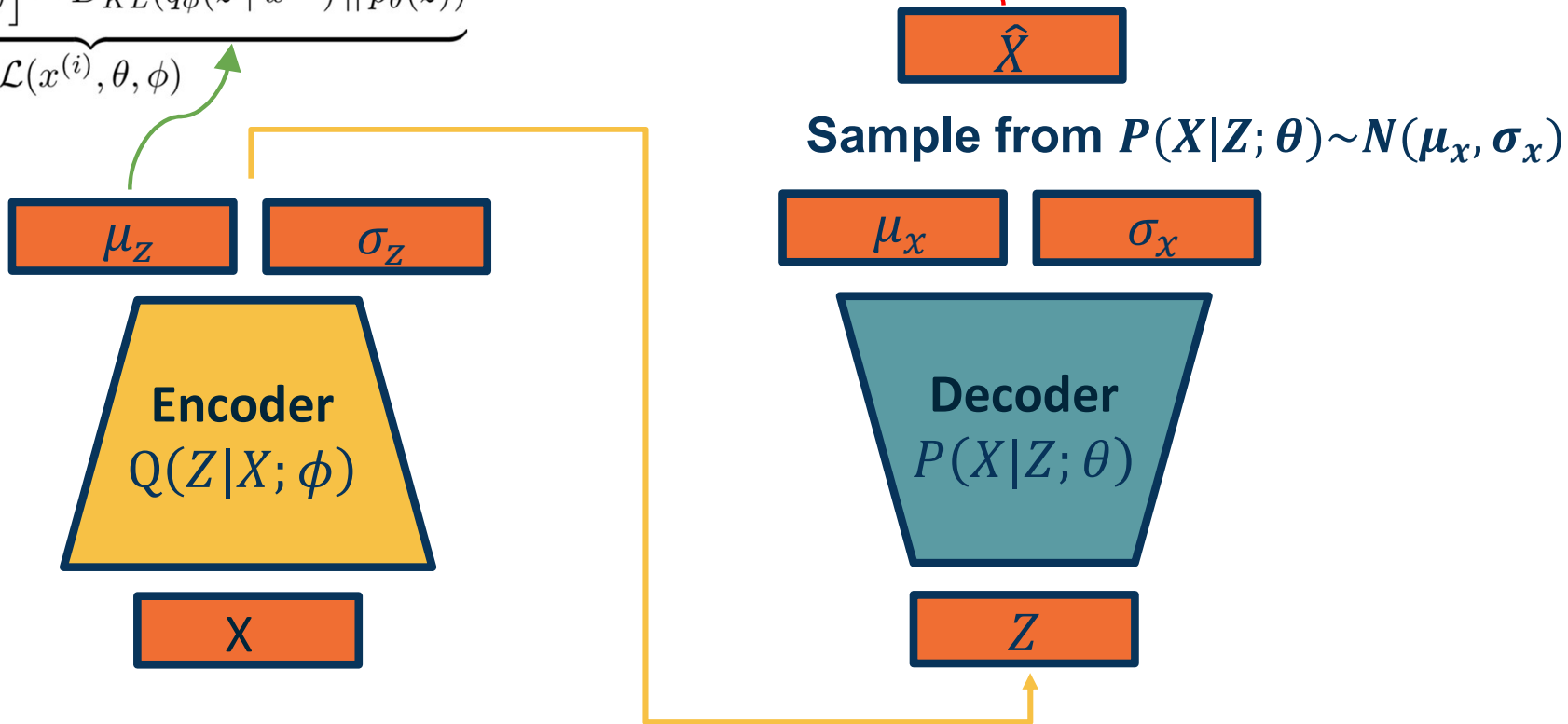
↑  
This KL term (between Gaussians for encoder and  $z$  prior) has nice closed-form solution!

↑  
 $p_{\theta}(z|x)$  intractable (saw earlier), can't compute this KL term :( But we know KL divergence always  $\geq 0$ .

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Maximize likelihood of original input being reconstructed



From CS231n, Fei-Fei Li, Justin Johnson, Serena Young



# Taxonomy of Generative Models

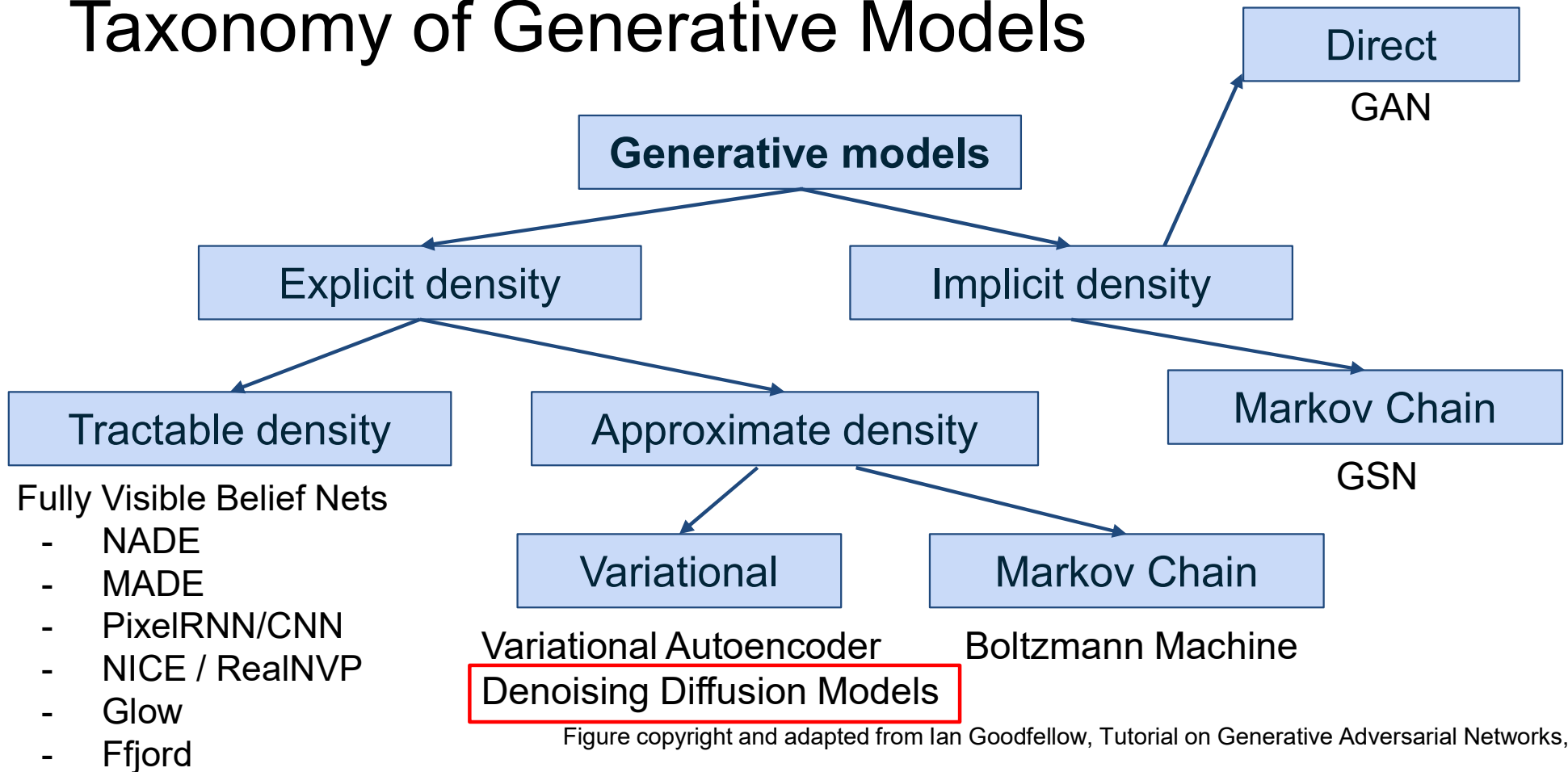


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

# Denoising Diffusion Probabilistic Models (DDPMs)

And Conditional Diffusion Models

TEXT DESCRIPTION

An astronaut    Teddy bears    A bowl of  
soup

riding a horse    lounging in a tropical resort  
in space    playing basketball with cats in  
space

in a photorealistic style    in the style of Andy  
Warhol    as a pencil drawing



DALL-E 2



<https://openai.com/dall-e-2/>



TEXT DESCRIPTION

An astronaut Teddy bears A bowl of soup

mixing sparkling chemicals as mad scientists shopping for groceries working on new AI research

as kids' crayon art on the moon in the 1980s underwater with 1990s technology



DALL-E 2



<https://openai.com/dall-e-2/>



<https://openai.com/dall-e-2/>

main 1 branch 0 tags

Go to file Add file Code

<b>pessier</b> Release under CreativeML Open RAIL M License ... 69ae4b3 on Aug 22 29 commits		
assets	Release under CreativeML Open RAIL M License	2 months ago
configs	stable diffusion	3 months ago
data	stable diffusion	3 months ago
ldm	stable diffusion	3 months ago
models	add configs for training unconditional/class-conditional Idms	11 months ago
scripts	Release under CreativeML Open RAIL M License	2 months ago
LICENSE	Release under CreativeML Open RAIL M License	2 months ago
README.md	Release under CreativeML Open RAIL M License	2 months ago
Stable_Diffusion_v1_Model_Card.md	Release under CreativeML Open RAIL M License	2 months ago
environment.yaml	Release under CreativeML Open RAIL M License	2 months ago
main.py	add configs for training unconditional/class-conditional Idms	11 months ago
notebook_helpers.py	add code	11 months ago
setup.py	add code	11 months ago

README.md

## Stable Diffusion

Stable Diffusion was made possible thanks to a collaboration with [Stability AI](#) and [Runway](#) and builds upon our previous work:

[High-Resolution Image Synthesis with Latent Diffusion Models](#)  
Robin Rombach\*, Andreas Blattmann\*, Dominik Lorenz, Patrick Esser, Björn Ommer  
[CVPR '22 Oral](#) | [GitHub](#) | [arXiv](#) | [Project page](#)

### About

A latent text-to-image diffusion model

[ommer-lab.com/research/latent-diffus...](#)

Readme  
View license  
33k stars  
321 watching  
5k forks

### Releases

No releases published

### Packages

No packages published

### Contributors 7



### Languages



https://github.com/CompVis/stable-diffusion

# Landscape Highlights of Diffusion Models (Nov 2022)

- basic principles
  - *Diffusion probabilistic models* ([Sohl-Dickstein et al., 2015](#))
  - *Noise-conditioned score network (NCSN)*; [Yang & Ermon, 2019](#))
  - *Denoising diffusion probabilistic models (DDPM)*; [Ho et al. 2020](#))
- conditional & high-res image generation
  - *Classifier-guided conditional generation* ([Dhariwal and Nichole, 2021](#))
  - *Classifier-free Diffusion Guidance* ([Ho and Salimans, 2022](#))
  - *Latent-space Diffusion (StableDiffusion)*; [Rombach and Blattmann et al., 2022](#))
- new applications
  - *Planning with Diffusion for Flexible Behavior Synthesis (Diffuser)*; [Janner et al., 2022](#))
  - *DreamFusion: Text-to-3D using 2D Diffusion* ([Poole and Jain et al., 2022](#))
  - *Make-A-Video: Text-to-Video Generation without Text-Video Data* ([Singer et al., 2022](#))

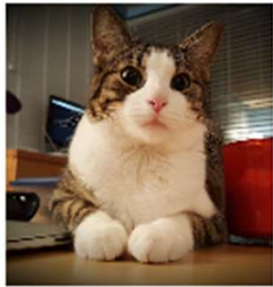
# Landscape Highlights of Diffusion Models (Nov 2022)

- *Diffusion probabilistic models* ([Sohl-Dickstein et al., 2015](#))
- *Noise-conditioned score network (NCSN)*; ([Yang & Ermon, 2019](#))
- *Denoising diffusion probabilistic models (DDPM)*; ([Ho et al. 2020](#))
- *Classifier-guided conditional generation* ([Dhariwal and Nichole, 2021](#))
- *Classifier-free Diffusion Guidance* ([Ho and Salimans, 2022](#))
- *Latent-space Diffusion (StableDiffusion)*; ([Rombach and Blattmann et al., 2022](#))
- *Planning with Diffusion for Flexible Behavior Synthesis (Diffuser)*; ([Janner et al., 2022](#))
- *DreamFusion: Text-to-3D using 2D Diffusion* ([Poole and Jain et al., 2022](#))
- *Make-A-Video: Text-to-Video Generation without Text-Video Data* ([Singer et al., 2022](#))

# The Denoising Diffusion Process

image from  
dataset

$x_0$

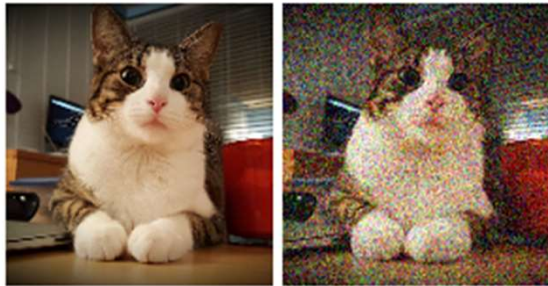


# The Denoising Diffusion Process

image from  
dataset

The “forward diffusion” process:  
add Gaussian noise each step

$x_0$  →  $x_1$  →



...

# The Denoising Diffusion Process

image from  
dataset

The “forward diffusion” process:  
add Gaussian noise each step

noise  $\mathcal{N}(0, I)$



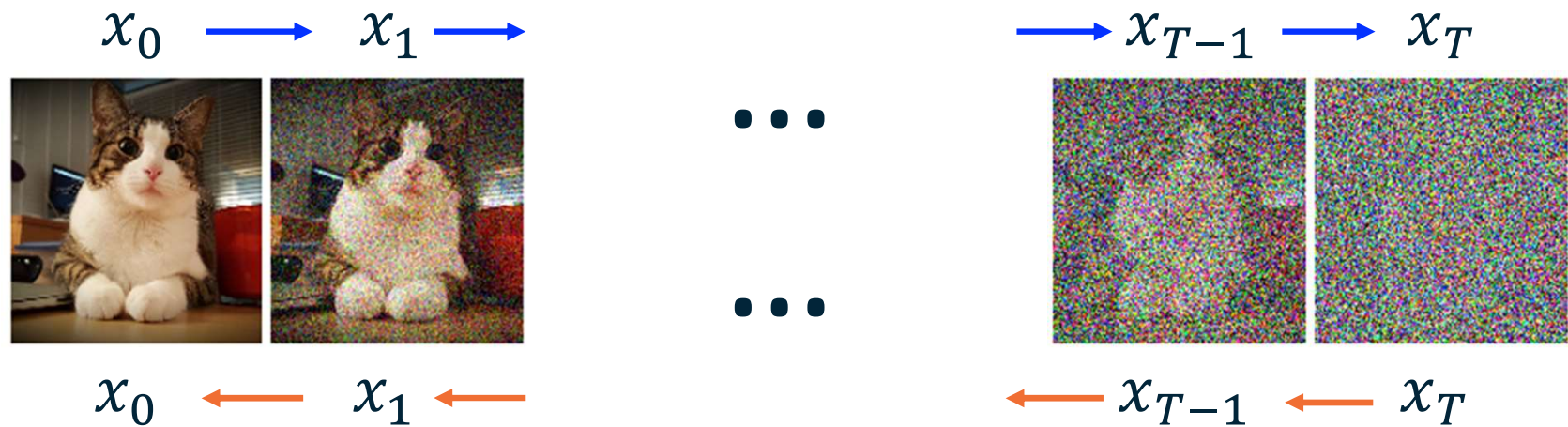


# The Denoising Diffusion Process

image from dataset

The “forward diffusion” process:  
add Gaussian noise each step

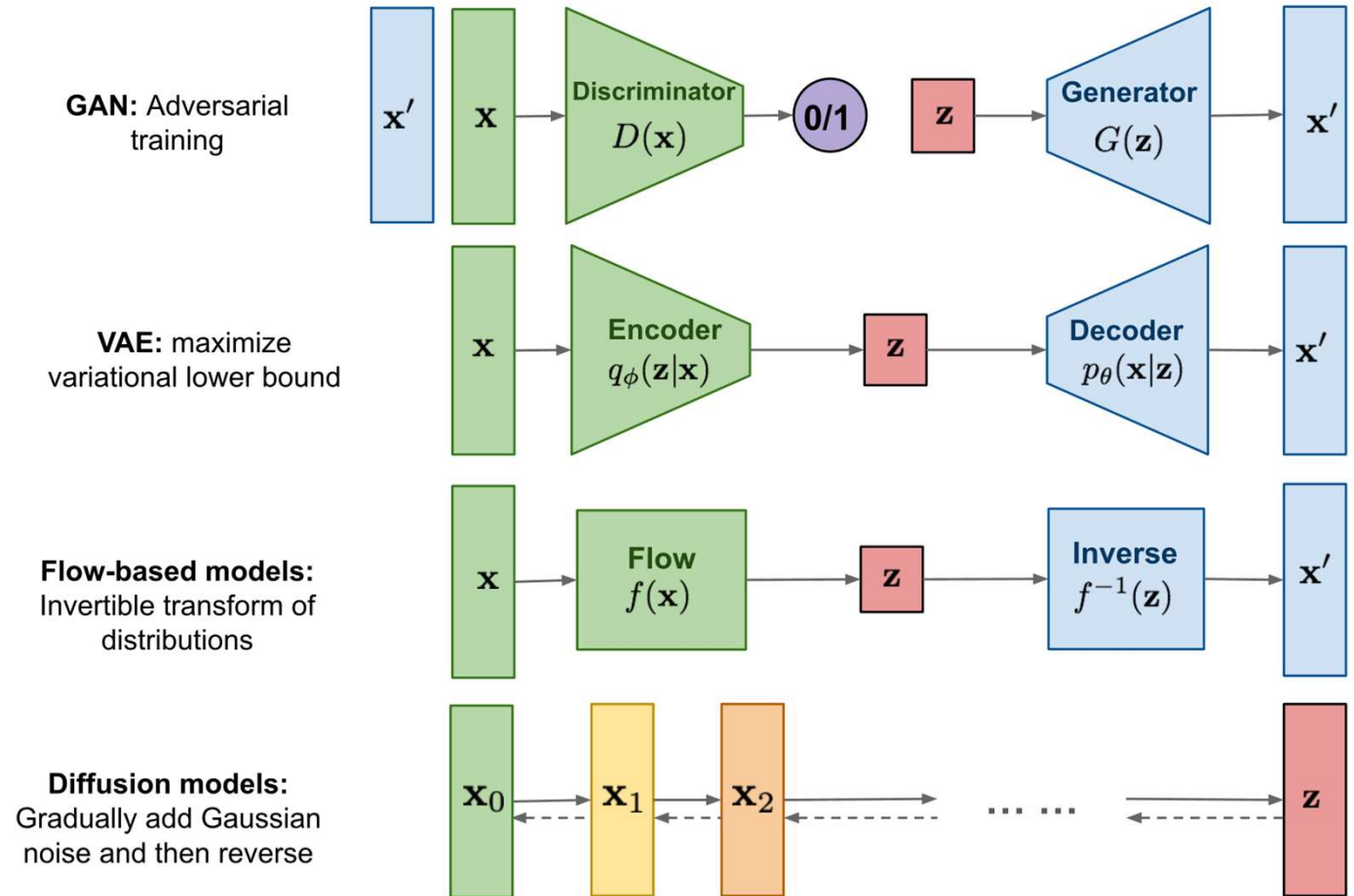
noise  $\mathcal{N}(0, I)$



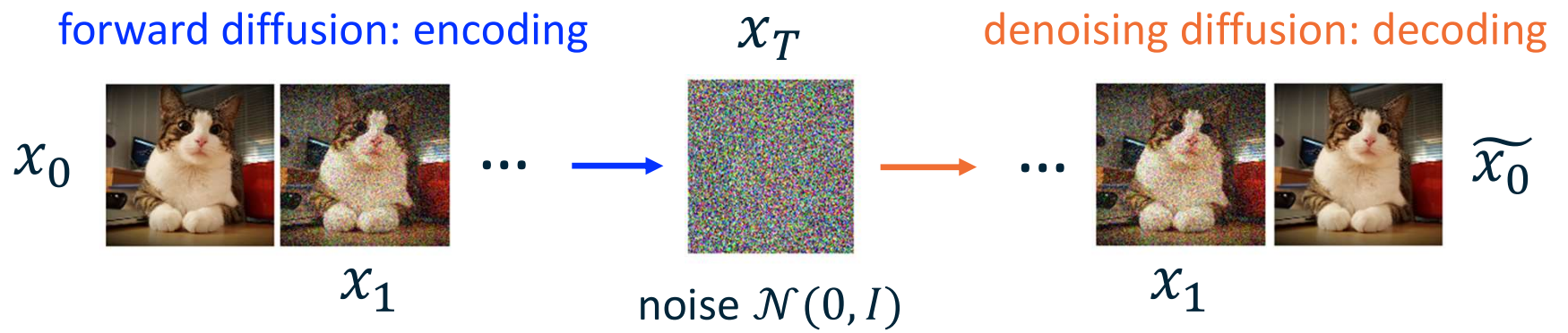
The “denoising diffusion” process:  
generate an image from noise by  
*denoising* the gaussian noises

Ties/inspiration from Annealed  
Importance Sampling in physics

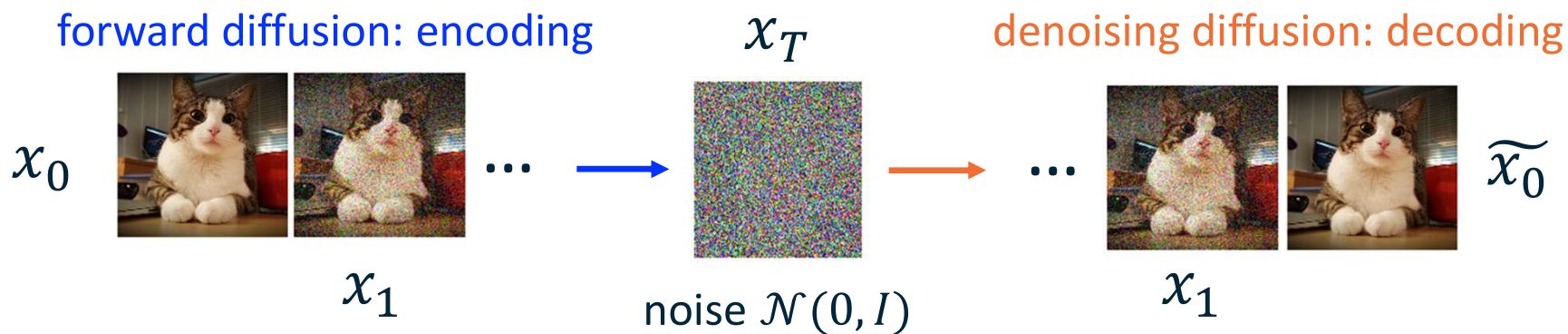
# Comparison



# Connection to VAEs



# Connection to VAEs



Known / predefined:

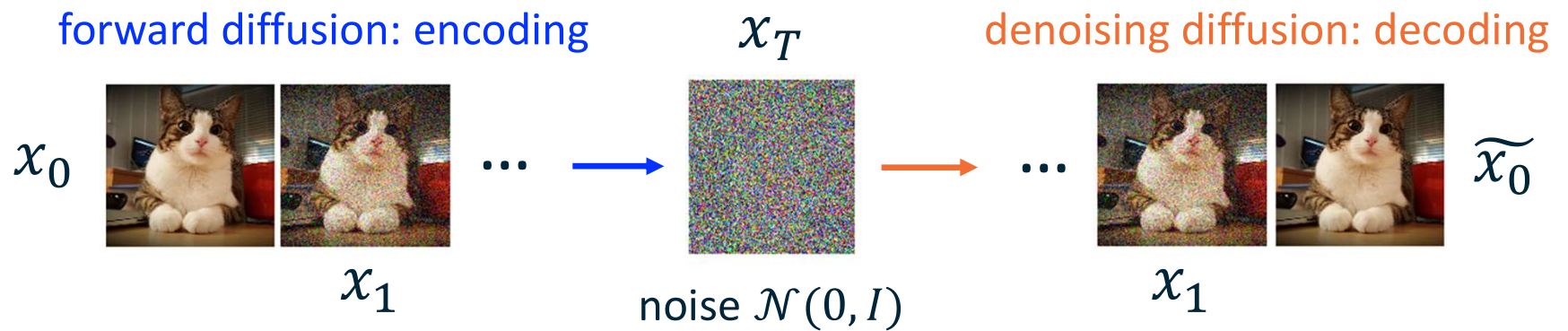
$$q(x_{1:T}|x_0)$$

Unknown / learned:

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$



# Connection to VAEs



Known / predefined:

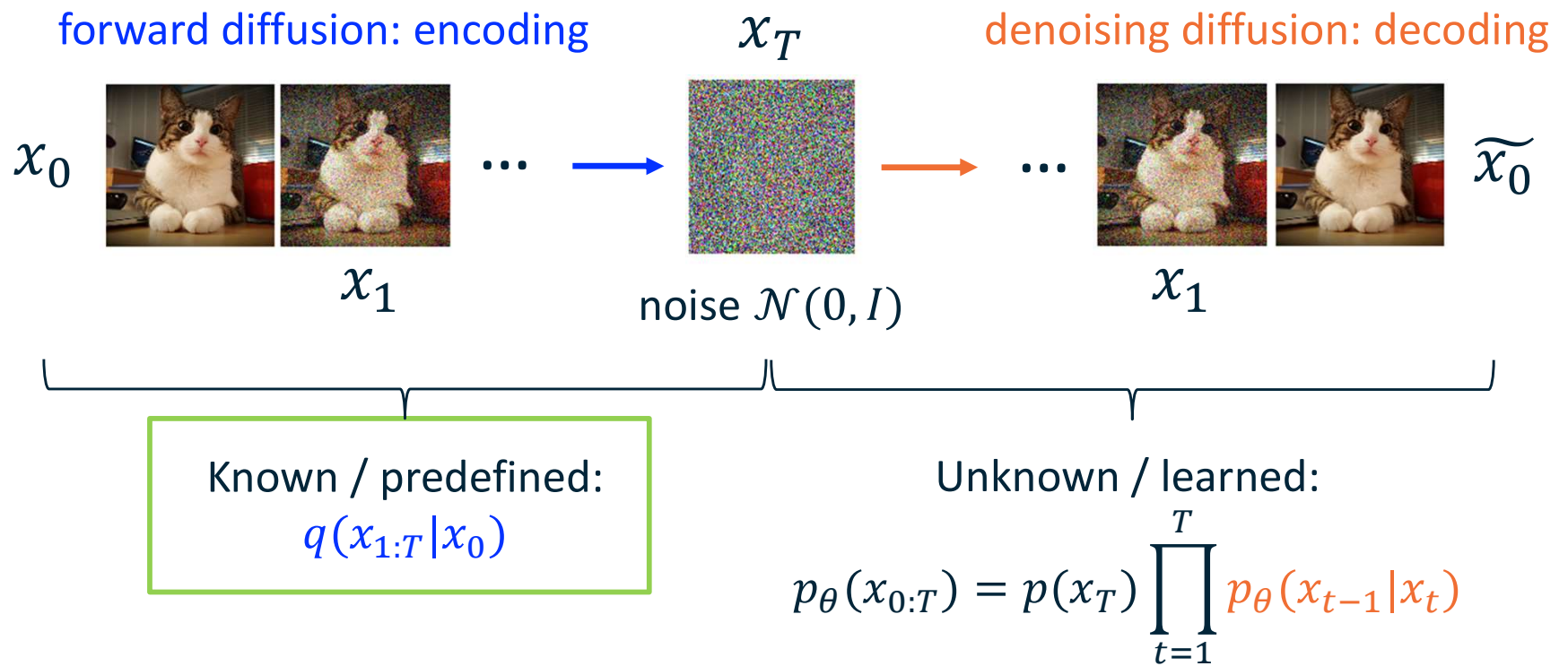
$$q(x_{1:T}|x_0)$$

Unknown / learned:

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$$

Similar to VAEs, use the denoising decoding process to generate new images.

# Connection to VAEs



# The Diffusion (Encoding) Process

The **known** forward process  $x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$

# The Diffusion (Encoding) Process

The **known** forward process  $x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \text{Probability Chain Rule (Markov Chain)}$$



# The Diffusion (Encoding) Process

The **known** forward process  $x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1 - \beta_t})x_{t-1}, \beta_t I) \quad \text{Conditional Gaussian}$$

# The Diffusion (Encoding) Process

The **known** forward process  $x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1 - \beta_t})x_{t-1}, \beta_t I) \quad \text{Conditional Gaussian}$$

Notation: A Gaussian distribution “for”  $x_t$

# The Diffusion (Encoding) Process

The **known** forward process  $x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1 - \beta_t})x_{t-1}, \beta_t I) \quad \text{Conditional Gaussian}$$

$\beta_t$  is the *variance schedule* at the diffusion step  $t$

# The Diffusion (Encoding) Process

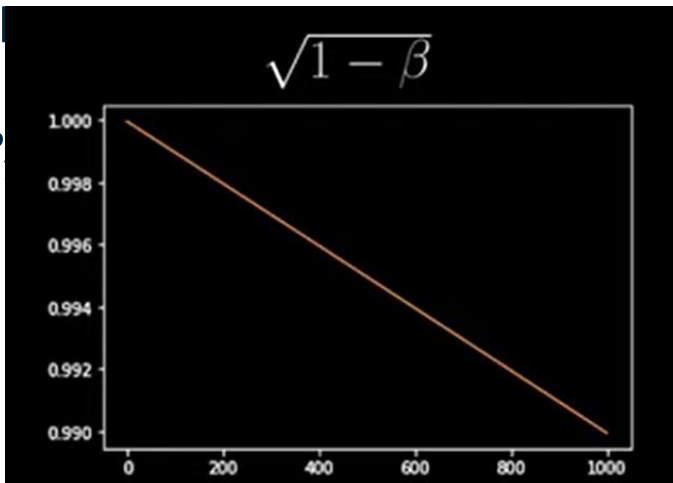
The **known** forward process  $x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1-\beta_t})x_{t-1}, \beta_t I) \quad \text{Conditional Gaussian}$$

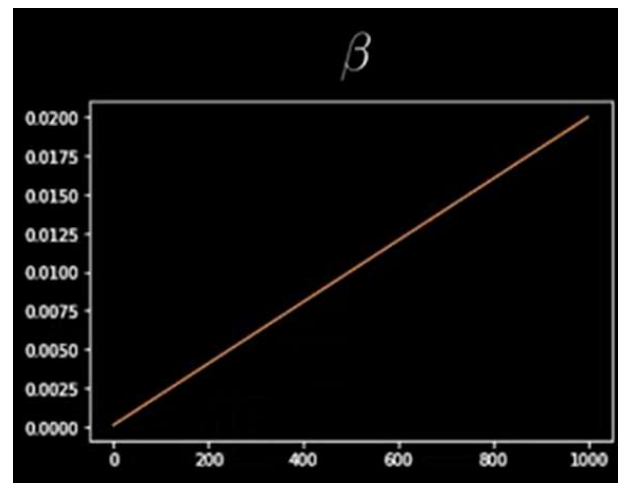
$\beta_t$  is the

$0 < \beta$



diffusion

value



$= 1000$

<https://www.youtube.com/watch?v=HoKDTa5jHvg&t=517s>

# The Diffusion (Encoding) Process

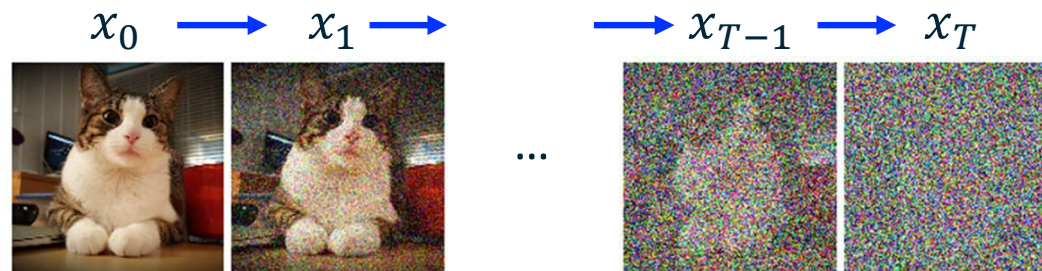
The **known** forward process  $x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1 - \beta_t})x_{t-1}, \beta_t I) \quad \text{Conditional Gaussian}$$

$\beta_t$  is the *variance schedule* at the diffusion step  $t$

$0 < \beta_1 < \beta_2 < \dots < \beta_T < 1$ , typical value range  $[0.0001, 0.02]$ , with  $T = 1000$



# The Diffusion (Encoding) Process

The **known** forward process  $x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$

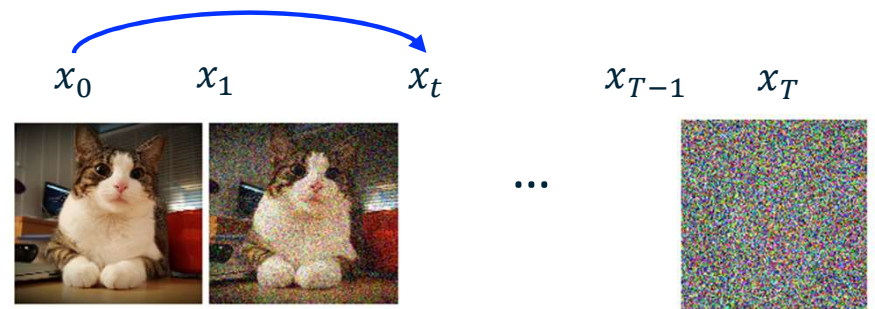
$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1 - \beta_t})x_{t-1}, \beta_t I) \quad \text{Conditional Gaussian}$$

**Nice property:** samples from an *arbitrary forward step* are also Gaussian-distributed!

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$

, where  $\alpha_t = (1 - \beta_t)$ ,  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$



<https://www.youtube.com/watch?v=HoKD1a5jHvg&t=517s>

$$\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t, \sqrt{1-\beta_t}x_{t-1}, \beta_t I)$$

$$= \sqrt{1-\beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon$$

$$= \sqrt{\alpha_t}x_{t-1} + \sqrt{1-\alpha_t}\epsilon$$

$$= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{1-\alpha_t\alpha_{t-1}}\epsilon$$

$$= \sqrt{\alpha_t\alpha_{t-1}\alpha_{t-2}}x_{t-3} + \sqrt{1-\alpha_t\alpha_{t-1}\alpha_{t-2}}\epsilon$$

$$= \sqrt{\alpha_t\alpha_{t-1}\dots\alpha_1\alpha_0}x_0 + \sqrt{1-\alpha_t\alpha_{t-1}\dots\alpha_1\alpha_0}\epsilon$$

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I) \leftarrow \boxed{= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon}$$

# The Diffusion (Encoding)

The **known** forward process  $x_0$  —

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \text{Probab}$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1-\beta_t})x_{t-1}, \beta_t I)$$

Conditional Gaussian

**Nice property:** samples from an *arbitrary forward step* are also Gaussian-distributed!

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I)$$

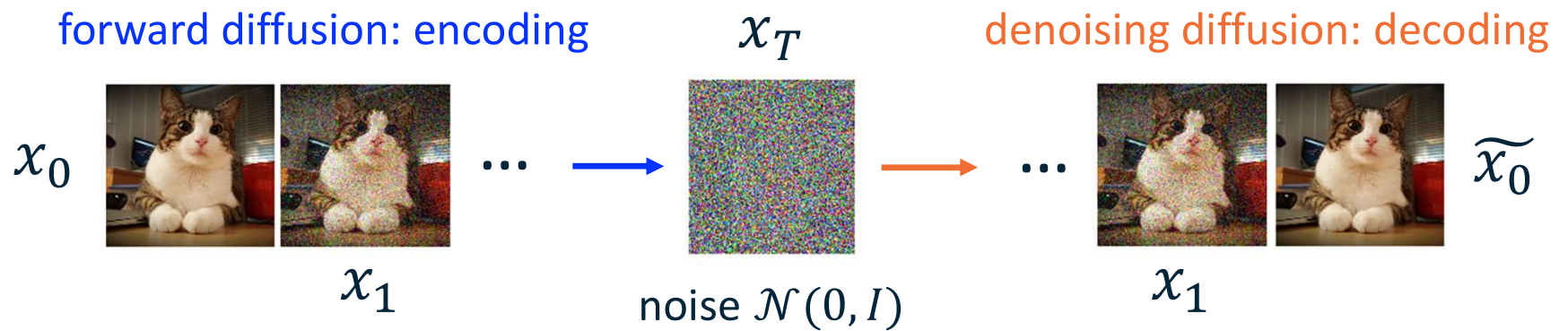
**Gaussian reparameterization trick** (recall from VAEs!):

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

(square root appears because reparameterization trick has just  $\sigma$ )

$$z = \mu + \epsilon * \sigma, \epsilon \sim N(0,1)$$

# The Diffusion and Denoising Process



Known / predefined:  
 $q(x_{1:T}|x_0)$

Unknown / learned:

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$



# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t) \quad \text{Probability Chain Rule (Markov Chain)}$$

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)) \quad \text{Conditional Gaussian}$$

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t)) \quad \text{Conditional Gaussian}$$

Want to learn time-  
dependent mean

Assume fixed / known variance  
(simplification)

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t)) \quad \text{Conditional Gaussian}$$

Want to learn time-  
dependent mean

Assume fixed / known variance  
(simplification)

How do we form a learning objective?

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t))$$

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t))$$

**High-level intuition:** derive a *ground truth denoising distribution*  $q(x_{t-1}|x_t, x_0)$  and train a neural net  $p_{\theta}(x_{t-1}|x_t)$  to match the distribution.

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_q(t))$$

**High-level intuition:** derive a *ground truth denoising distribution*  $q(x_{t-1}|x_t, x_0)$  and train a neural net  $p_\theta(x_{t-1}|x_t)$  to match the distribution.

**The learning objective:**  $\operatorname{argmin}_\theta D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))$



# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_q(t))$$

**High-level intuition:** derive a *ground truth denoising distribution*  $q(x_{t-1}|x_t, x_0)$  and train a neural net  $p_\theta(x_{t-1}|x_t)$  to match the distribution.

**The learning objective:**  $\operatorname{argmin}_\theta D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))$

What does it look like?  $q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_q(t), \Sigma_q(t))$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I)$$

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_q(t))$$

**High-level intuition:** derive a *ground truth denoising distribution*  $q(x_{t-1}|x_t, x_0)$  and train a neural net  $p_\theta(x_{t-1}|x_t)$  to match the distribution.

**The learning objective:**  $\operatorname{argmin}_\theta D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))$

What does it look like?  $q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_q(t), \Sigma_q(t))$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I) \longleftarrow \text{Recall: Gaussian reparameterization trick}$$

The “ground truth” noise that brought  $x_{t-1}$  to  $x_t$

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_q(t))$$

**High-level intuition:** derive a *ground truth denoising distribution*  $q(x_{t-1}|x_t, x_0)$  and train a neural net  $p_\theta(x_{t-1}|x_t)$  to match the distribution.

**The learning objective:**  $\operatorname{argmin}_\theta D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))$

What does it look like?  $q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_q(t), \Sigma_q(t))$

Assuming identical variance  $\Sigma_q(t)$ , we have:

$$\operatorname{argmin}_\theta D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t)) = \operatorname{argmin}_\theta w ||\mu_q(t) - \mu_\theta(x_t, t)||^2$$

Should be variance-dependent, but constant works better in practice

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_q(t))$$

**High-level intuition:** derive a *ground truth denoising distribution*  $q(x_{t-1}|x_t, x_0)$  and train a neural net  $p_\theta(x_{t-1}|x_t)$  to match the distribution.

**The learning objective:**  $\operatorname{argmin}_\theta D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))$

What does it look like?  $q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_q(t), \Sigma_q(t))$

Assuming identical variance  $\Sigma_q(t)$ , we have:

$$\operatorname{argmin}_\theta D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t)) = \operatorname{argmin}_\theta w ||\mu_q(t) - \mu_\theta(x_t, t)||^2$$

**Simplified learning objective:**  $\operatorname{argmin}_\theta ||\epsilon - \epsilon_\theta(x_t, t)||^2$

Predict the one-step noise that was added (and remove it)!

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t)) \quad \text{Conditional Gaussian}$$

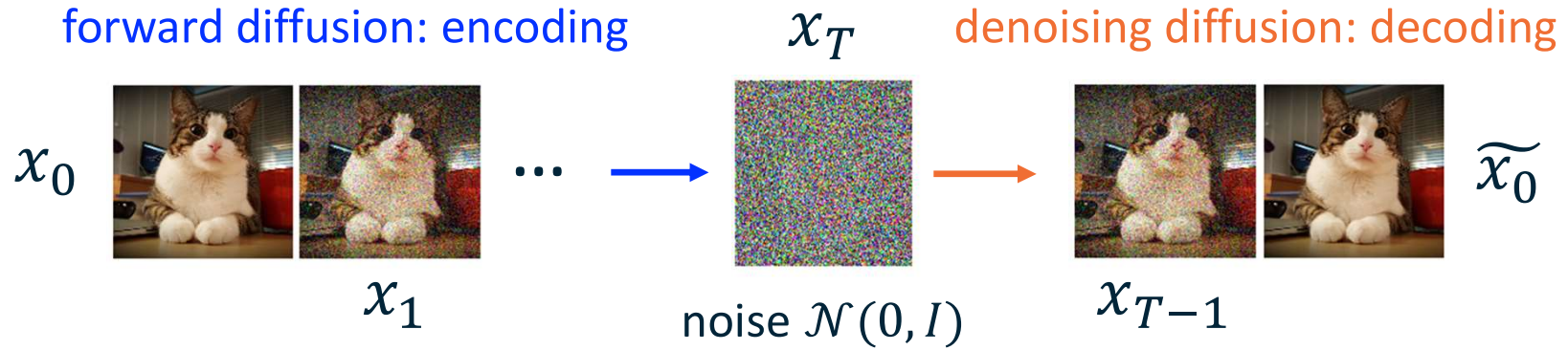
Assume fixed / known variance

How did we arrive at the learning objective?

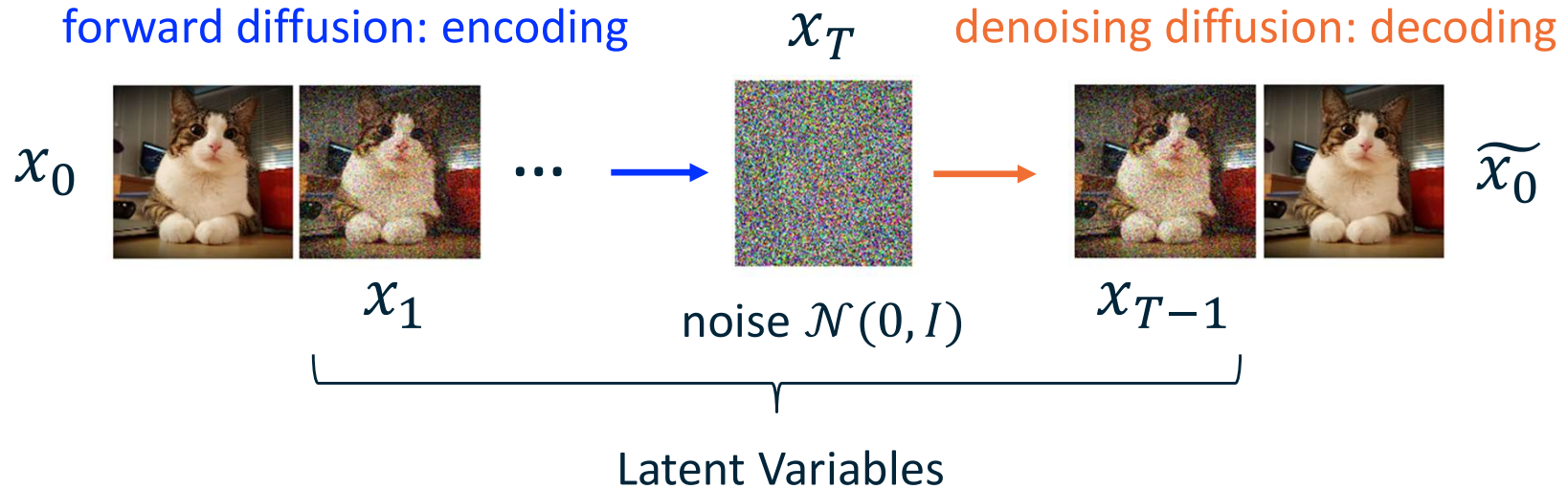
Let's go back to the basics of variational models ...

**(Quick) Derivation!**

# Connection to VAEs



# Connection to VAEs







$p(x) = \int p(x|z)p(z)dz$     Intractable to estimate!

## Variational Inference

$p(x) = \int p(x|z)p(z)dz$      Intractable to estimate!

$$\log p(x) = \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x) || p(z|x))$$
$$\geq \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right]$$

Evidence Lower Bound (ELBO) – **From last lecture on VAEs**

Variational  
Inference



Simplify to  
KL

$p(x) = \int p(x|z)p(z)dz$  Intractable to estimate!

$$\begin{aligned} \log p(x) &= \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x) || p(z|x)) \\ &\geq \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \end{aligned} \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq \mathbb{E}_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

Variational  
Inference



Simplify to  
KL

$p(x) = \int p(x|z)p(z)dz$  Intractable to estimate!

$$\log p(x) = \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x) || p(z|x))$$
$$\geq \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right]$$

Evidence Lower Bound (ELBO)

$$\log p(x_0) \geq \mathbb{E}_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= \mathbb{E}_q \left[ \log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{\prod_{t=1}^T q(x_t|x_{t-1})} \right]$$

← reverse denoising (orange arrow)  
← forward diffusion (blue arrow)

Variational  
Inference



Simplify to  
KL

$p(x) = \int p(x|z)p(z)dz$  Intractable to estimate!

$$\begin{aligned} \log p(x) &= \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x) || p(z|x)) \\ &\geq \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \end{aligned} \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq \mathbb{E}_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

Variational  
Inference



Simplify to  
KL

$$p(x) = \int p(x|z)p(z)dz \quad \text{Intractable to estimate!}$$

$$\begin{aligned} \log p(x) &= \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x)||p(z|x)) \\ &\geq \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \end{aligned} \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq \mathbb{E}_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= -\mathbb{E}_q [D_{KL}(q(x_T|x_0)||p(x_T))] - \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1)$$

Variational Inference



Simplify to KL

$p(x) = \int p(x|z)p(z)dz$  Intractable to estimate!

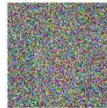
$$\log p(x) = E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x)||p(z|x))$$
$$\geq E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq E_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= -E_q[D_{KL}(q(x_T|x_0)||p(x_T))] - \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1)$$

fixed



Easy to optimize / sometimes omitted

Variational Inference



Simplify to KL

$p(x) = \int p(x|z)p(z)dz$  Intractable to estimate!

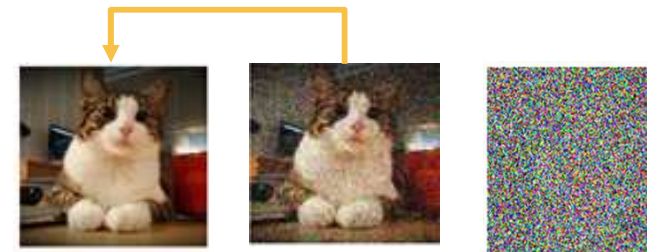
$$\log p(x) = E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x)||p(z|x))$$
$$\geq E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq E_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= -E_q [D_{KL}(q(x_T|x_0)||p(x_T))] - \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1)$$

Maximize the agreement between the predicted reverse diffusion distribution  $p_\theta$  and the “ground truth” reverse diffusion distribution  $q$







$$p(x) = \int p(x|z)p(z)dz \quad \text{Intractable to estimate!}$$

$$\log p(x) = \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x) || p(z|x))$$

$$\geq \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq \mathbb{E}_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= -\mathbb{E}_q [D_{KL}(q(x_T|x_0) || p(x_T))] - \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1)$$



$$p(x) = \int p(x|z)p(z)dz \quad \text{Intractable to estimate!}$$

$$\log p(x) = E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x)||p(z|x))$$

$$\geq E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq E_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= -E_q[D_{KL}(q(x_T|x_0)||p(x_T))] - \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1)$$

$$q(x_{t-1}|x_t) = q(x_{t-1}|x_t, x_0) \quad \text{(markov assumption)}$$

$$= \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)} \quad \text{(Bayes rule)}$$

$$= \frac{\mathcal{N}(x_t; \sqrt{\bar{a}_t}x_{t-1}, \beta_t I) \mathcal{N}(x_{t-1}; \sqrt{\bar{a}_{t-1}}x_{t-1}, (1-\bar{a}_{t-1})I)}{\mathcal{N}(x_t; \sqrt{\bar{a}_t}x_0, (1-\bar{a}_{t-1})I)}$$

$$\propto \mathcal{N}\left(x_{t-1}; \frac{\sqrt{\bar{a}_t}(1-\bar{a}_{t-1})x_t + \sqrt{\bar{a}_{t-1}}(1-\bar{a}_t)x_0}{1-\sqrt{\bar{a}_t}}, \Sigma_q(t)\right) \quad \text{(Property of Gaussian)}$$



$$p(x) = \int p(x|z)p(z)dz \quad \text{Intractable to estimate!}$$

$$\log p(x) = \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x) || p(z|x))$$

$$\geq \mathbb{E}_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq \mathbb{E}_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= -\mathbb{E}_q [D_{KL}(q(x_T|x_0) || p(x_T))] - \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1)$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_q(t), \Sigma_q(t))$$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I)$$

Proof using bayes rule and gaussian reparameterization trick



$$p(x) = \int p(x|z)p(z)dz \quad \text{Intractable to estimate!}$$

$$\log p(x) = E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x)||p(z|x))$$

$$\geq E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq E_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= -E_q[D_{KL}(q(x_T|x_0)||p(x_T))] - \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1)$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_q(t), \Sigma_q(t))$$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I)$$

Proof using bayes rule and gaussian reparameterization trick

The “ground truth” noise that brought  $x_{t-1}$  to  $x_t$



$$p(x) = \int p(x|z)p(z)dz \quad \text{Intractable to estimate!}$$

$$\log p(x) = E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x)||p(z|x))$$

$$\geq E_q \left[ \log \frac{p(x|z)p(z)}{q(z|x)} \right] \quad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq E_q \left[ \log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \quad x = x_0, z = x_{1:T}$$

... (derivation omitted, see Sohl-Dickstein *et al.*, 2015 Appendix B)

$$= -E_q[D_{KL}(q(x_T|x_0)||p(x_T))] - \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1)$$

Minimize the difference of distribution means (assuming identical variance)

$$\operatorname{argmin}_\theta w \|\mu_q(t) - \mu_\theta(x_t, t)\|^2$$



# Learning the Denoising Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t)) \quad \text{Conditional Gaussian}$$

Learning objective:  $\operatorname{argmin}_{\theta} \|\mu_q(t) - \mu_{\theta}(x_t, t)\|^2$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I)$$



# Learning the Denoising Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t)) \quad \text{Conditional Gaussian}$$

Learning objective:  $\operatorname{argmin}_{\theta} \|\mu_q(t) - \mu_{\theta}(x_t, t)\|^2$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I)$$

Do we actually need to learn the entire  $\mu_{\theta}(x_t, t)$ ?



# Learning the Denoising Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t)) \quad \text{Conditional Gaussian}$$

Learning objective:  $\operatorname{argmin}_{\theta} \|\mu_q(t) - \mu_{\theta}(x_t, t)\|^2$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I)$$

known during inference     
 Unknown during inference     
 Recall: this is the “ground truth” noise that brought  $x_{t-1}$  to  $x_t$





# Learning the Denoising Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma(t)) \quad \text{Conditional Gaussian}$$

Learning objective:  $\operatorname{argmin}_\theta \|\mu_q(t) - \mu_\theta(x_t, t)\|^2$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I)$$

known during inference     
 Unknown during inference     
 Recall: this is the “ground truth” noise that brought  $x_{t-1}$  to  $x_t$

Idea: just learn  $\epsilon$  with  $\epsilon_\theta(x_t, t)$ !



# Learning the Denoising Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t)) \quad \text{Conditional Gaussian}$$

Simplified learning objective:  $\operatorname{argmin}_{\theta} \|\epsilon - \epsilon_{\theta}(x_t, t)\|^2$



# Learning the Denoising Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t)) \quad \text{Conditional Gaussian}$$

Simplified learning objective:  $\operatorname{argmin}_{\theta} \|\epsilon - \epsilon_{\theta}(x_t, t)\|^2$

Recall: the simplified  $t$ -step forward sample:

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$$



# Learning the Denoising Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t)) \quad \text{Conditional Gaussian}$$

Simplified learning objective:  $\operatorname{argmin}_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$

Recall: the simplified  $t$ -step forward sample:

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$$



# Learning the Denoising Process

The **learned** denoising process  $x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t)) \quad \text{Conditional Gaussian}$$

Simplified learning objective:  $\operatorname{argmin}_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$

# The Denoising (Decoding) Process

The **learned** denoising process  $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t)) \quad \text{Conditional Gaussian}$$

We know how to learn

Assume fixed / known variance

$$\text{Inference time: } \mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(x_t, t) \right)$$



# The Denoising (Decoding) Process

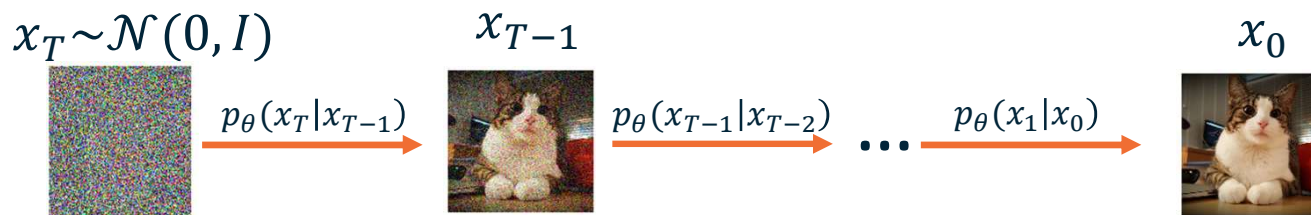
The **learned** denoising process  $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t)) \quad \text{Conditional Gaussian}$$

We know how to learn

Assume fixed / known variance



Generate new images!

# The Denoising Diffusion Algorithm

---

## Algorithm 1 Training

---

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



# The Denoising Diffusion Algorithm

---

## 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  
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$
  - 6: **until** converged
- 

---

## Algorithm 2 Sampling

---

- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 2: **for**  $t = T, \dots, 1$  **do**
  - 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
  - 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**  $\mathbf{x}_0$
-

# The Denoising Diffusion Algorithm

---

## Algorithm 1 Training

---

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

---

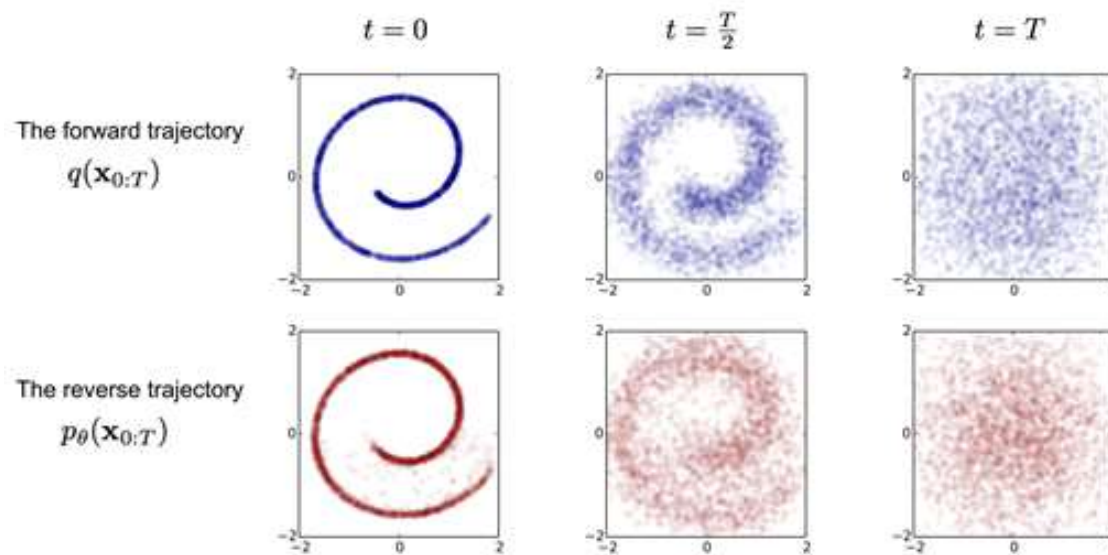
## Algorithm 2 Sampling

---

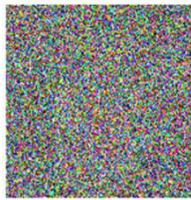
- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 2: **for**  $t = T, \dots, 1$  **do**
  - 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
  - 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}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


$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t))$$

# Visualizing the Diffusion Process on 2D data



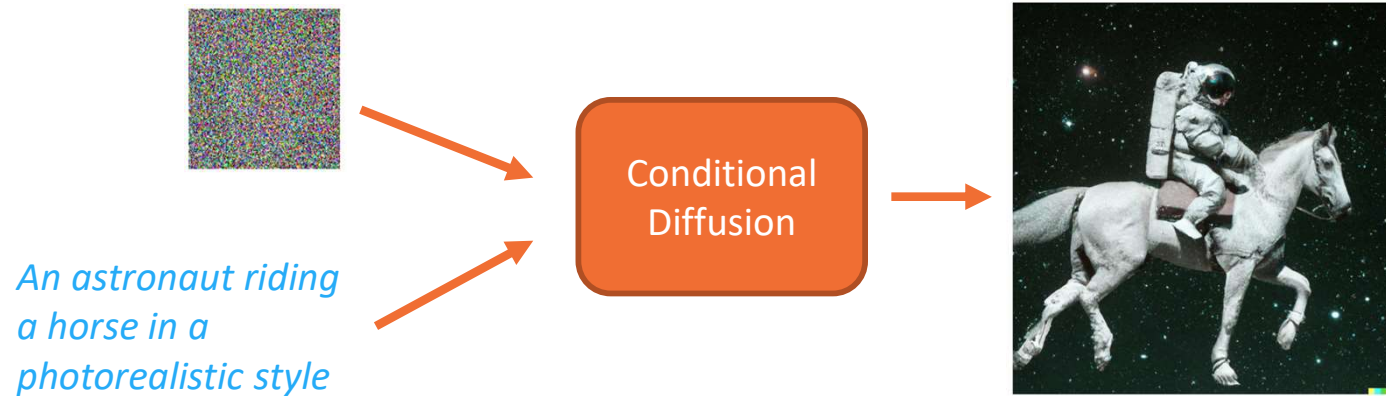
# Conditional Diffusion Models



*An astronaut riding  
a horse in a  
photorealistic style*



# Conditional Diffusion Models



Simple idea: just condition the model on some text labels  $y$ !

$$\epsilon_{\theta}(x_t, y, t)$$