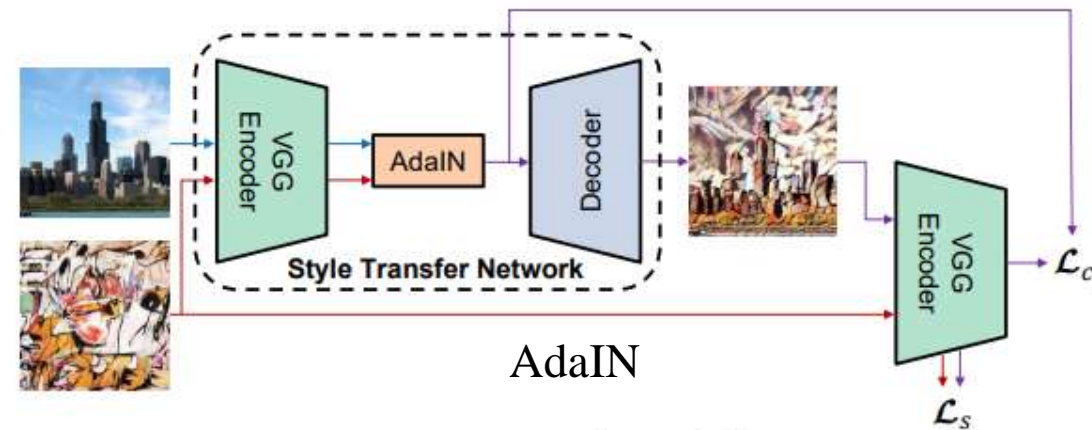
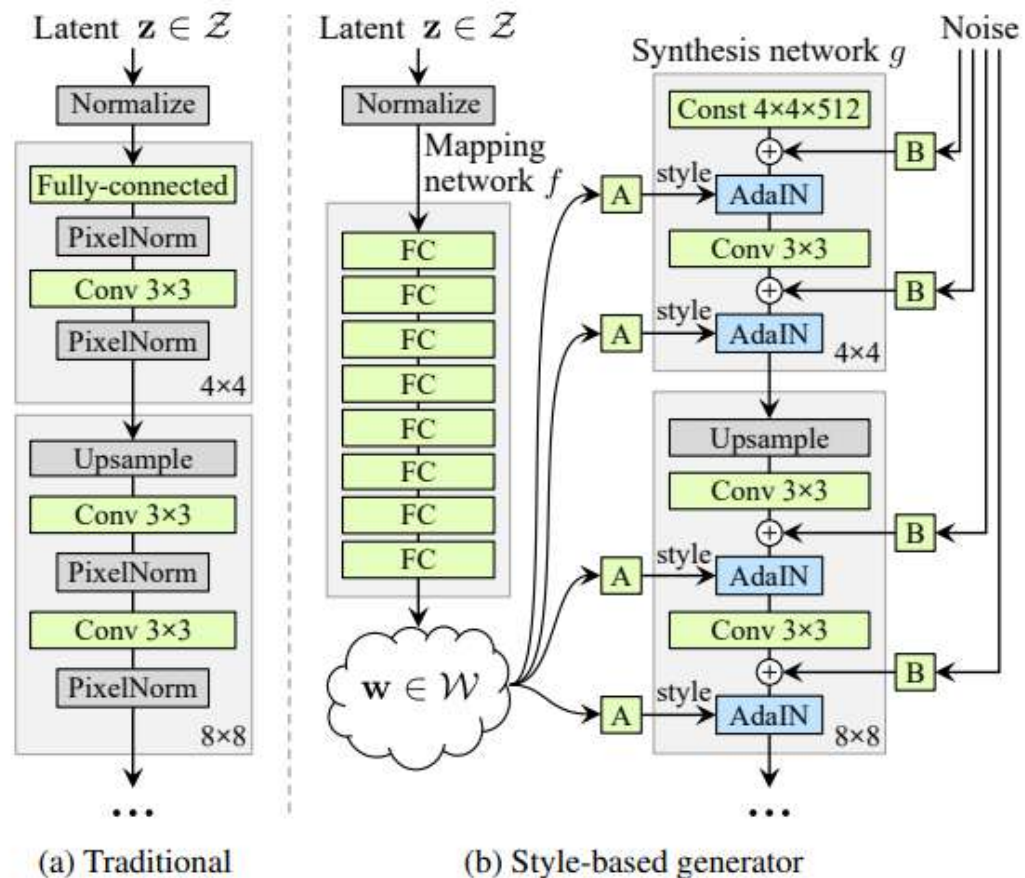


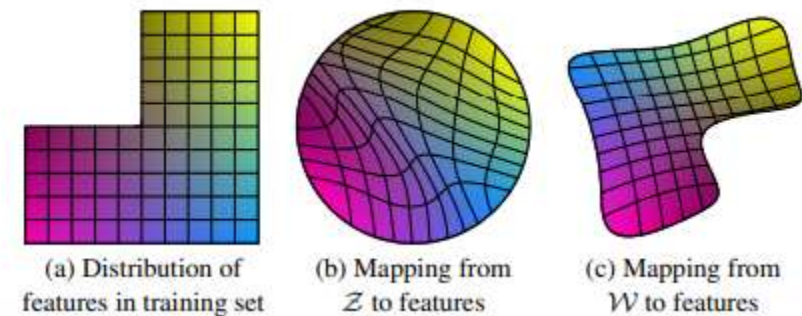
Recap.

STYLEGAN: A Style-Based Generator Architecture for Generative Adversarial Networks

- SoTA GAN before tokenizers.



$$\text{AdaIN}(x, y) = \sigma(y) \left(\frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$

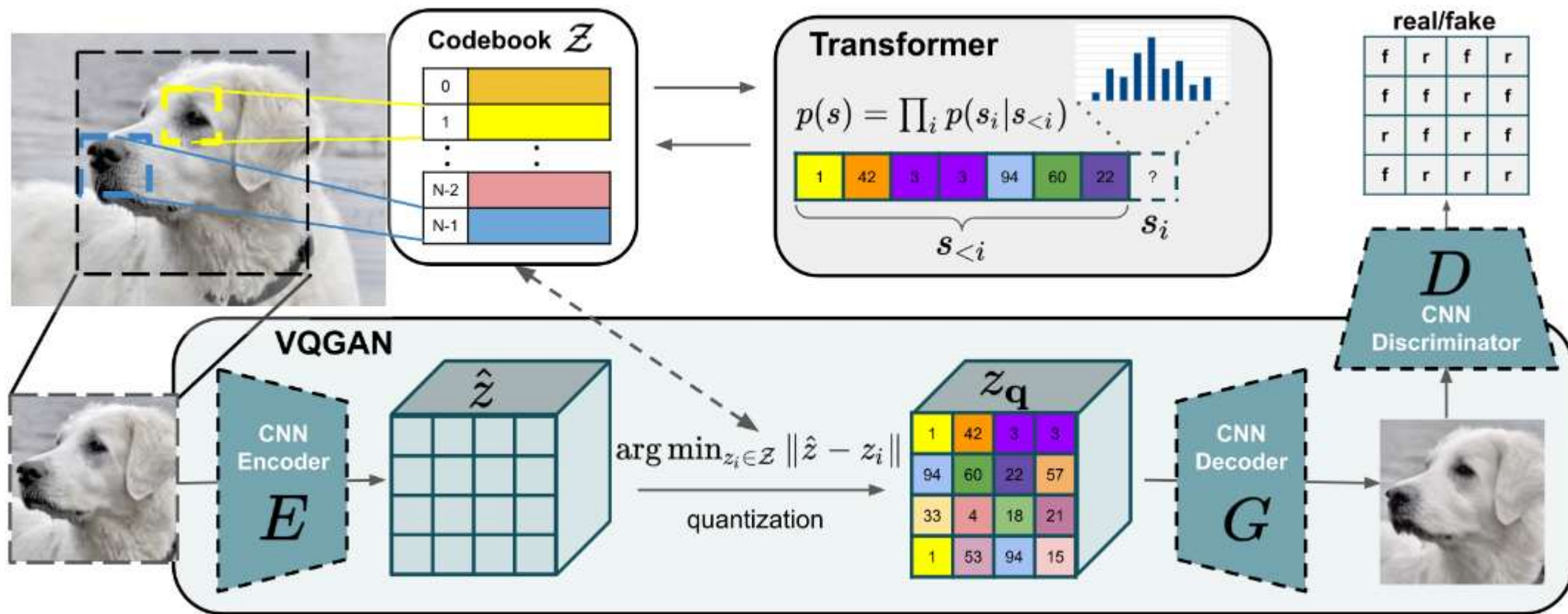


STYLEGANv3: Alias-Free Generative Adversarial Networks

- Excellent interpolation.



VQGAN



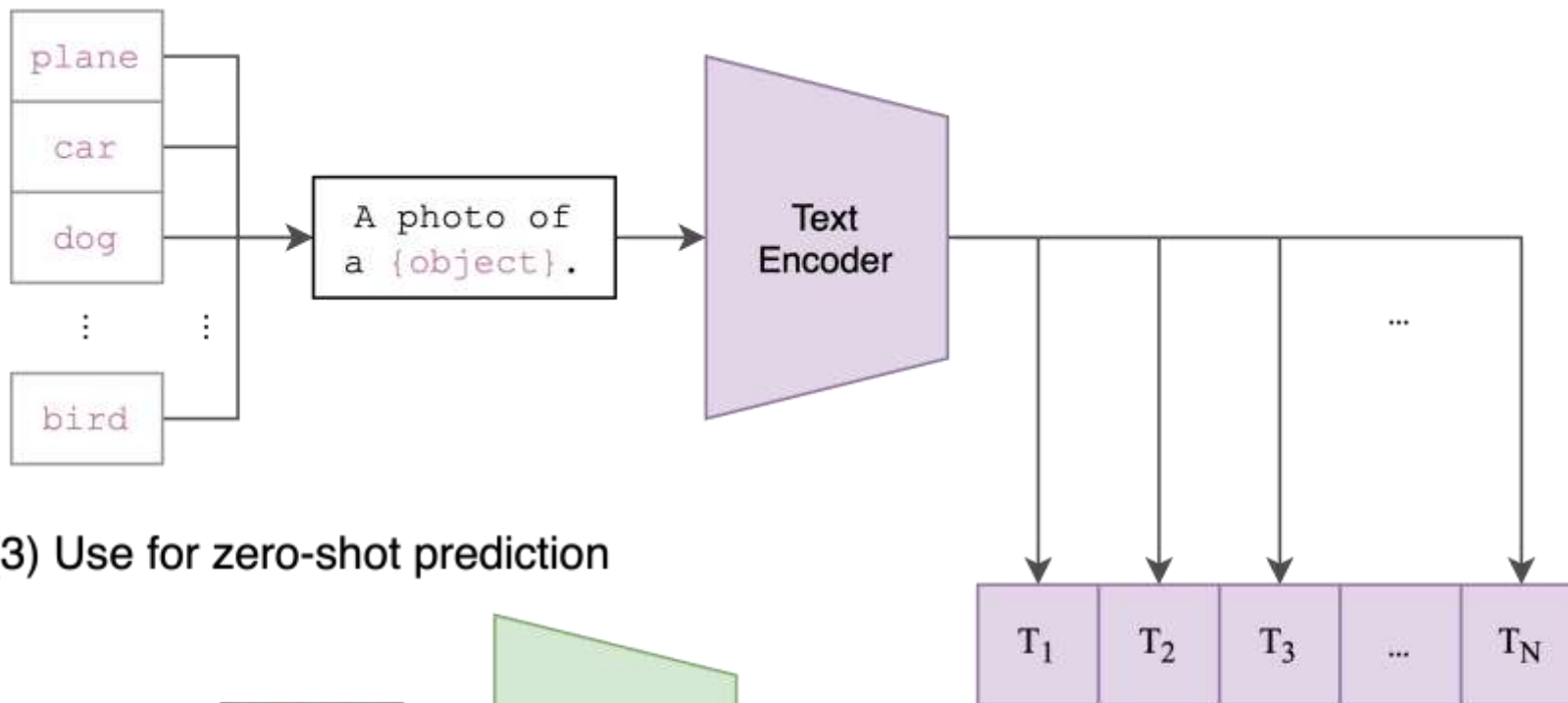
CLIP

[Radford et al., 2021]

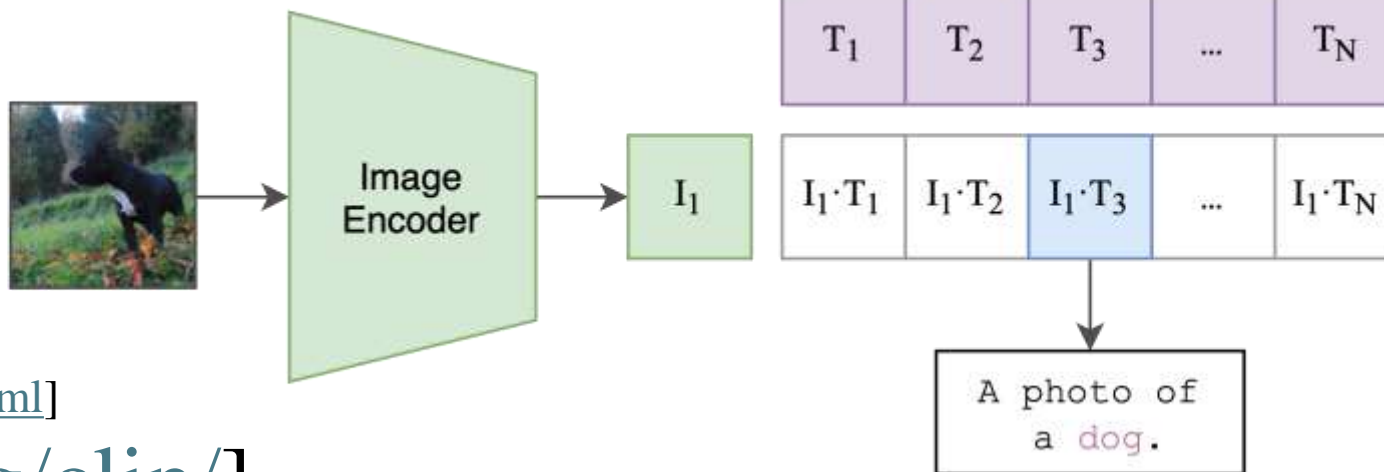
<https://arxiv.org/pdf/2103.00020.pdf>

2. Adaptor: Just ask

(2) Create dataset classifier from label text



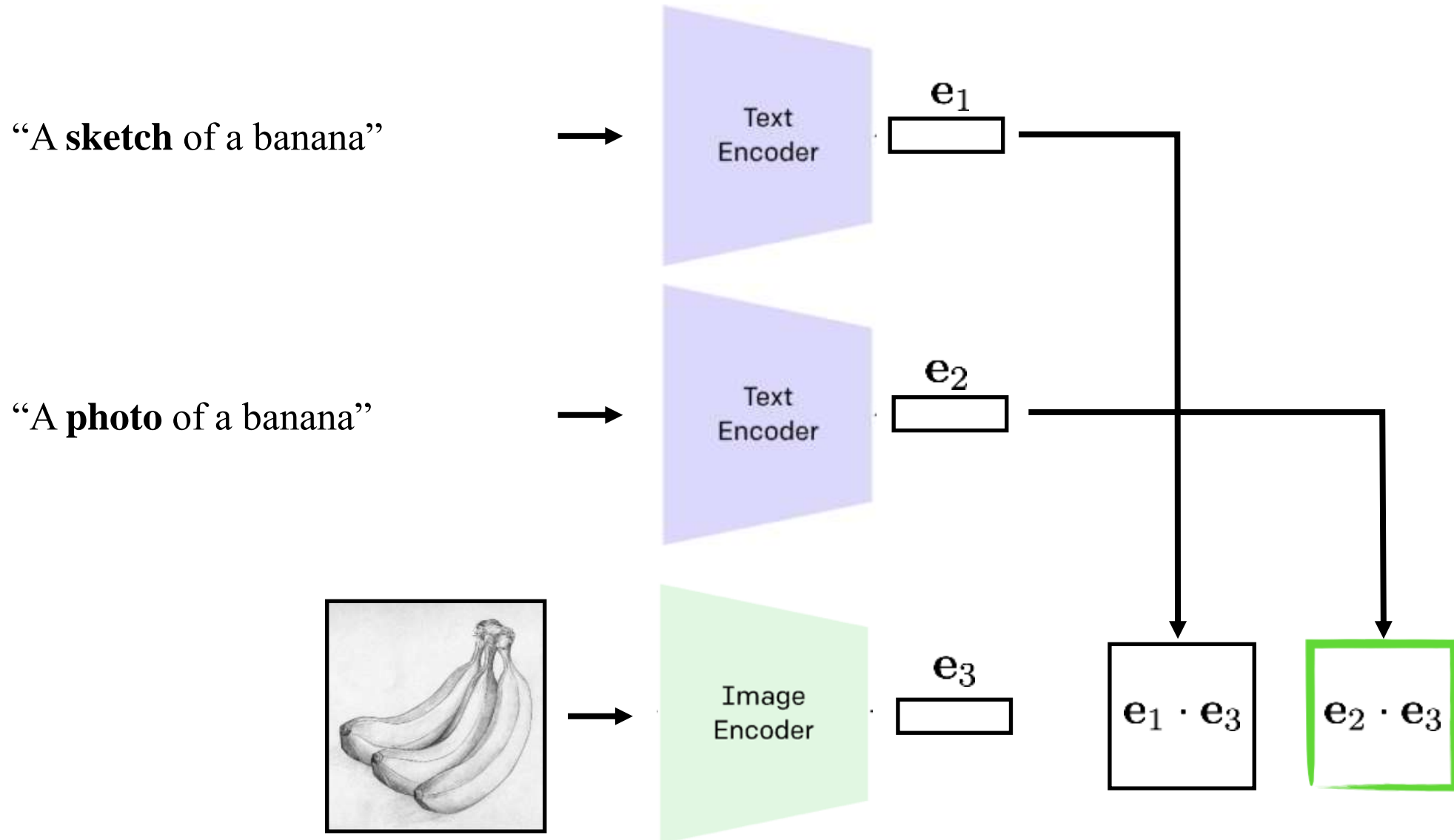
(3) Use for zero-shot prediction



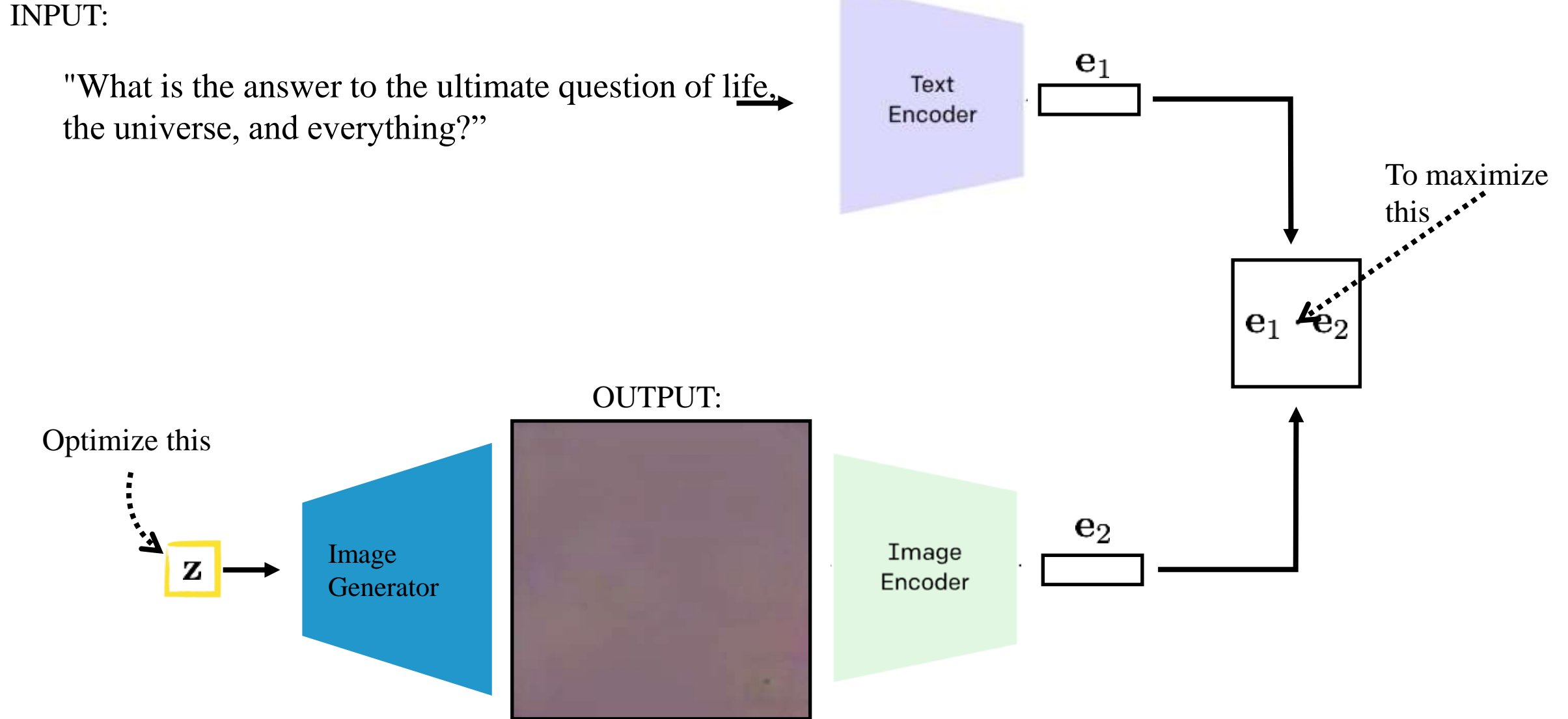
[\[https://evjang.com/2021/10/23/generalization.html\]](https://evjang.com/2021/10/23/generalization.html)

[\[https://openai.com/blog/clip/\]](https://openai.com/blog/clip/)

New capabilities by just asking



New capabilities by plugging pretrained models together: CLIP+GAN



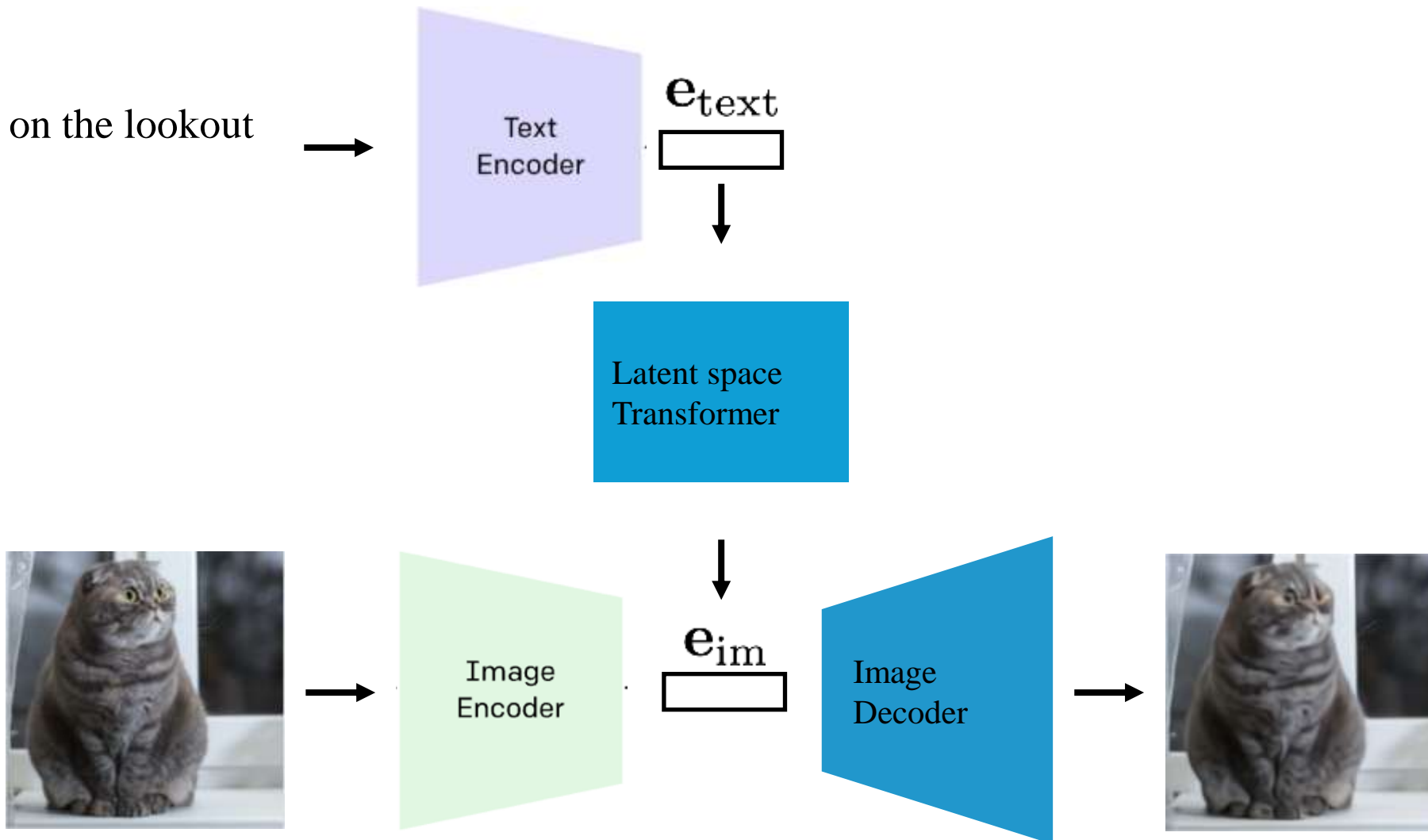
DALL-E [Ramesh et al. 2021]

<https://arxiv.org/pdf/2102.12092.pdf>

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

INPUT:

“A wide-eyed cat on the lookout for food”



StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery



“Emma Stone”

“Mohawk hairstyle”

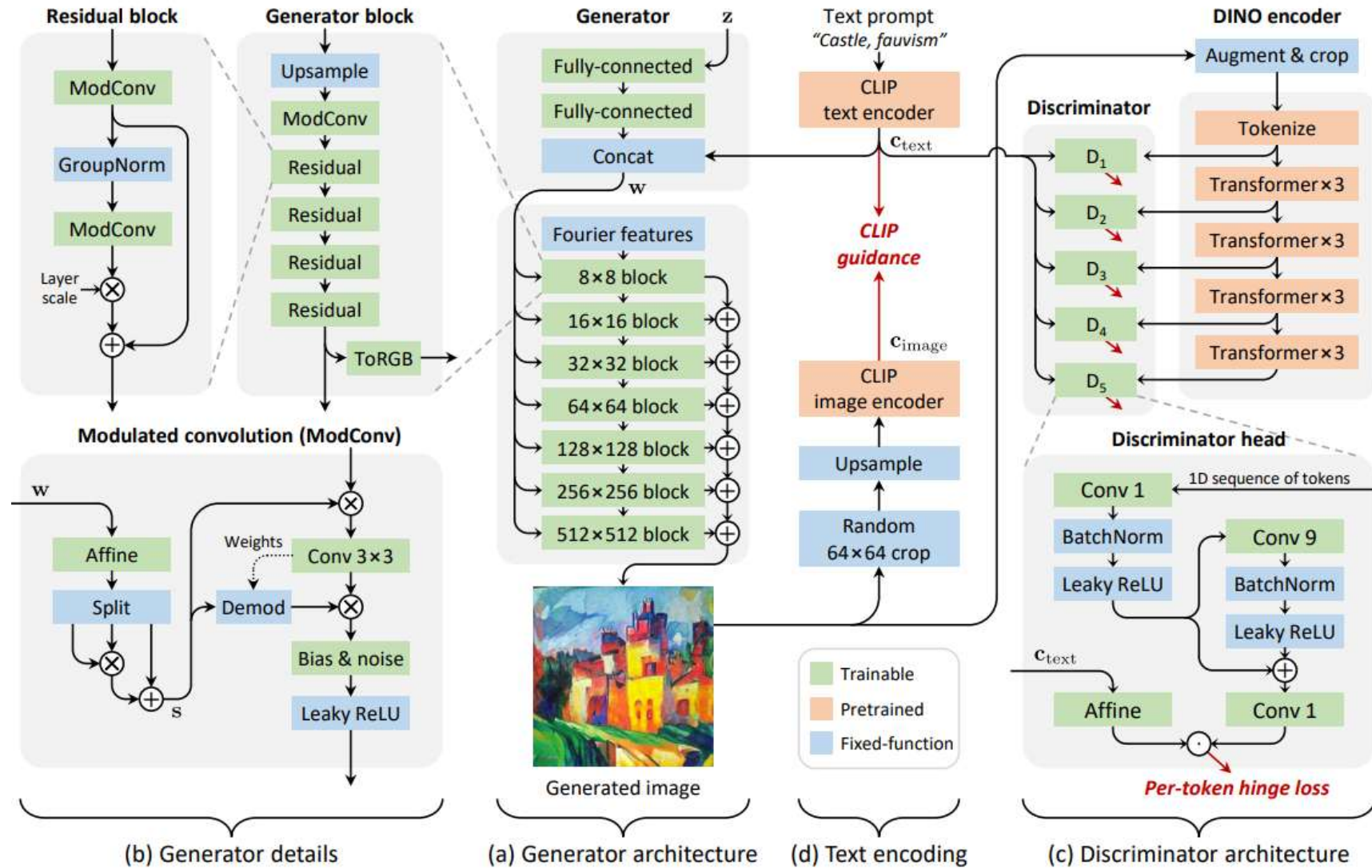
“Without makeup”

“Cute cat”

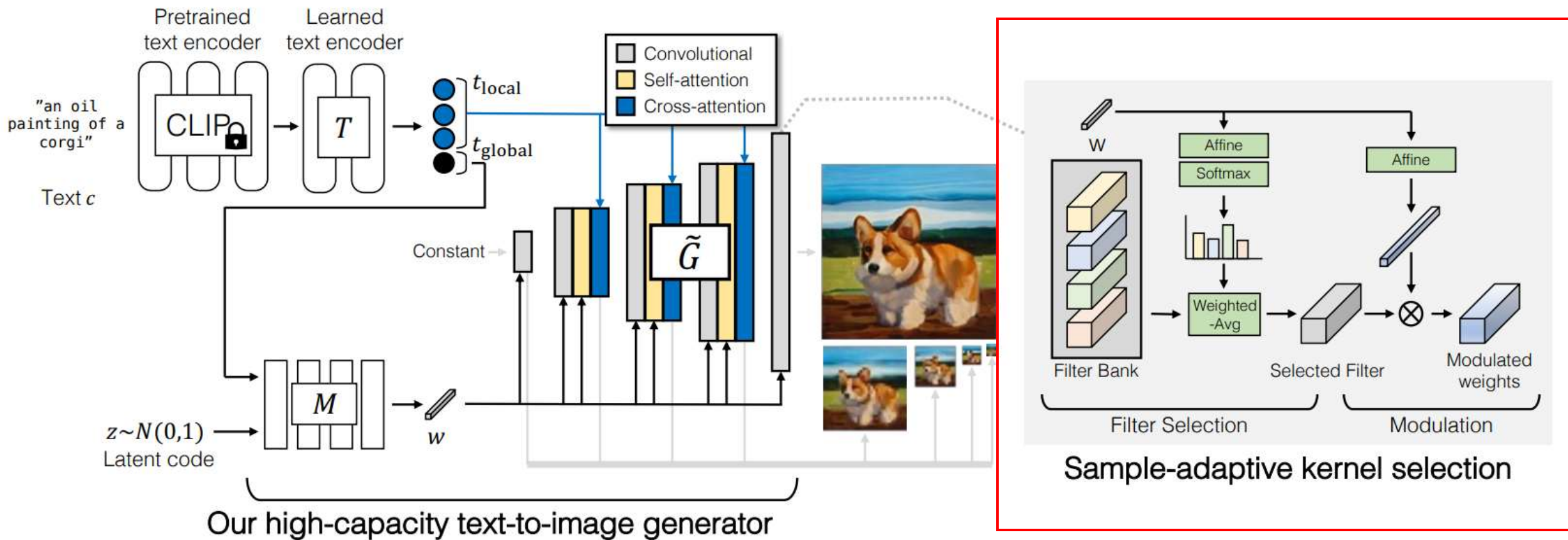
“Lion”

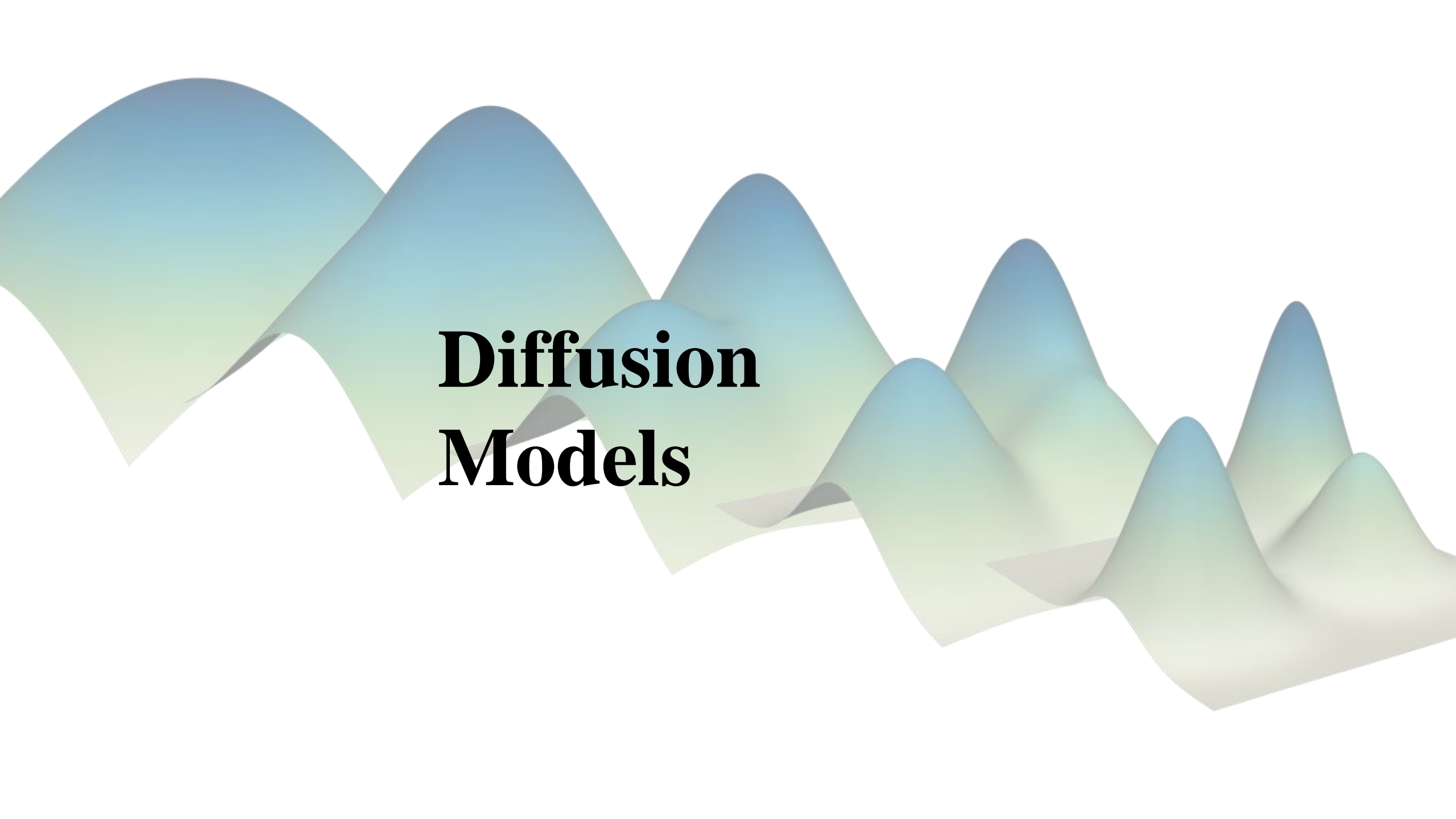
“Gothic church”

StyleGAN-T: Unlocking the Power of GANs for Fast Large-Scale Text-to-Image Synthesis



GigaGAN: Scaling up GANs for Text-to-Image Synthesis





**Diffusion
Models**

Overview

- Diffusion Models
- Energy-based Models and Score Matching

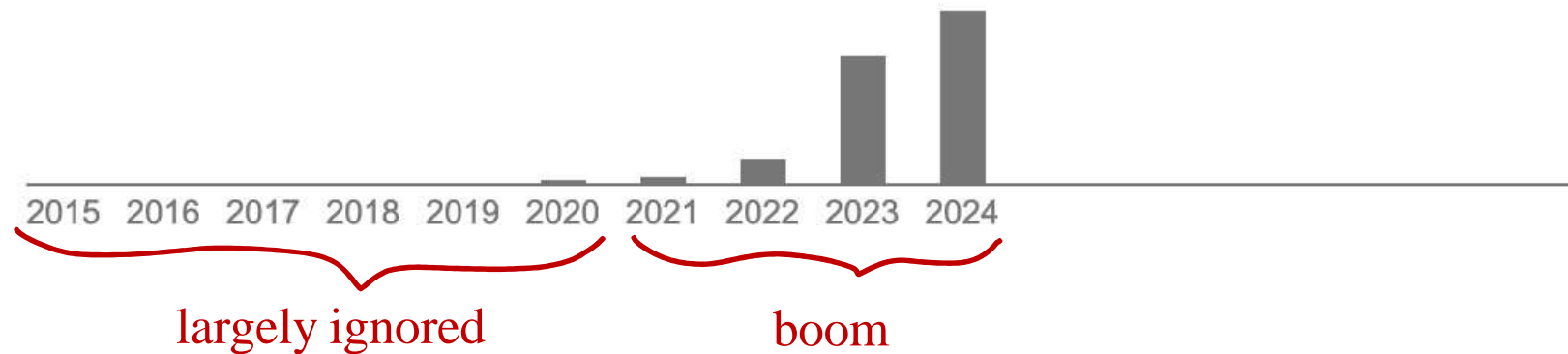
Deep unsupervised learning using nonequilibrium thermodynamics

Authors Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, Surya Ganguli

Publication date 2015/3/12

Journal International Conference on Machine Learning

Total citations Cited by 5630



Diffusion Models

Diffusion Models

- Forward process
 - add noise to data
- Reverse process
 - learn to denoise
- Training objective
 - from Hierarchical VAE to L2 loss
- Noise Conditional Network
 - represent distributions

... in a nutshell

noise

data



x_T

...

x_t

x_{t-1}

...

x_0



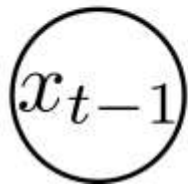
... in a nutshell

noise

data



...

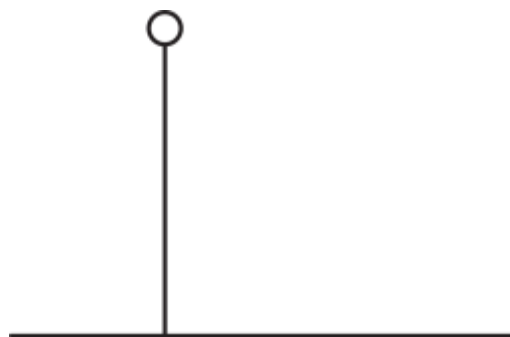


...

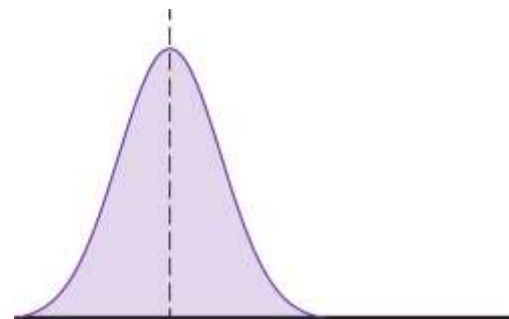


What is noise?

- Adding Gaussian noise \Leftrightarrow sampling $x \sim \mathcal{N}(x | x_0, \sigma)$



$$p(x) = \delta(x - x_0)$$

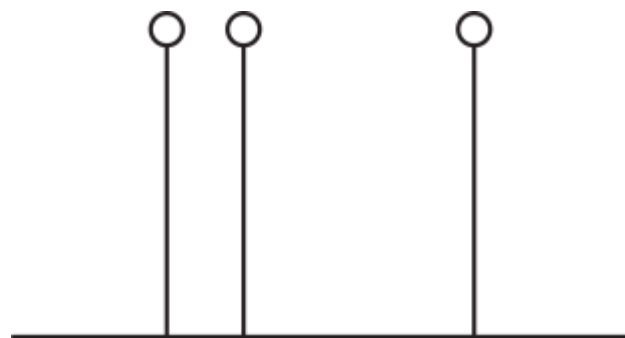


$$p(x) = \mathcal{N}(x | x_0, \sigma)$$

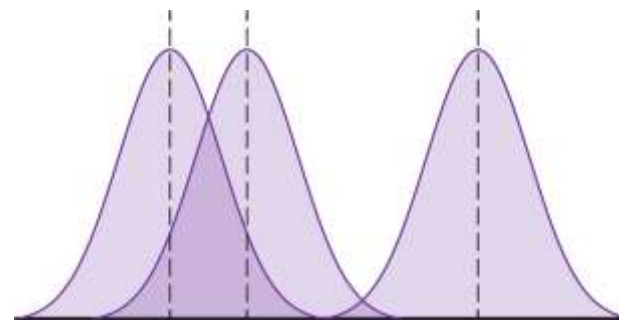


What is noise?

- Adding Gaussian noise \Leftrightarrow sampling $x \sim \mathcal{N}(x \mid x_0, \sigma)$



$p_{\text{data}}(x)$



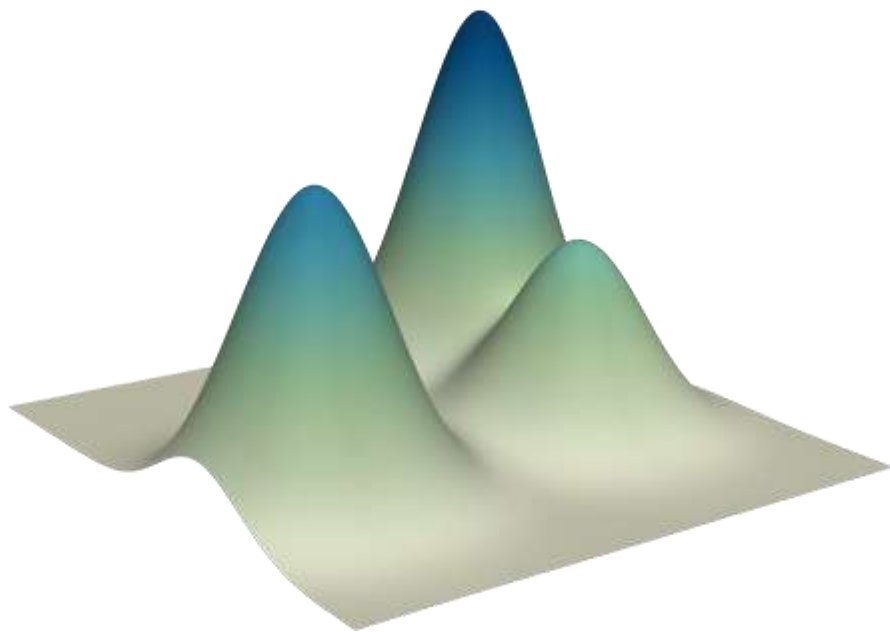
$p_{\text{data}}(x) * \mathcal{N}(x \mid 0, \sigma)$

convolution
(of pdf)

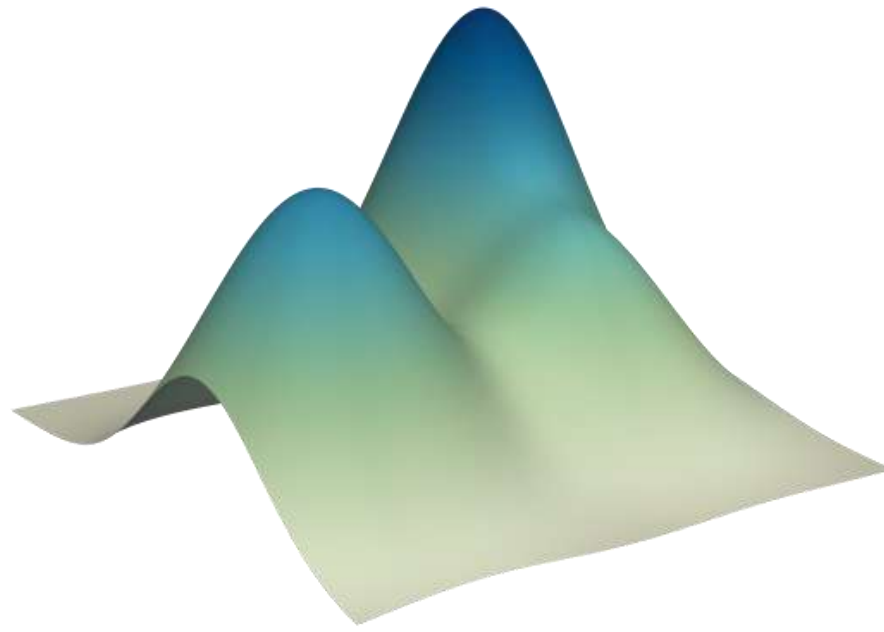
What is noise?

- Adding Gaussian noise \Leftrightarrow sampling

$$x \sim \mathcal{N}(x \mid x_0, \sigma)$$



$$p_{\text{data}}(x)$$

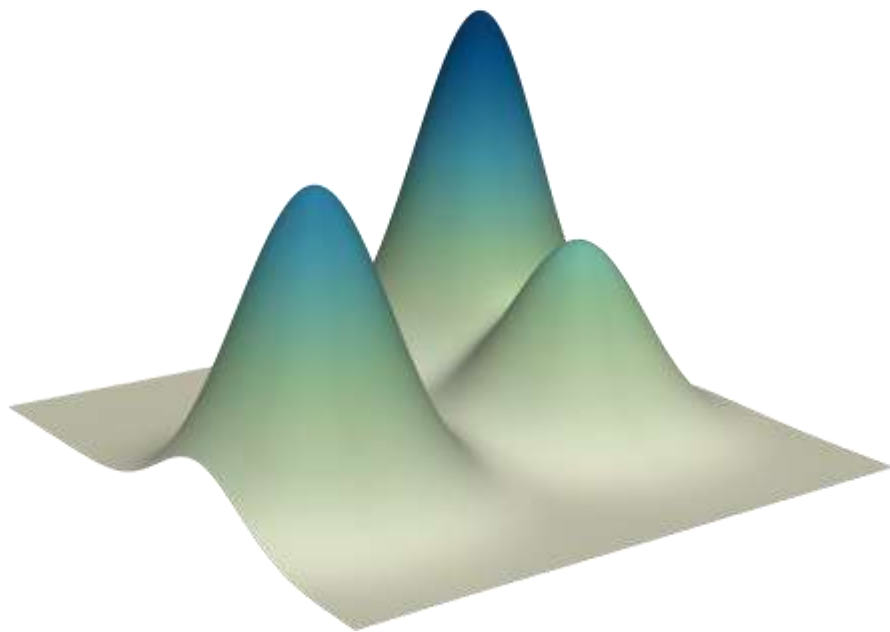


$$p_{\text{data}}(x) * \mathcal{N}(x \mid 0, \sigma)$$

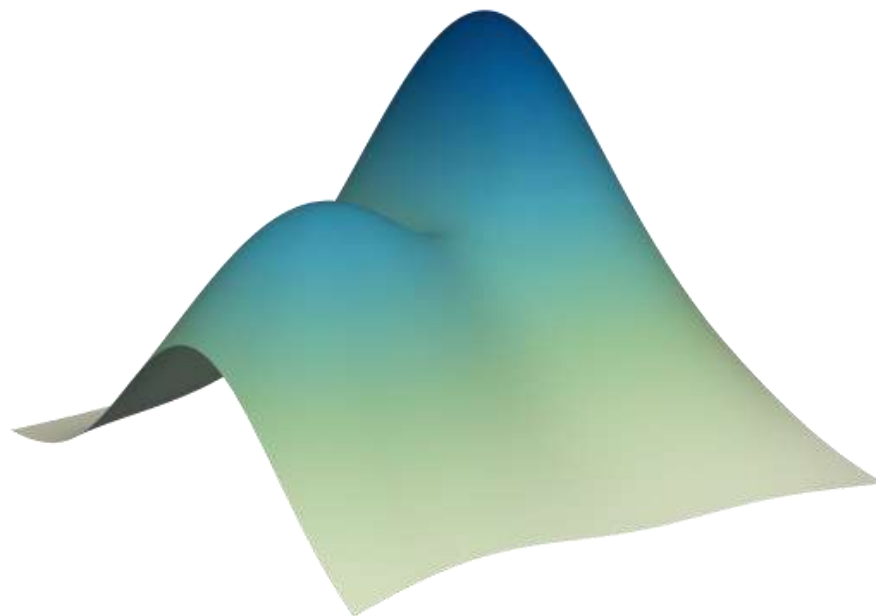
What is noise?

- Adding Gaussian noise \Leftrightarrow sampling

$$x \sim \mathcal{N}(x | x_0, \sigma)$$



$$p_{\text{data}}(x)$$

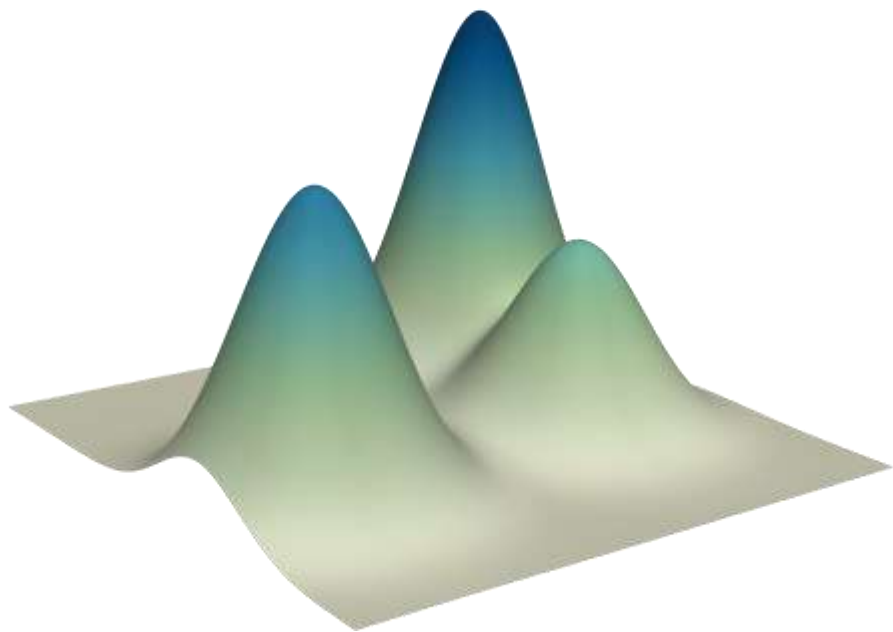


$$p_{\text{data}}(x) * \mathcal{N}(x | 0, \sigma)$$

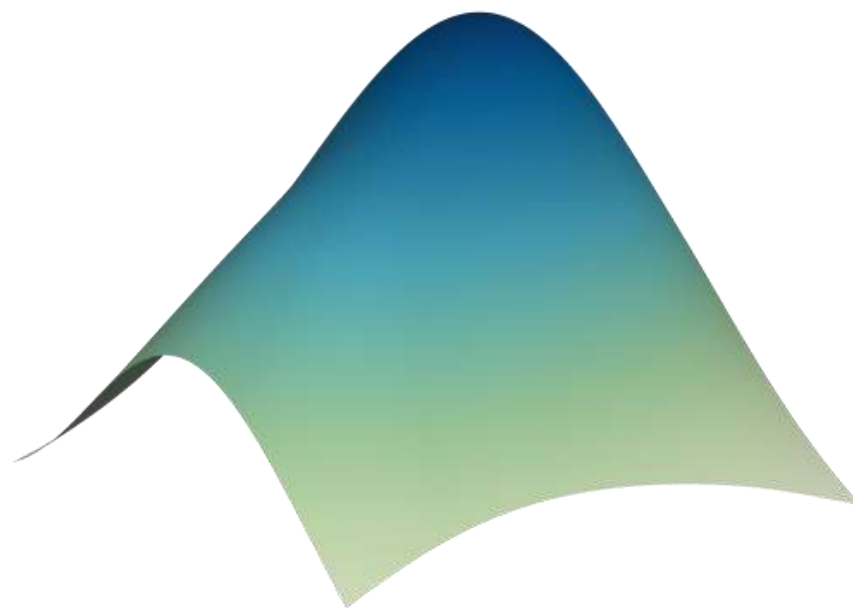
What is noise?

- Adding Gaussian noise \Leftrightarrow sampling

$$x \sim \mathcal{N}(x | x_0, \sigma)$$



$$p_{\text{data}}(x)$$

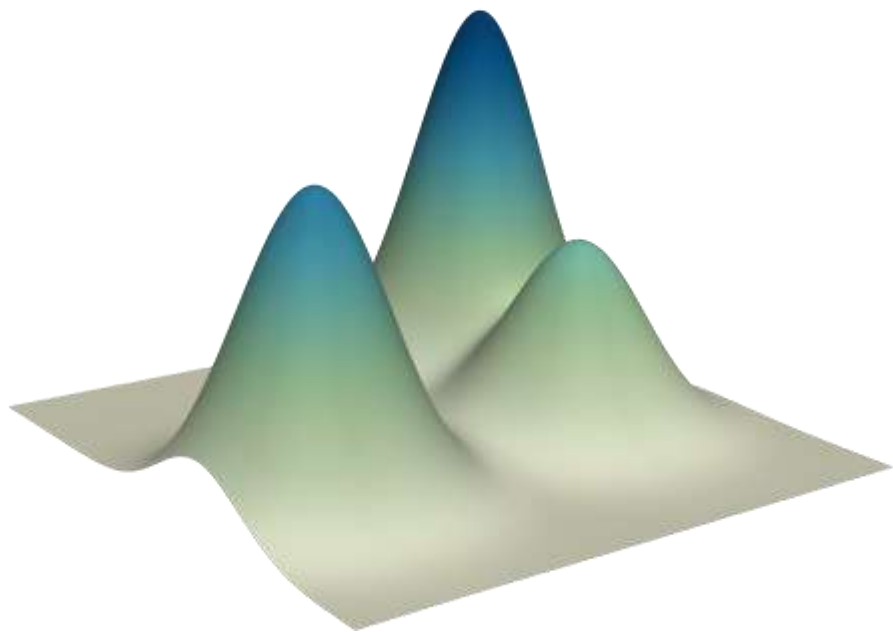


$$p_{\text{data}}(x) * \mathcal{N}(x | 0, \sigma)$$

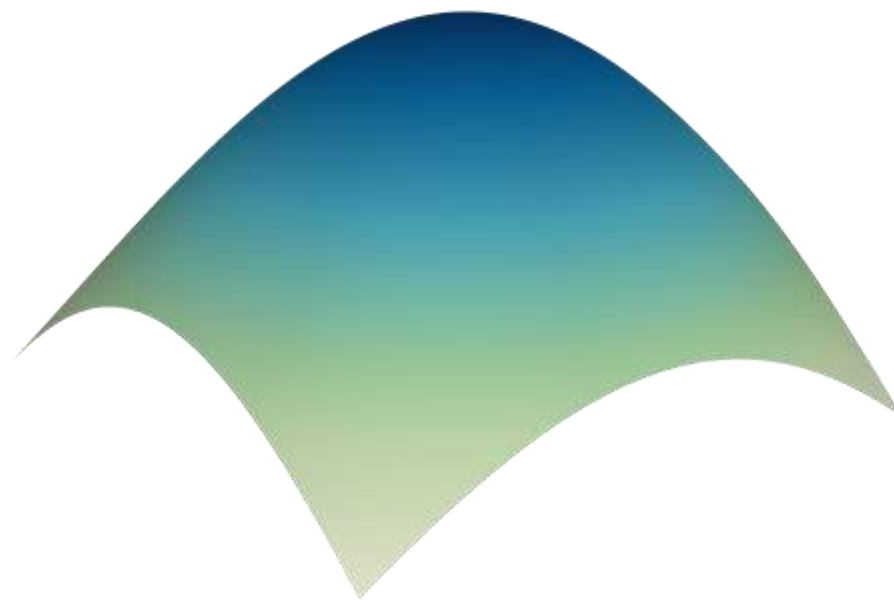
What is noise?

- Adding Gaussian noise \Leftrightarrow sampling

$$x \sim \mathcal{N}(x \mid x_0, \sigma)$$

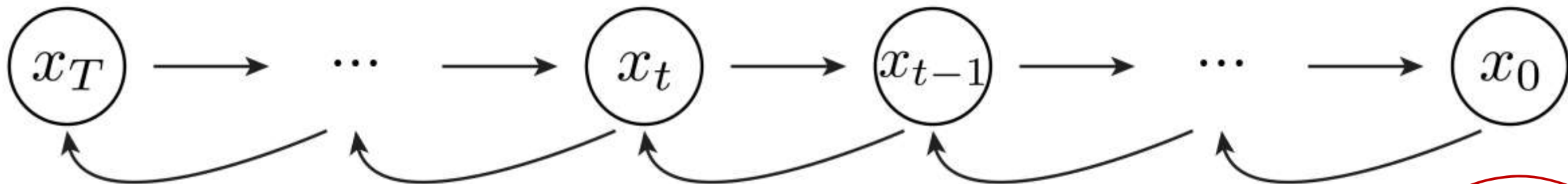


$$p_{\text{data}}(x)$$

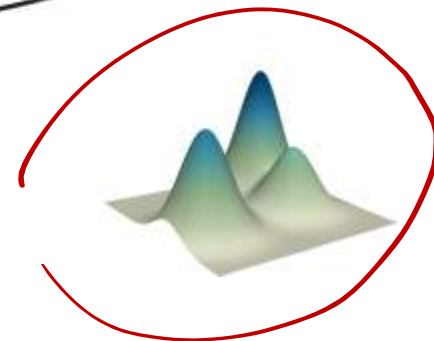
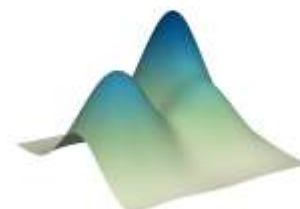
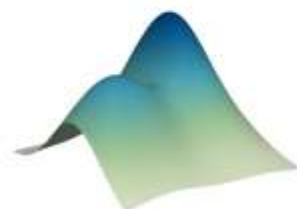


$$p_{\text{data}}(x) * \mathcal{N}(x \mid 0, \sigma)$$

What is noise?

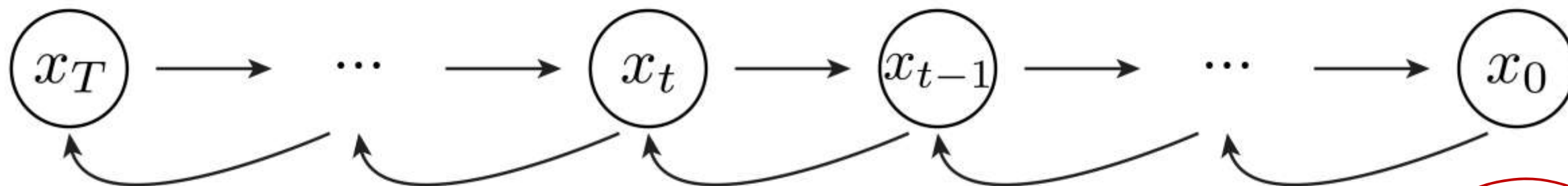


noise
distribution

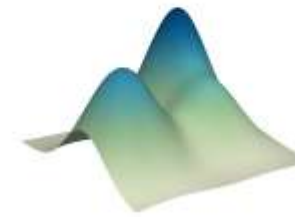
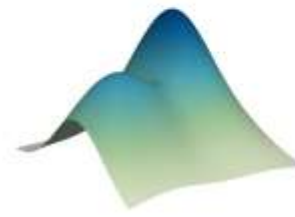
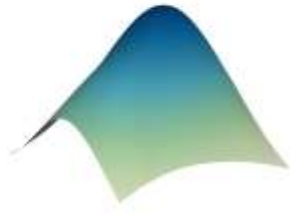


data
distribution

What is noise?



latent
distribution

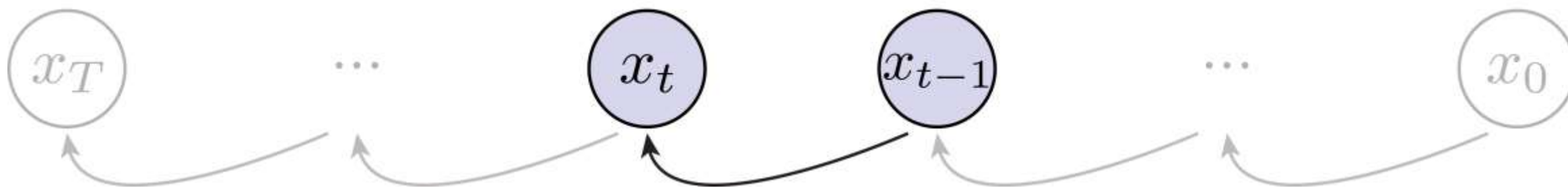


data
distribution

Diffusion Models

- Forward process
 - add noise to data
- Reverse process
 - learn to denoise
- Training objective
 - from Hierarchical VAE to L2 loss
- Noise Conditional Network
 - represent distributions

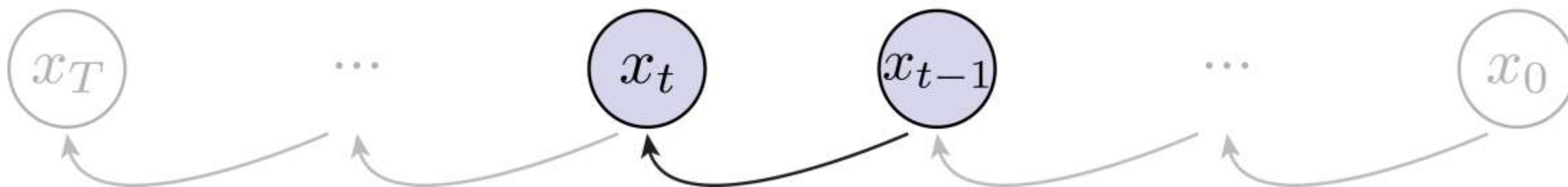
Forward Process



$$x_t = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

coefficients:
variance preserving

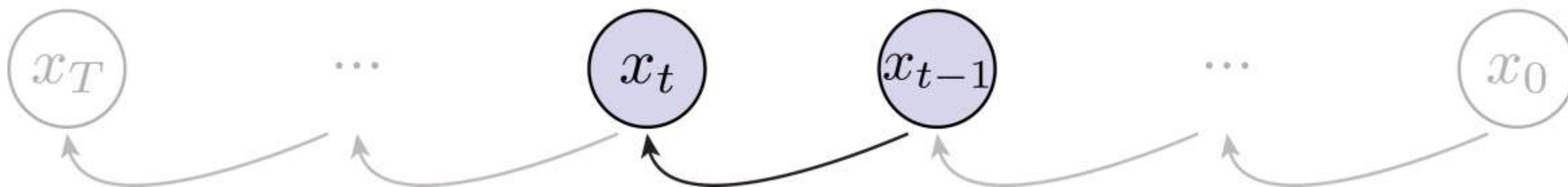
Forward Process



$$x_t = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

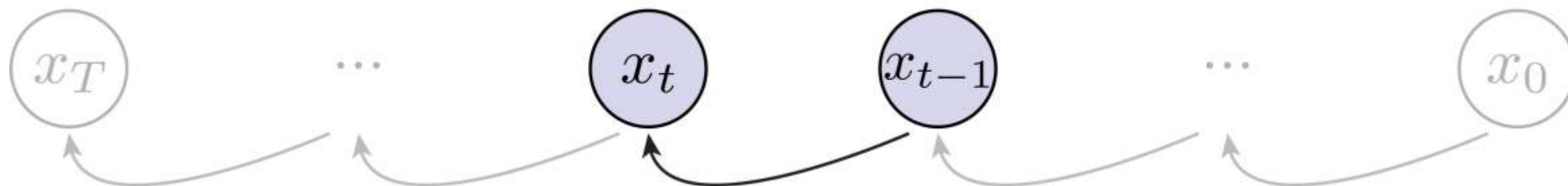
t : “schedule”,
key to Diffusion Models’ success

Forward Process



$$x_t = \underbrace{\sqrt{1 - \beta_t}}_{\text{mean of } x_t} x_{t-1} + \underbrace{\sqrt{\beta_t}}_{\text{std of } x_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

Forward Process



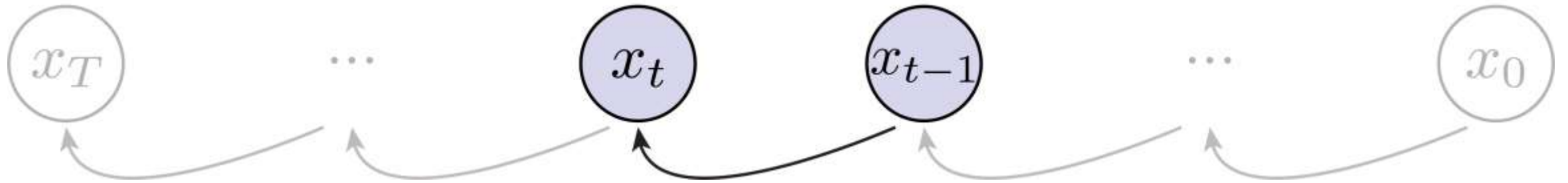
$$q(x_t | x_{t-1})$$

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

mean of x_t

var of x_t

Forward Process



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

identity matrix

- sampling is i.i.d.
- dim = dim of data

$$q(x_t|x_{t-1})$$

 x_t

$$= \sqrt{1 - \beta_t}$$

 x_{t-1}

$$+ \sqrt{\beta_t}$$



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

 $\beta_1, \beta_2, \dots, \beta_T$

$$q(x_t|x_0)$$

 x_0

+



+

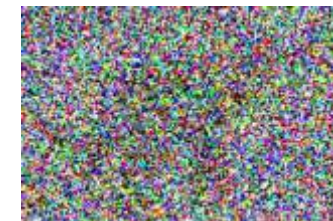


...

...



+

 x_t

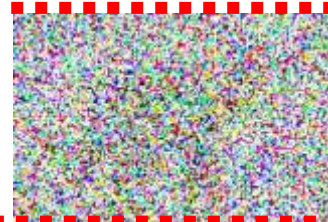


x_1

$$= \sqrt{1 - \beta_1} x_0 + \sqrt{\beta_1} z_1$$



x_0



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

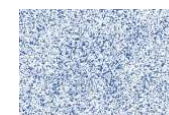


x_2

$$= \sqrt{1 - \beta_2} x_1 + \sqrt{\beta_2} z_2$$



x_1



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

Ind.

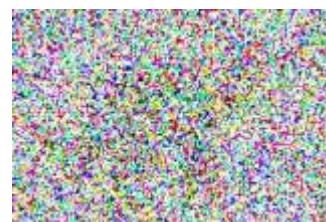


x_2

$$= \sqrt{1 - \beta_2} \sqrt{1 - \beta_1} x_0 + \sqrt{1 - \beta_2} \sqrt{\beta_1} z_1 + \sqrt{\beta_2} z_2$$



x_0





x_2

$$= \sqrt{1 - \beta_2} \sqrt{1 - \beta_1}$$



x_0



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



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

$$+ \sqrt{1 - \beta_2} \sqrt{\beta_1} \text{ [white noise] } + \sqrt{\beta_2} \text{ [blue noise]}$$

$$+ \sqrt{1 - (1 - \beta_2)(1 - \beta_1)} \text{ [yellow noise]}$$

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

$$q(x_t|x_0)$$

$$\beta_1, \beta_2, \dots, \beta_T$$



$$= \sqrt{1 - \beta_1}$$



$$+ \sqrt{\beta_1}$$



$$= \sqrt{1 - \beta_2}$$



$$+ \sqrt{\beta_2}$$



⋮

⋮

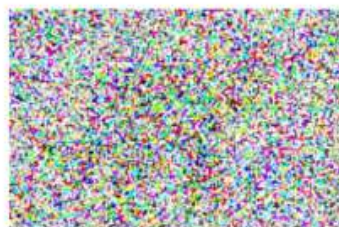
⋮



$$= \sqrt{1 - \beta_t}$$



$$+ \sqrt{\beta_t}$$

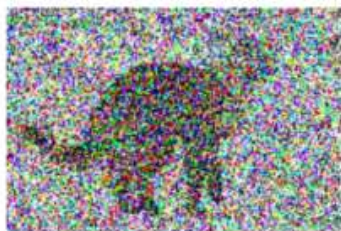


$$\sim \mathcal{N}(\mathbf{0}, I)$$

$$\alpha_t = 1 - \beta_t$$

$$\bar{\alpha}_t = \alpha_1 \alpha_2 \dots \alpha_t$$

⋮



$$= \frac{\sqrt{1 - \beta_1} \dots \sqrt{1 - \beta_t}}{\sqrt{\bar{\alpha}_t}}$$



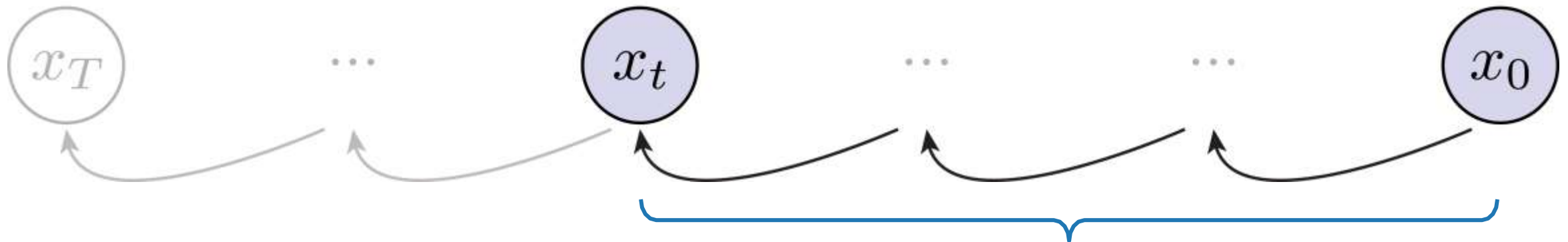
+

$$\frac{\sqrt{1 - (1 - \beta_1) \dots (1 - \beta_t)}}{\sqrt{1 - \bar{\alpha}_t}}$$

$$\sqrt{1 - \bar{\alpha}_t}$$



Forward Process



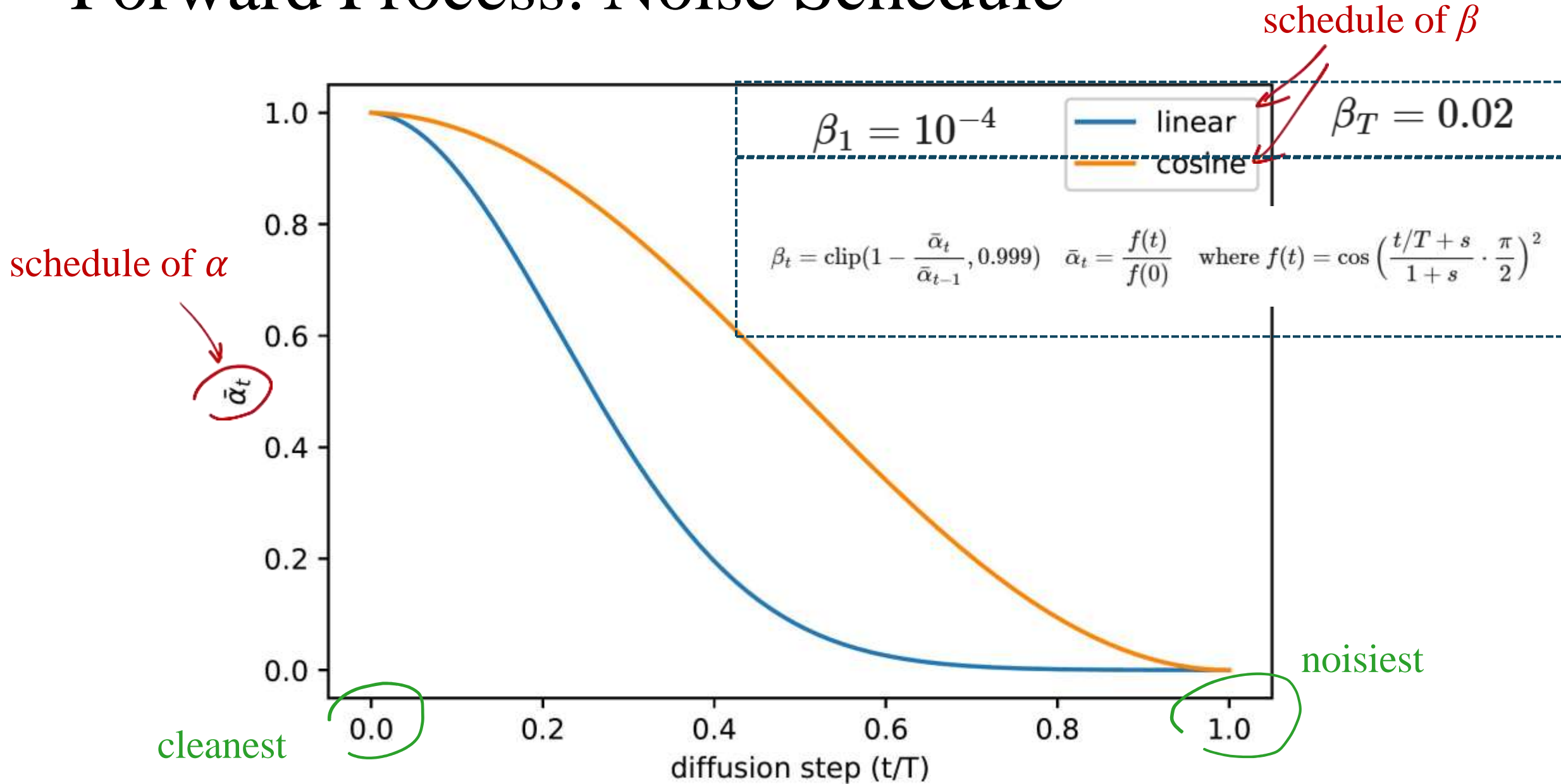
- sampling without simulation
- x_t from x_0 in closed form

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

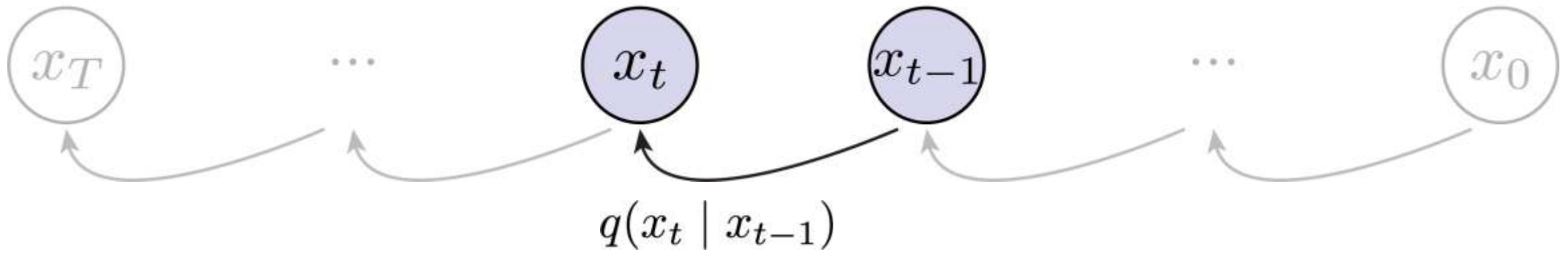
coefficients
given by β

$$\alpha_t := 1 - \beta_t$$
$$\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$$

Forward Process: Noise Schedule



Forward Process



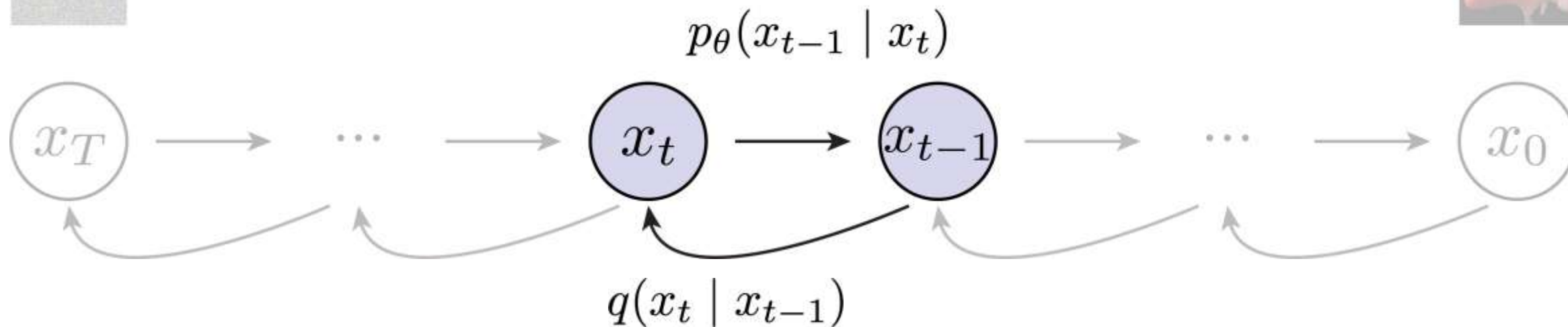
tl; dr:

- pre-defined conditional distributions
- Gaussian w/ controllable mean/std
- divide and conquer

Diffusion Models

- Forward process
 - add noise to data
- Reverse process
 - learn to denoise
- Training objective
 - from Hierarchical VAE to L2 loss
- Noise Conditional Network
 - represent distributions

Reverse Process



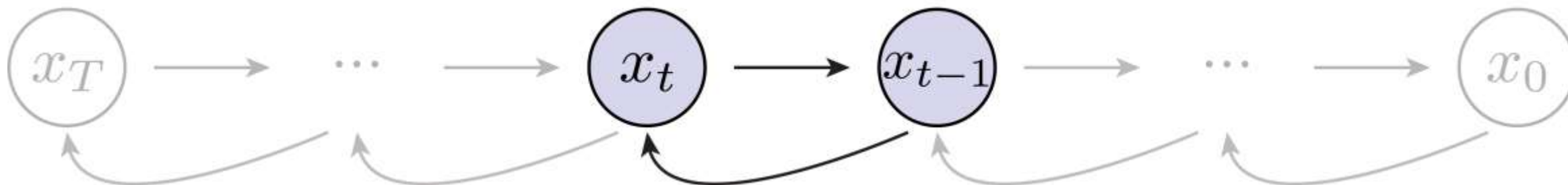
Reverse Process



parameterized
by a network

reverse the
time steps

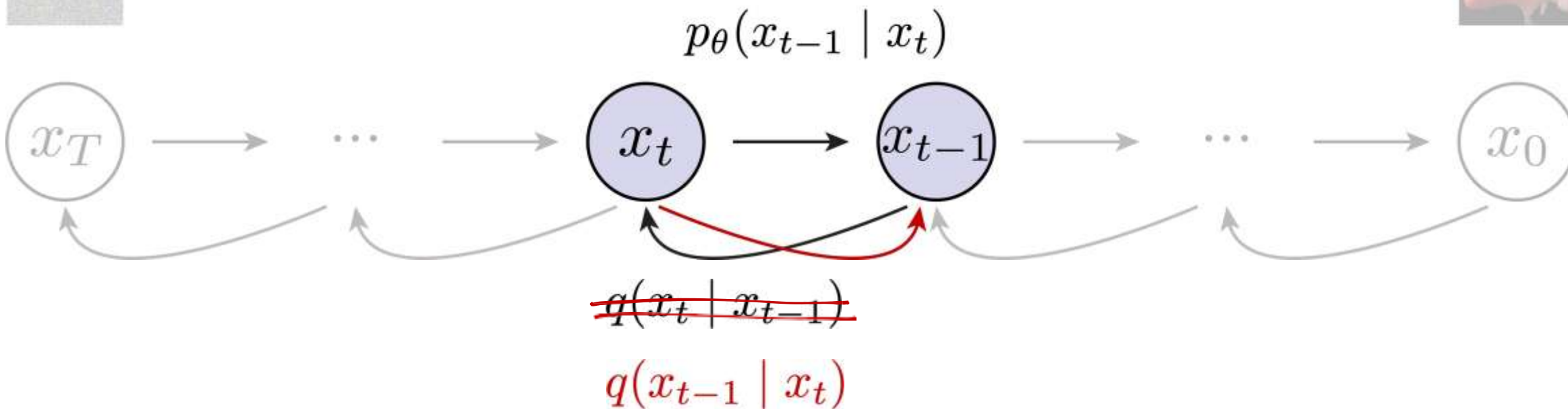
$$p_{\theta}(\underline{x_{t-1}} \mid \underline{x_t})$$



$$q(\underline{x_t \mid x_{t-1}})$$

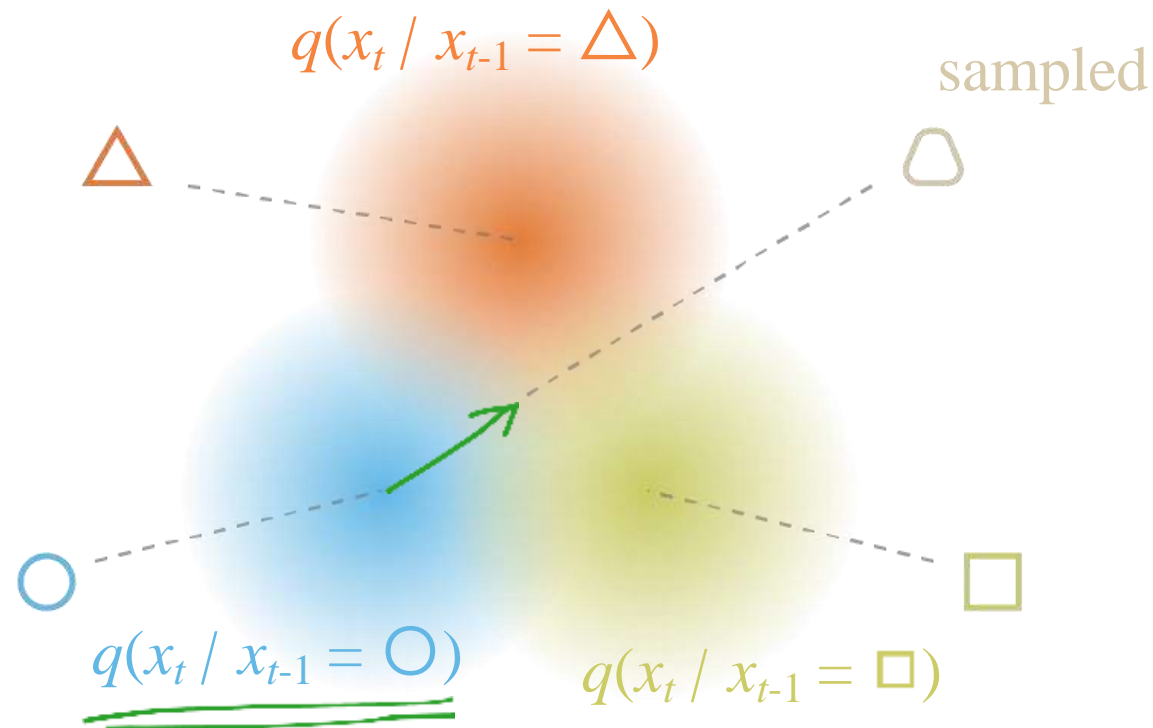
- known
- but not our target

Reverse Process



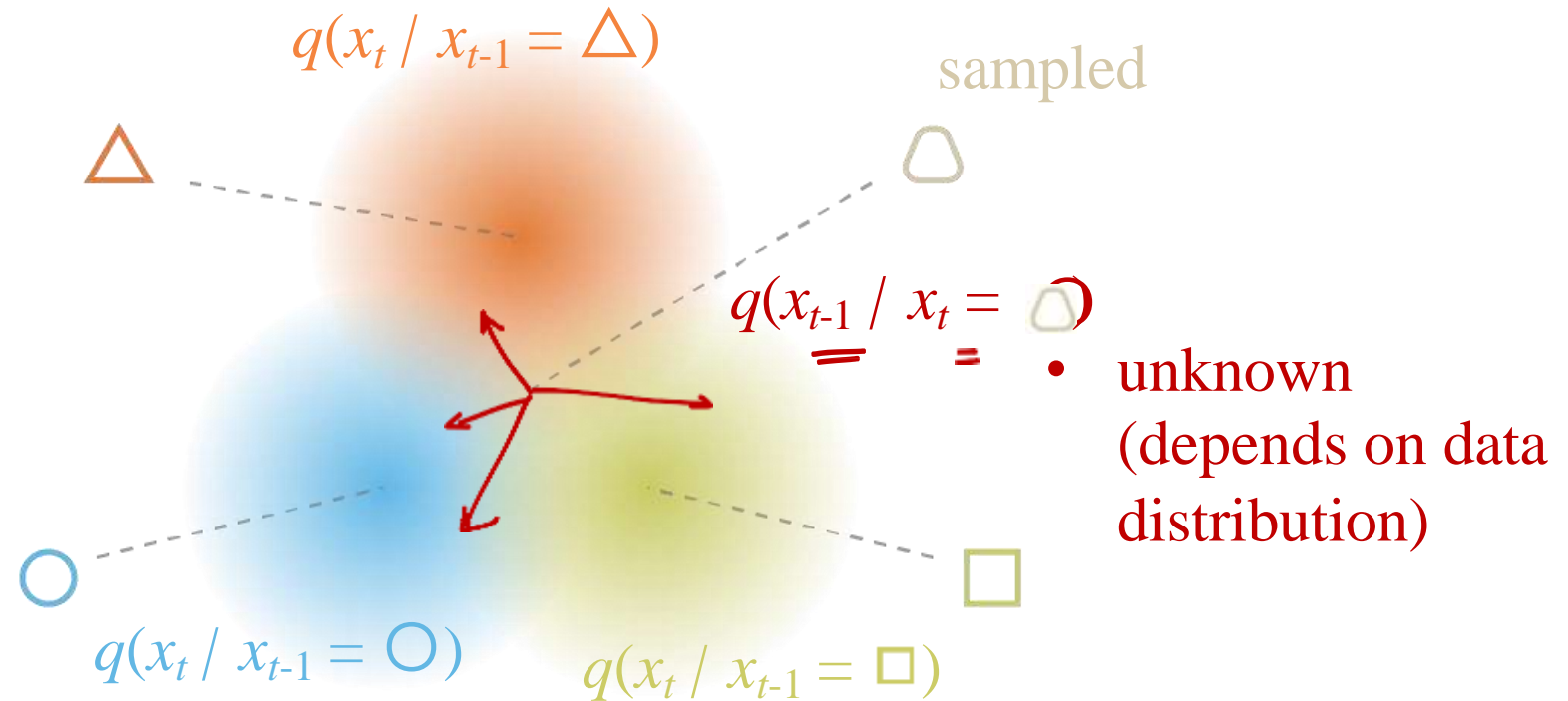
- our target
- but unknown

Why are the reverse conditionals unknown?

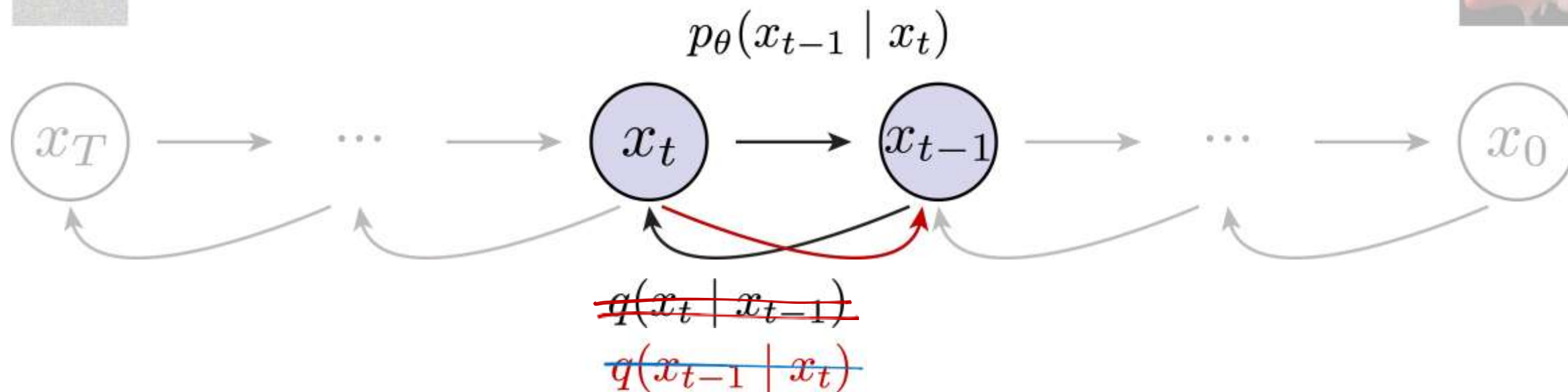


- known (Gaussian)

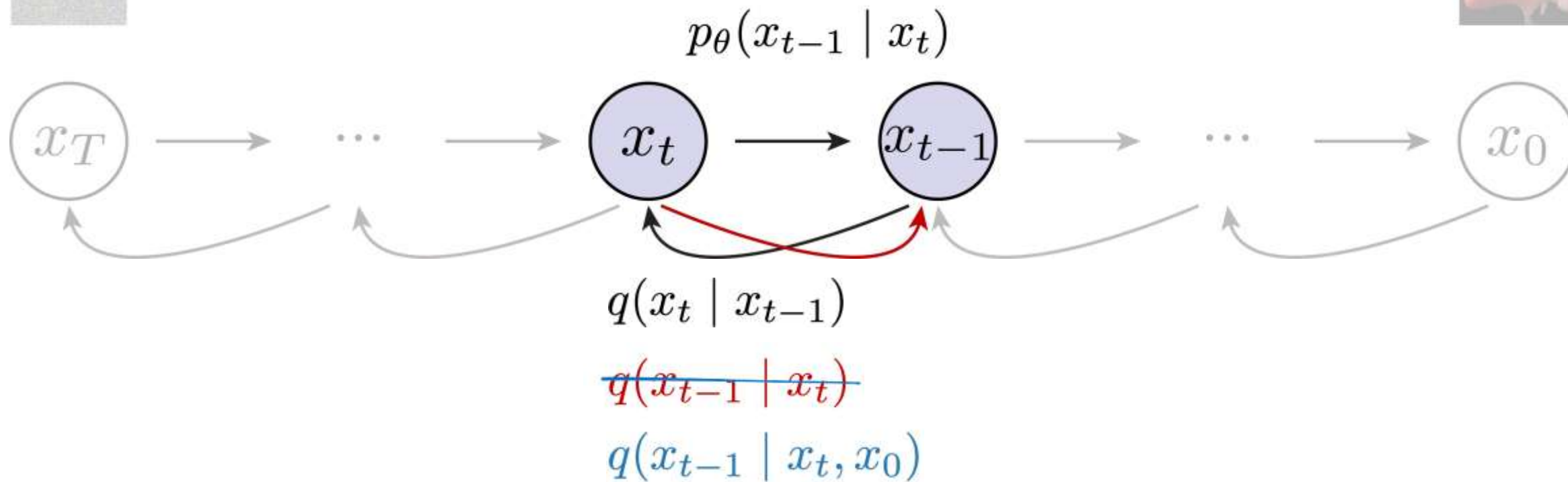
Why are the reverse conditionals unknown?



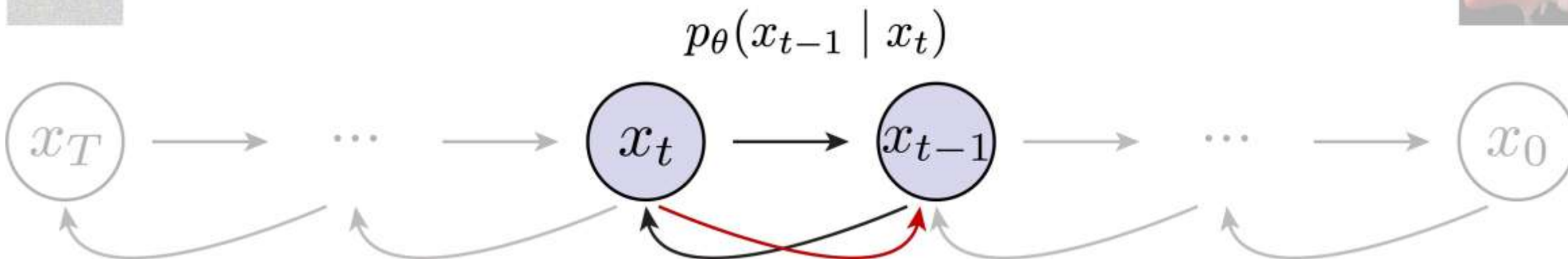
Reverse Process



Reverse Process



Reverse Process



$$p_{\theta}(x_{t-1} | x_t)$$

~~$$q(x_t | x_{t-1})$$~~

~~$$q(x_{t-1} | x_t)$$~~

$$q(x_{t-1} | x_t, x_0)$$

- known
- Gaussian

Reverse Process

$$q(x_t|x_0) \quad \text{[noisy cat]} = \sqrt{\bar{\alpha}_t} \quad \text{[cat]} + \sqrt{1 - \bar{\alpha}_t} \quad \text{[noise]}$$

$$q(x_{t-1}|x_0) \quad \text{[less noisy cat]} = \sqrt{\bar{\alpha}_{t-1}} \quad \text{[cat]} + \sqrt{1 - \bar{\alpha}_{t-1}} \quad \text{[noise]}$$

$$q(x_t|x_{t-1}) \quad \text{[noisy cat]} = \sqrt{1 - \beta_t} \quad \text{[less noisy cat]} + \sqrt{\beta_t} \quad \text{[noise]}$$

$$q(x_{t-1}|x_t, x_0)$$

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

- known
- Gaussian
- known
- Gaussian
- known
- Gaussian

$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \quad (34)$$

$$= \log \int \frac{p(\mathbf{x}_{0:T}) q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \quad (35)$$

$$= \log \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (36)$$

$$\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (37)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad (38)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_\theta(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_T|\mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad (39)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_\theta(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=1}^{T-1} p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_T|\mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad (40)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \prod_{t=1}^{T-1} \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad (41)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^{T-1} \log \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad (42)$$

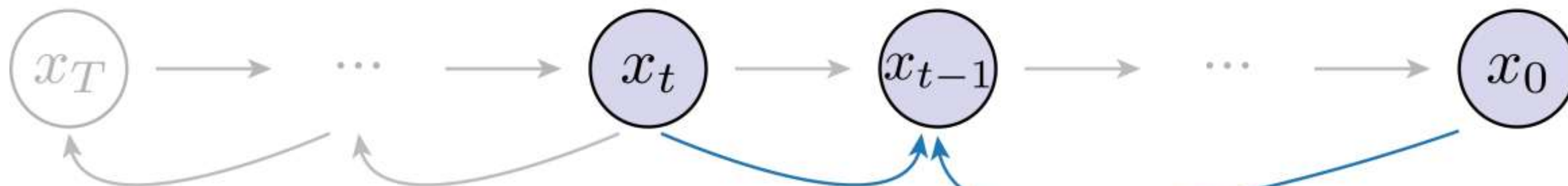
$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad (43)$$

$$= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{T-1}, \mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{x}_{t+1}|\mathbf{x}_0)} \left[\log \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad (44)$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_{T-1}) \parallel p(\mathbf{x}_T))]}_{\text{prior matching term}} \quad (45)$$

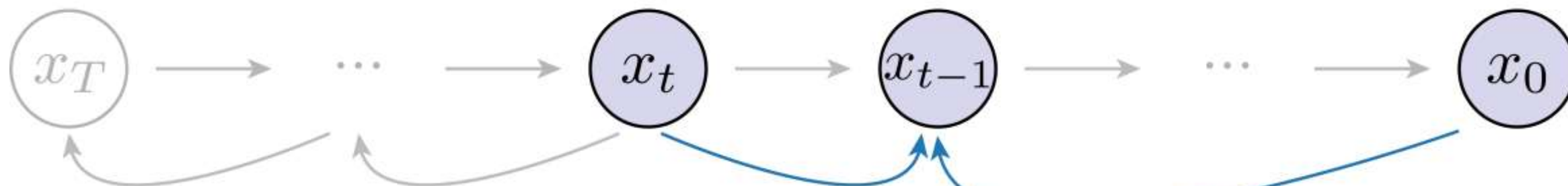
$$- \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) \parallel p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}))]}_{\text{consistency term}}$$

Reverse Process



$$q(x_{t-1} | x_t, x_0)$$
$$= \mathcal{N}(x_{t-1} | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \mathbf{I})$$

Reverse Process

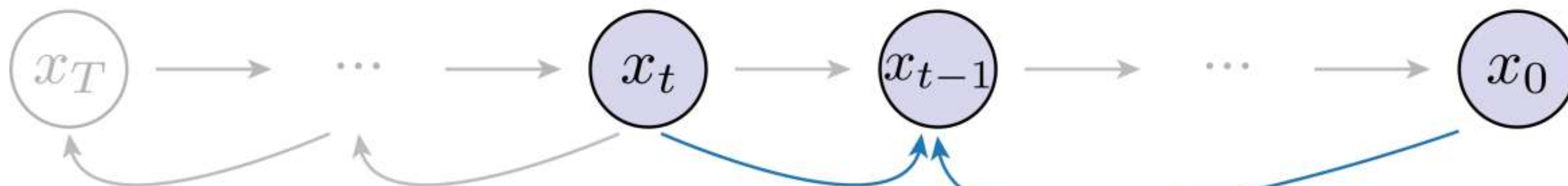


$$q(x_{t-1} | x_t, x_0)$$
$$= \mathcal{N}(x_{t-1} | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \mathbf{I})$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

var

Reverse Process



$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1} | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \mathbf{I})$$

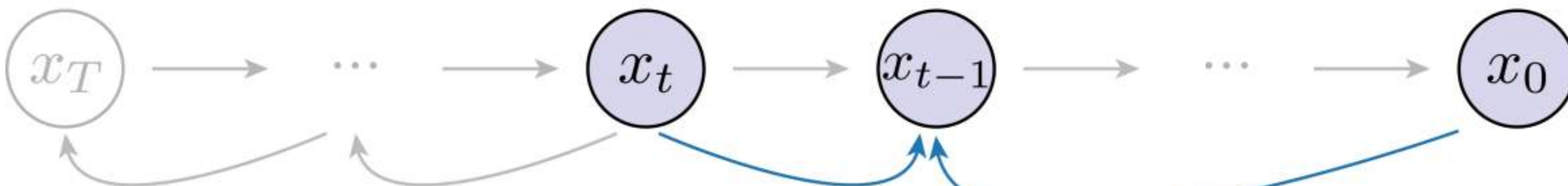
$$\tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} x_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t$$

mean

$$\tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

var

Reverse Process



$$q(x_{t-1} | x_t, x_0)$$

$$= \mathcal{N}(x_{t-1} | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \mathbf{I})$$

$$\tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} x_0 + \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t$$

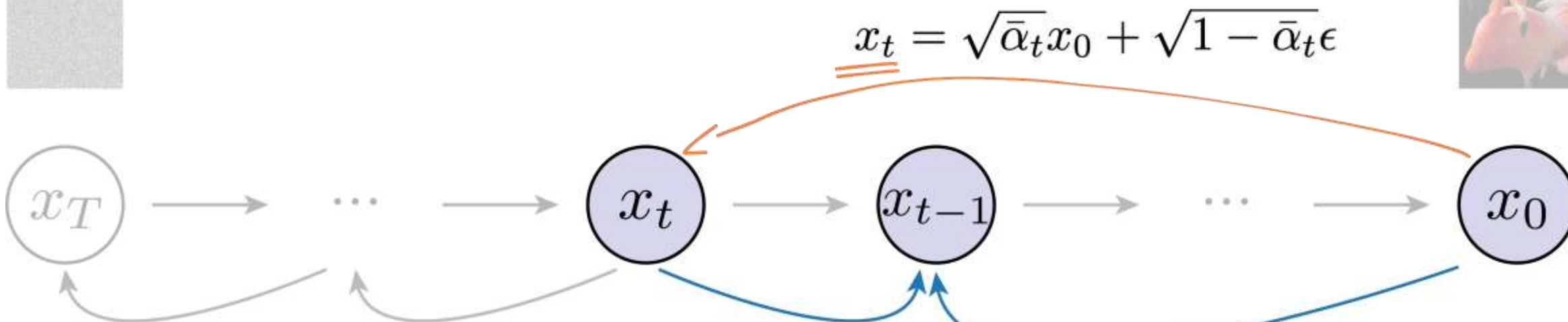
mean

linear combination

$$\tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

var

Reverse Process



$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1} | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \mathbf{I})$$

$$\tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\alpha_t} \beta_t}{1 - \alpha_t} x_0 + \frac{\sqrt{\alpha_t} (1 - \alpha_{t-1})}{1 - \alpha_t} x_t$$

$$\tilde{\beta}_t := \frac{1 - \alpha_{t-1}}{1 - \alpha_t} \beta_t$$

mean

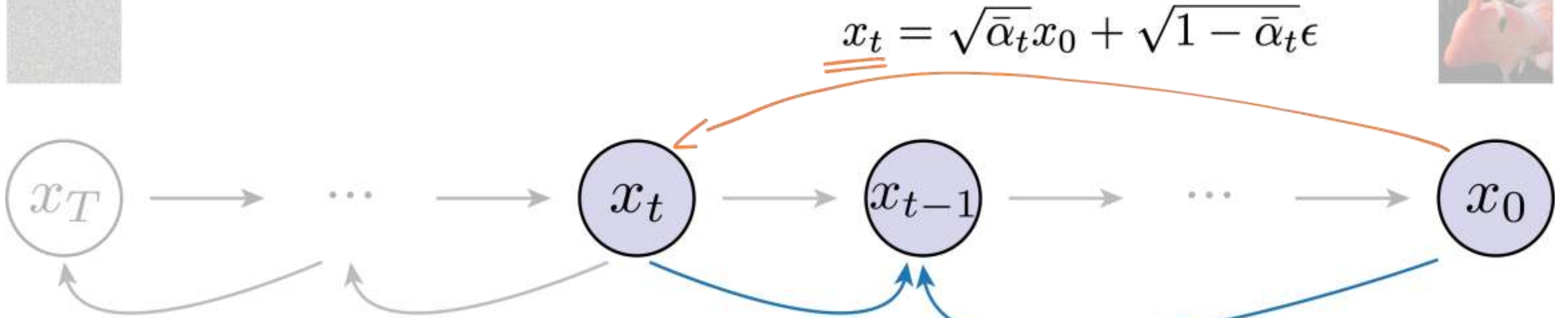
var

$$= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon \right) \text{ "noise"}$$

Reverse Process

tl; dr:

- outcome of the dependency graph
- some linear combinations



$$q(x_{t-1} | x_t, x_0)$$

$$= \mathcal{N}(x_{t-1} | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \mathbf{I})$$

$$\tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\alpha_t} \beta_t}{1 - \alpha_t} x_0 + \frac{\sqrt{\alpha_t} (1 - \alpha_{t-1})}{1 - \alpha_t} x_t$$

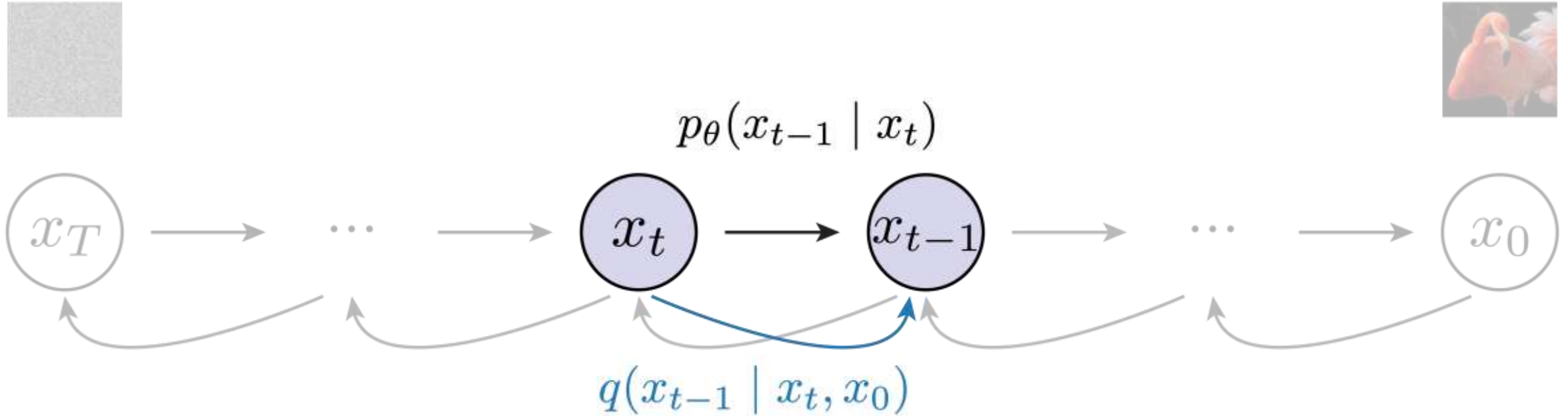
mean

$$= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon \right) \text{ "noise"}$$

$$\tilde{\beta}_t := \frac{1 - \alpha_{t-1}}{1 - \alpha_t} \beta_t$$

var

Reverse Process



- **tl; dr:** a known Gaussian
- we want to learn it by p_θ
- we can represent p_θ by a Gaussian
- minimize KL divergence

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \quad (71)$$

$$= \frac{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1-\alpha_t)\mathbf{I})\mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1-\bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I})} \quad (72)$$

$$\propto \exp \left\{ - \left[\frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{2(1-\alpha_t)} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{2(1-\bar{\alpha}_{t-1})} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{2(1-\bar{\alpha}_t)} \right] \right\} \quad (73)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{1-\alpha_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{1-\bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{1-\bar{\alpha}_t} \right] \right\} \quad (74)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{(-2\sqrt{\alpha_t}\mathbf{x}_t\mathbf{x}_{t-1} + \alpha_t\mathbf{x}_{t-1}^2)}{1-\alpha_t} + \frac{(\mathbf{x}_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}\mathbf{x}_0)}{1-\bar{\alpha}_{t-1}} + C(\mathbf{x}_t, \mathbf{x}_0) \right] \right\} \quad (75)$$

$$\propto \exp \left\{ - \frac{1}{2} \left[- \frac{2\sqrt{\alpha_t}\mathbf{x}_t\mathbf{x}_{t-1}}{1-\alpha_t} + \frac{\alpha_t\mathbf{x}_{t-1}^2}{1-\alpha_t} + \frac{\mathbf{x}_{t-1}^2}{1-\bar{\alpha}_{t-1}} - \frac{2\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right] \right\} \quad (76)$$

$$= \exp \left\{ - \frac{1}{2} \left[\left(\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (77)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t(1-\bar{\alpha}_{t-1}) + 1-\alpha_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (78)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t - \bar{\alpha}_t + 1-\alpha_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (79)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (80)$$

$$= \exp \left\{ - \frac{1}{2} \left(\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right)}{\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}} \mathbf{x}_{t-1} \right] \right\} \quad (81)$$

$$= \exp \left\{ - \frac{1}{2} \left(\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) (1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \quad (82)$$

$$= \exp \left\{ - \frac{1}{2} \left(\frac{1}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\mathbf{x}_0}{1-\bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \quad (83)$$

$$\propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\mathbf{x}_0}{1-\bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}}_{\Sigma_q(t)} \mathbf{I}) \quad (84)$$

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$



Gaussian



Mean

Variance

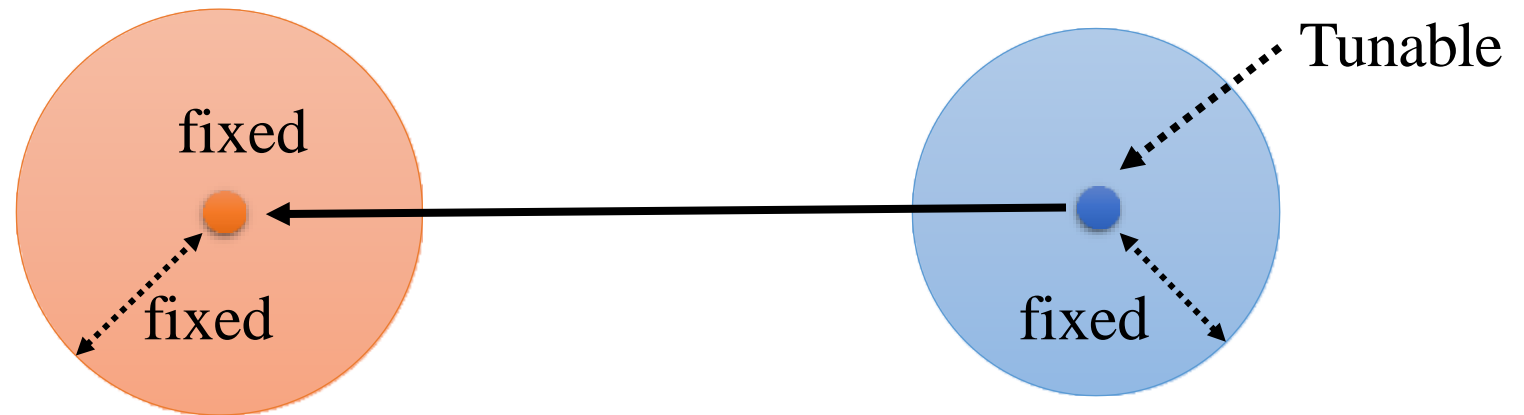
$$\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t x_0 + \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

$$\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t I$$

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0) || P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0) || P(x_{t-1}|x_t))]$$

How to
minimize KL
divergence?



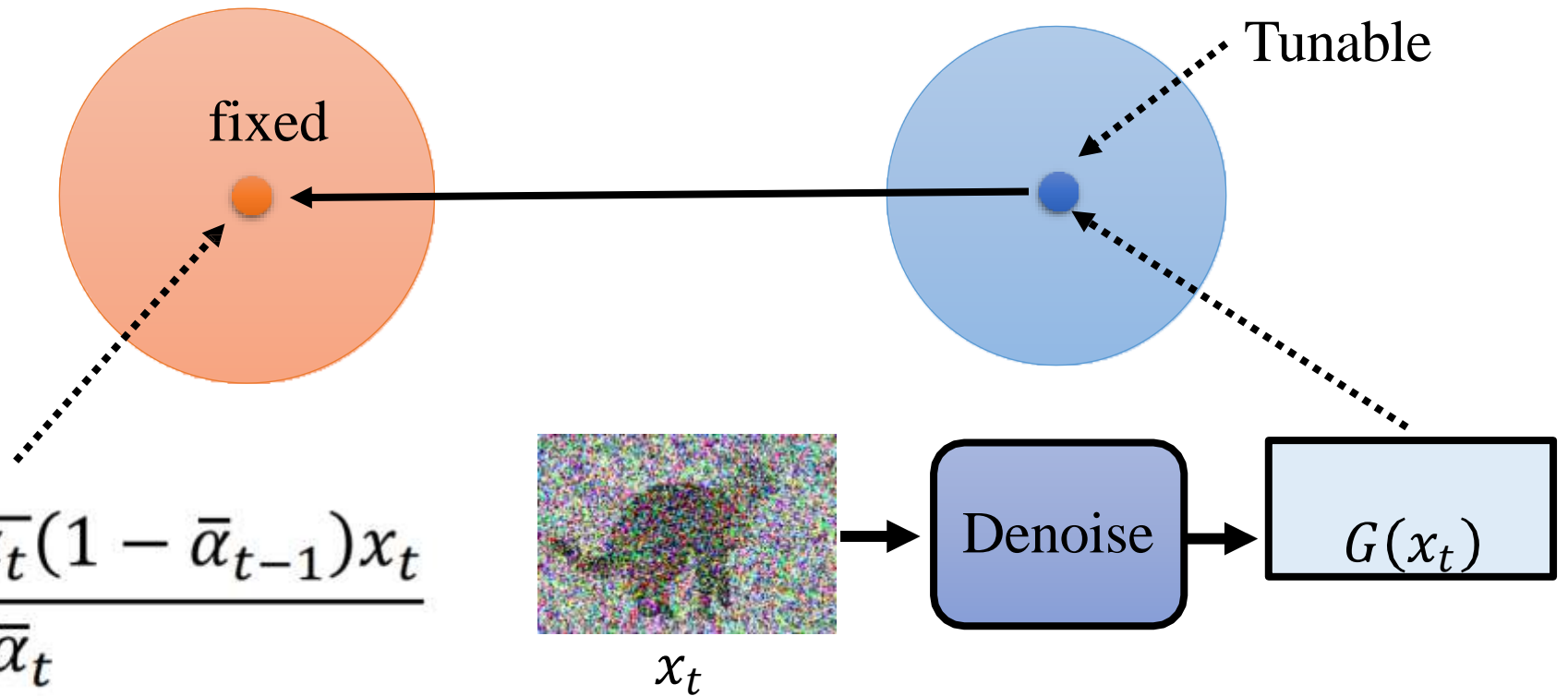
Recall that the KL Divergence between two Gaussian distributions is:

$$D_{\text{KL}}(\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x) || \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)) = \frac{1}{2} \left[\log \frac{|\boldsymbol{\Sigma}_y|}{|\boldsymbol{\Sigma}_x|} - d + \text{tr}(\boldsymbol{\Sigma}_y^{-1} \boldsymbol{\Sigma}_x) + (\boldsymbol{\mu}_y - \boldsymbol{\mu}_x)^T \boldsymbol{\Sigma}_y^{-1} (\boldsymbol{\mu}_y - \boldsymbol{\mu}_x) \right]$$

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0) || P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0) || P(x_{t-1}|x_t))]$$

How to minimize KL divergence?

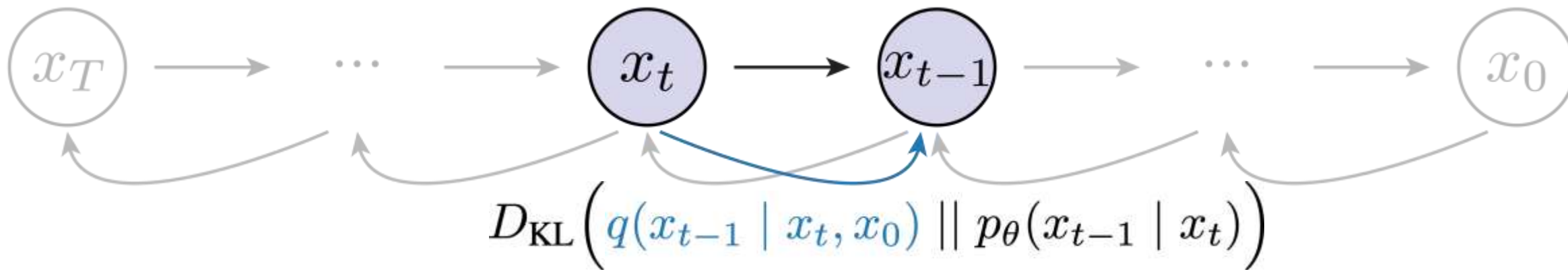


$$\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t x_0 + \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

Reverse Process

D_{KL} of two Gaussians is like L2 loss:

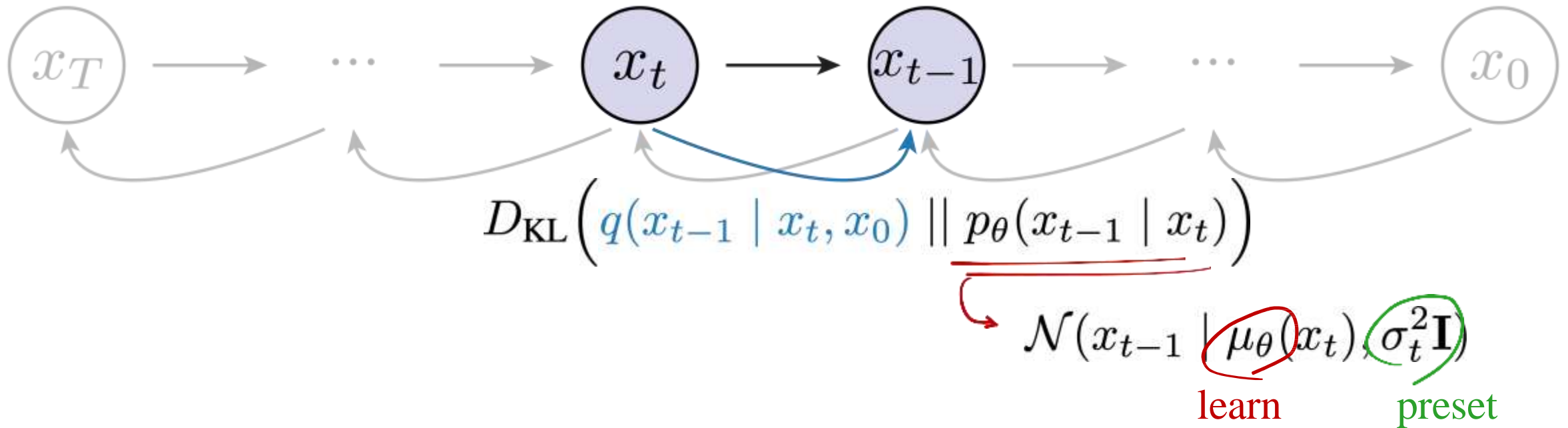
$$D_{\text{KL}}(\mathcal{N}_1 \parallel \mathcal{N}_2) = \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1^2 + \|\mu_1 - \mu_2\|^2}{2\sigma_2^2} - \frac{1}{2}$$



Reverse Process

D_{KL} of two Gaussians is like L2 loss:

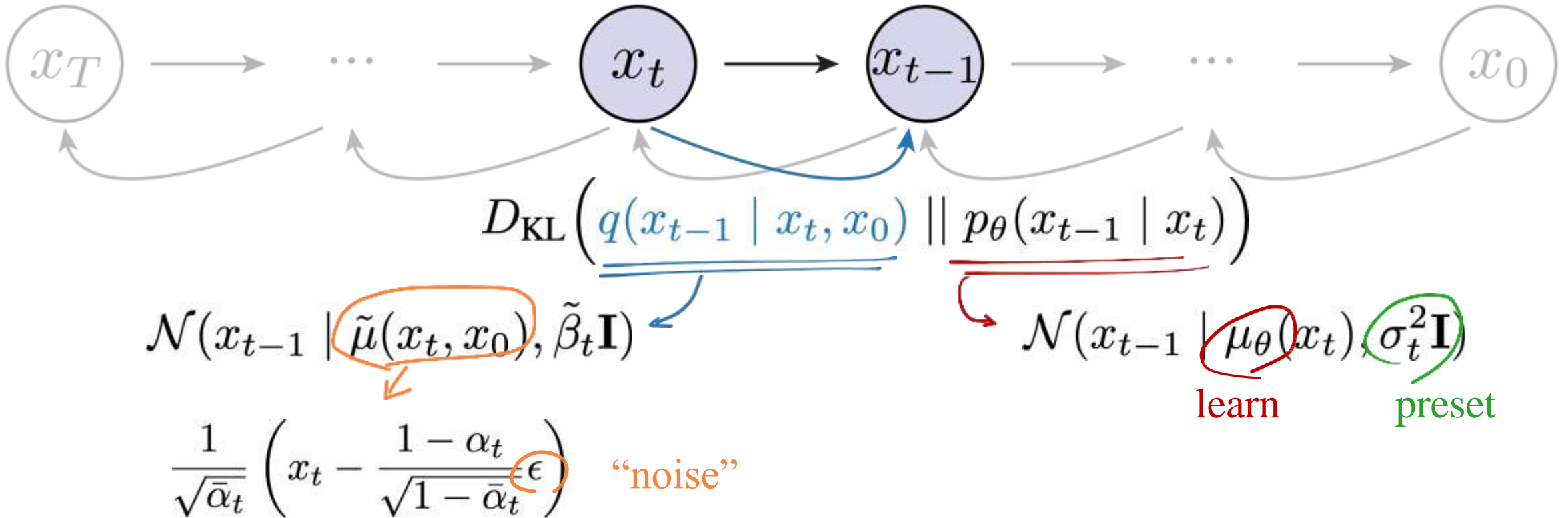
$$D_{\text{KL}}(\mathcal{N}_1 \parallel \mathcal{N}_2) = \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1^2 + \|\mu_1 - \mu_2\|^2}{2\sigma_2^2} - \frac{1}{2}$$



Reverse Process

D_{KL} of two Gaussians is like L2 loss:

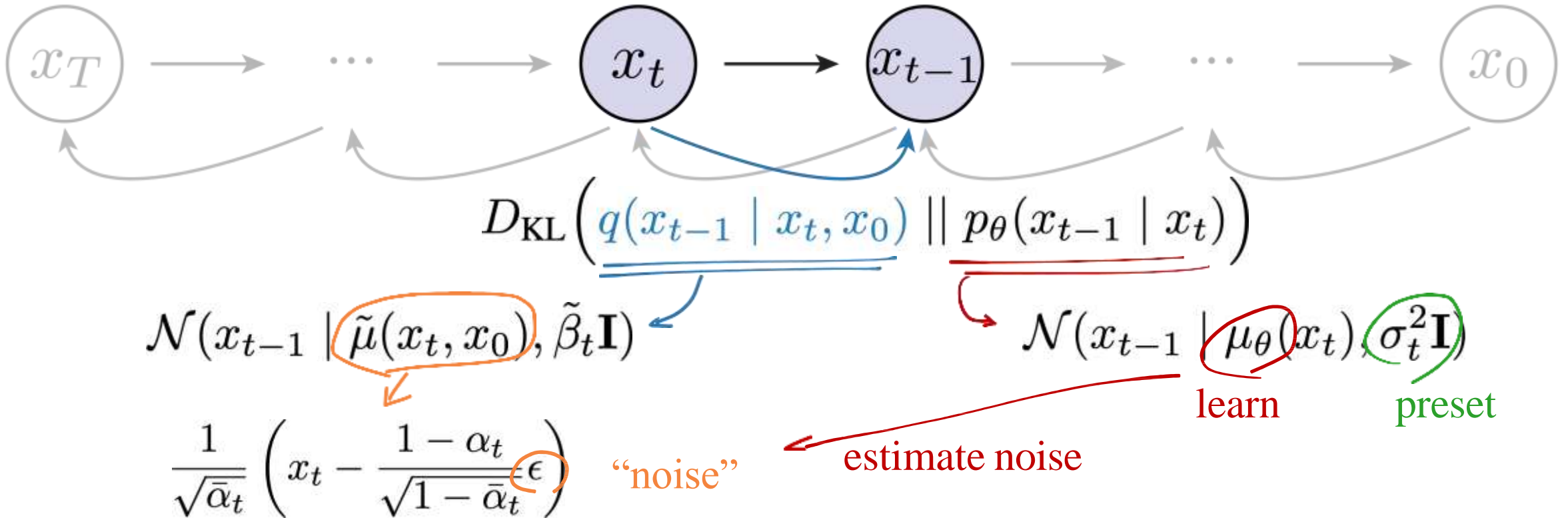
$$D_{KL}(\mathcal{N}_1 \parallel \mathcal{N}_2) = \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1^2 + \|\mu_1 - \mu_2\|^2}{2\sigma_2^2} - \frac{1}{2}$$



Reverse Process

D_{KL} of two Gaussians is like L2 loss:

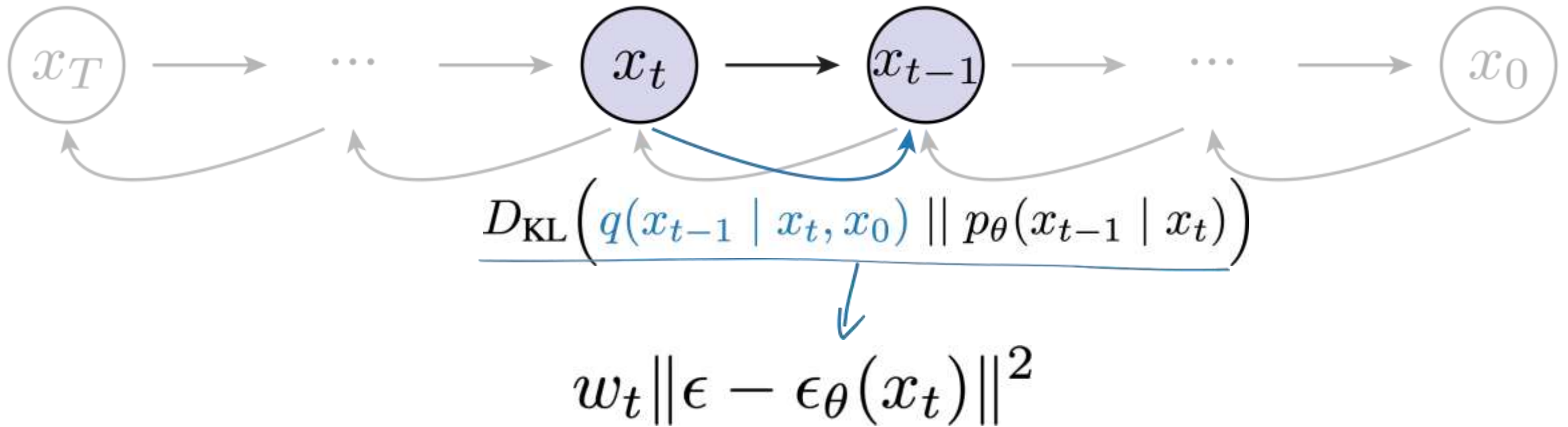
$$D_{KL}(\mathcal{N}_1 \parallel \mathcal{N}_2) = \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1^2 + \|\mu_1 - \mu_2\|^2}{2\sigma_2^2} - \frac{1}{2}$$



Reverse Process

D_{KL} of two Gaussians is like L2 loss:

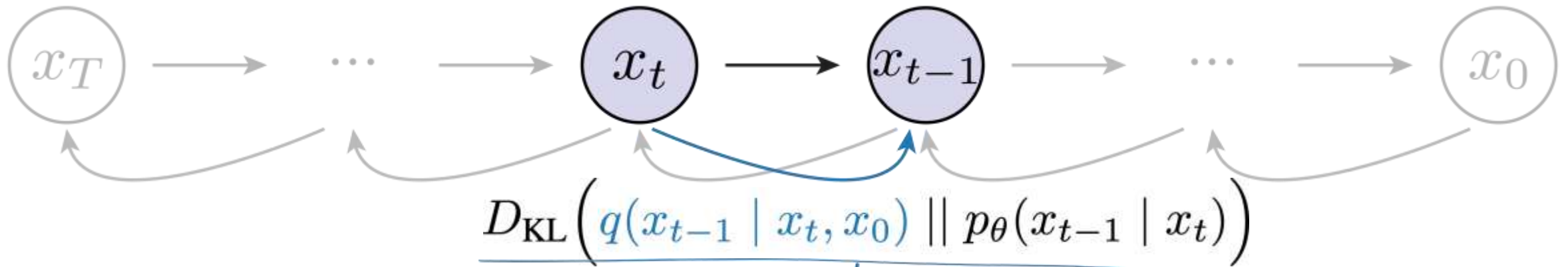
$$D_{\text{KL}}(\mathcal{N}_1 \parallel \mathcal{N}_2) = \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1^2 + \|\mu_1 - \mu_2\|^2}{2\sigma_2^2} - \frac{1}{2}$$



Reverse Process

D_{KL} of two Gaussians is like L2 loss:

$$D_{\text{KL}}(\mathcal{N}_1 \parallel \mathcal{N}_2) = \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1^2 + \|\mu_1 - \mu_2\|^2}{2\sigma_2^2} - \frac{1}{2}$$



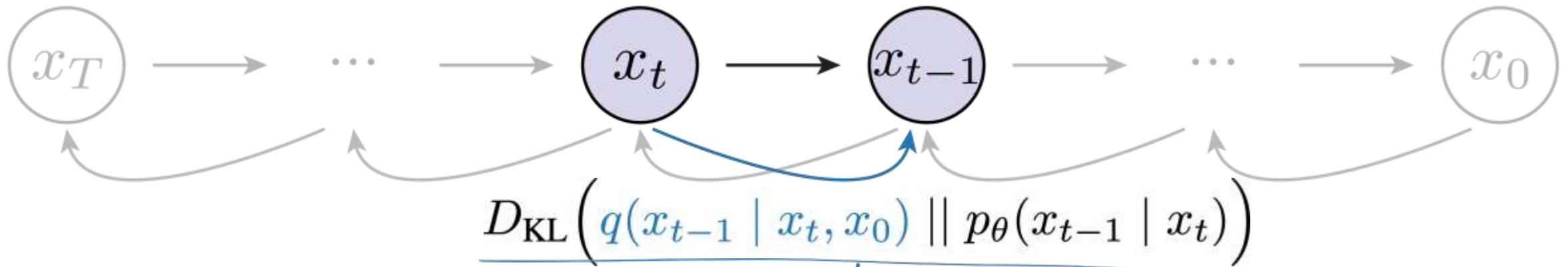
$$w_t \|\epsilon - \epsilon_{\theta}(x_t)\|^2$$

- a network to predict noise
- input: noisy image

Reverse Process

D_{KL} of two Gaussians is like L2 loss:

$$D_{\text{KL}}(\mathcal{N}_1 \parallel \mathcal{N}_2) = \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1^2 + \|\mu_1 - \mu_2\|^2}{2\sigma_2^2} - \frac{1}{2}$$



$$w_t \|\epsilon - \epsilon_{\theta}(x_t)\|^2$$

- weights due to α_t, β_t
- but set as 1 (**critical**)
- a network to predict noise
- input: noisy image

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0) || P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0) || P(x_{t-1}|x_t)))]$$



Sample x_0



Sample x_t



x_t

$$= \sqrt{\bar{\alpha}_t}$$



x_0

$$+ \sqrt{1 - \bar{\alpha}_t}$$



ϵ

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

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$



x_0

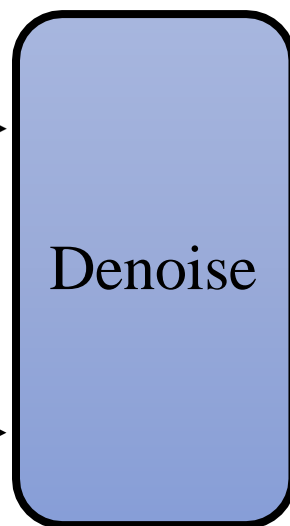


x_t



Sample x_t

t



$$\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t x_0 + \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

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



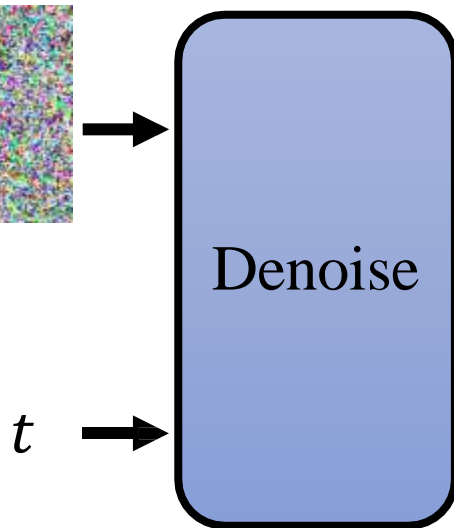
x_0



x_t



Sample x_t



$$\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t \boxed{x_0} + \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

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

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

$$\frac{x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon}{\sqrt{\bar{\alpha}_t}} = x_0$$

$$= \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t \boxed{\frac{x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon}{\sqrt{\bar{\alpha}_t}}} + \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

$$= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon \right)$$



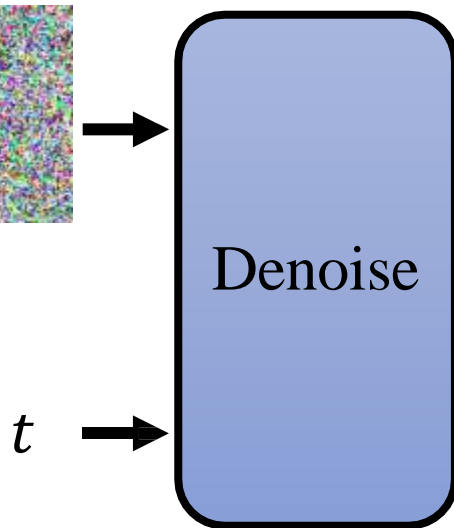
x_0



x_t



Sample x_t



$$? \longleftrightarrow \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \begin{matrix} \boxed{\varepsilon} \\ \vdots \end{matrix} \right)$$

To predict

Algorithm 2 Sampling

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

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

$$\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \varepsilon}{\sqrt{\bar{\alpha}_t}} = x_0$$

- 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: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

tl; dr

- some dependency graphs
- some linear combinations
- D_{KL}
- L2 loss of noise

Recap.

... in a nutshell

noise

data



x_T

...

x_t

x_{t-1}

...

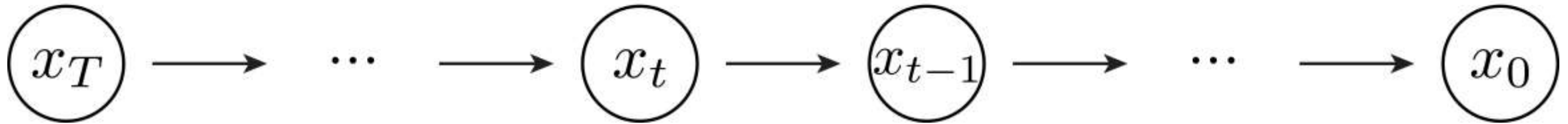
x_0



... in a nutshell

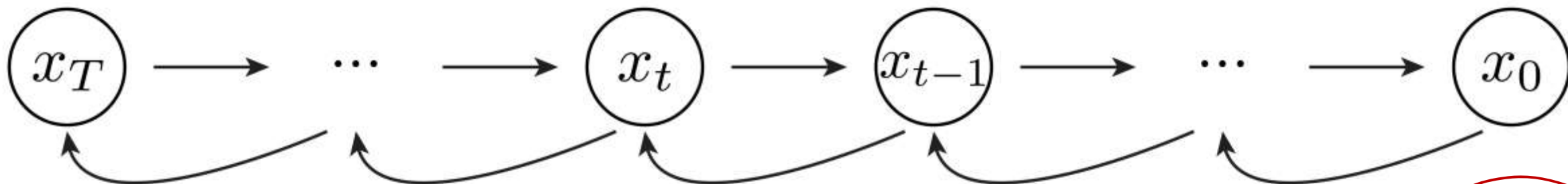
noise

data

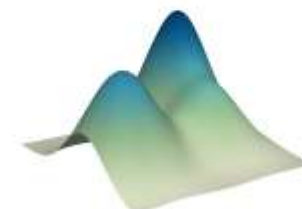
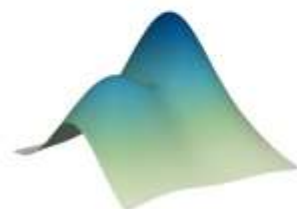


Reverse process: denoise

What is noise?

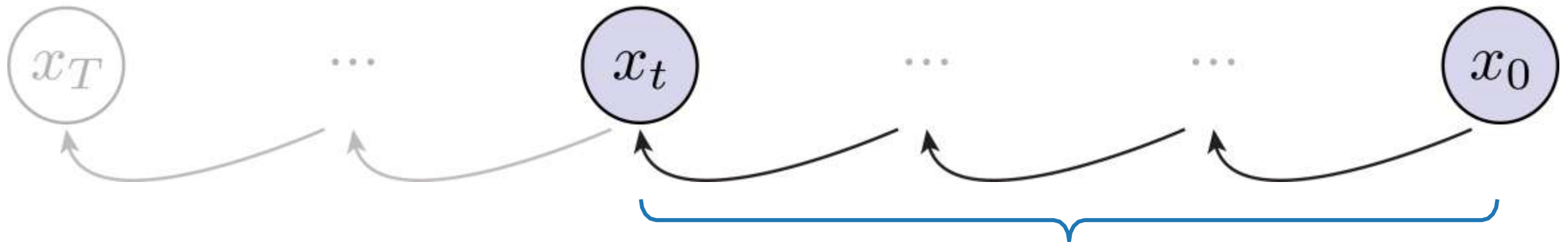


latent
distribution



data
distribution

Forward Process



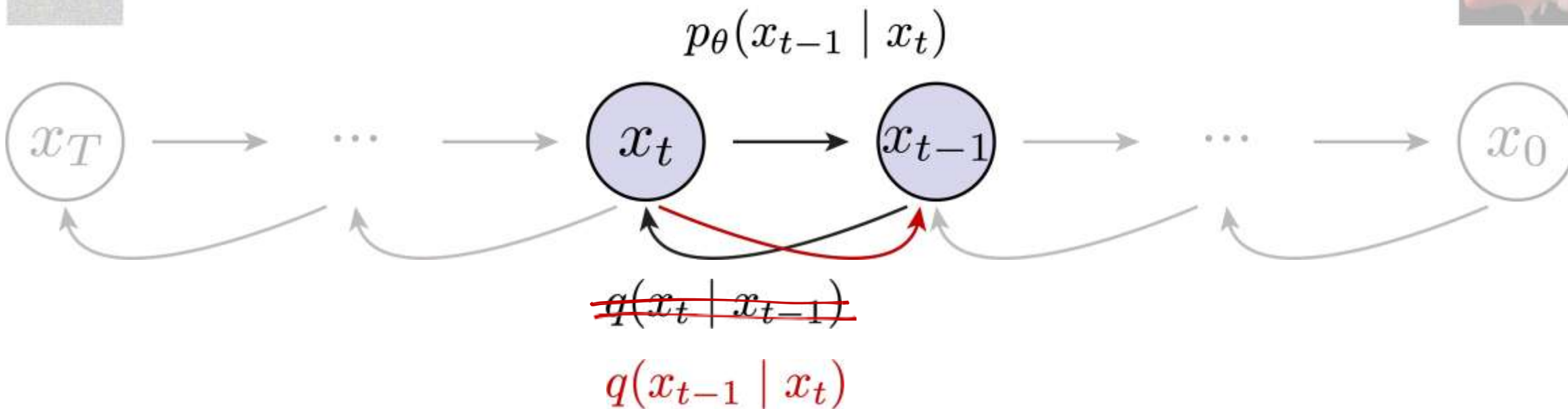
- sampling without simulation
- x_t from x_0 in closed form

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

coefficients
given by β

$$\alpha_t := 1 - \beta_t$$
$$\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$$

Reverse Process



- our target
- but unknown

Reverse Process

$$q(x_t|x_0) \quad \text{[noisy cat]} = \sqrt{\bar{\alpha}_t} \quad \text{[clear cat]} + \sqrt{1 - \bar{\alpha}_t} \quad \text{[noise]}$$

$$q(x_{t-1}|x_0) \quad \text{[less noisy cat]} = \sqrt{\bar{\alpha}_{t-1}} \quad \text{[clear cat]} + \sqrt{1 - \bar{\alpha}_{t-1}} \quad \text{[noise]}$$

$$q(x_t|x_{t-1}) \quad \text{[noisy cat]} = \sqrt{1 - \beta_t} \quad \text{[less noisy cat]} + \sqrt{\beta_t} \quad \text{[noise]}$$

$$q(x_{t-1}|x_t, x_0)$$

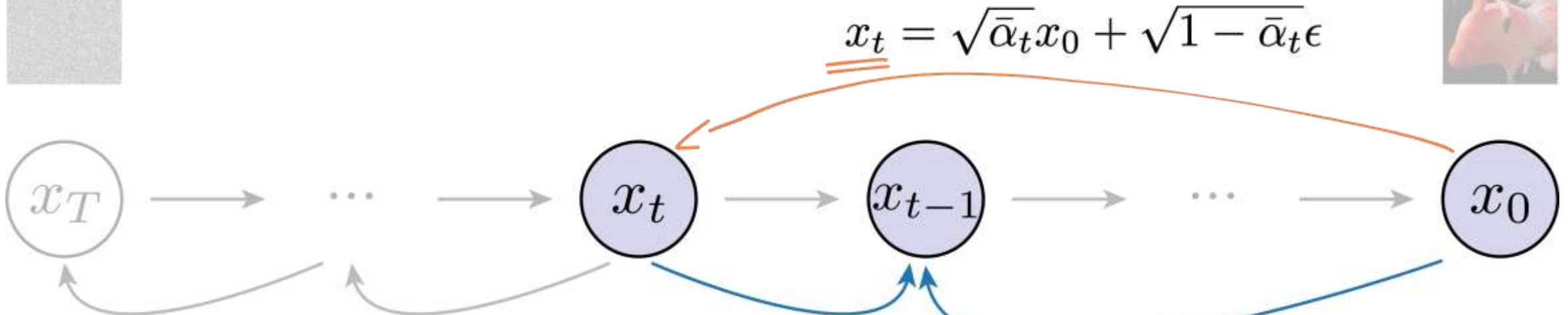
$$= \frac{q(x_{t-1}, x_t, x_0)}{q(x_t, x_0)} = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)q(x_0)}{q(x_t|x_0)q(x_0)} = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}$$

- known
- Gaussian
- known
- Gaussian
- known
- Gaussian

Reverse Process

tl; dr:

- outcome of the dependency graph
- some linear combinations



$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1} | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \mathbf{I})$$

$$\tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\alpha_t} \beta_t}{1 - \alpha_t} x_0 + \frac{\sqrt{\alpha_t} (1 - \alpha_{t-1})}{1 - \bar{\alpha}_t} x_t$$

mean

$$= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right) \text{ "noise"}$$

$$\tilde{\beta}_t := \frac{1 - \alpha_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

var



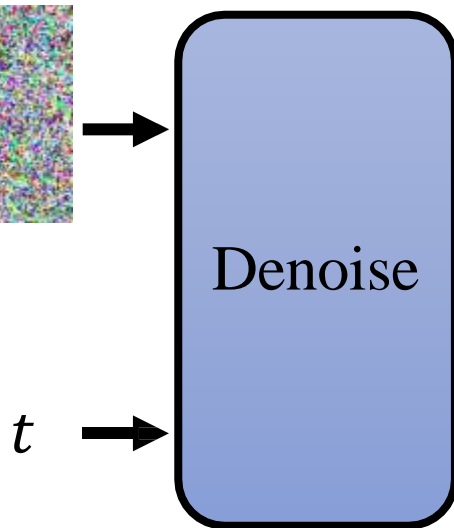
x_0



x_t



Sample x_t



$$\frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \begin{matrix} \boxed{\varepsilon} \\ \vdots \end{matrix} \right)$$

To predict

Algorithm 2 Sampling

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

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

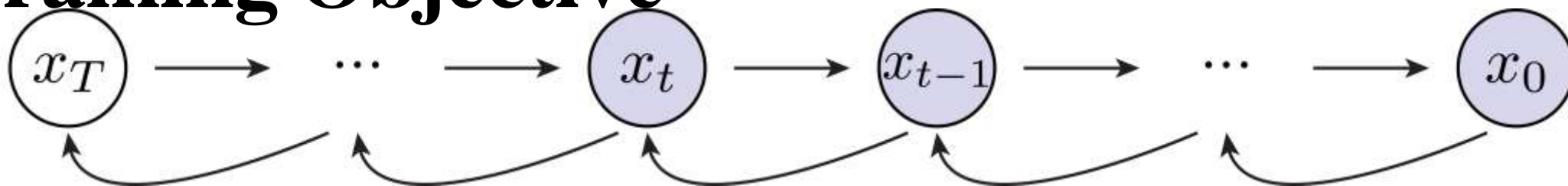
$$\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \varepsilon}{\sqrt{\bar{\alpha}_t}} = x_0$$

- 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: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Diffusion Models

- Forward process
 - add noise to data
- Reverse process
 - learn to denoise
- Training objective
 - from Hierarchical VAE to L2 loss
- Noise Conditional Network
 - represent a distribution

Training Objective



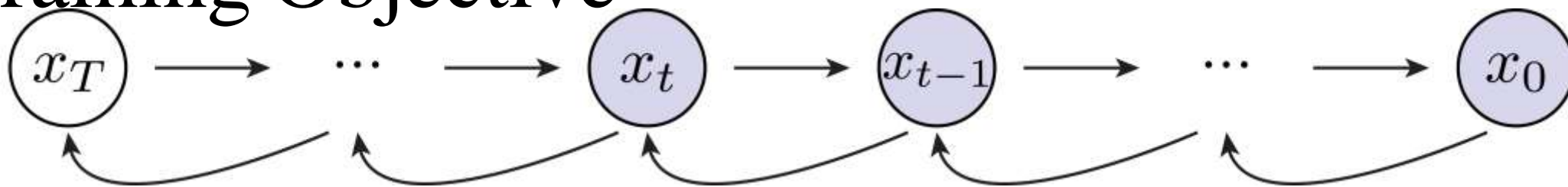
$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

$$\mathcal{L}_T := D_{\text{KL}}\left(q(x_T | x_0) \parallel p_{\theta}(x_T)\right)$$

$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t)\right)$$

$$\mathcal{L}_0 := -\log p_{\theta}(x_0 | x_1)$$

Training Objective



- variational lower bound
- like ELBO

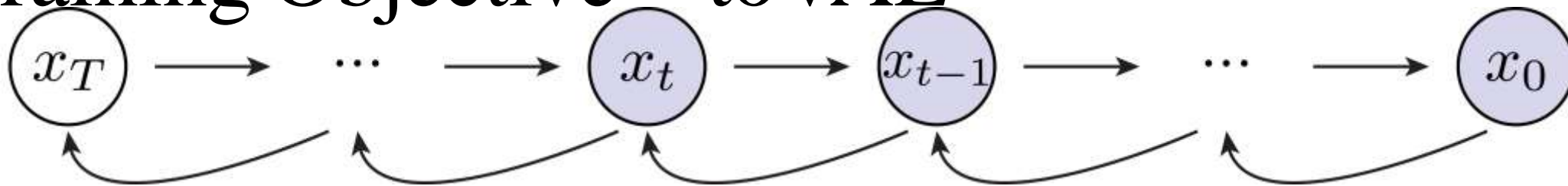
$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

$$\mathcal{L}_T := D_{\text{KL}}\left(q(x_T | x_0) \parallel p_{\theta}(x_T)\right)$$

$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t)\right)$$

$$\mathcal{L}_0 := -\log p_{\theta}(x_0 | x_1)$$

Training Objective – to VAE



- variational lower bound
- like ELBO

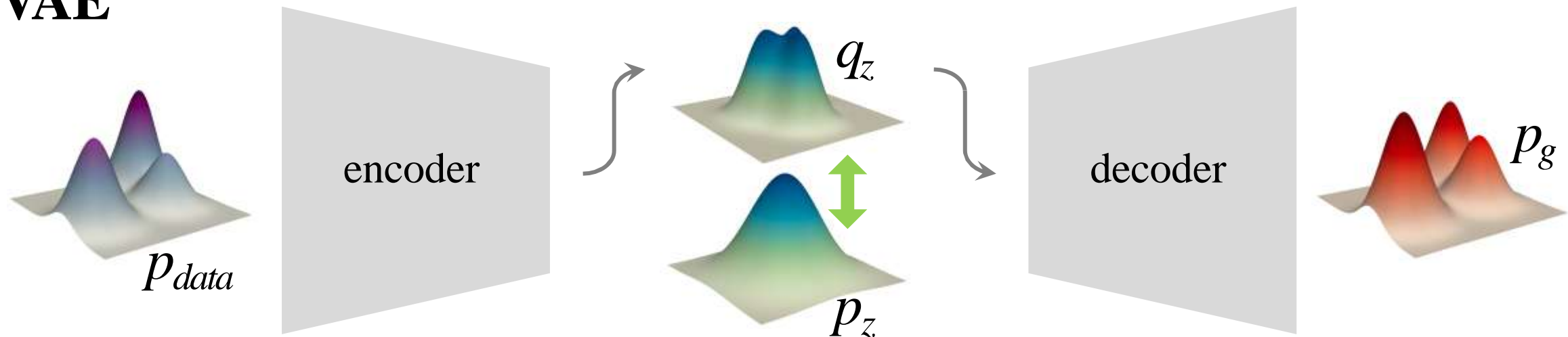
$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

$$\mathcal{L}_T := D_{\text{KL}}\left(q(\overset{z}{\cancel{x_T}} \mid x_0) \parallel p_{\theta}(\overset{z}{\cancel{x_T}})\right) \quad \text{it's ELBO if one step}$$

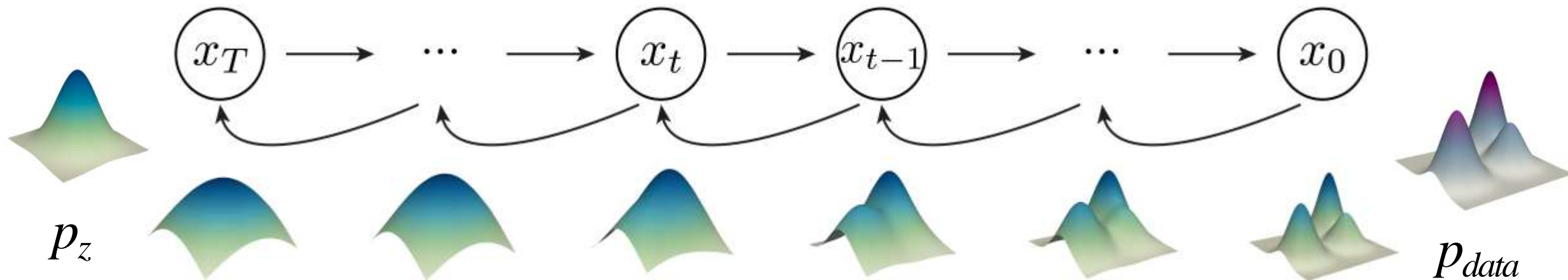
~~$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} \mid x_t, x_0) \parallel p_{\theta}(x_{t-1} \mid x_t)\right)$$~~

$$\mathcal{L}_0 := -\log p_{\theta}(x_0 \mid \overset{z}{\cancel{x_1}})$$

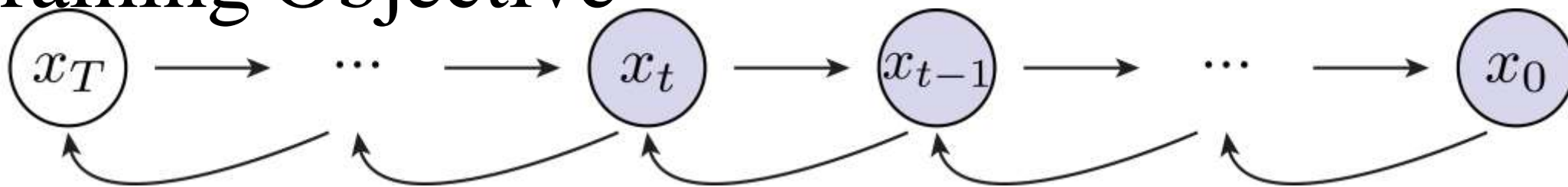
VAE



DM/Flow



Training Objective



$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

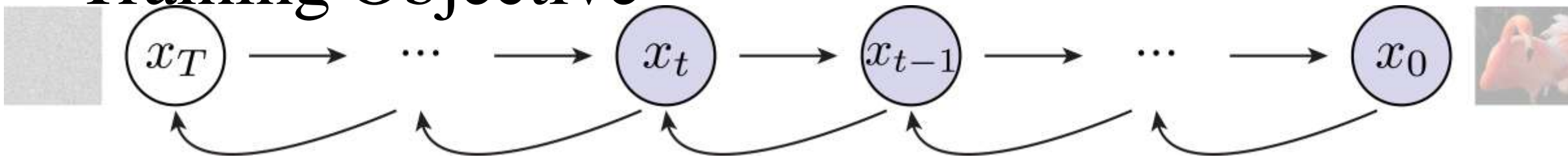
no parameter, unlike VAE's q_ϕ

$$\mathcal{L}_T := D_{\text{KL}} \left(\underbrace{q(x_T | x_0)}_{\text{constant}} \parallel \underbrace{p_\theta(x_T)}_{\text{Gaussian}} \right)$$

$$\mathcal{L}_{t-1} := D_{\text{KL}} \left(q(x_{t-1} | x_t, x_0) \parallel p_\theta(x_{t-1} | x_t) \right)$$

$$\mathcal{L}_0 := -\log p_\theta(x_0 | x_1)$$

Training Objective



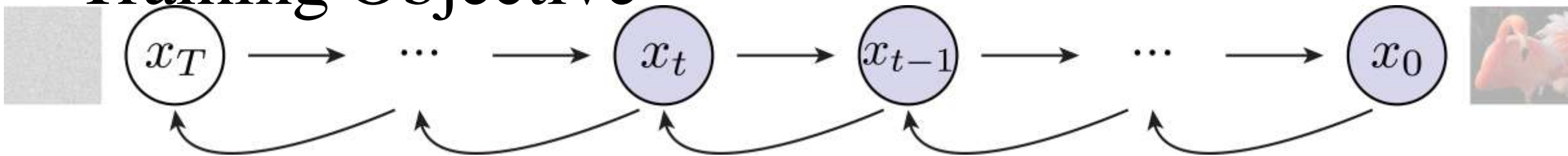
$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

$$\mathcal{L}_T := D_{\text{KL}}\left(q(x_T | x_0) \parallel p_{\theta}(x_T)\right)$$

$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t)\right)$$

reconstruction loss
like VAE $\mathcal{L}_0 := -\log p_{\theta}(x_0 | x_1)$

Training Objective



$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

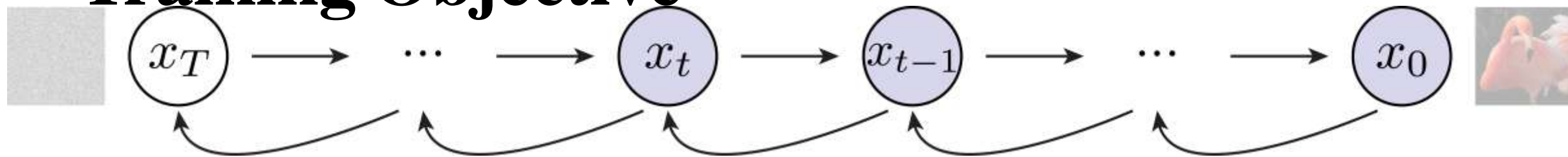
$$\mathcal{L}_T := D_{\text{KL}}\left(q(x_T | x_0) \parallel p_{\theta}(x_T)\right)$$

L2 loss on
noise

$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t)\right)$$

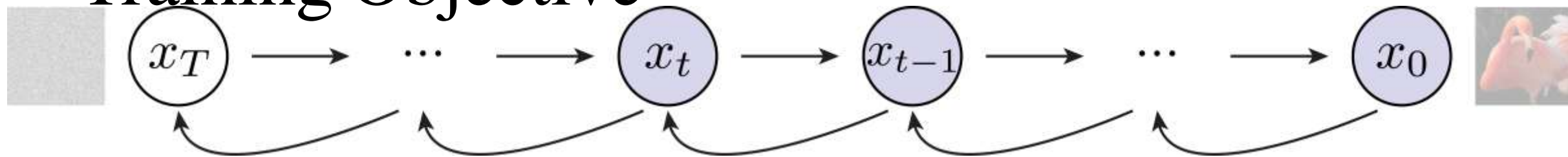
$$\mathcal{L}_0 := -\log p_{\theta}(x_0 | x_1)$$

Training Objective



$$\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \left[w_t \|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right]$$

Training Objective



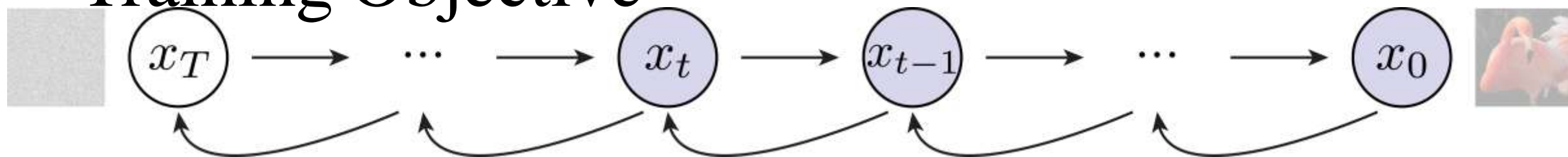
$$\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \left[w_t \|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right]$$

over p_{data}

over $[1, T]$

over $\mathcal{N}(0, \mathbf{I})$

Training Objective



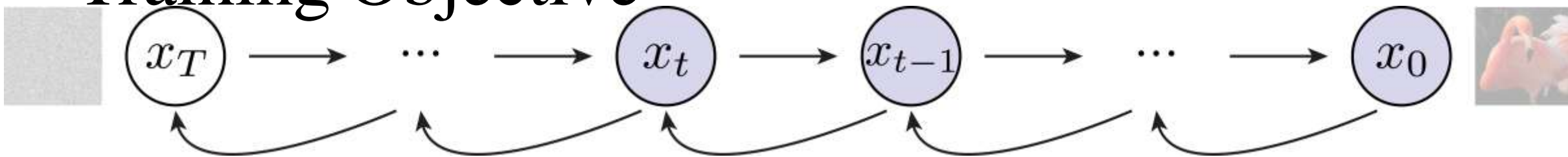
$$\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \left[w_t \|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right]$$

set as 1 (critical)

Objective	IS	FID
L , learned diagonal Σ	–	–
L , fixed isotropic Σ	7.67 ± 0.13	13.51
$\ \tilde{\epsilon} - \epsilon_\theta\ ^2$ (L_{simple})	9.46 ± 0.11	3.17

[Ho et al. 2020]; see more in [Salimans & Ho, 2022]

Training Objective



$$\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \left[\left\| w_t \left\| \epsilon - \underbrace{\epsilon_\theta}_{\text{network to predict noise}}(x_t, \underbrace{t}_{\text{conditioned on noise level (critical)}}) \right\|^2 \right]$$

Diffusion Models

- Forward process
 - add noise to data
- Reverse process
 - learn to denoise
- Training objective
 - from Hierarchical VAE to L2 loss
- Noise Conditional Network
 - represent a distribution

Noise Conditional Network

- Diffusion models decompose a distribution into **many** simpler ones.
- We need the same # networks to fit **all** of them.
- We can **combine** all into one “powerful” network.
- This network is conditioned on noise level t .
- **Noise Conditional Network** [Song & Ermon 2019]: things made work

Noise Conditional Network

How to represent $p_{\theta}(x_{t-1} | x_t)$

- network input: x_t
- network output: μ and σ of a distribution
- parametrize μ by: $\epsilon_{\theta}(x_t, t)$

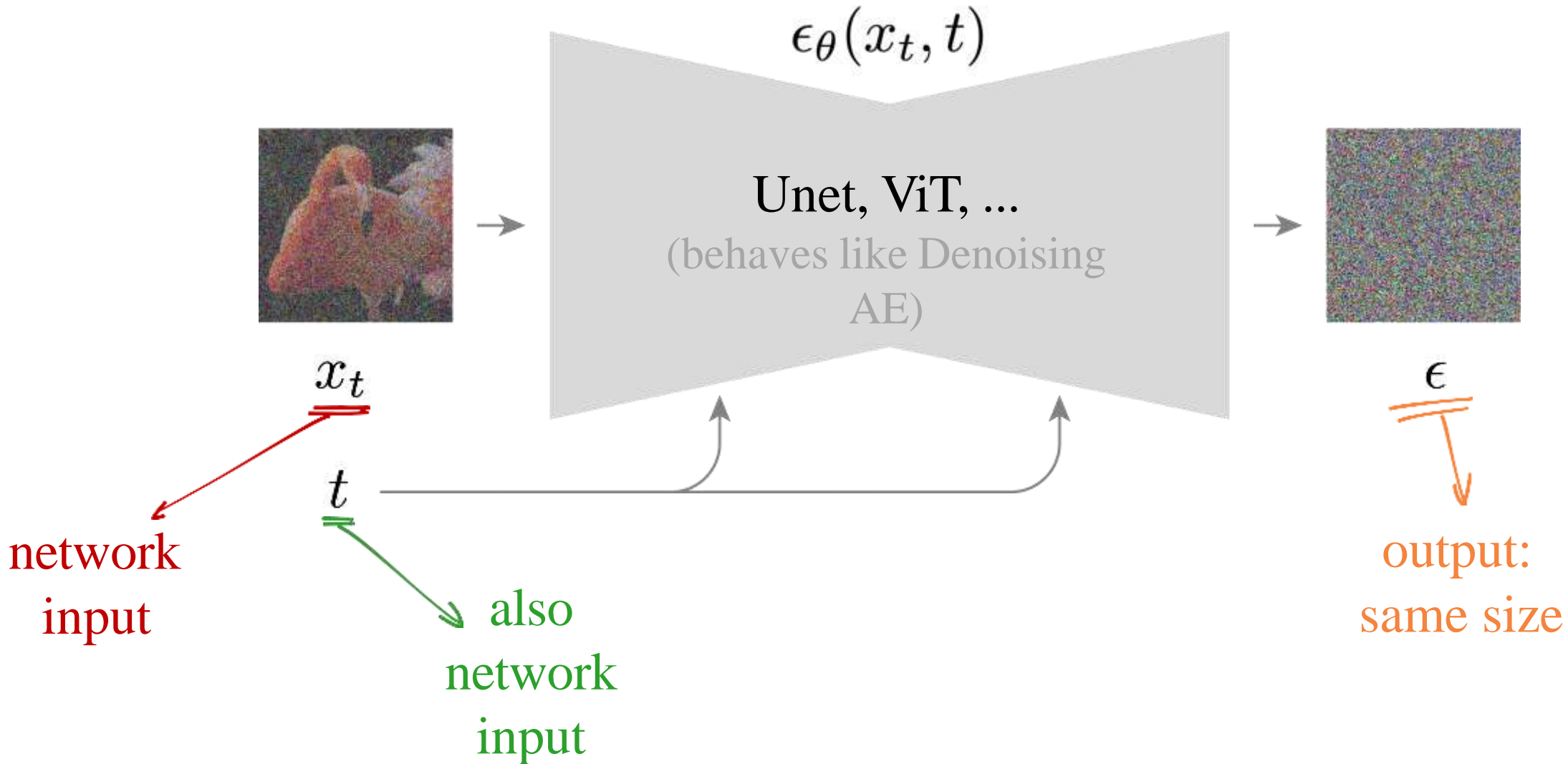
noisy image:

- condition
- network input

noise level:

- condition
- network input

Noise Conditional Network



Diffusion algorithm annotated.

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

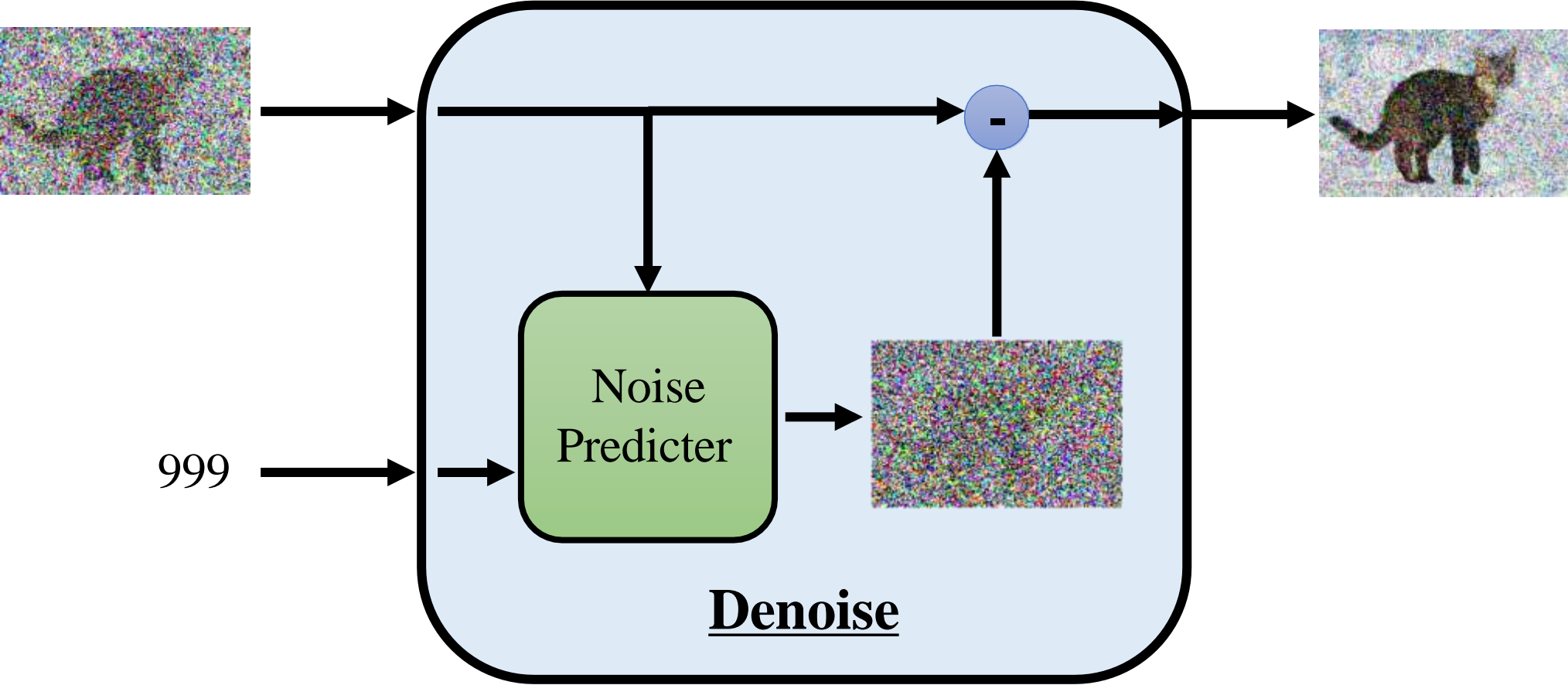
Algorithm 2 Sampling

- 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

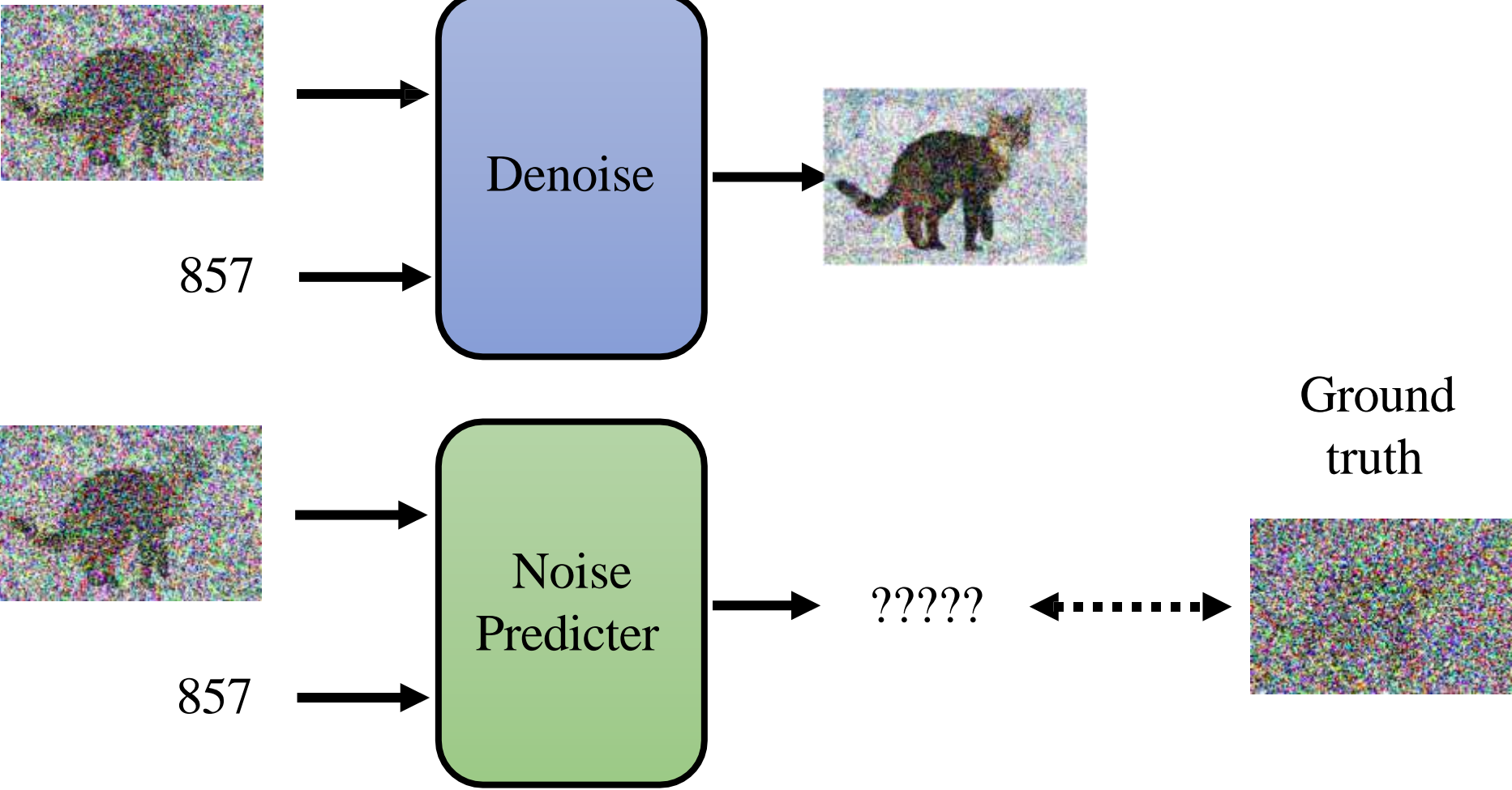
estimated μ

sampling from
estimated distribution

Noise Conditional Network



Noise Conditional Network



Diffusion algorithm annotated.

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

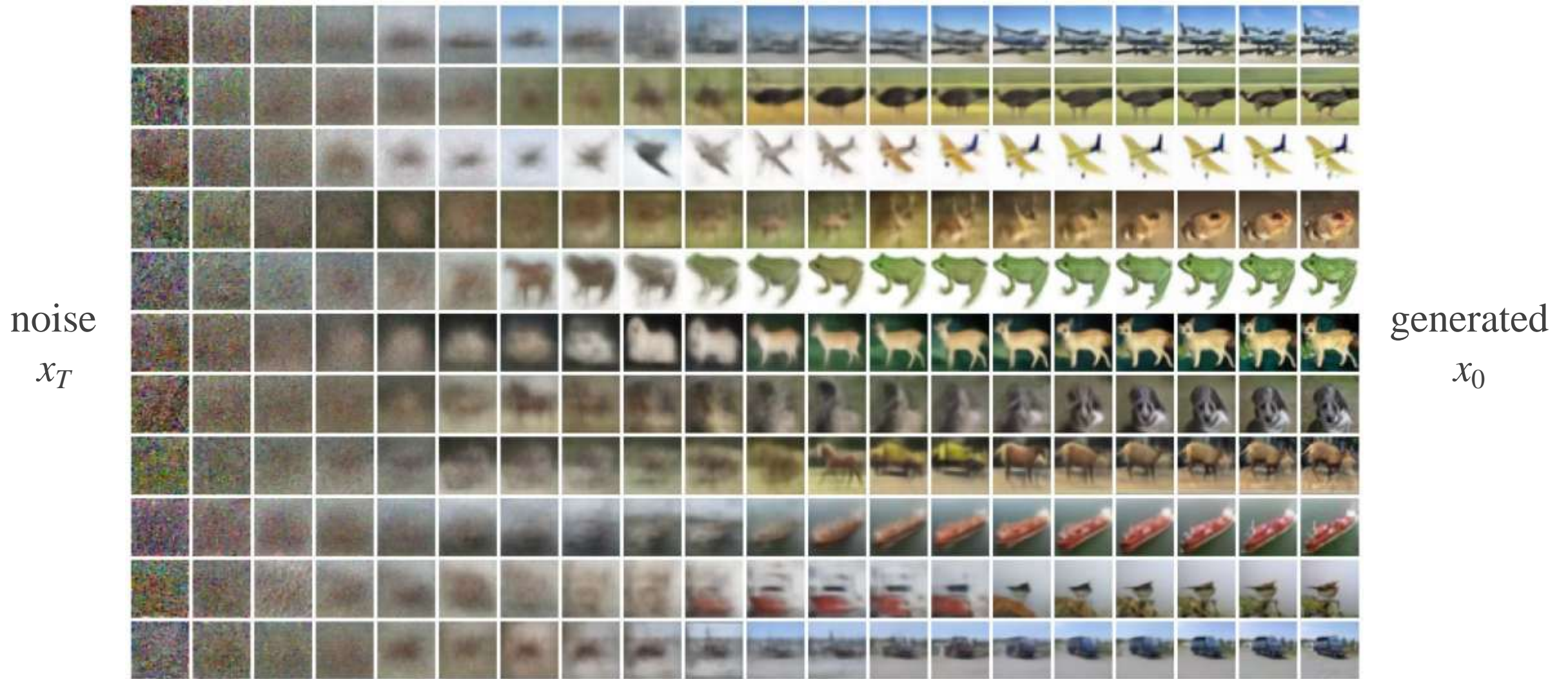
Algorithm 2 Sampling

- 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

tl; dr: noising and denoising

- Turns out to be extremely simple
- Being “simple and effective” moves the needle

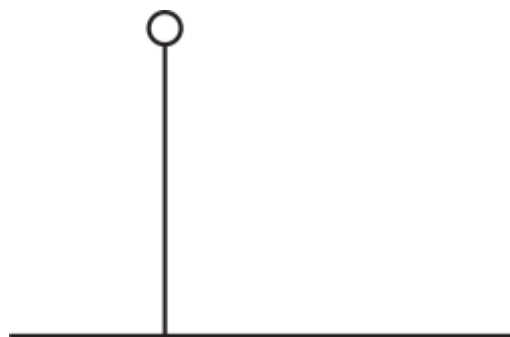
Example: Unconditional Generation on CIFAR-10



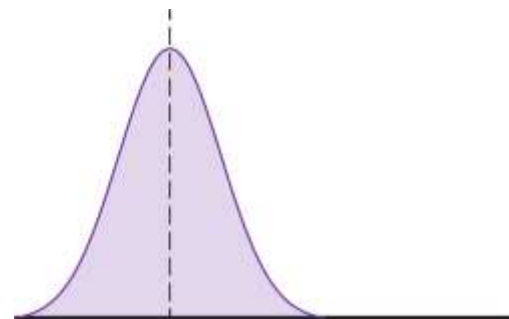
Example: shared intermediate latents



Recap.
• Adding Gaussian noise \Leftrightarrow sampling $x \sim \mathcal{N}(x \mid x_0, \sigma)$



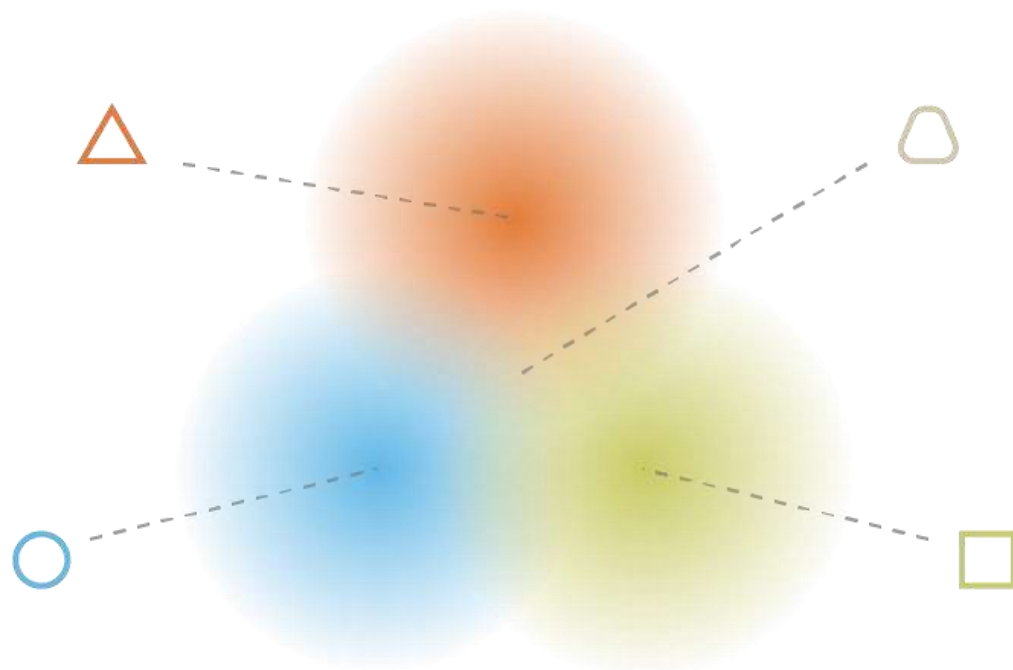
$$p(x) = \delta(x - x_0)$$



$$p(x) = \mathcal{N}(x \mid x_0, \sigma)$$



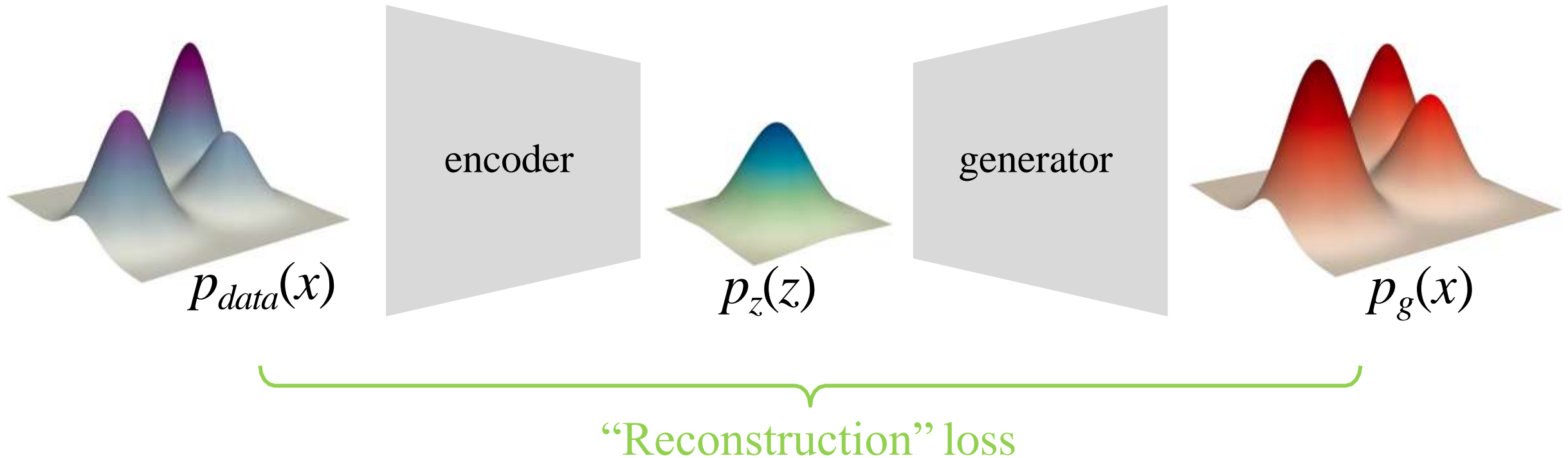
Problems



Recap: Variational Autoencoder (VAE)

Autoencoding distributions:

“Encoding” data distribution p_{data} into latent distribution p_z



Summary

Forward process

- add noise to data

Reverse process

- learn to denoise

Training objective

- from Hierarchical VAE to L2 loss

Noise Conditional Network

- represent distributions