

Diffusion Probabilistic Models

Bartek Sobieski

MI2.AI Research Seminar – Winter 2023/24

Outline

- 1. What is a generative model?**
- 2. Background knowledge**
- 3. Diffusion models - three equivalent interpretations**
- 4. Guidance - conditional generation**
- 5. DDIM - hidden gems in trained diffusion models**
- 6. Applications:**
 - a. Diffusion autoencoders**
 - b. Latent diffusion models**
 - c. Stable diffusion**
- 7. Summary**

Jascha Sohl-Dickstein



J Sohl-Dickstein, EA Weiss, N Maheswaranathan, S Ganguli.
Deep unsupervised learning using nonequilibrium thermodynamics.
International Conference on Machine Learning (2015).



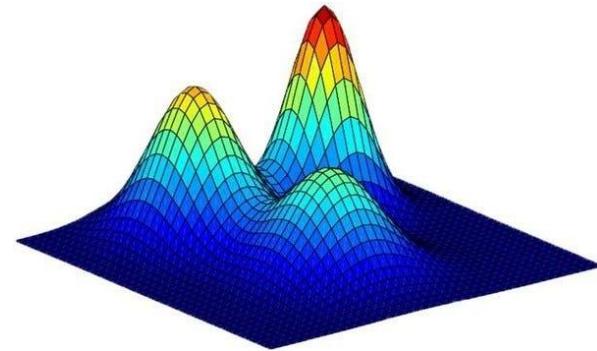
Jascha is a principal scientist in Google DeepMind. He is most (in)famous for [inventing diffusion models](#). His recent work has focused on theory of [overparameterized neural networks](#), meta-training of [learned optimizers](#), and [understanding the capabilities of large language models](#). Before working at Google, Jascha was a visiting scholar in [Surya Ganguli's lab](#) at Stanford University, and an academic resident at the [Khan Academy](#) education nonprofit. He earned his PhD in 2012 in the [Redwood Center for Theoretical Neuroscience](#) at UC Berkeley, in [Bruno Olshausen's](#) lab. Prior to his PhD, he worked [sending rovers to Mars](#).

What is a generative model?

Generative models



Observed samples
 \mathbf{x}



True distribution
 $p(\mathbf{x})$

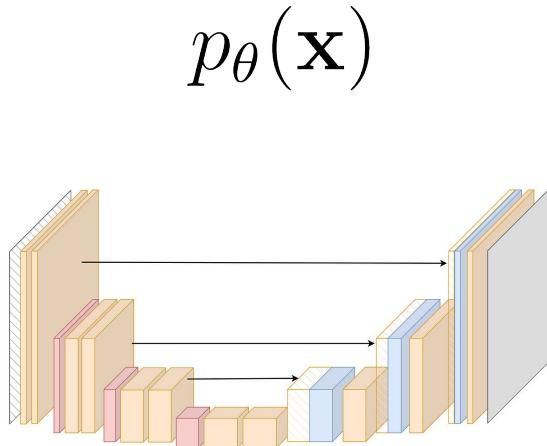
Generative models learn to *model* the true distribution using the observed samples

Generative models

Approximate
true distribution



Synthetic samples

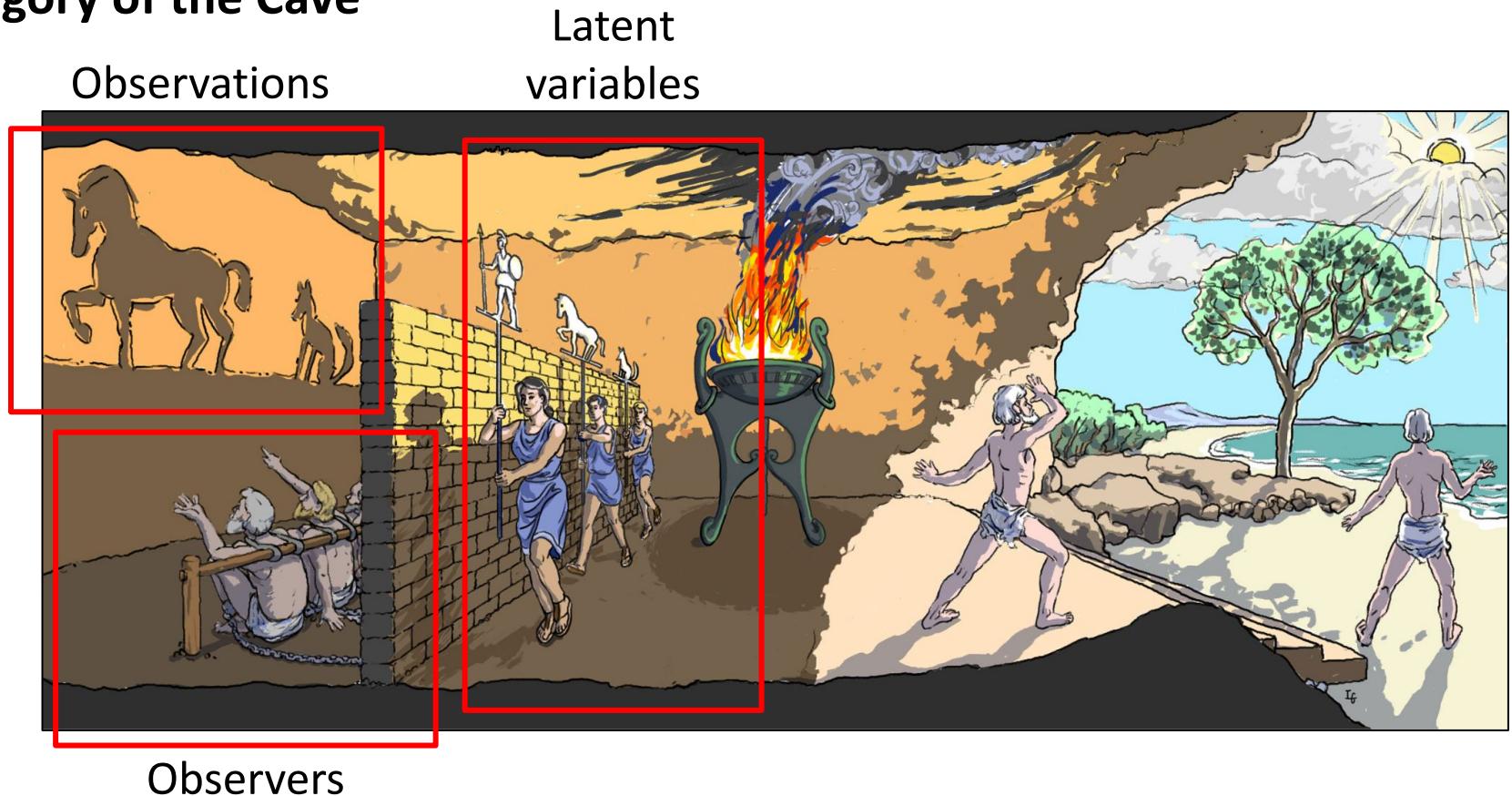


Background knowledge

Allegory of the Cave



Allegory of the Cave



Evidence Lower Bound

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

Bayes rule

$$p(\mathbf{x}, \mathbf{z}) \xrightarrow{\text{Marginalizing}} p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

Evidence Lower Bound

$$\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

ELBO

Evidence Lower Bound

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

ELBO

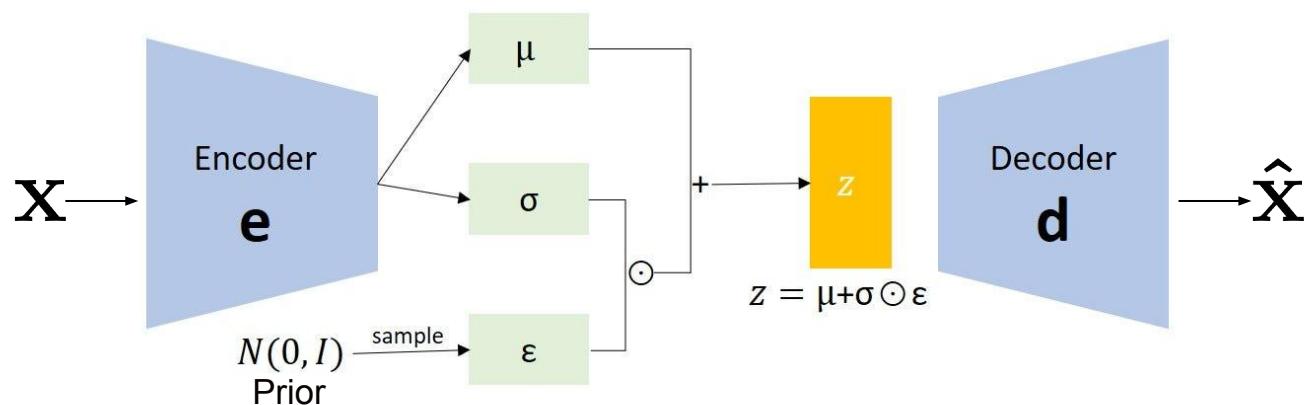
Evidence Lower Bound

$$\begin{aligned}\log p(\mathbf{x}) &= \mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] + \mathbf{D}_{\mathbf{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x})) \\ &\geq \mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]\end{aligned}$$

Variational autoencoders

$$\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] = \mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))$$

ELBO
Decoder
Encoder
Prior



Variational autoencoders

- Latent dimension is smaller than data dimension
- Encoder is learned and the encoding process consists of a single forward pass through it
- Latent variables come from a standard normal distribution

Diffusion models

Diffusion models

- Latent dimension is **exactly equal to** data dimension
- Encoder is **not learned** and the encoding process consists of **many steps**
- Latent variables come from a standard normal distribution

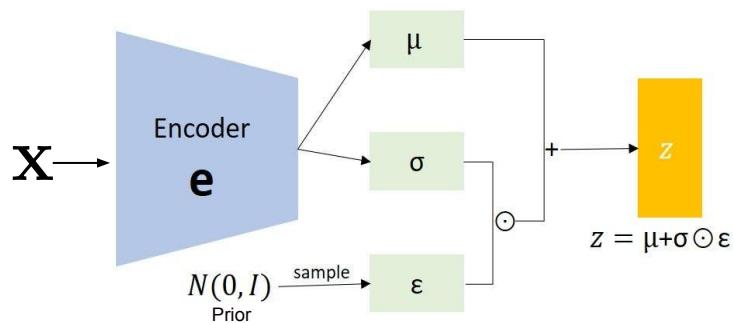
Variational autoencoders

- Latent dimension is **smaller than** data dimension
- Encoder **is learned** and the encoding process consists of a **single forward pass** through it
- Latent variables come from a standard normal distribution

Encoding in diffusion models

Variational autoencoders

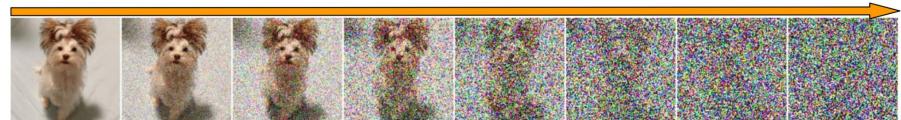
$$q_{\phi}(\mathbf{z}|\mathbf{x})$$



Diffusion models

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

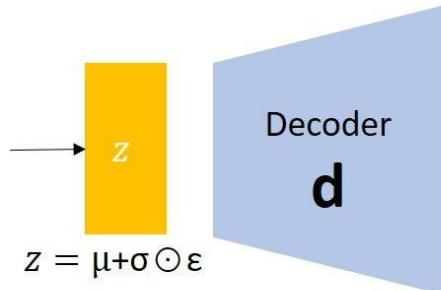
$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$



Decoding in diffusion models

Variational autoencoders

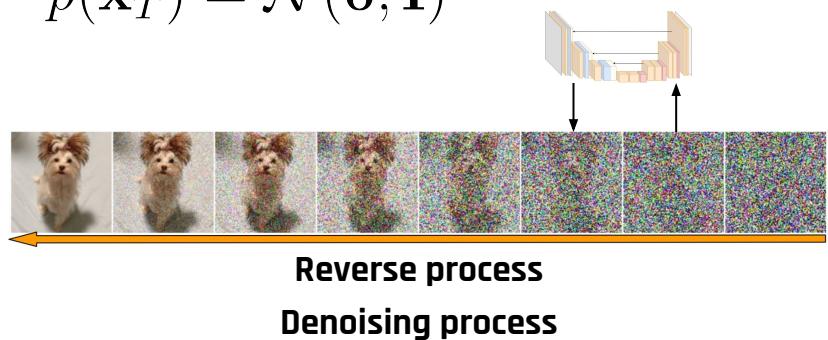
$$p_{\theta}(\mathbf{x}|\mathbf{z})$$



Diffusion models

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

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



ELBO in diffusion models

ELBO

$$\mathbf{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] = \mathbf{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] - \mathbf{D}_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T)) - \sum_{t=2}^T \mathbf{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\mathbf{D}_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))] = 0$$

Reconstruction term**Prior matching term****Denoising matching term**

Three equivalent interpretations

Denoising matching term objective

Approximate
reverse process step

$$\operatorname{argmin}_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

Ground truth
reverse process step

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \propto \mathcal{N} \left(\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I} \right)$$
$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$
$$\mu_q(\mathbf{x}_t, \mathbf{x}_0)$$
$$\sigma_q^2(t)$$
$$\Sigma_q(t)$$

Predicting the mean

$$\operatorname{argmin}_{\theta} \mathbf{D}_{\mathbf{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \operatorname{argmin}_{\theta} \frac{1}{2\sigma_q^2(t)} [\|\mu_{\theta} - \mu_q\|_2^2]$$

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}$$

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$$

Predicting the mean via approximate clean image

|

$$\begin{aligned} \operatorname{argmin}_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \\ = \operatorname{argmin}_{\theta} \frac{1}{2\sigma_q^2(t)} [\|\mu_{\theta} - \mu_q\|_2^2] \\ = \operatorname{argmin}_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)^2} [\|\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2] \end{aligned}$$

Reparameterization trick

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_0 \quad \boldsymbol{\epsilon}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0}{\sqrt{\bar{\alpha}_t}}$$

Predicting the mean via approximate source noise

$$\operatorname{argmin}_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

$$= \operatorname{argmin}_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} [\|\boldsymbol{\epsilon}_0 - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)\|_2^2]$$

Tweedie's formula

For a random variable $z \sim \mathcal{N}(\mu_z, \Sigma_z)$

$$\mathbf{E} [\mu_z | z] = z + \Sigma_z \nabla_z \log p(z)$$

**Score function**

Predicting the mean via approximate score function

$$\operatorname{argmin}_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

$$= \operatorname{argmin}_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{\alpha_t} [\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla \log p(\mathbf{x}_t)\|_2^2]$$

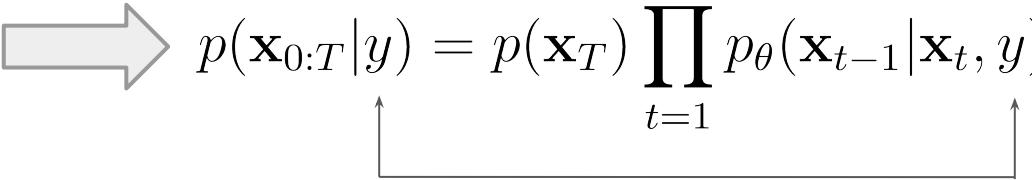
$$\nabla \log p(\mathbf{x}_t) = -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_0$$

Guidance - conditional generation

Guidance

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

Conditioning signal



Classifier guidance

**Unconditional
score function**

$$\nabla \log p(\mathbf{x}_t | y) = \boxed{\nabla \log p(\mathbf{x}_t)} + \nabla \log p(y | \mathbf{x}_t)$$

**Adversarial
gradient**

Classifier guidance

$$\nabla \log p(\mathbf{x}_t | y) = \boxed{\nabla \log p(\mathbf{x}_t)} + \gamma \boxed{\nabla \log p(y | \mathbf{x}_t)}$$

Unconditional score function

Adversarial gradient

Classifier-free guidance

$$\nabla \log p(\mathbf{x}_t|y) = \gamma \boxed{\nabla \log p(\mathbf{x}_t|y)} + (1 - \gamma) \boxed{\nabla \log p(\mathbf{x}_t)}$$

**Conditional
score function**

**Unconditional
score function**

DDIM - hidden gems in trained diffusion models

Non-markovian forward processes

$$q_\sigma(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I} \right)$$


Parameterized by σ

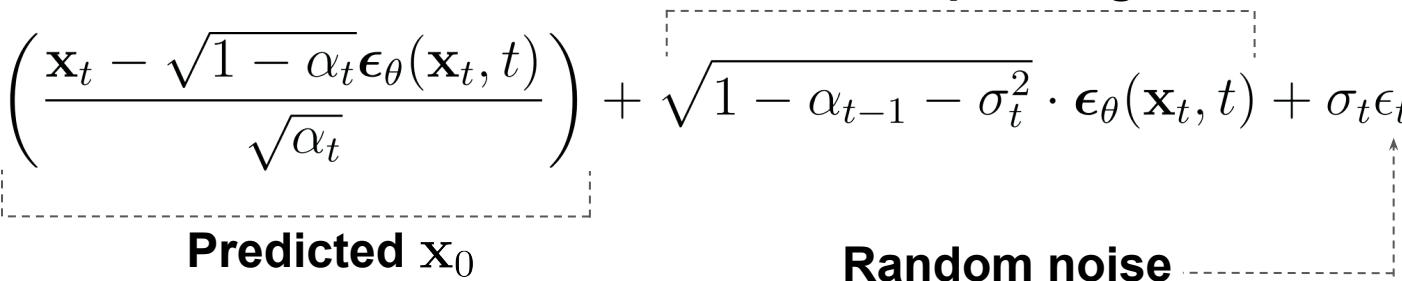
Denoising Diffusion Implicit Models

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\mathbf{x}_t, t) + \sigma_t \epsilon_t$$

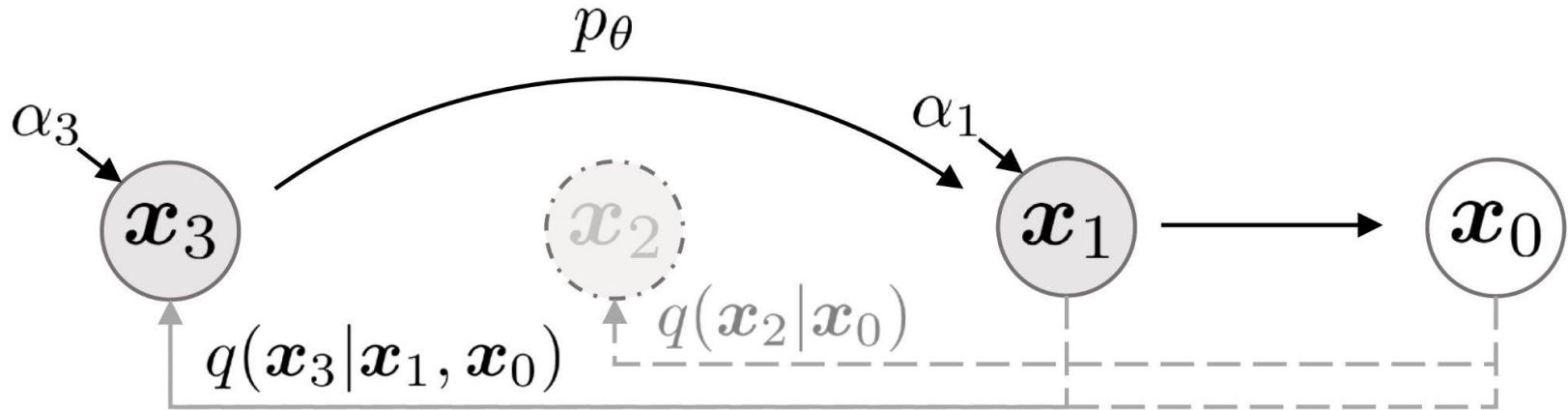
Predicted \mathbf{x}_0

Direction pointing to \mathbf{x}_t

Random noise



Accelerated generation



Connection to neural ODEs

$$\frac{d\bar{\mathbf{x}}(t)}{dt} = \frac{d\sigma(t)}{dt} \boldsymbol{\epsilon}_{\theta} \left(\frac{\bar{\mathbf{x}}(t)}{\sqrt{\sigma^2(t) + 1}}, t \right)$$

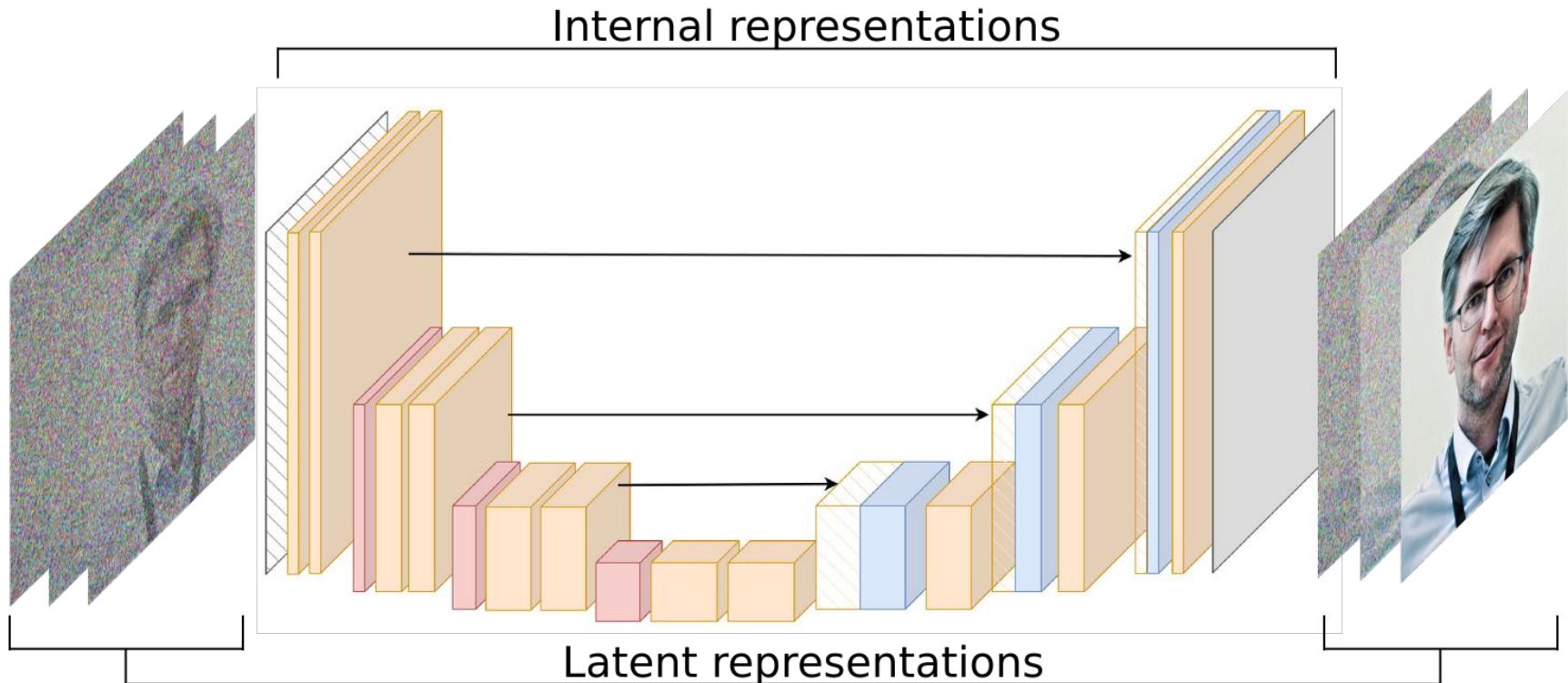
$$\bar{\mathbf{x}}(t) = \frac{\mathbf{x}(t)}{\sqrt{\alpha(t)}}$$

$$\sigma(t) = \sqrt{\frac{1 - \alpha(t)}{\alpha(t)}}$$

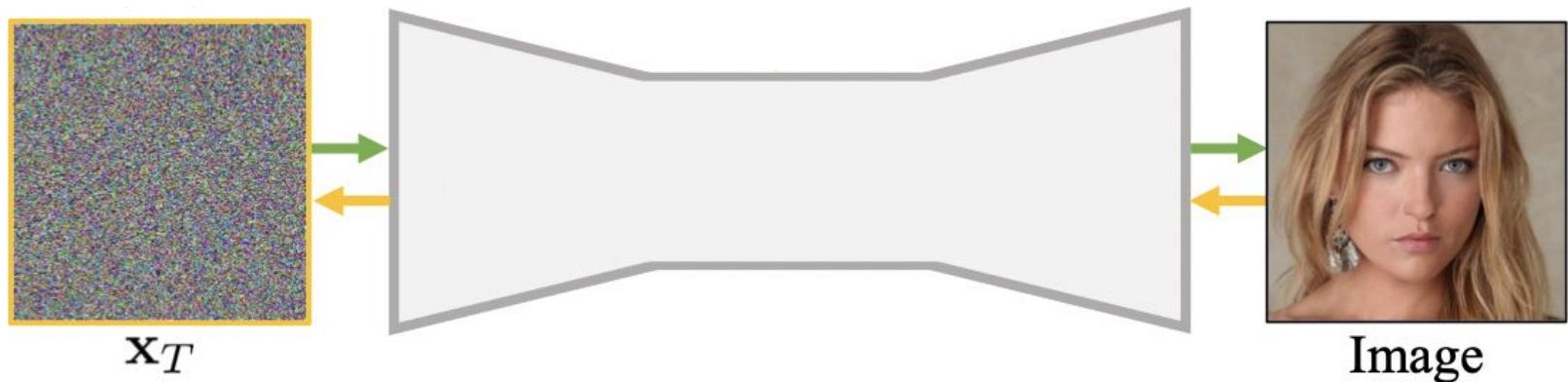
Applications

Enhancing the latent space

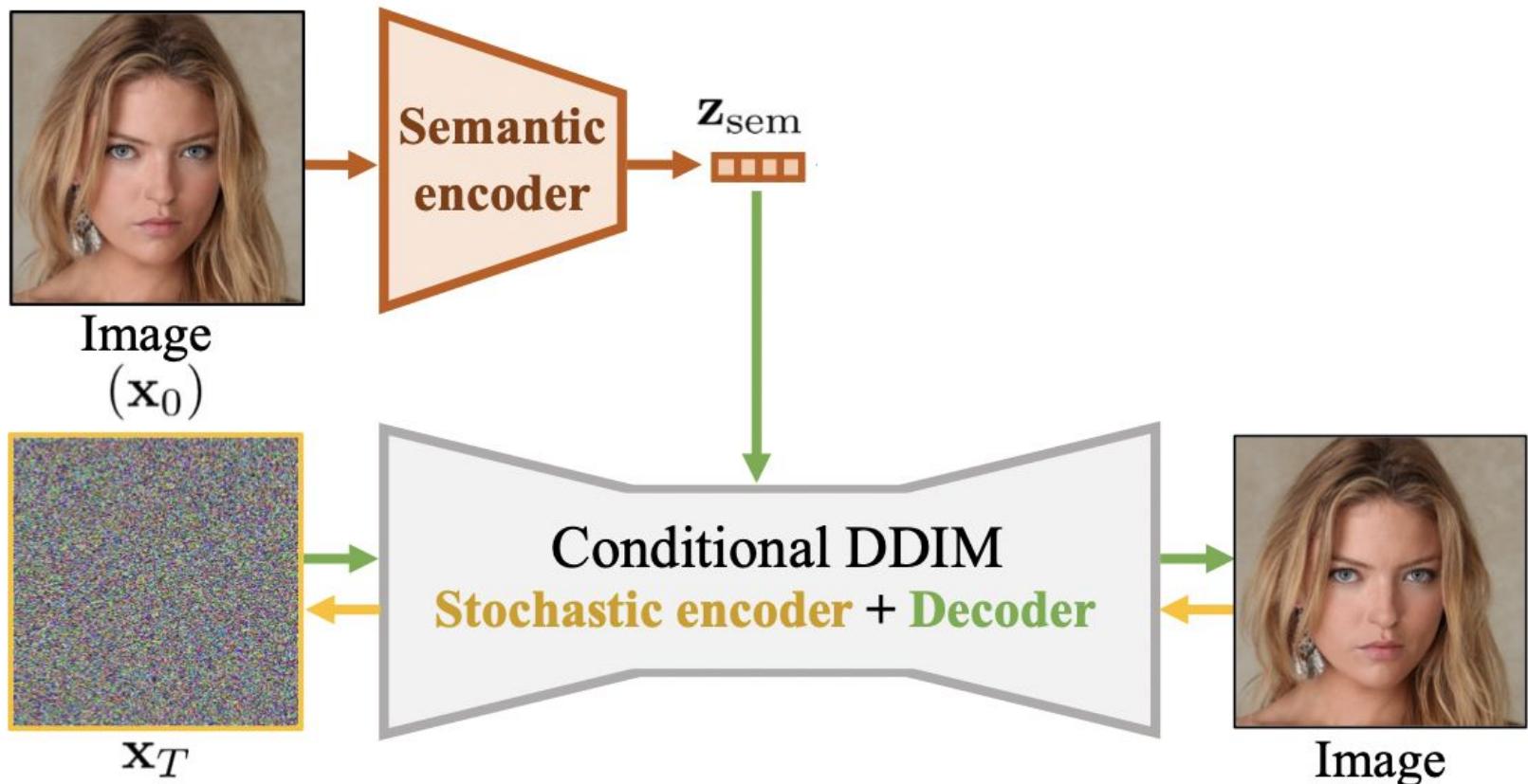
Enhancing the latent space



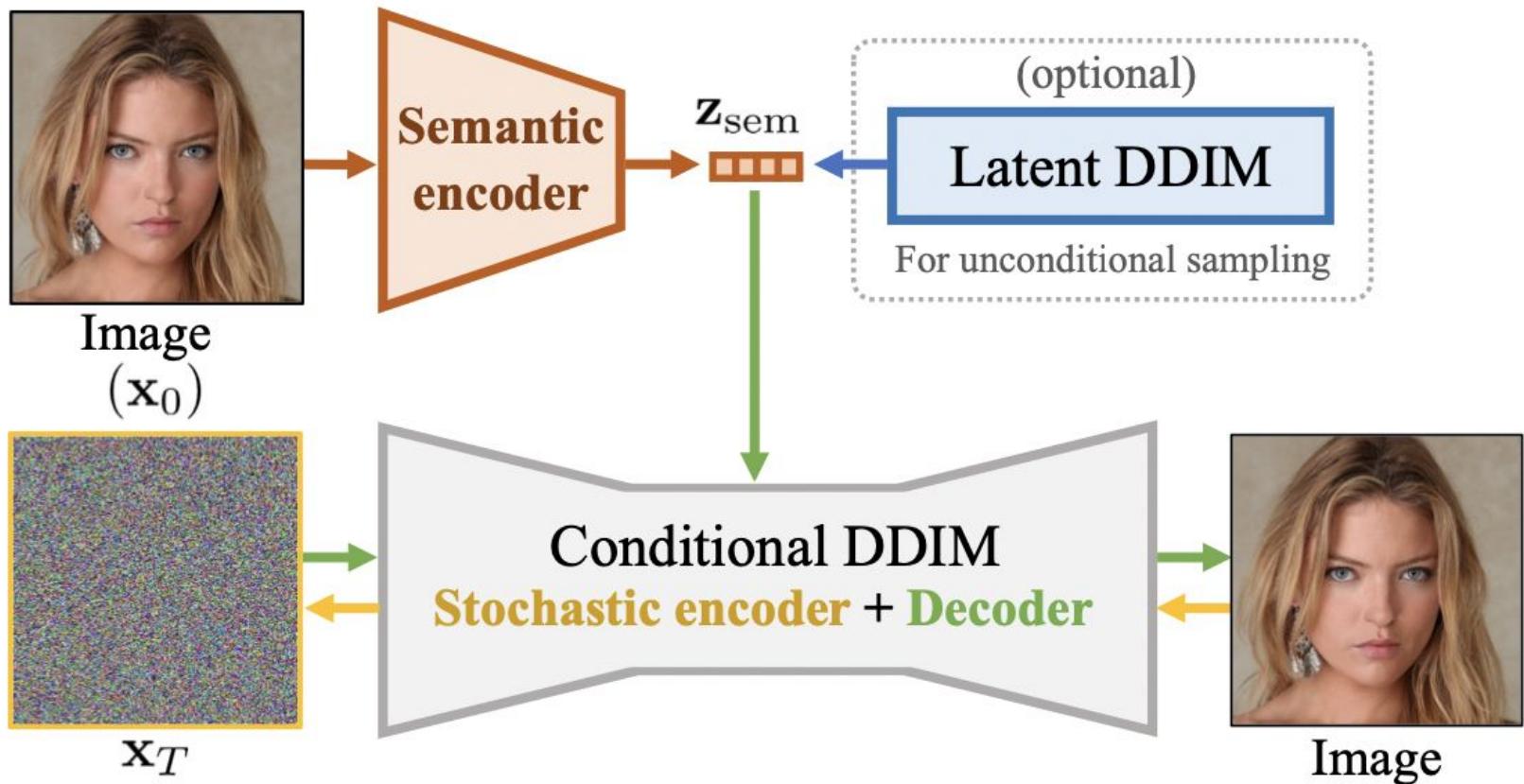
Enhancing the latent space



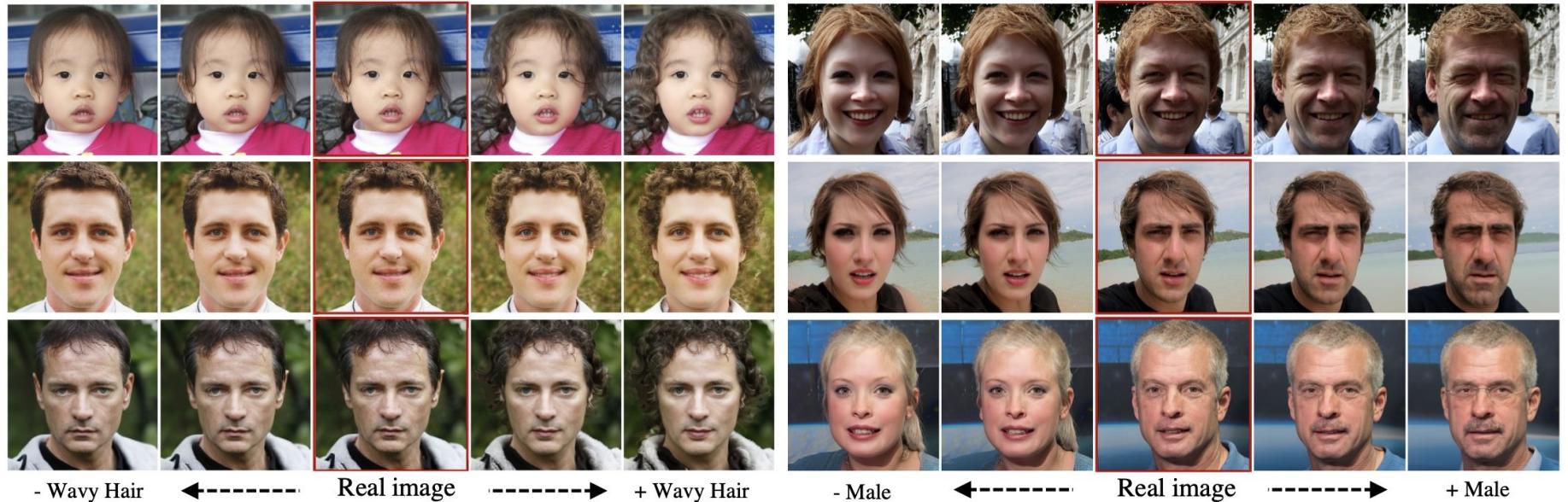
Diffusion Autoencoders



Diffusion Autoencoders

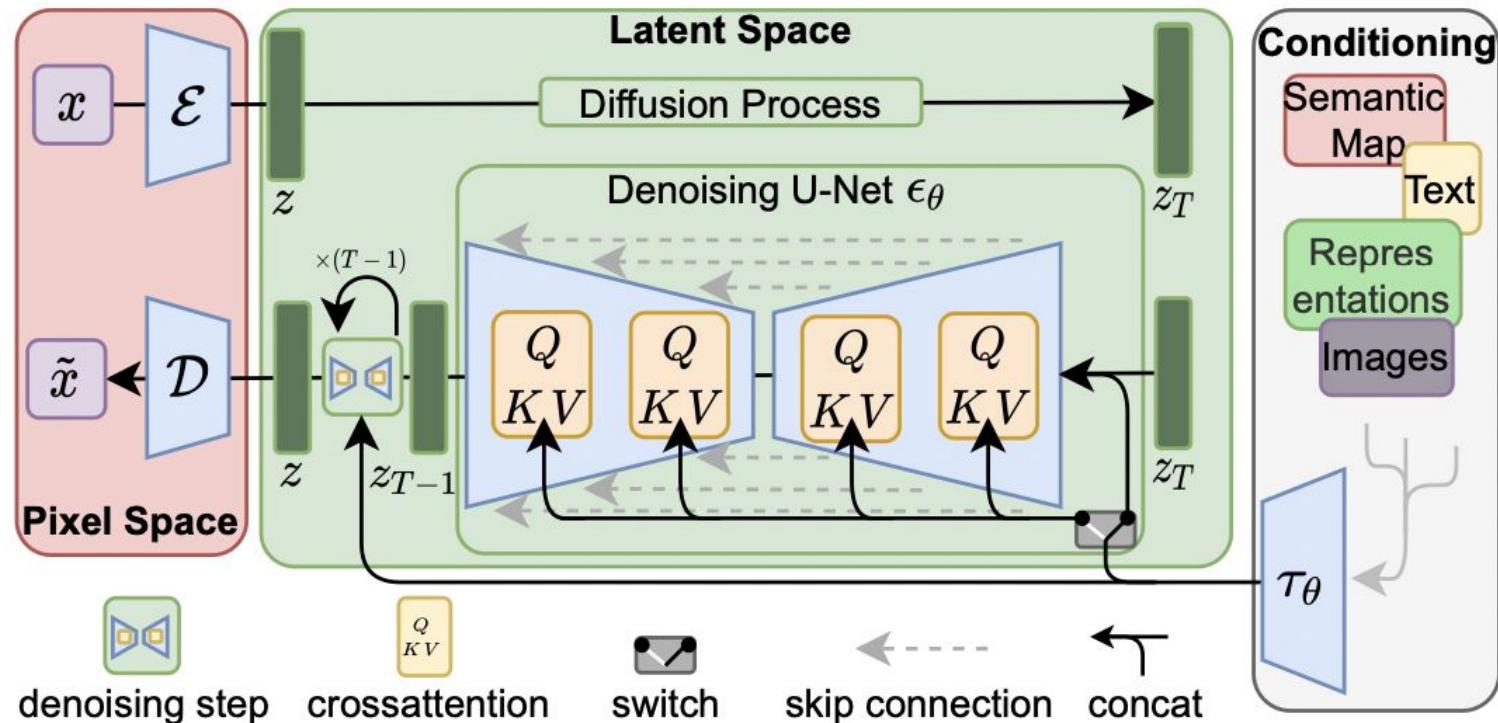


Intriguing properties



Improving computational efficiency

Latent diffusion models



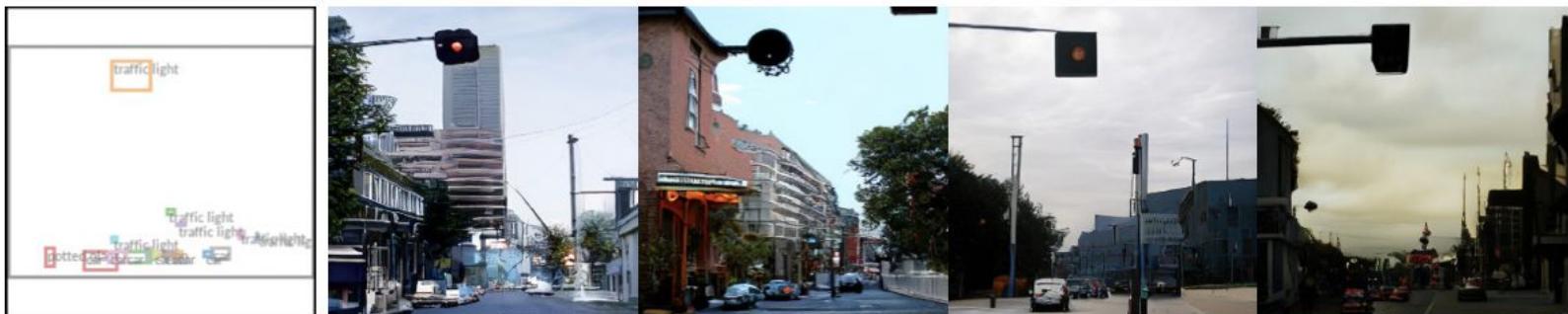
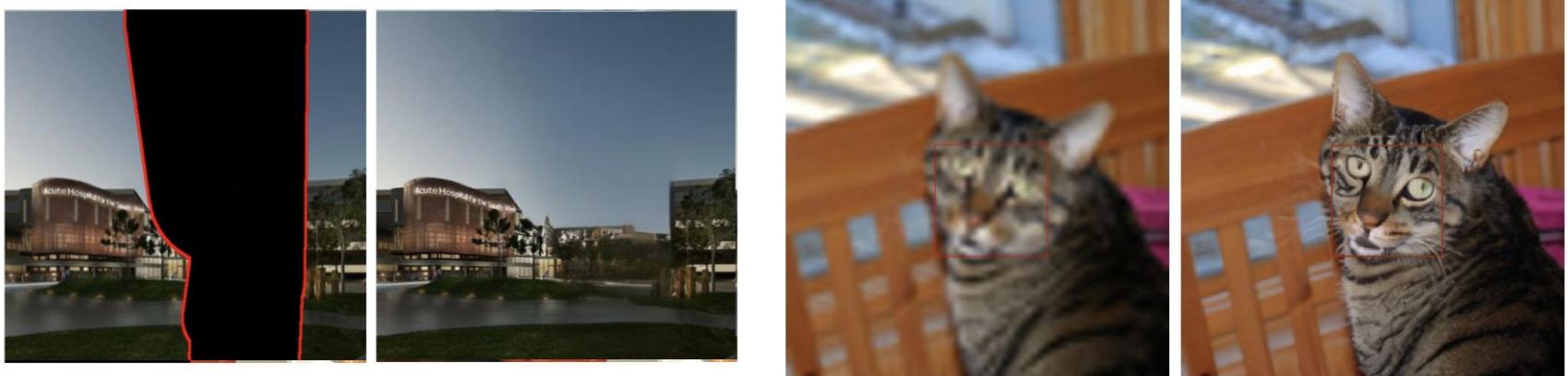
Latent diffusion models



'An oil painting of a latent space.'



Latent diffusion models



Questions

References

1. Calvin Luo, *Understanding Diffusion Models: A Unified Perspective*, 2022
2. Dhariwal, Prafulla, and Alexander Nichol, Diffusion models beat gans on image synthesis, NeurIPS 2021
3. Ho, Jonathan, and Tim Salimans, *Classifier-free diffusion guidance*, 2022
4. Song, Jiaming, Chenlin Meng, and Stefano Ermon, *Denoising diffusion implicit models*, ICLR 2021
5. Preechakul, Konpat, et al., *Diffusion autoencoders: Toward a meaningful and decodable representation*, CVPR 2022
6. Rombach, Robin, et al, *High-resolution image synthesis with latent diffusion models*, CVPR 2022.

**Thank you
for your attention**