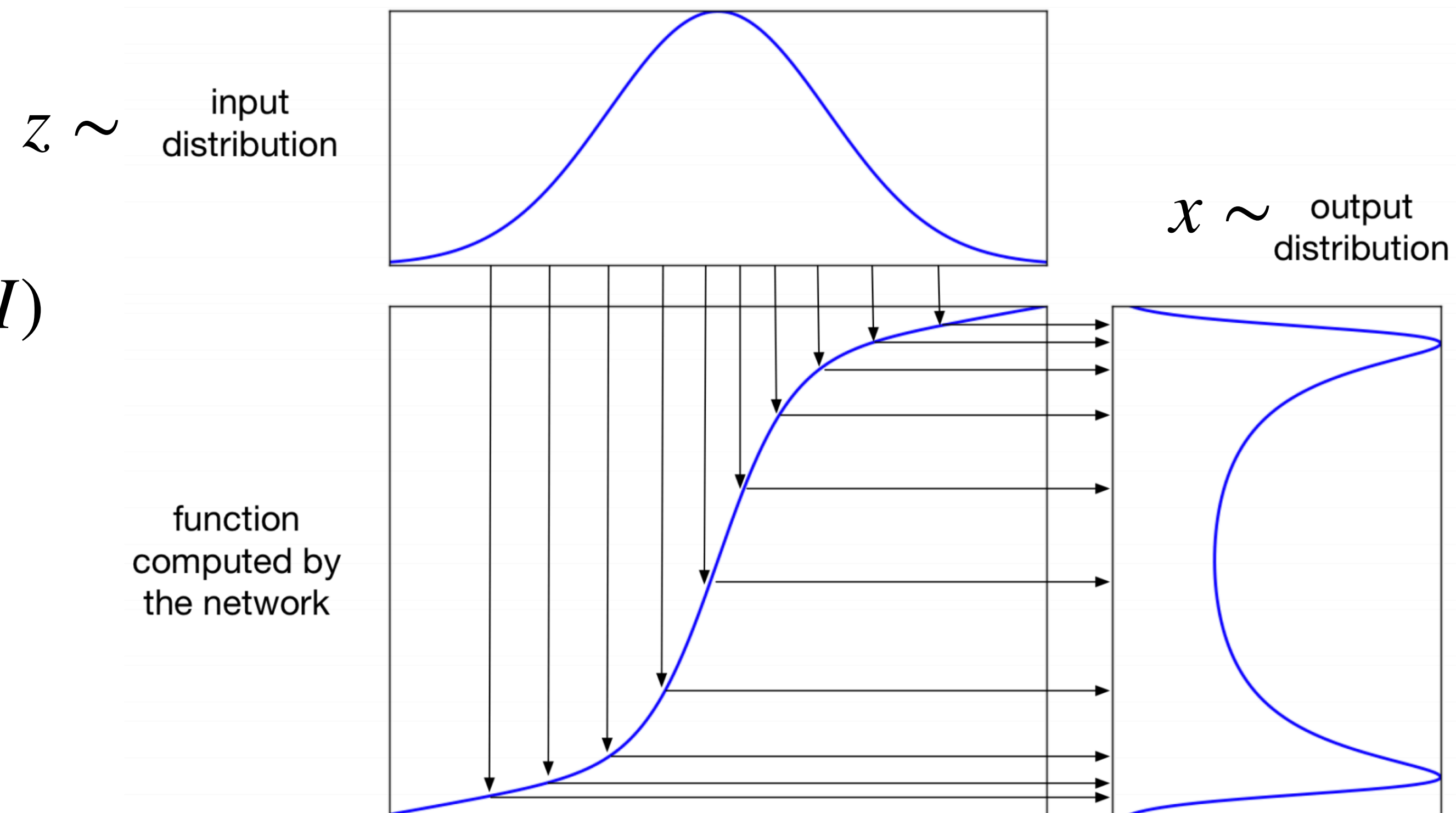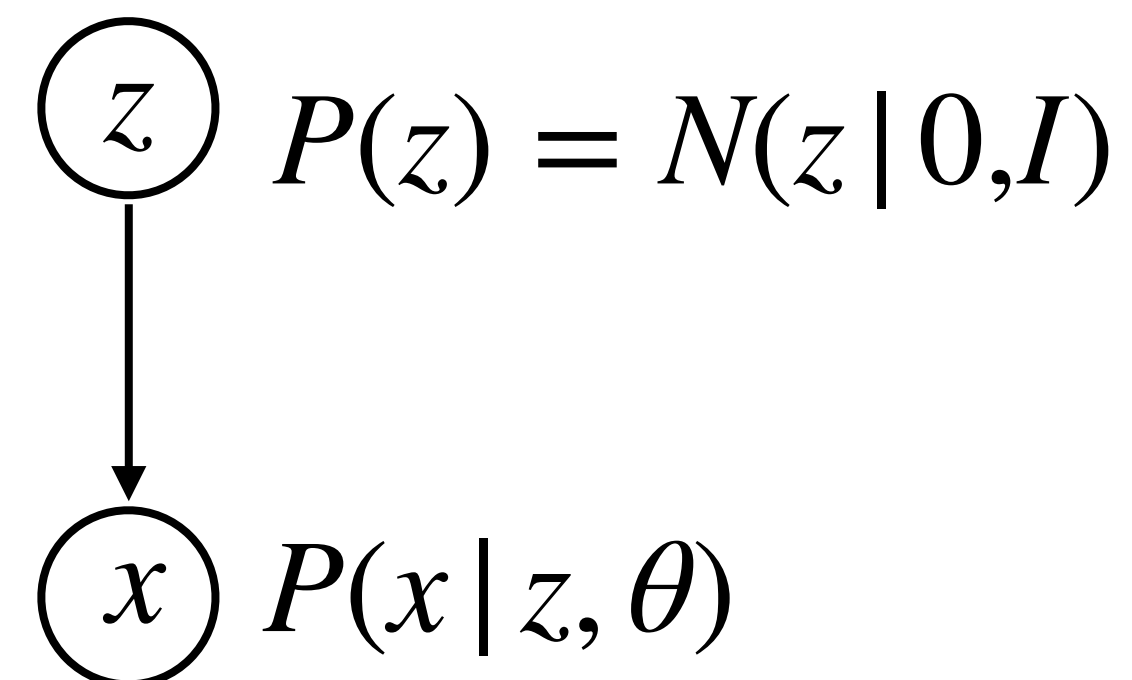# 6.7900 Machine Learning (Fall 2024)

## Lecture 22: generative models — diffusion

# Preface: from simple samples to complex objects

‣ We can realize complex distributions over (e.g.) images by drawing samples from a fixed simple distribution and then mapping such samples through a complicated function (a neural network)

$z \sim$ input distribution

$x \sim$ output distribution

$z$   $P(z) = N(z|0,I)$

$x$   $P(x|z,\theta)$

function computed by the network

# Preface: latents and how they are used

‣ The latent "degrees of freedom" z are useful if they are strongly coupled with x values; specifically, they should help us realize x's more easily conditionally on z

‣ In VAEs, we have an encoder and a decoder and both models represent this coupling between x and z.



inferential model or "encoder"

$z$   $Q(z|x, \phi)$

$x$   $Q_D(x)$   fixed data distribution

$\approx$

$z$   $P(z)$   fixed latent distribution

$x$   $P(x|z, \theta)$

generative model or "decoder"

‣ The ELBO criterion adjusts these models so they agree as distributions

# Preface: latents and how they are used

- The latent "degrees of freedom" z are useful if they are strongly coupled with x values; specifically, they should help us realize x's more easily conditionally on z

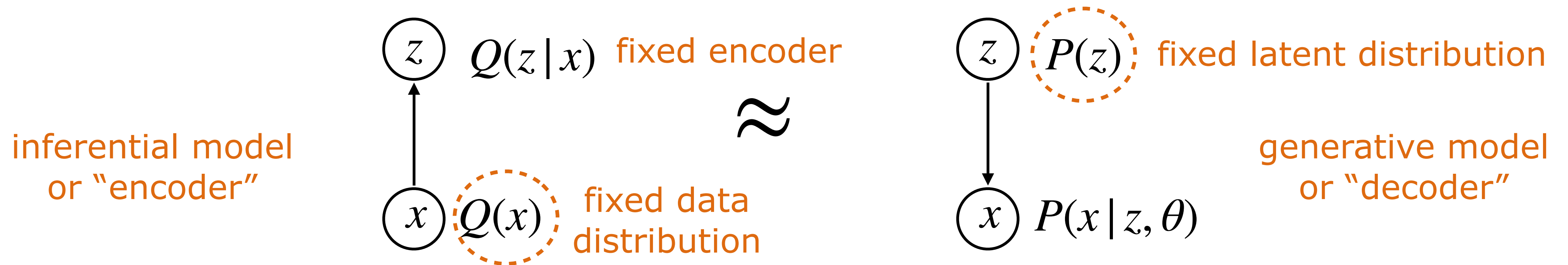- In VAEs, we have an encoder and a decoder and both models represent this coupling between x and z.



inferential model or "encoder"

$z$   $Q(z|x, \phi)$

$x$   $Q_D(x)$   fixed data distribution

$\approx$

$z$   $P(z)$   fixed latent distribution

$x$   $P(x|z, \theta)$

generative model or "decoder"

- The ELBO criterion adjusts these models so they agree as distributions, including marginals

$$\int Q_D(x)Q(z|x, \phi)dx \approx P(z) \qquad \int P(z)P(x|z, \theta)dz \approx Q_D(x)$$

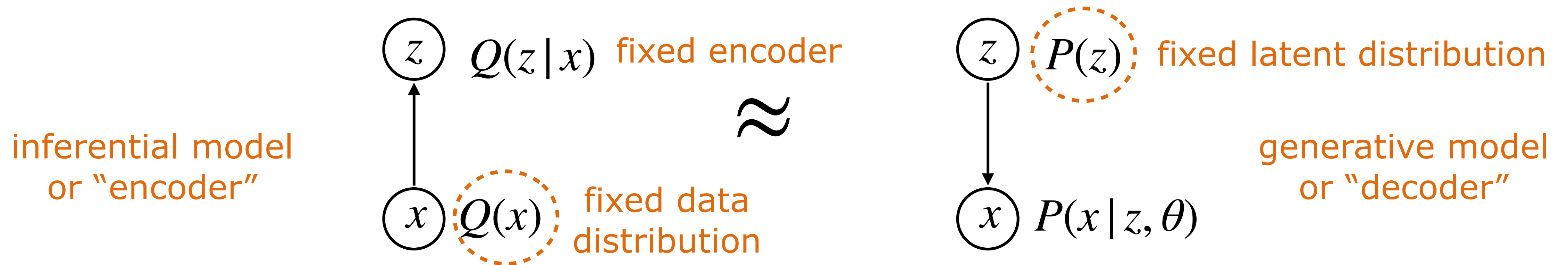# Preface: latents and how they are used

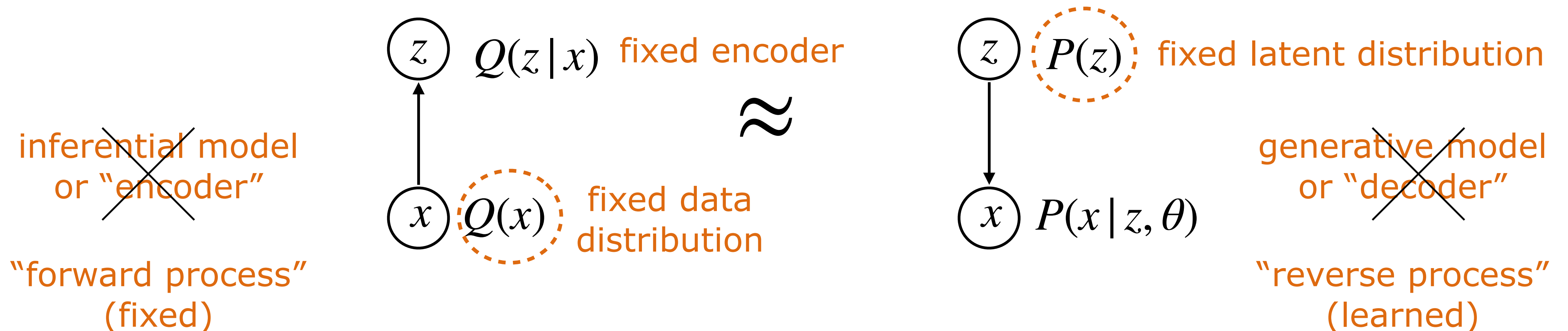‣ We can actually always just **fix the encoder** and only learn the decoder



$z$  $Q(z|x)$  fixed encoder

$\approx$

$z$  $P(z)$  fixed latent distribution

inferential model or "encoder"

$x$  $Q(x)$  fixed data distribution

$x$  $P(x|z,\theta)$

generative model or "decoder"

# Preface: latents and how they are used

‣ We can actually always just **fix the encoder** and only learn the decoder

$z$  $Q(z|x)$  fixed encoder

inferential model
or "encoder"

$x$  $Q(x)$  fixed data
distribution

$\approx$

$z$  $P(z)$  fixed latent distribution

generative model
or "decoder"

$x$  $P(x|z,\theta)$

‣ The fixed encoder should still agree with the decoder's fixed latent marginal (otherwise they'd be permanently inconsistent)
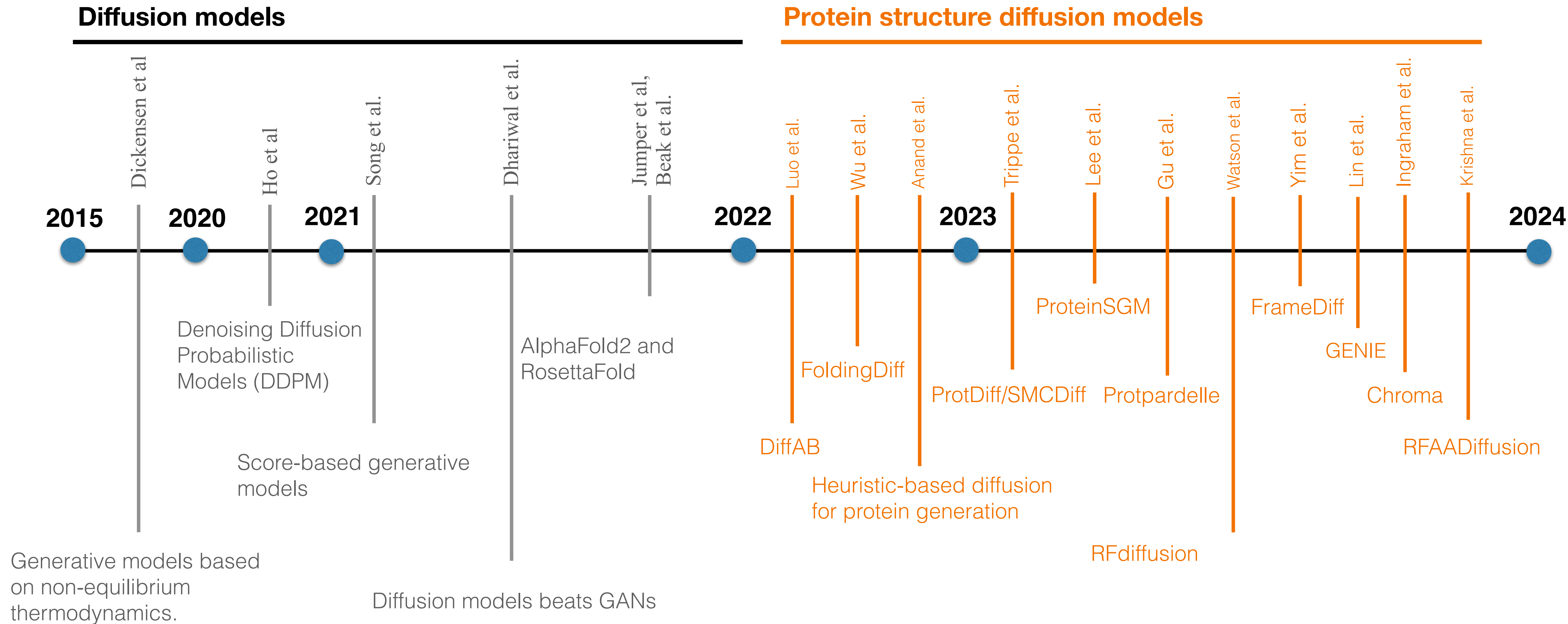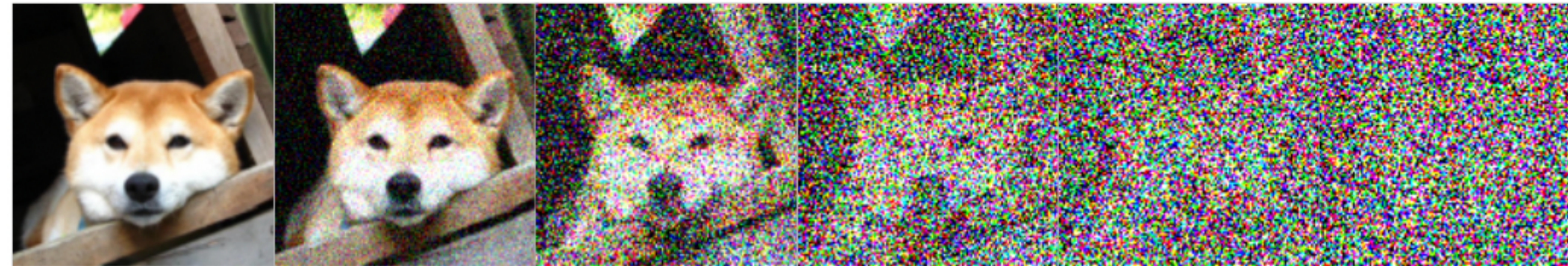
$$\int Q_D(x)Q(z|x)dx \approx P(z)$$

# Preface: latents and how they are used

‣ We can actually always just **fix the encoder** and only learn the decoder



$z$  $Q(z|x)$  fixed encoder

$z$  $P(z)$  fixed latent distribution

$\approx$

inferential model
or "encoder"

$x$  $Q(x)$  fixed data
distribution

generative model
or "decoder"

"forward process"
(fixed)

$x$  $P(x|z,\theta)$

"reverse process"
(learned)

‣ The fixed encoder should still agree with the decoder's fixed latent marginal
(otherwise they'd be permanently inconsistent)

$$\int Q_D(x)Q(z|x)dx \approx P(z)$$

# A slice of the generative "landscape"

# Ex: rapid adoption of diffusion models



[Yim et al 2024, "Diffusion models in protein structure and docking"]
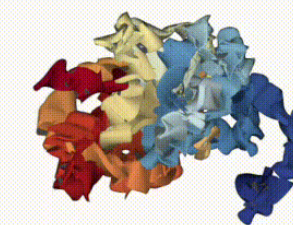
# Diffusion models

- De-noising diffusion models over images (e.g., Ho et al., Song et al.)
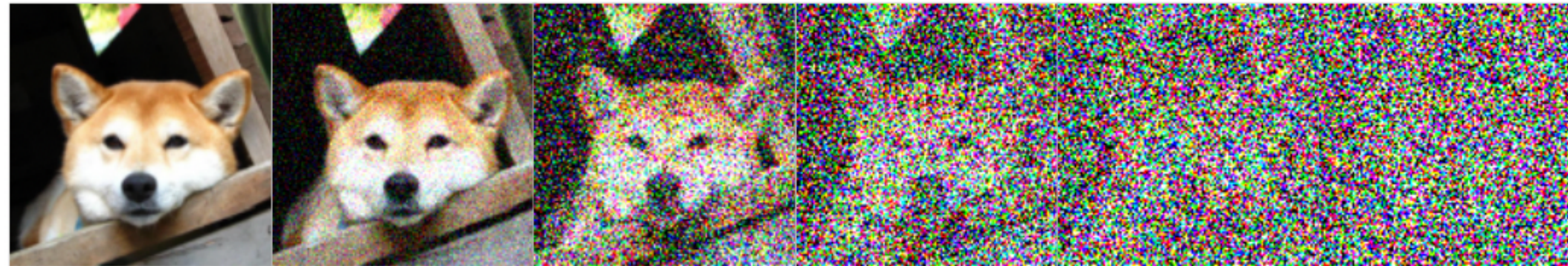


[image from Rissanen et al 2022]

$$x \longleftarrow z \sim N(0,I)$$

- 3D structures of molecules (e.g., proteins)



- Discrete diffusion for language, sequence modeling, etc.

# Diffusion motivation

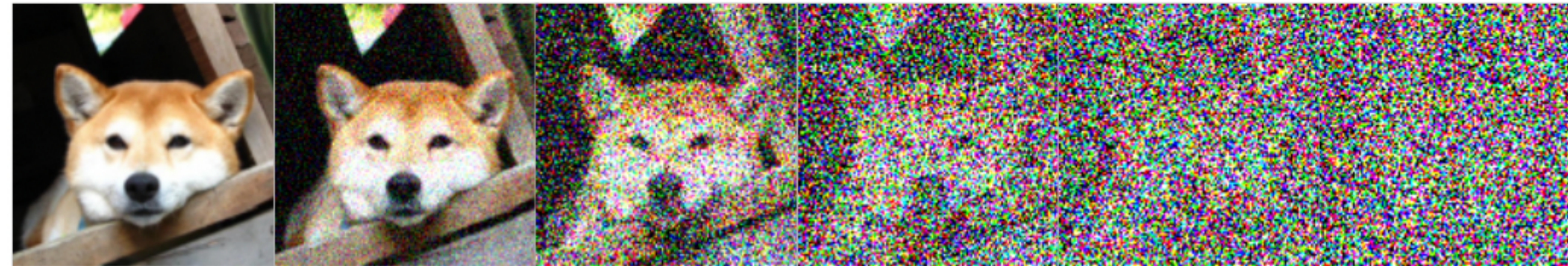‣ De-noising diffusion models over images (e.g., Ho et al., Song et al.)



$x$

[image from Rissanen et al 2022]

‣ It's really helpful to give a generative model a stack of noisy versions of the same image… different features are present at different noise levels

# Diffusion motivation

‣ De-noising diffusion models over images (e.g., Ho et al., Song et al.)
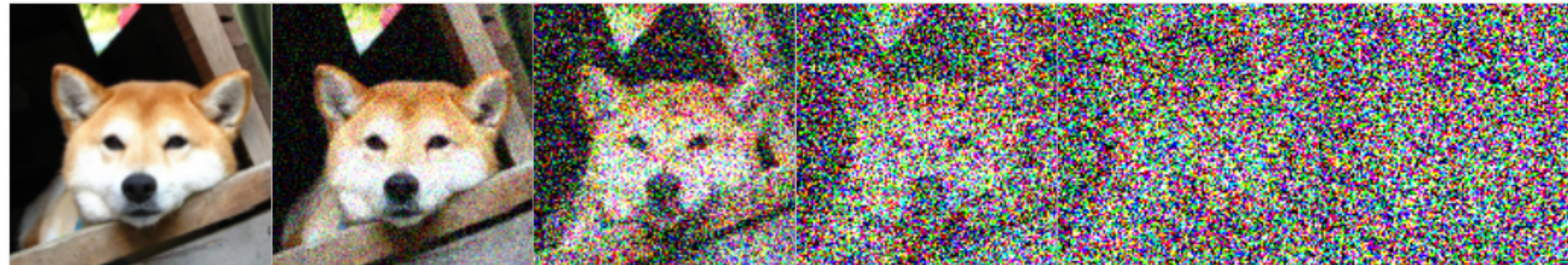


[image from
Rissanen et al 2022]

$x$

our model

‣ It's really helpful to give a generative model a stack of noisy versions of the same image… different features are present at different noise levels

‣ The goal is then to learn to successively uncover these lost features (de-noise images back into cleaner instances)

# ...Diffusion motivation



- De-noising diffusion models over images (e.g., Ho et al., Song et al.)

[image from Rissanen et al 2022]

our model

- It's really helpful to give a generative model a stack of noisy versions of the same image... different features are present at different noise levels
- The goal is then to learn to successively uncover these features (de-noise images back into cleaner instances)
- This is hard since there are potentially many images that could have similar noisy versions

# Diffusion motivation

‣ De-noising diffusion models over images (e.g., Ho et al., Song et al.)



[image from Rissanen et al 2022]

$$x \longleftarrow \qquad\qquad z \sim N(0,I)$$

‣ It's really helpful to give a generative model a stack of noisy versions of the same image… different features are present at different noise levels

‣ The goal is then to learn to successively uncover these features (de-noise images back into cleaner instances)

‣ This is hard since there are potentially many images that could have similar noisy versions

‣ Once we have our de-noising model, we can draw a noisy sample, and iterate to generate a new clean instance (works the same with other types of objects, e.g., molecule structures)

# **Mechanisms of generation: diffusion**

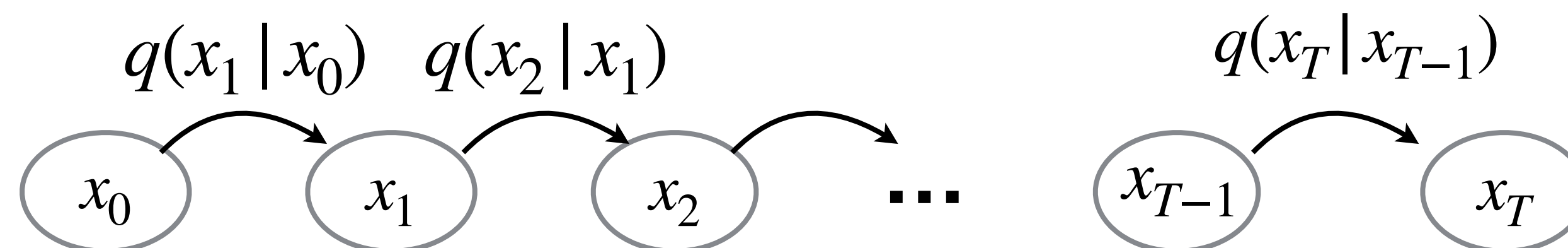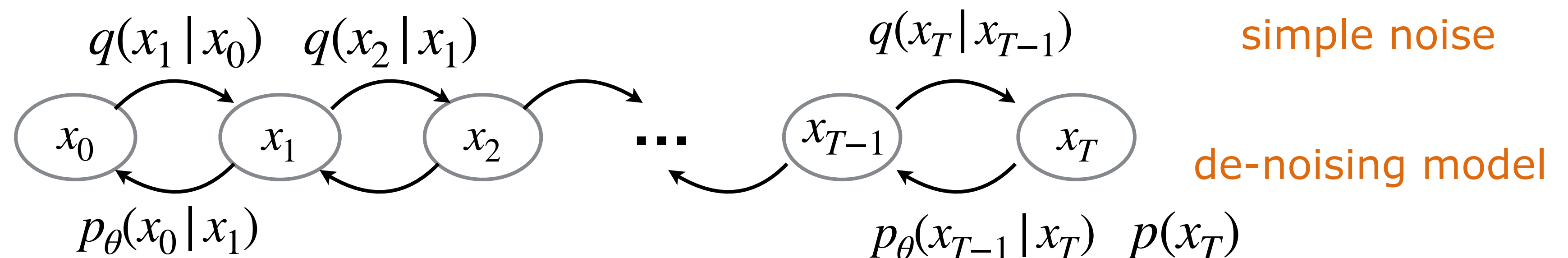‣ De-noising diffusion models over images (e.g., Ho et al., Song et al.)
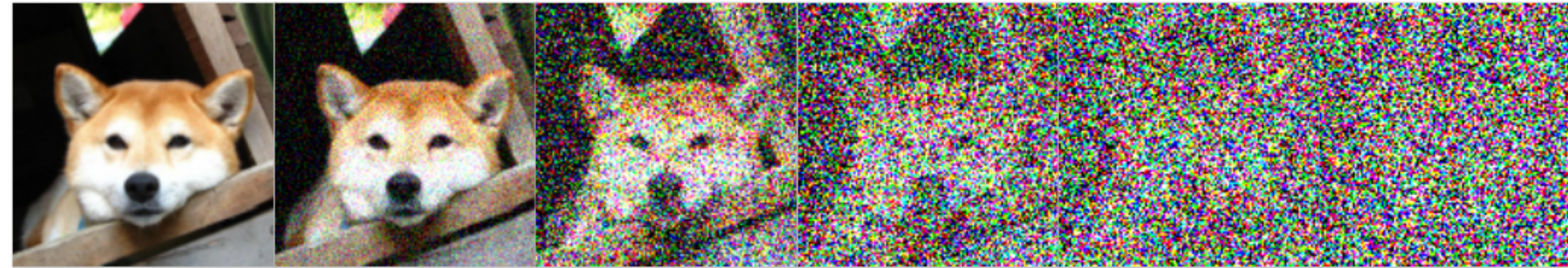


[image from
Rissanen et al 2022]

$$x_0 \longrightarrow z \sim N(0,I)$$

‣ A <u>forward process</u> "simplifies" objects (images) by adding more and more noise (each noise addition removes high frequency features from the image)

$$q(x_1|x_0) \quad q(x_2|x_1) \qquad\qquad q(x_T|x_{T-1})$$

simple noise

$$x_0 \longrightarrow x_1 \longrightarrow x_2 \quad \dots \quad x_{T-1} \longrightarrow x_T$$

# **Mechanisms of generation: diffusion**

‣ De-noising diffusion models over images (e.g., Ho et al., Song et al.)



[image from Rissanen et al 2022]

$$x_0 \longleftarrow z \sim N(0,I)$$

‣ A <u>forward process</u> "simplifies" objects (images) by adding more and more noise (each noise addition removes high frequency features from the image)

‣ The <u>reverse process</u> tries to reverse each step, i.e., it learns to predict how to de-noise images back into high quality versions
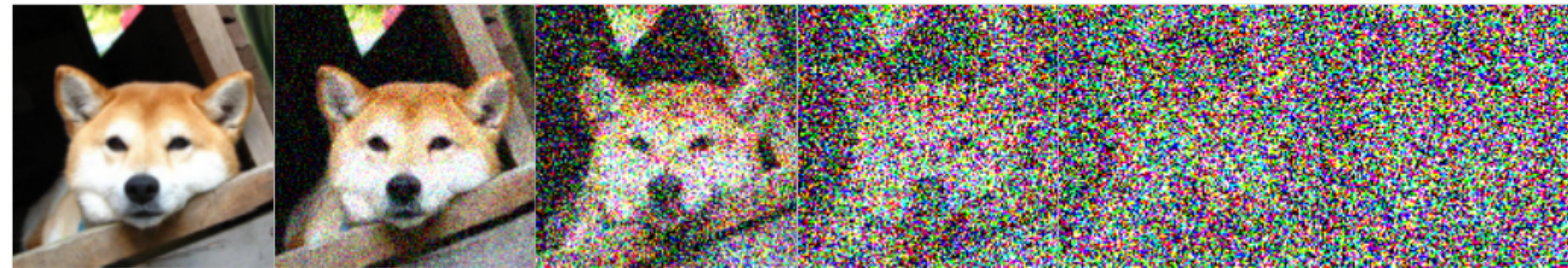


$q(x_1|x_0)$  $q(x_2|x_1)$  $q(x_T|x_{T-1})$  simple noise

$x_0$  $x_1$  $x_2$  $\cdots$  $x_{T-1}$  $x_T$  de-noising model

$p_\theta(x_0|x_1)$  $p_\theta(x_{T-1}|x_T)$  $p(x_T)$

# Basic diffusion models



$q(x_{1:T}|x_0) =$   $q(x_1|x_0)$   $q(x_2|x_1)$   $q(x_T|x_{T-1})$

$$x_0 \quad x_1 \quad x_2 \quad \ldots \quad x_{T-1} \quad x_T$$

$p_\theta(x_{0:T}) =$   $p_\theta(x_0|x_1)$   $p_\theta(x_{T-1}|x_T)$   $p(x_T)$

- We have a *fixed* (not learned) "forward process" $q(x_t|x_{t-1})$ which adds Gaussian noise (+ shrinkage)

- We use a reverse de-noising process $p_\theta(x_{t-1}|x_t)$ to (statistically) invert those steps

- This is a latent variable model where z = latent trajectories $(x_1, \ldots, x_T)$
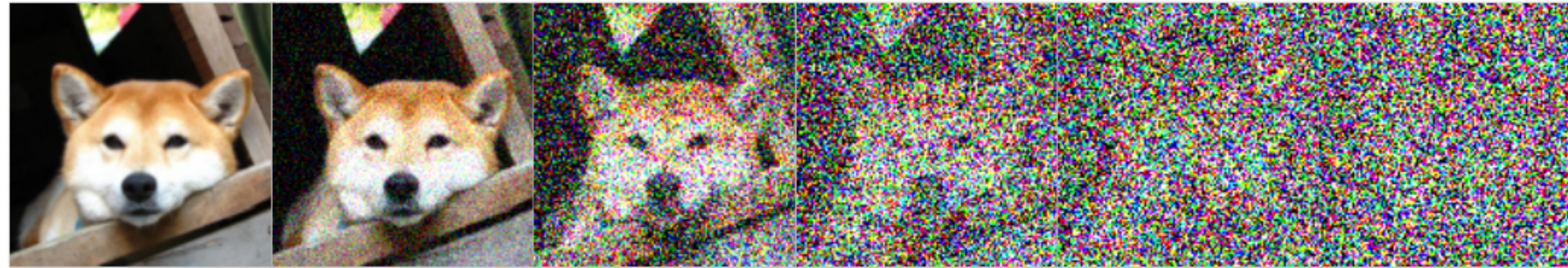
# Basic diffusion models



$$q(x_{1:T}|x_0) =$$
$$q(x_1|x_0) \quad q(x_2|x_1) \qquad\qquad q(x_T|x_{T-1})$$

$x_0$   $x_1$   $x_2$   ...   $x_{T-1}$   $x_T$

$$p_\theta(x_{0:T}) =$$
$$p_\theta(x_0|x_1) \qquad\qquad p_\theta(x_{T-1}|x_T) \quad p(x_T)$$

‣ We have a *fixed* (not learned) "forward process" $q(x_t|x_{t-1})$ which adds Gaussian noise (+ shrinkage)

‣ We use a reverse de-noising process $p_\theta(x_{t-1}|x_t)$ to (statistically) invert those steps

‣ This is a latent variable model where z = latent trajectories $(x_1, ..., x_T)$

‣ In principle, we can learn the reverse de-noising process by maximizing the ELBO criterion

$$\log p_\theta(x_0) \geq E_q \left[ \log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right]$$
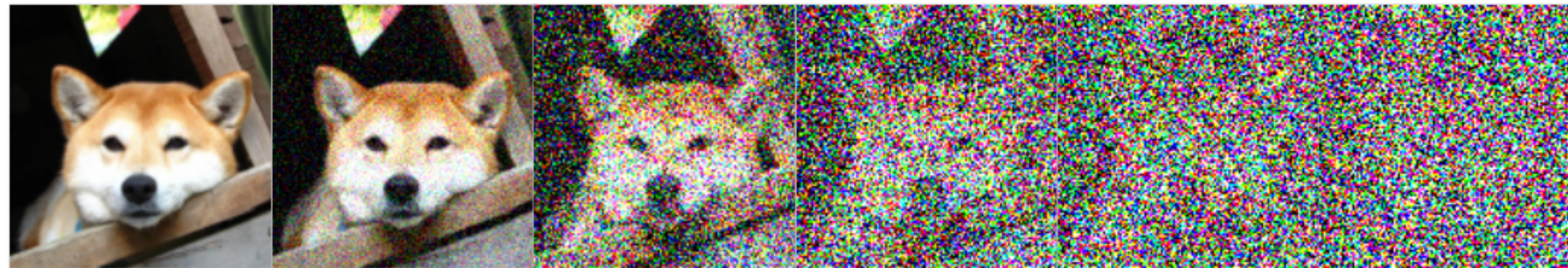
# **Forward process**



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

$$q(x_1|x_0) \quad q(x_2|x_1) \qquad\qquad q(x_T|x_{T-1})$$

$$x_0 \quad x_1 \quad x_2 \quad \cdots \quad x_{T-1} \quad x_T$$

$$q(x_{1:T}|x_0) = \prod_{t=1}^{T} N(x_t | \sqrt{1-\beta_t} x_{t-1}, \beta_t I)$$

‣ We start with $x_0$ and add a bit of Gaussian noise at each step (with shrinkage). So the distribution of $x_t$ at step t will have to be Gaussian as well

# Forward process: diffusion kernel



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

$$q(x_1 | x_0) \quad q(x_2 | x_1) \quad\quad\quad q(x_T | x_{T-1})$$

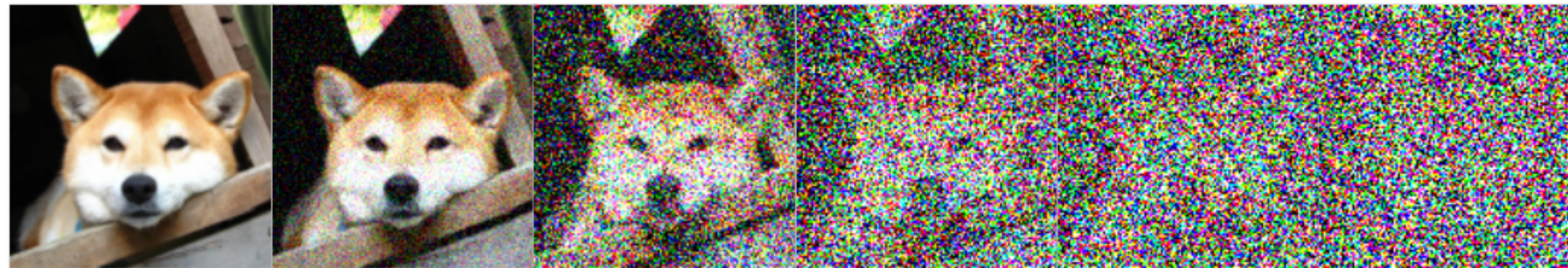$x_0 \quad x_1 \quad x_2 \quad \ldots \quad x_{T-1} \quad x_T$

$$q(x_{1:T} | x_0) = \prod_{t=1}^{T} N(x_t | \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

- We start with $x_0$ and add a bit of Gaussian noise at each step (with shrinkage). So the distribution of $x_t$ at step t will have to be Gaussian as well

$$q_t(x_t | x_0) = N\big(x_t | \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I\big) \quad \text{(diffusion kernel)} \quad\quad \bar{\alpha}_t = \prod_{s=1}^{t} (1 - \beta_s)$$

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

# Forward process: diffusion kernel



$$q(x_{1:T}|x_0) = \qquad \overset{q(x_1|x_0)}{\qquad} \overset{q(x_2|x_1)}{\qquad} \qquad \overset{q(x_T|x_{T-1})}{\qquad}$$
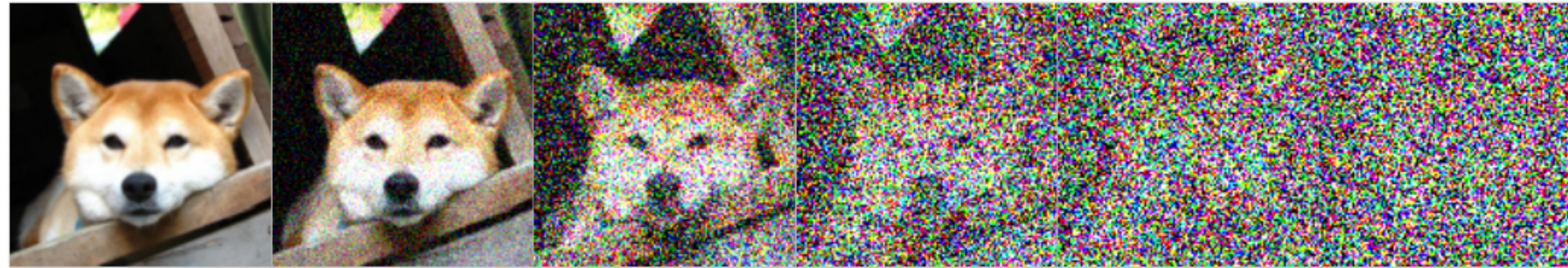
$$x_0 \to x_1 \to x_2 \cdots x_{T-1} \to x_T$$

$$q(x_{1:T}|x_0) = \prod_{t=1}^{T} N(x_t|\sqrt{1-\beta_t}x_{t-1}, \beta_t I)$$

‣ We start with $x_0$ and add a bit of Gaussian noise at each step (with shrinkage). So the distribution of $x_t$ at step t will have to be Gaussian as well

$$q_t(x_t|x_0) = N\big(x_t|\sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I\big) \qquad \text{(diffusion kernel)} \qquad \bar{\alpha}_t = \prod_{s=1}^{t}(1-\beta_s)$$

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, \ \ \epsilon \sim N(0,I) \qquad \text{For large T} \quad x_T \approx \epsilon, \ \ \epsilon \sim N(0,I)$$
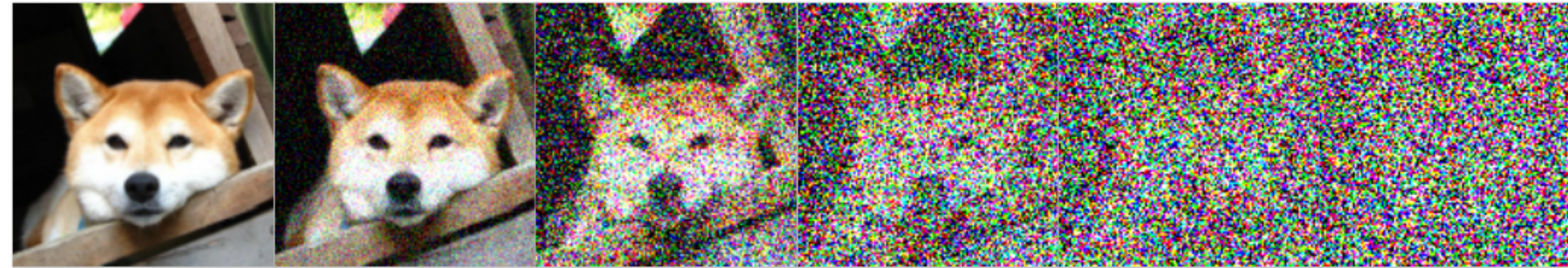
# Forward process: marginal at t



$$q(x_{1:T}|x_0) = \quad \overset{q(x_1|x_0)}{\curvearrowright} \quad \overset{q(x_2|x_1)}{\curvearrowright} \quad \cdots \quad \overset{q(x_T|x_{T-1})}{\curvearrowright}$$

$x_0 \rightarrow x_1 \rightarrow x_2 \quad \cdots \quad x_{T-1} \rightarrow x_T$

‣ The diffusion kernel tells us how $x_t$ is distributed conditioned on $x_0$

$$q_t(x_t|x_0) = N\big(x_t|\sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I\big)$$

‣ If we sample data $x_0 \sim q(x_0)$ (e.g., uniform at random), the noisy image $x_t$ at step t has a complex marginal distribution
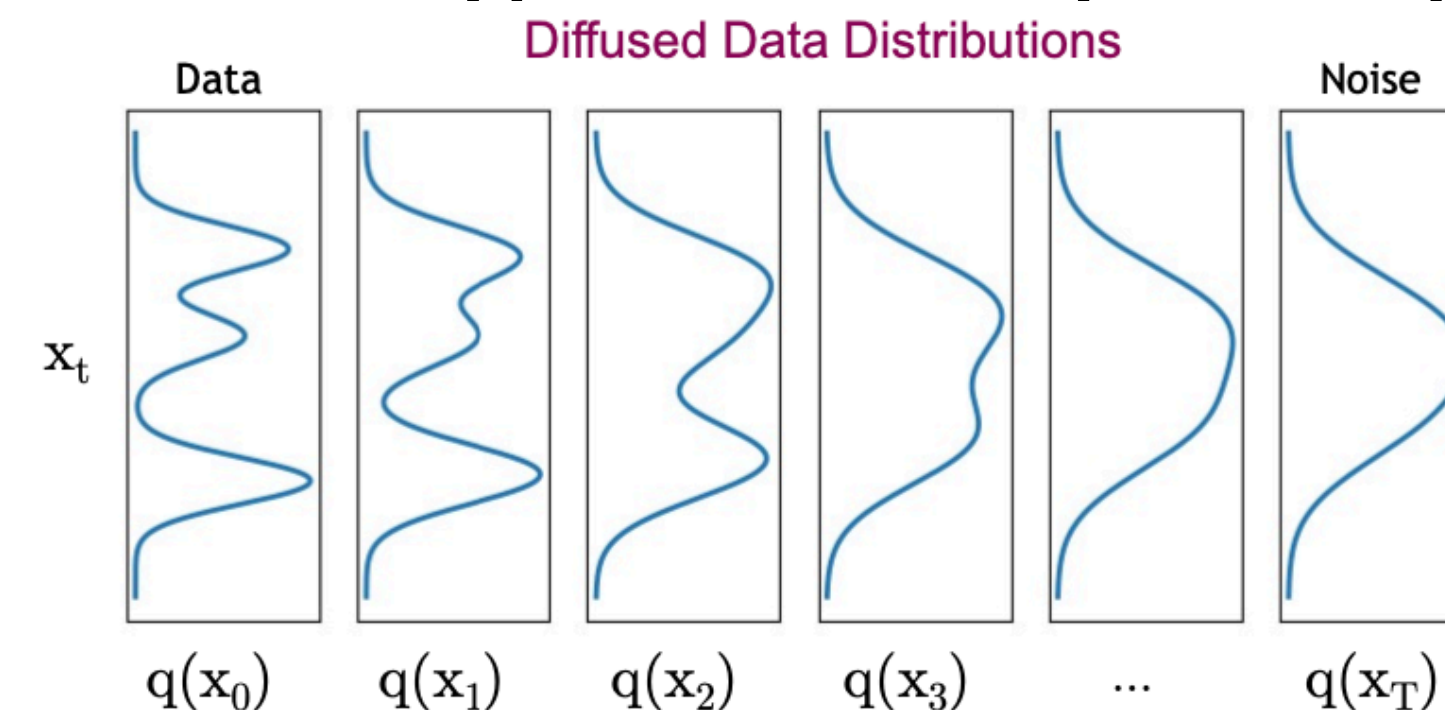
# Forward process: marginal at t



$$q(x_{1:T} \mid x_0) = \qquad \overset{q(x_1 \mid x_0)}{\longrightarrow} \quad \overset{q(x_2 \mid x_1)}{\longrightarrow} \qquad \qquad \overset{q(x_T \mid x_{T-1})}{\longrightarrow}$$

$x_0 \to x_1 \to x_2 \to \ldots \quad x_{T-1} \to x_T$

‣ The diffusion kernel tells us how $x_t$ is distributed conditioned on $x_0$

$$q_t(x_t \mid x_0) = N\big(x_t \mid \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I\big)$$
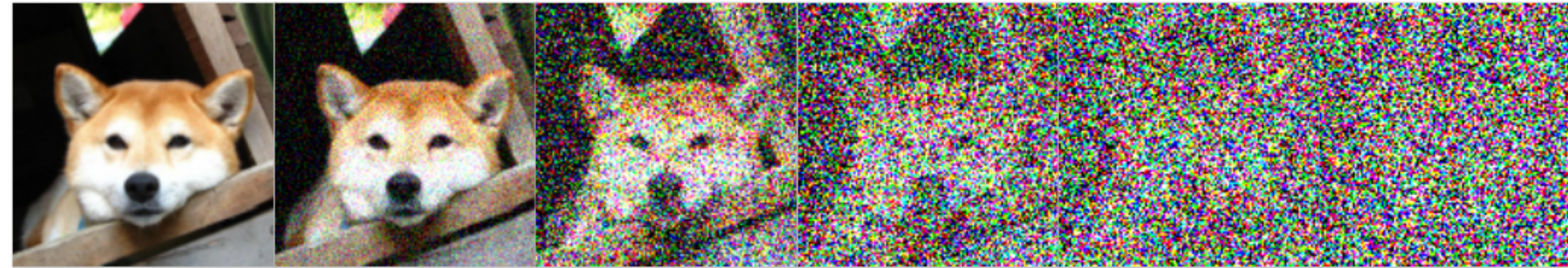
‣ If we sample data $x_0 \sim q(x_0)$ (e.g., uniform at random), the noisy image $x_t$ at step t has a complex marginal distribution

$$q_t(x_t) = \int q_t(x_t \mid x_0) q(x_0) dx_0 \approx \frac{1}{n} \sum_{i=1}^{n} q_t(x_t \mid x_0^i)$$



Diffused Data Distributions

Data                      Noise

$x_t$

$q(x_0) \quad q(x_1) \quad q(x_2) \quad q(x_3) \quad \ldots \quad q(x_T)$

[Vahdat et al 2022]

# Reverse process — our generative model



$$p_\theta(x_{0:T}) = \quad p_\theta(x_0 | x_1) \qquad\qquad p_\theta(x_{T-1} | x_T) \quad p(x_T)$$

‣ We use a reverse generative process to try to de-noise images back into clean versions, starting with $x_T$

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

# Reverse process — our generative model



$$p_\theta(x_{0:T}) = \qquad p_\theta(x_0 \mid x_1) \qquad\qquad p_\theta(x_{T-1} \mid x_T) \quad p(x_T)$$
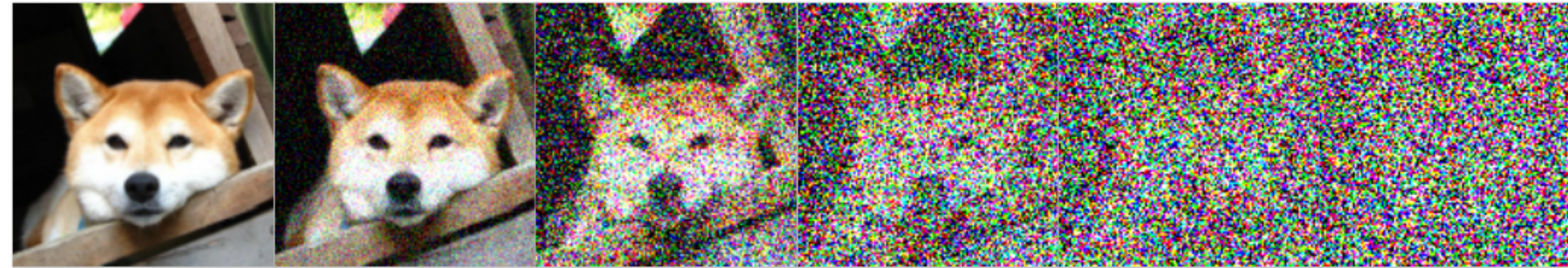
‣ We use a reverse generative process to try to de-noise images back into clean versions, starting with $x_T$

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^{T} N(x_{t-1} \mid \mu_\theta(x_t, t), \beta_t I)$$

‣ We can set $p(x_T) = N(x_T \mid 0, I)$ since the forward process always ends with $x_T \sim N(0, I)$

# Reverse process — our generative model



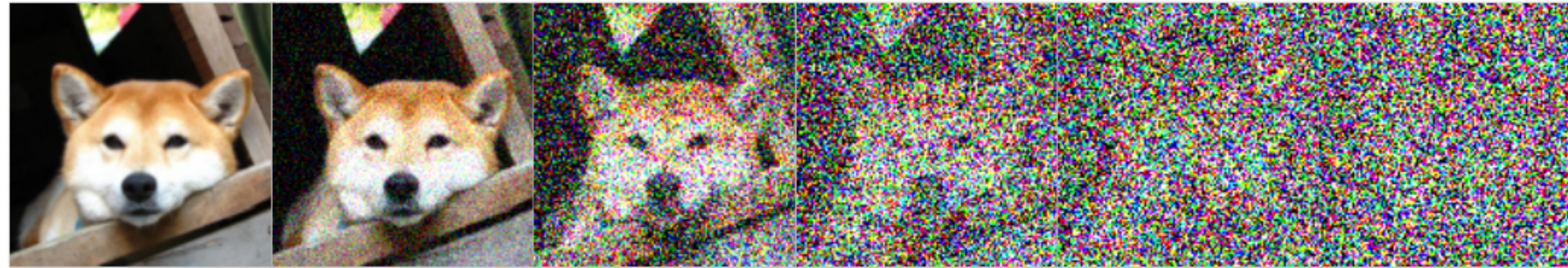$$p_\theta(x_{0:T}) = \quad p_\theta(x_0|x_1) \qquad\qquad p_\theta(x_{T-1}|x_T) \quad p(x_T)$$

‣ We use a reverse generative process to try to de-noise images back into clean versions, starting with $x_T$

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

‣ We can set $p(x_T) = N(x_T|0,I)$ since the forward process always ends with $x_T \sim N(0,I)$

‣ Learning the de-noising "vector field" $\mu_\theta(x_t, t)$ is the key part!

# Reparameterization of the reverse process



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

$$q(x_1 | x_0) \quad q(x_2 | x_1) \quad\quad\quad q(x_T | x_{T-1})$$

$$x_0 \quad x_1 \quad x_2 \quad \ldots \quad x_{T-1} \quad x_T$$

$$p_\theta(x_{0:T}) =$$

$$p_\theta(x_0 | x_1) \quad\quad\quad p_\theta(x_{T-1} | x_T) \quad p(x_T)$$

$$\alpha_t = (1 - \beta_t)$$

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^{T} N(x_{t-1} | \mu_\theta(x_t, t), \beta_t I) \qquad \bar{\alpha}_t = \prod_{s=1}^{t} (1 - \beta_s)$$

‣ Recall that based on the forward process: $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \ \epsilon \sim N(0, I)$

# Reparameterization of the reverse process



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

$$q(x_1|x_0) \quad q(x_2|x_1) \qquad\qquad q(x_T|x_{T-1})$$

$$x_0 \quad x_1 \quad x_2 \quad \cdots \quad x_{T-1} \quad x_T$$

$$p_\theta(x_{0:T}) = \qquad p_\theta(x_0|x_1) \qquad\qquad p_\theta(x_{T-1}|x_T) \quad p(x_T)$$

$$\alpha_t = (1 - \beta_t)$$

$$p_\theta(x_{0:T}) = p(x_T)\prod_{t=1}^{T} N(x_{t-1}|\mu_\theta(x_t, t), \beta_t I) \qquad\qquad \bar{\alpha}_t = \prod_{s=1}^{t}(1 - \beta_s)$$

- Recall that based on the forward process: $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \ \epsilon \sim N(0,I)$
- We can predict the added noise $\epsilon$ instead (at the same scale for all t).

# Reparameterization of the reverse process



$$q(x_{1:T} | x_0) = \qquad \overset{q(x_1 | x_0)}{\qquad} \overset{q(x_2 | x_1)}{\qquad} \qquad \overset{q(x_T | x_{T-1})}{\qquad}$$

$$\boxed{x_0} \ \boxed{x_1} \ \boxed{x_2} \ \cdots \ \boxed{x_{T-1}} \ \boxed{x_T}$$

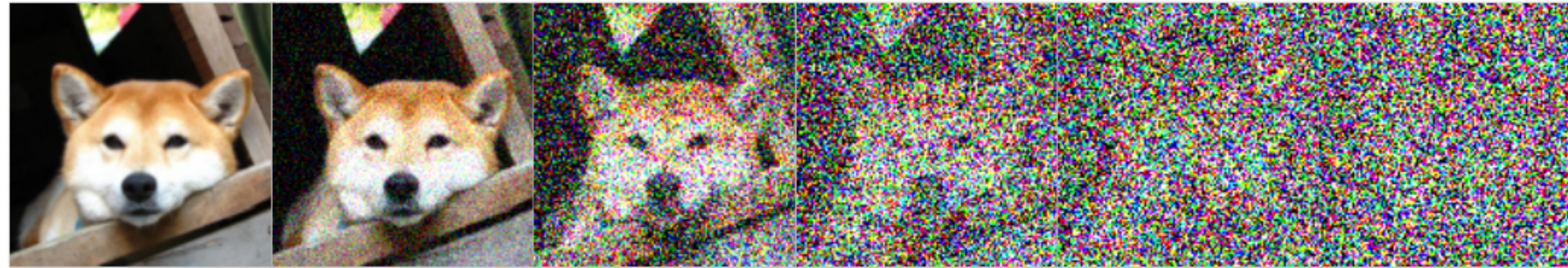$$p_\theta(x_{0:T}) = \qquad p_\theta(x_0 | x_1) \qquad\qquad p_\theta(x_{T-1} | x_T) \quad p(x_T)$$

$$\alpha_t = (1 - \beta_t)$$

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^{T} N(x_{t-1} | \mu_\theta(x_t, t), \beta_t I) \qquad \bar{\alpha}_t = \prod_{s=1}^{t} (1 - \beta_s)$$

‣ Recall that based on the forward process: $\quad x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \ \ \epsilon \sim N(0, I)$

‣ We can predict the added noise $\epsilon$ instead (at the same scale for all t). It allows us to calculate an estimate of $x_0$

$$\hat{x}_0 = \frac{1}{\sqrt{\bar{\alpha}}} (x_t - \sqrt{1 - \bar{\alpha}_t} \, \epsilon)$$

# Reparameterization of the reverse process



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

$$q(x_1 | x_0) \quad q(x_2 | x_1) \qquad\qquad q(x_T | x_{T-1})$$

$$x_0 \quad x_1 \quad x_2 \quad \cdots \quad x_{T-1} \quad x_T$$

$$p_\theta(x_{0:T}) = \qquad p_\theta(x_0 | x_1) \qquad\qquad p_\theta(x_{T-1} | x_T) \quad p(x_T)$$

$$\alpha_t = (1 - \beta_t)$$

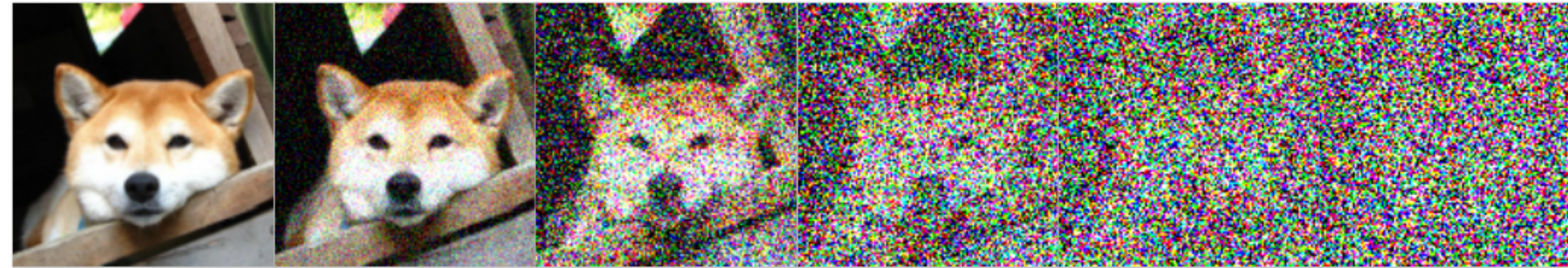$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^{T} N(x_{t-1} | \mu_\theta(x_t, t), \beta_t I) \qquad \bar{\alpha}_t = \prod_{s=1}^{t} (1 - \beta_s)$$

· Recall that based on the forward process: $\quad x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim N(0, I)$

· We can predict the added noise $\epsilon$ instead (at the same scale for all t). It allows us to calculate an estimate of $x_0$ and what the de-noising step should be

$$\hat{x}_0 = \frac{1}{\sqrt{\bar{\alpha}}} (x_t - \sqrt{1 - \bar{\alpha}_t} \, \epsilon) \qquad q(x_{t-1} | x_t, \hat{x}_0) = \frac{q(x_t | x_{t-1}) q(x_{t-1} | \hat{x}_0)}{q(x_t | \hat{x}_0)} \qquad \text{all the terms are Gaussian}$$

# Reparameterization of the reverse process



$$q(x_{1:T} | x_0) = \quad q(x_1 | x_0) \quad q(x_2 | x_1) \qquad\qquad q(x_T | x_{T-1})$$

$$x_0 \quad x_1 \quad x_2 \quad \ldots \quad x_{T-1} \quad x_T$$

$$p_\theta(x_{0:T}) = \quad p_\theta(x_0 | x_1) \qquad\qquad p_\theta(x_{T-1} | x_T) \quad p(x_T)$$

$$\alpha_t = (1 - \beta_t)$$

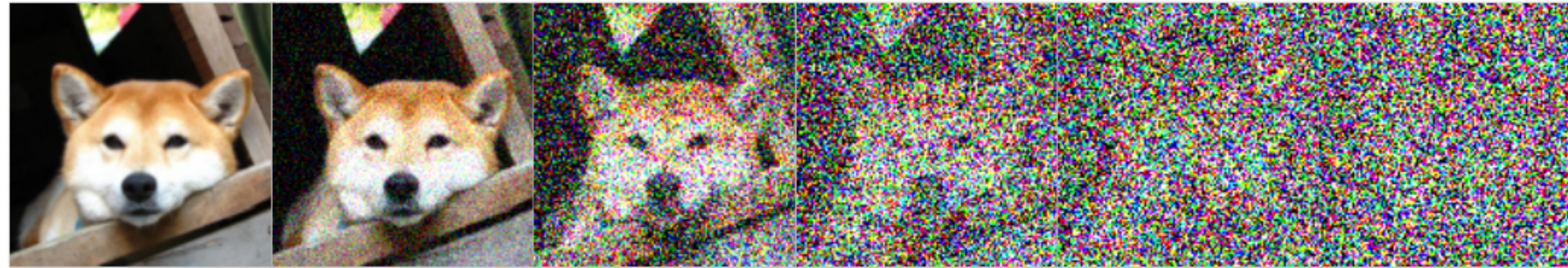$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^{T} N(x_{t-1} | \mu_\theta(x_t, t), \beta_t I) \qquad\qquad \bar{\alpha}_t = \prod_{s=1}^{t} (1 - \beta_s)$$

‣ Recall that based on the forward process: $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \ \epsilon \sim N(0, I)$

‣ We can predict the added noise $\epsilon$ instead (at the same scale for all t). It allows us to calculate an estimate of $x_0$ and what the de-noising step should be

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

# Basic training, sampling



$$q(x_{1:T} | x_0) = \quad q(x_1 | x_0) \quad q(x_2 | x_1) \qquad\qquad q(x_T | x_{T-1})$$

$$p_\theta(x_{0:T}) = \quad p_\theta(x_0 | x_1) \qquad\qquad p_\theta(x_{T-1} | x_T) \quad p(x_T)$$

---

**Algorithm 1** Training

---

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

---

# Basic training, sampling



$$q(x_{1:T} \mid x_0) = \quad q(x_1 \mid x_0) \quad q(x_2 \mid x_1) \qquad\qquad q(x_T \mid x_{T-1})$$

$$\boxed{x_0} \quad \boxed{x_1} \quad \boxed{x_2} \quad \cdots \quad \boxed{x_{T-1}} \quad \boxed{x_T}$$

$$p_\theta(x_{0:T}) = \quad p_\theta(x_0 \mid x_1) \qquad\qquad p_\theta(x_{T-1} \mid x_T) \quad p(x_T)$$

<span style="color:orange">strictly, ELBO would imply a different weighting</span>

---

**Algorithm 1** Training

---

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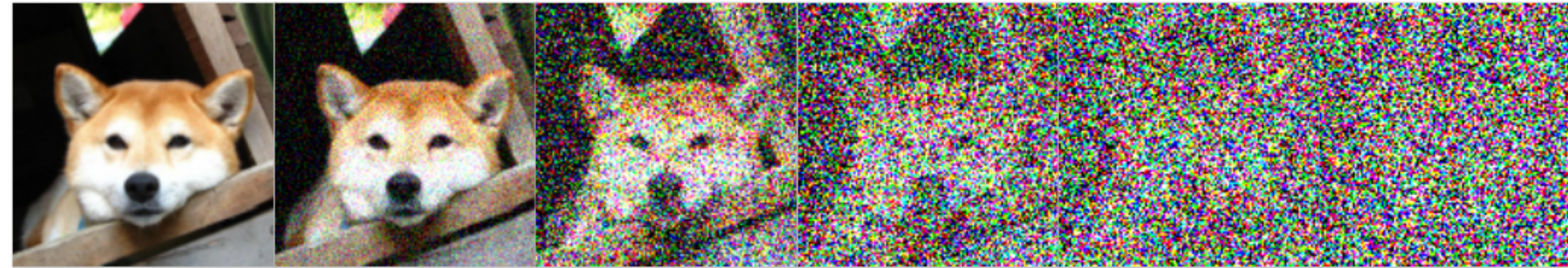6: **until** converged

---

# Basic training, sampling



$$q(x_{1:T}|x_0) = \quad q(x_1|x_0) \quad q(x_2|x_1) \qquad\qquad q(x_T|x_{T-1})$$

$$p_\theta(x_{0:T}) = \quad p_\theta(x_0|x_1) \qquad\qquad p_\theta(x_{T-1}|x_T) \quad p(x_T)$$

strictly, ELBO would
imply a different
weighting

---

**Algorithm 1** Training

1: **repeat**
2:   $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:   $t \sim \mathrm{Uniform}(\{1, \ldots, T\})$
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** $t = T, \ldots, 1$ **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** $\mathbf{x}_0$

# Score based (continuous) diffusion models



$$x_0 \qquad x_1 \qquad\qquad\qquad x_{T-1} \qquad x_T$$

‣ What if we used more steps in between [0,T], adding less noise at each?

# Score based (continuous) diffusion models



$x_0$      $x_1$          $x_{T-1}$    $x_T$

‣ What if we used more steps in between [0,T], adding less noise at each? We can write the noise variance now as $\beta(t)\Delta t$ (approaching zero)

$$x_t = \sqrt{1 - \beta(t)\Delta t}\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \;\; \epsilon \sim N(0,I)$$

# Score based (continuous) diffusion models



$x_0$      $x_1$      $x_{T-1}$    $x_T$

‣ What if we used more steps in between [0,T], adding less noise at each? We can write the noise variance now as $\beta(t)\Delta t$ (approaching zero)

$$x_t = \sqrt{1 - \beta(t)\Delta t}\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \;\; \epsilon \sim N(0,I)$$

$$x_t - x_{t-1} = (\sqrt{1 - \beta(t)\Delta t} - 1)\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \;\; \epsilon \sim N(0,I)$$

# Score based (continuous) diffusion models



$x_0$  $x_1$  $x_{T-1}$  $x_T$

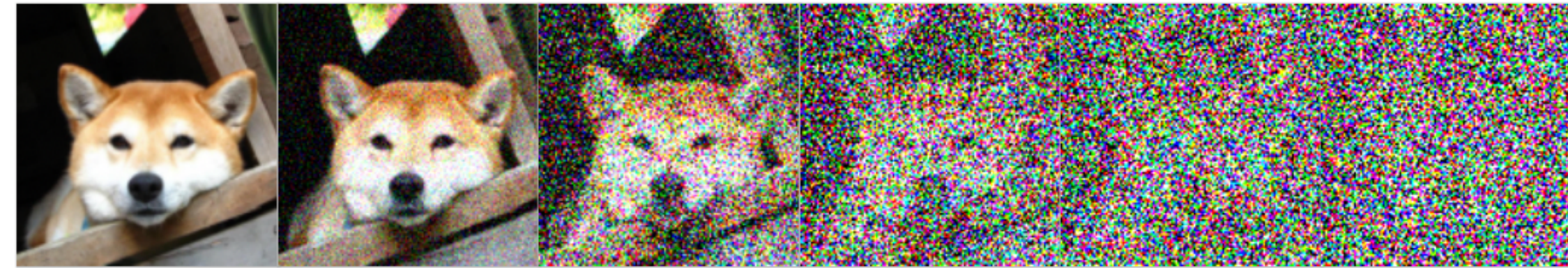‣ What if we used more steps in between [0,T], adding less noise at each? We can write the noise variance now as $\beta(t)\Delta t$ (approaching zero)

$$x_t = \sqrt{1 - \beta(t)\Delta t}\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \; \epsilon \sim N(0,I)$$

$$x_t - x_{t-1} = (\sqrt{1 - \beta(t)\Delta t} - 1)\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \; \epsilon \sim N(0,I)$$

$$\Delta x_t \approx -\frac{1}{2}\beta(t)\Delta t\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \; \epsilon \sim N(0,I)$$

# Score based (continuous) diffusion models
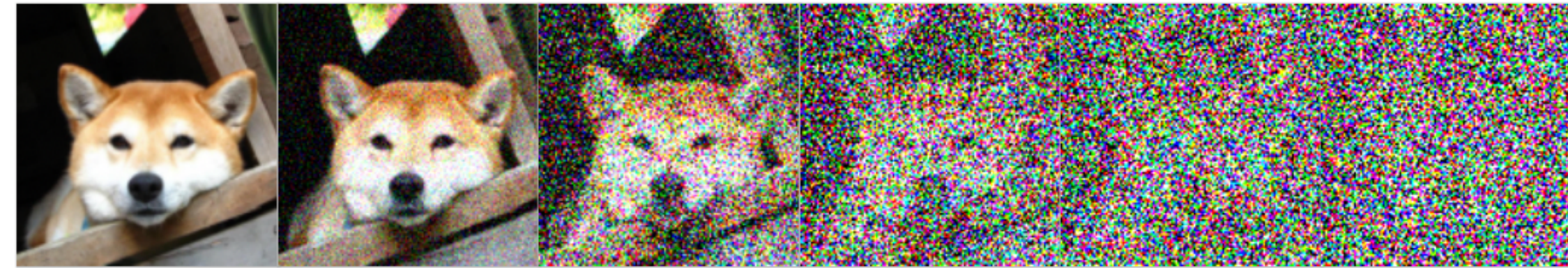


$x_0 \qquad x_1 \qquad\qquad\qquad x_{T-1} \qquad x_T$

- What if we used more steps in between [0,T], adding less noise at each? We can write the noise variance now as $\beta(t)\Delta t$ (approaching zero)

$$x_t = \sqrt{1 - \beta(t)\Delta t}\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \ \ \epsilon \sim N(0,I)$$

$$x_t - x_{t-1} = (\sqrt{1 - \beta(t)\Delta t} - 1)\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \ \ \epsilon \sim N(0,I)$$

$$\Delta x_t \approx -\frac{1}{2}\beta(t)\Delta t\, x_{t-1} + \sqrt{\beta(t)\Delta t}\, \epsilon, \ \ \epsilon \sim N(0,I)$$

- Taking the limit of $\Delta t$ we arrive at a stochastic differential equation (SDE)

Brownian motion

$$dx_t = -\frac{1}{2}\beta(t)\, x_t\, dt + \sqrt{\beta(t)}\, dw_t$$

# Forward diffusion process



Forward diffusion process (fixed)

$q(\mathbf{x}_0)$          $q(\mathbf{x}_T)$

$\mathbf{x}_0$   $\cdots$   $\mathbf{x}_t$   $\cdots$   $\mathbf{x}_T$

[image from Vahdat et al 2022]

‣ Forward process:   $dx_t = -\dfrac{1}{2}\beta(t)\,x_t\,dt + \sqrt{\beta(t)}\,dw_t$

# Reverse diffusion process



Forward diffusion process (fixed)

$q(\mathbf{x}_0)$     $q(\mathbf{x}_T)$

$\mathbf{x}_0$   ...   $\mathbf{x}_t$   ...   $\mathbf{x}_T$

[image from Vahdat et al 2022]

‣ Forward process:   $dx_t = -\dfrac{1}{2}\beta(t)\,x_t\,dt + \sqrt{\beta(t)}\,dw_t$

‣ Key result [Anderson 82]: reverse process is also a diffusion process

also Brownian motion

$$dx_t = \left[-\frac{1}{2}\beta(t)x_t - \beta(t)\,\underline{\nabla_{x_t}\log q_t(x_t)}\right]dt + \sqrt{\beta(t)}\,d\tilde{w}_t$$

score function

(dt now negative)

# Reverse diffusion process

- We would like to use the reverse process to sample new images    $(x_T \sim N(0,I)$ as before)

$$dx_t = \left[ -\frac{1}{2}\beta(t)x_t - \beta(t) \underbrace{\nabla_{x_t} \log q_t(x_t)}_{\text{score function}} \right] dt + \sqrt{\beta(t)}\, d\tilde{w}_t$$

<span style="color:orange">score function</span>

- To do so we should learn a neural model $s_\theta(x, t)$ to approx. the score function
- How to learn $s_\theta(x, t)$?

# **Reverse diffusion process**

‣ We would like to use the reverse process to sample new images $(x_T \sim N(0,I)$ as before$)$

$$dx_t = \left[ -\frac{1}{2}\beta(t)x_t - \beta(t)\underline{\nabla_{x_t}\log q_t(x_t)} \right] dt + \sqrt{\beta(t)}\, d\tilde{w}_t$$

<span style="color:orange">score function</span>

‣ To do so we should learn a neural model $s_\theta(x,t)$ to approx. the score function

‣ How to learn $s_\theta(x,t)$? We could try to (similar to before)

$$x_0 \sim q(x_0), \;\; t \sim U(0,T)$$

$$x_t \sim q_t(x_t|x_0) \;\; \text{diffusion kernel (it is still just a Gaussian)}$$

$$\theta \leftarrow \theta - \eta\,\nabla_\theta \| s_\theta(x_t,t) - \underline{\nabla_{x_t}\log q_t(x_t)} \|^2$$

$$q_t(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \gamma_t\mathbf{x}_0, \sigma_t^2\mathbf{I})$$

$$\gamma_t = e^{-\frac{1}{2}\int_0^t \beta(s)ds}$$

$$\sigma_t^2 = 1 - e^{-\int_0^t \beta(s)ds}$$

[Vahdat et al 2022]

‣ But $q_t(x_t)$ is a complex distribution!

# Reverse diffusion process

- We would like to use the reverse process to sample new images $(x_T \sim N(0,I)$ as before$)$

$$dx_t = \left[ -\frac{1}{2}\beta(t)x_t - \beta(t)\underbrace{\nabla_{x_t}\log q_t(x_t)}_{\text{score function}} \right] dt + \sqrt{\beta(t)}\, d\tilde{w}_t$$

- To do so we should learn a neural model $s_\theta(x, t)$ to approx. the score function
- How to learn $s_\theta(x, t)$? We can do

$$x_0 \sim q(x_0), \;\; t \sim U(0,T)$$

$x_t \sim q_t(x_t | x_0)$  diffusion kernel (it is still just a Gaussian)

$$\theta \leftarrow \theta - \eta\, \nabla_\theta \left\| s_\theta(x_t, t) - \underline{\nabla_{x_t}\log q_t(x_t | x_0)} \right\|^2$$

$$q_t(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \gamma_t \mathbf{x}_0, \sigma_t^2 \mathbf{I})$$

$$\gamma_t = e^{-\frac{1}{2}\int_0^t \beta(s)ds}$$

$$\sigma_t^2 = 1 - e^{-\int_0^t \beta(s)ds}$$

[Vahdat et al 2022]

- $\nabla_{x_t}\log q_t(x_t | x_0)$ (score of a Gaussian) can be easily calculated and has the same expected value (over $x_0$ for a given $x_t$) when $x_0 \sim q(x_0)$ !

# Score expectation: derivation

▸ $\nabla_{x_t} \log q_t(x_t | x_0)$ is a score of a Gaussian that can be easily calculated

▸ It has the same expected value as the complex score function $\nabla_{x_t} \log q_t(x_t)$ we want and thus can be used as a (noisy) learning target

$$E_{x_0 \sim q(x_0|x_t)}\{ \nabla_{x_t} \log q_t(x_t | x_0) \} = \int q(x_0 | x_t) \nabla_{x_t} \log q_t(x_t | x_0) \, dx_0$$

$$= \int \frac{q(x_0)q_t(x_t | x_0)}{q_t(x_t)} \nabla_{x_t} \log q_t(x_t | x_0) \, dx_0$$

$$= \int \frac{q(x_0)}{q_t(x_t)} \nabla_{x_t} q_t(x_t | x_0) \, dx_0$$

$$= \frac{1}{q_t(x_t)} \nabla_{x_t} \int q_t(x_t | x_0) q(x_0) \, dx_0$$

$$= \frac{1}{q_t(x_t)} \nabla_{x_t} q_t(x_t) = \nabla_{x_t} \log q_t(x_t)$$

# Reverse diffusion process

‣ We would like to use the reverse process to sample new images  $(x_T \sim N(0,I)$ as before)

$$dx_t = \left[ -\frac{1}{2}\beta(t)x_t - \beta(t)\underline{\nabla_{x_t}\log q_t(x_t)} \right] dt + \sqrt{\beta(t)}\, d\tilde{w}_t$$

<span style="color:orange">score function</span>

‣ To do so we should learn a neural model $s_\theta(x,t)$ to approx. the score function

‣ How to learn $s_\theta(x,t)$? We could try to (similar to before)

$$x_0 \sim q(x_0), \; t \sim U(0,T)$$

$$x_t \sim q_t(x_t|x_0) \;\; \text{diffusion kernel (it is still just a Gaussian)}$$

$$\theta \leftarrow \theta - \eta \nabla_\theta \|s_\theta(x_t,t) - \underline{\nabla_{x_t}\log q_t(x_t)}\|^2$$

$$q_t(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \gamma_t\mathbf{x}_0, \sigma_t^2\mathbf{I})$$

$$\gamma_t = e^{-\frac{1}{2}\int_0^t \beta(s)ds}$$

$$\sigma_t^2 = 1 - e^{-\int_0^t \beta(s)ds}$$

[Vahdat et al 2022]

# Additional reading

- **Some early diffusion model papers:**
- Sohl-Dickstein et al., "Deep Unsupervised Learning using Nonequilibrium Thermodynamics", http://proceedings.mlr.press/v37/sohl-dickstein15.pdf
- Ho et al., "Denoising Diffusion Probabilistic Models", https://arxiv.org/abs/2006.11239
- Song et al., "Generative Modeling by Estimating Gradients of the Data Distribution", https://arxiv.org/abs/1907.05600
- **Tutorials (excerpts used in this lecture)**
- A. Vahdat, K. Kreis, R. Gao, "CVPR 2022 Tutorial — Denoising Diffusion-based Generative Modeling: Foundations and Applications", https://cvpr2022-tutorial-diffusion-models.github.io/