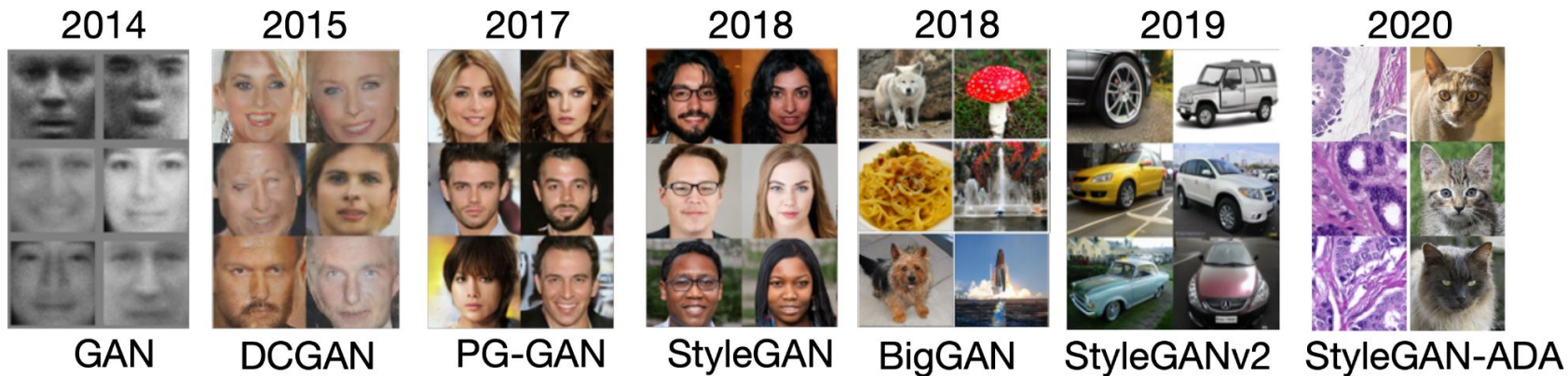


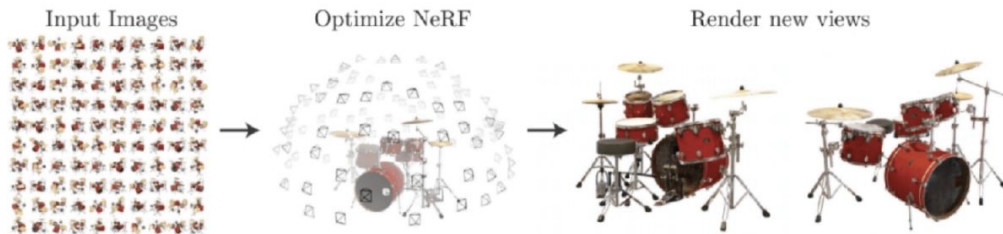
1. GANs and StyleGANs

2. AutoEncoders

Progress for Image Generation



2020: NeRF (Neural Radiance Fields)



2021: OpenAI DALLE (VQ-VAE) Arm chair in shape of avocado



Progress for Image Generation

2020? 2022?

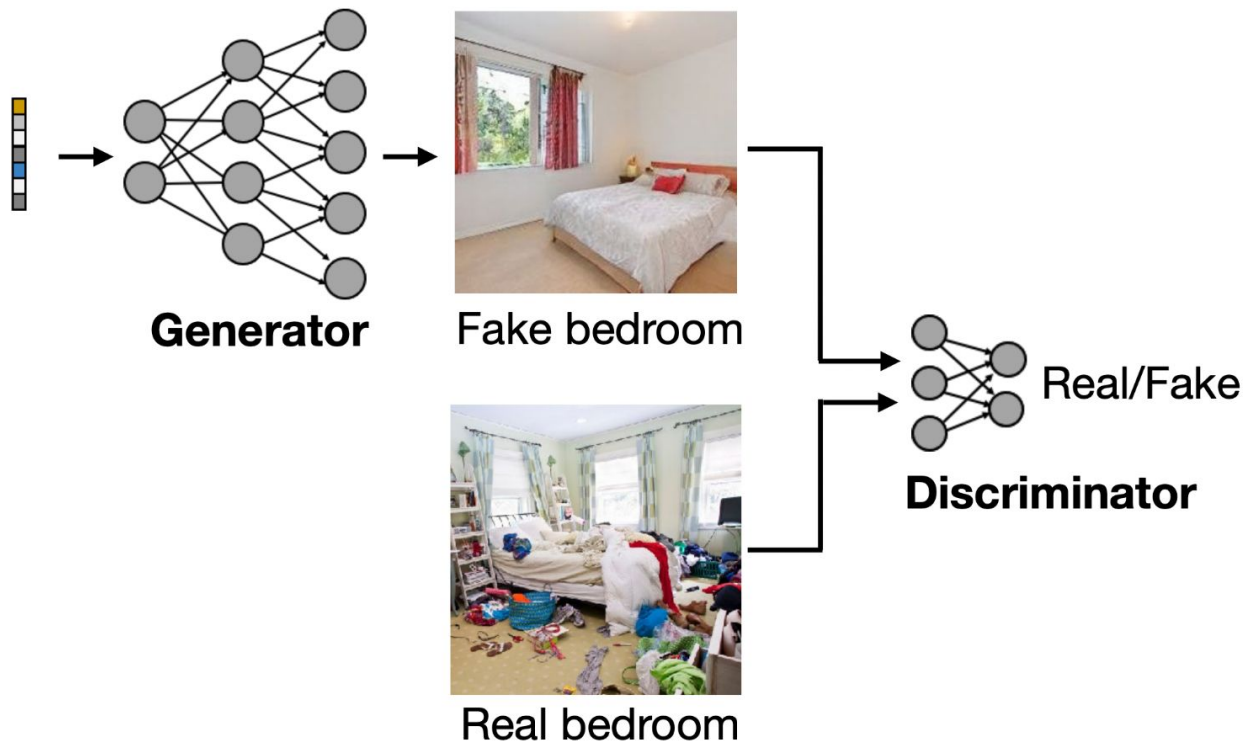
An astronaut riding a horse in a photorealistic style.



Diffusion Models

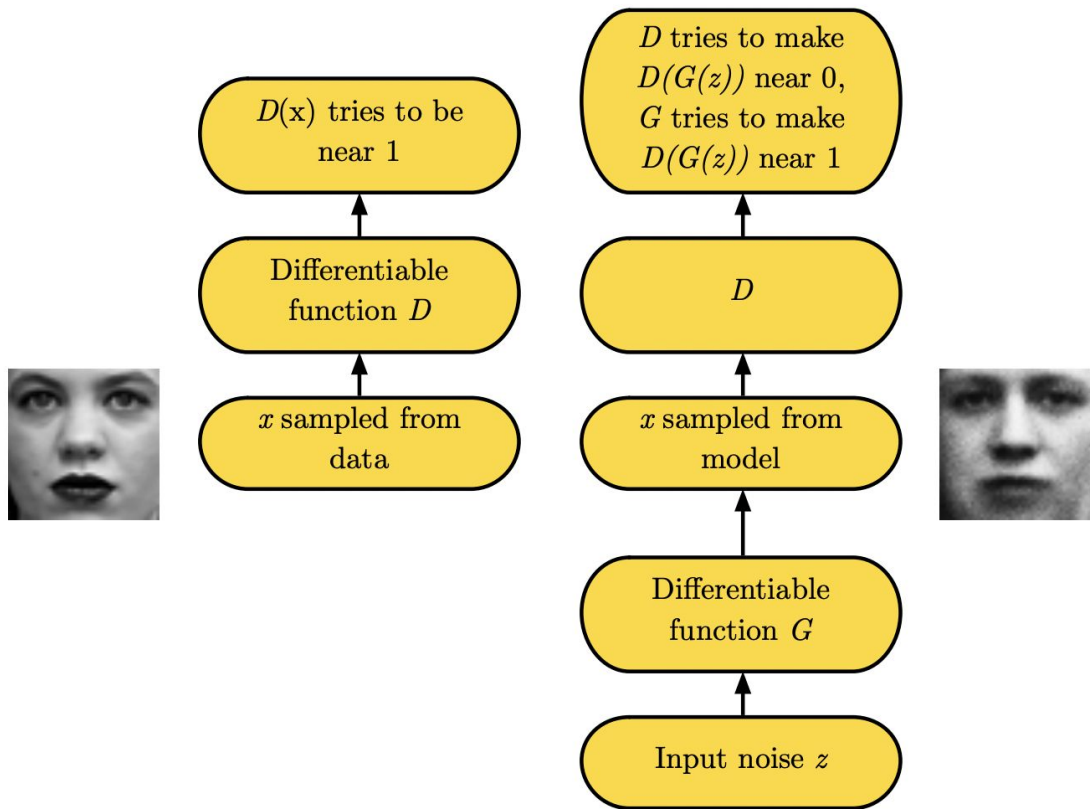
OpenAI DALL-E2 <https://openai.com/dall-e-2/>
Google Imagen <https://imagen.research.google/>
Meta Make-A-Video <https://makeavideo.studio/>
GPT-4 ???

Generative Adversarial Networks (GANs)

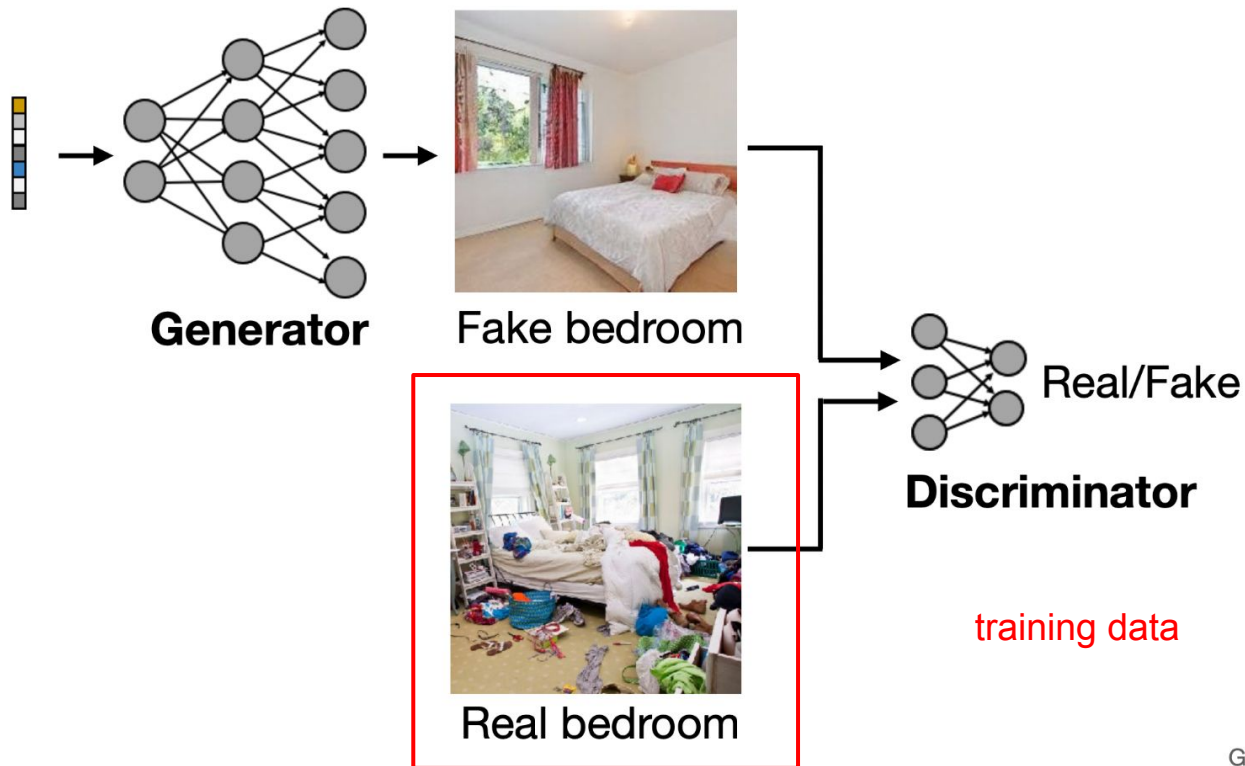


Generative Adversarial Networks (GANs)

- The basic idea of GANs is to set up a game between two players.
 - Generator
 - Creates samples that are intended to come from the same distribution as the training data
 - The counterfeiter: Trained to fool the discriminator
 - Discriminator
 - Examines samples to determine whether they are real or fake
 - The police: Trained to distinguish between the generated or the real (training data)
- Formally, GANs are a structured probabilistic model containing latent variables \mathbf{z} and observed variables \mathbf{x} .

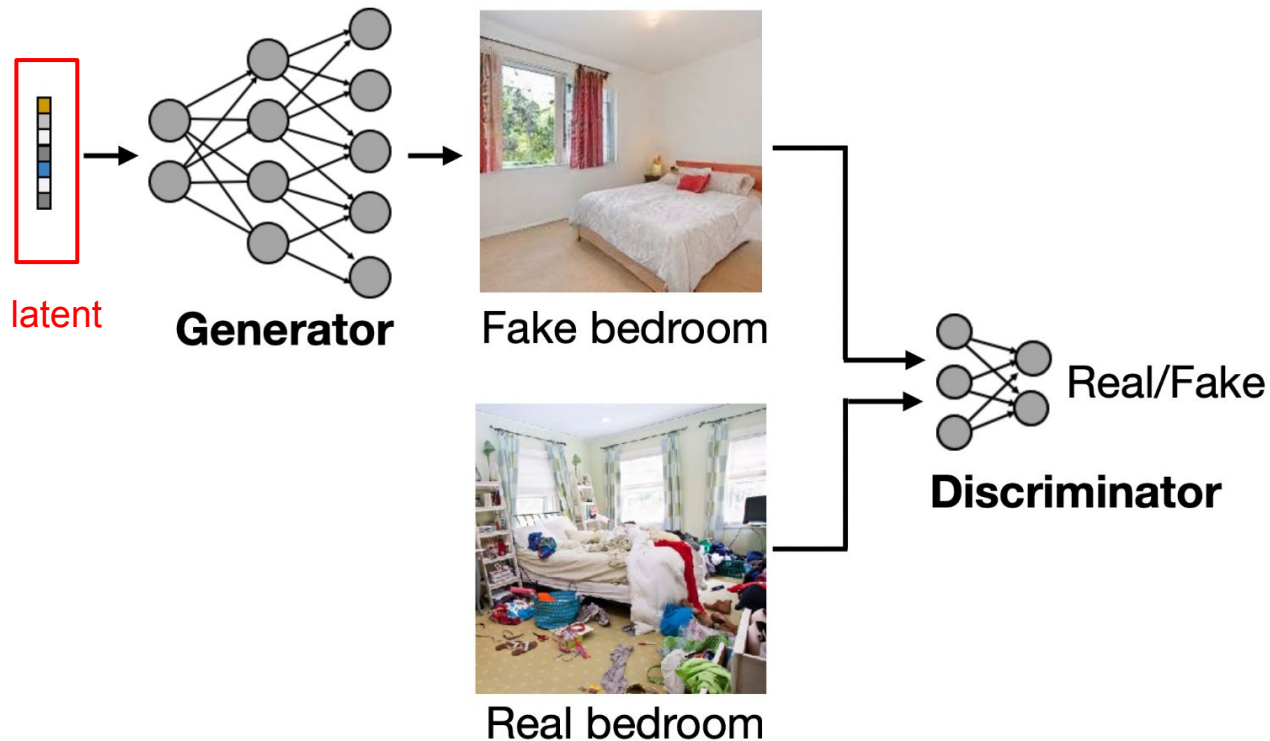


Generative Adversarial Networks (GANs)

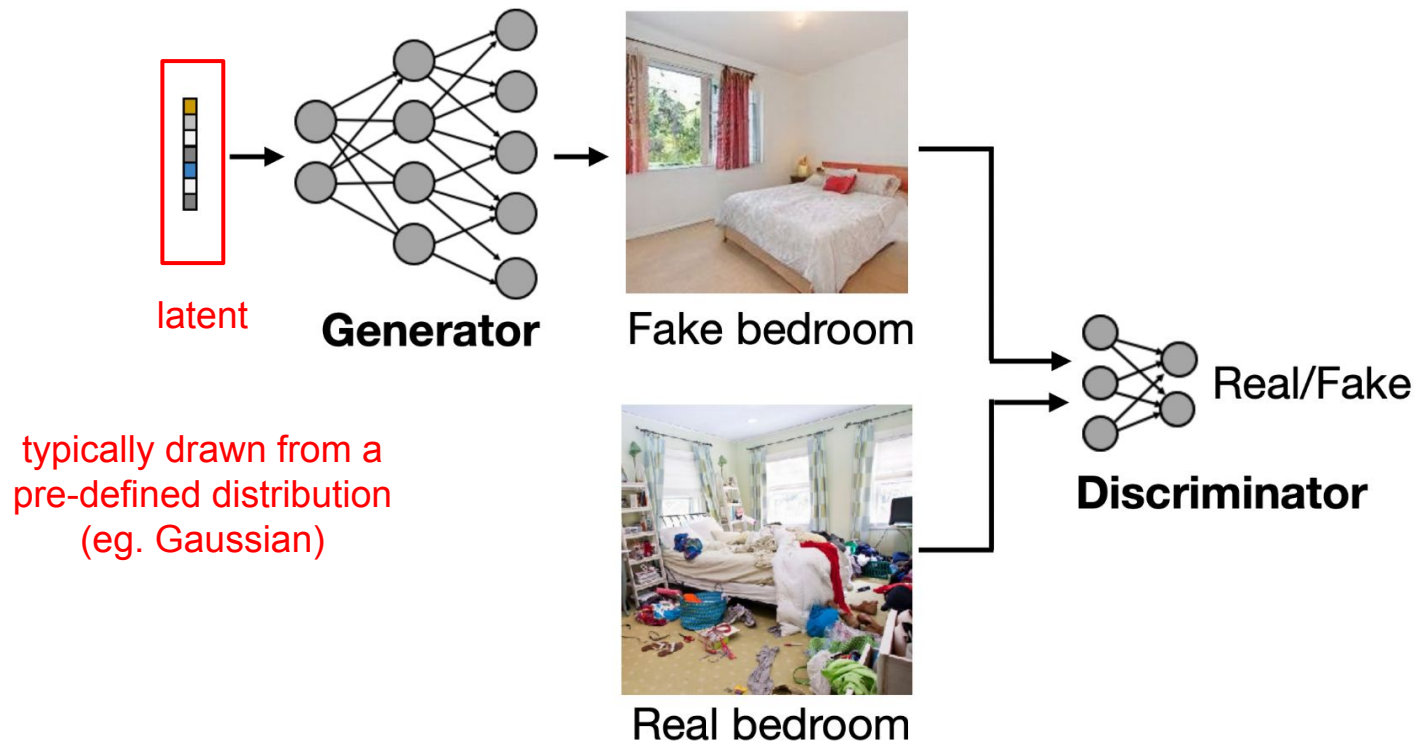


Goodfellow et al. NeurIPS'14

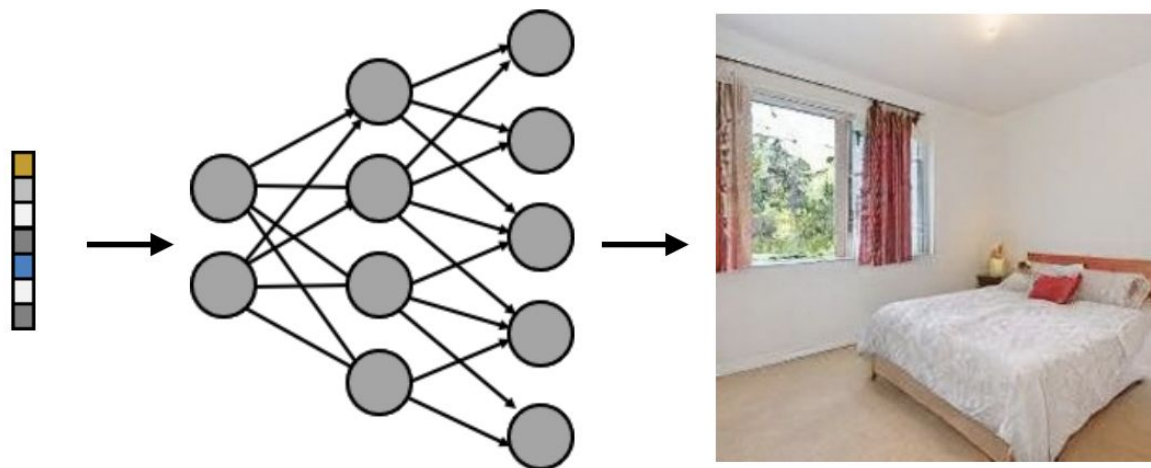
Generative Adversarial Networks (GANs)



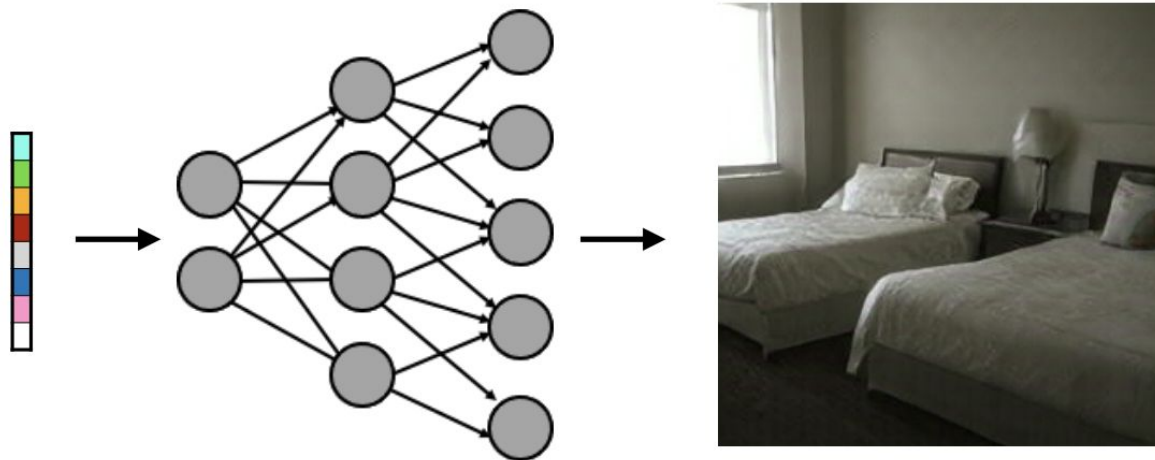
Generative Adversarial Networks (GANs)



Generative Adversarial Networks (GANs)

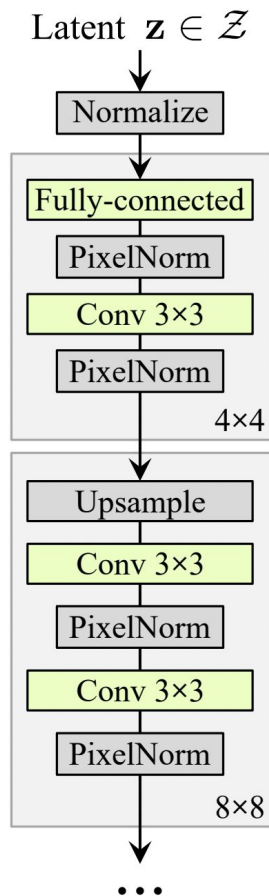


Generative Adversarial Networks (GANs)

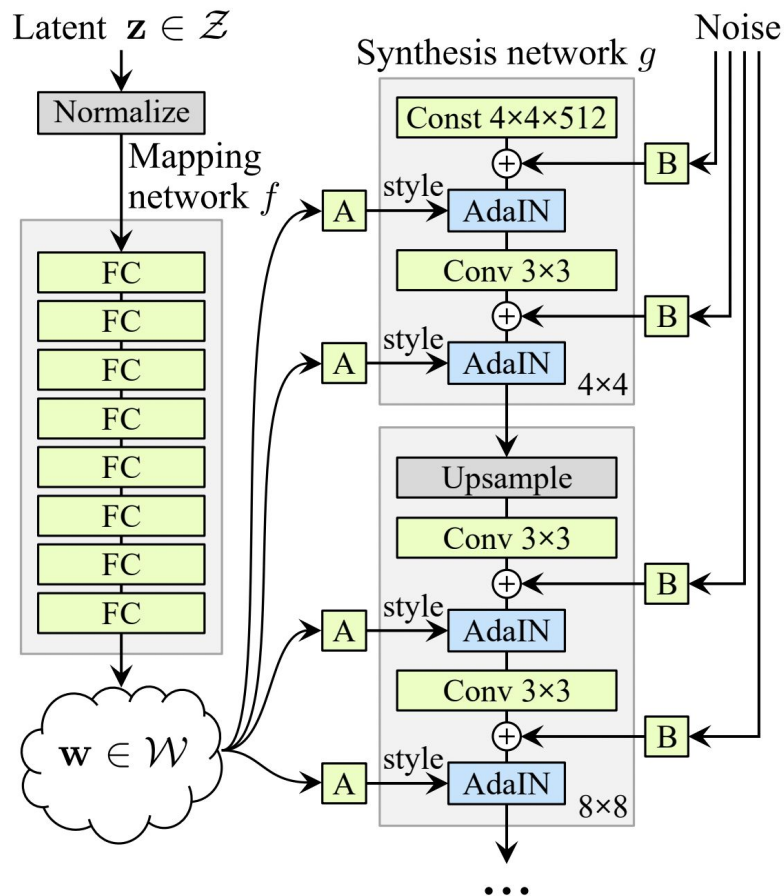


StyleGANs

Generator



(a) Traditional



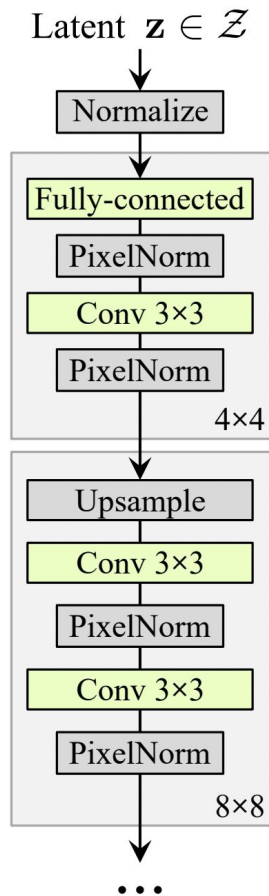
(b) Style-based generator

StyleGANs

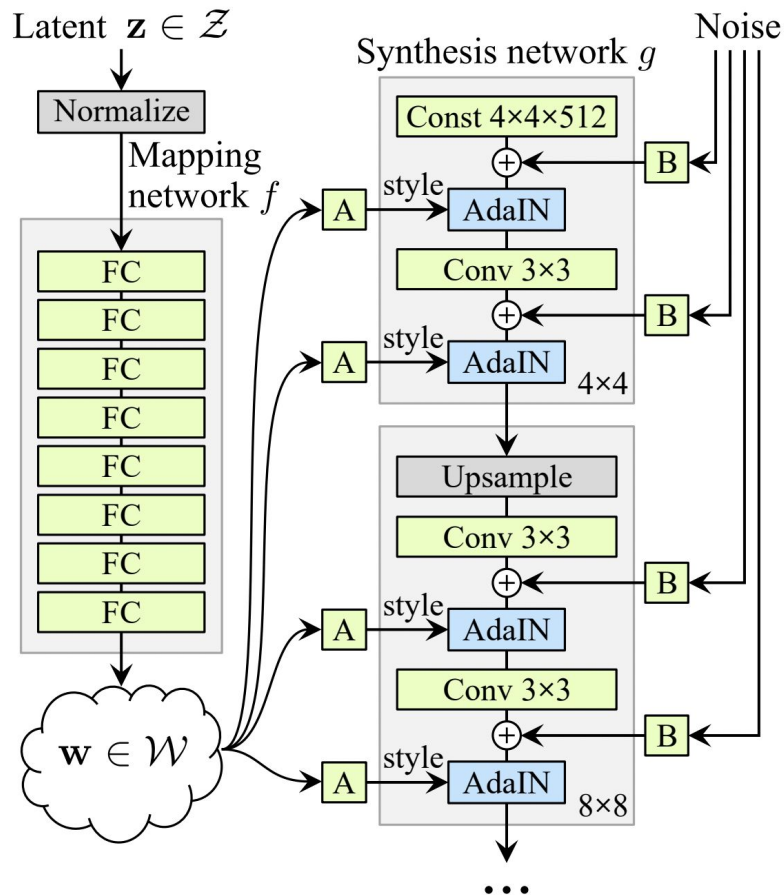
AdaIN

adaptive instance normalization

$$\text{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i},$$



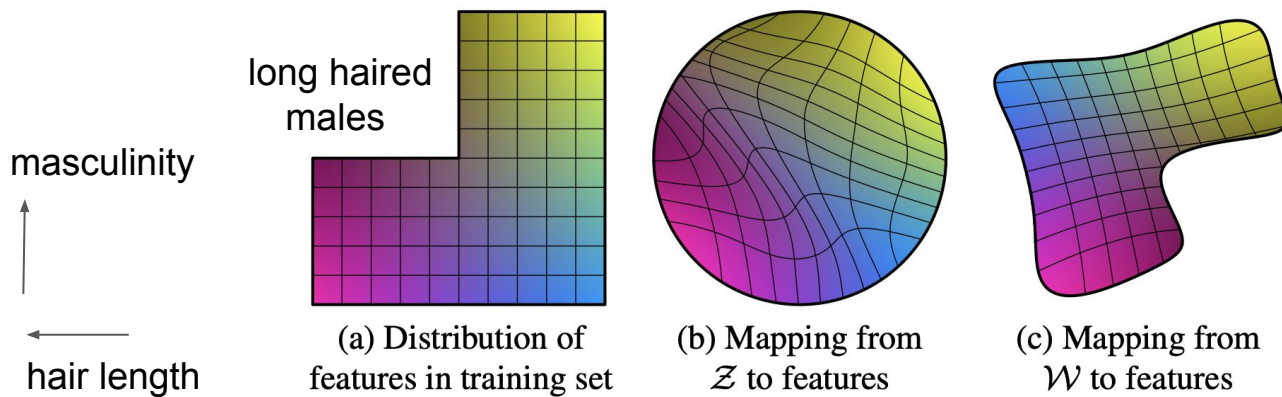
(a) Traditional



(b) Style-based generator

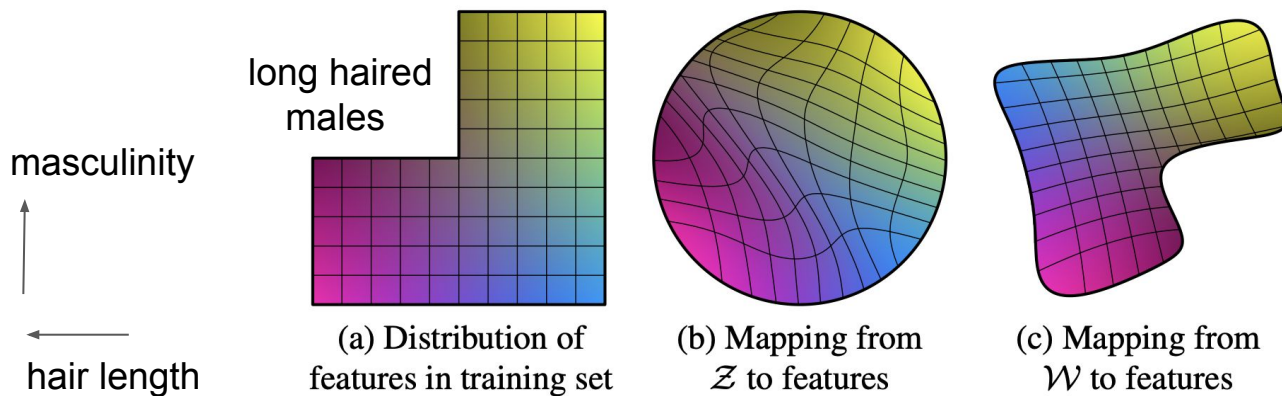
StyleGANs

- StyleGAN embeds the input latent code \mathbf{z} into an intermediate latent space \mathbf{w}
 - $\mathbf{w} = F(\mathbf{z})$
- Now it is \mathbf{w} , not \mathbf{z} , that controls the *style* of the generated images



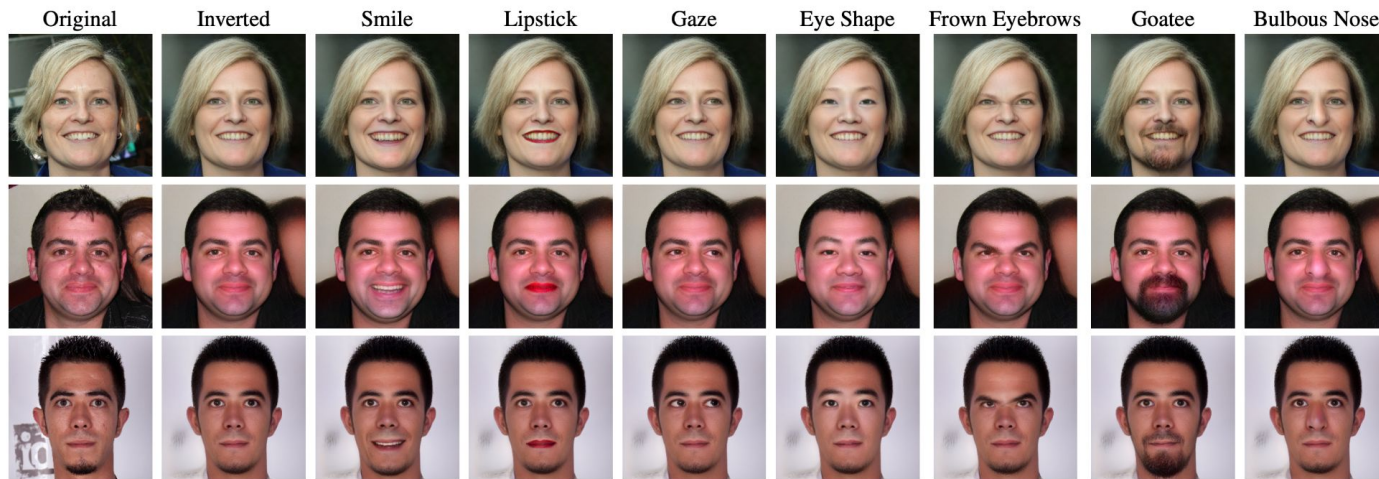
StyleGANs

- StyleGAN embeds the input latent code \mathbf{z} into an intermediate latent space \mathbf{w}
 - $\mathbf{w} = \mathbf{F}(\mathbf{z})$
- Now it is \mathbf{w} , not \mathbf{z} , that controls the *style* of the generated images
 - The mapping \mathbf{F} can “unwarp” \mathbf{z} to \mathbf{w} so that the factors of variation become more linear.
 - Style: factors of variation of the domain of interest

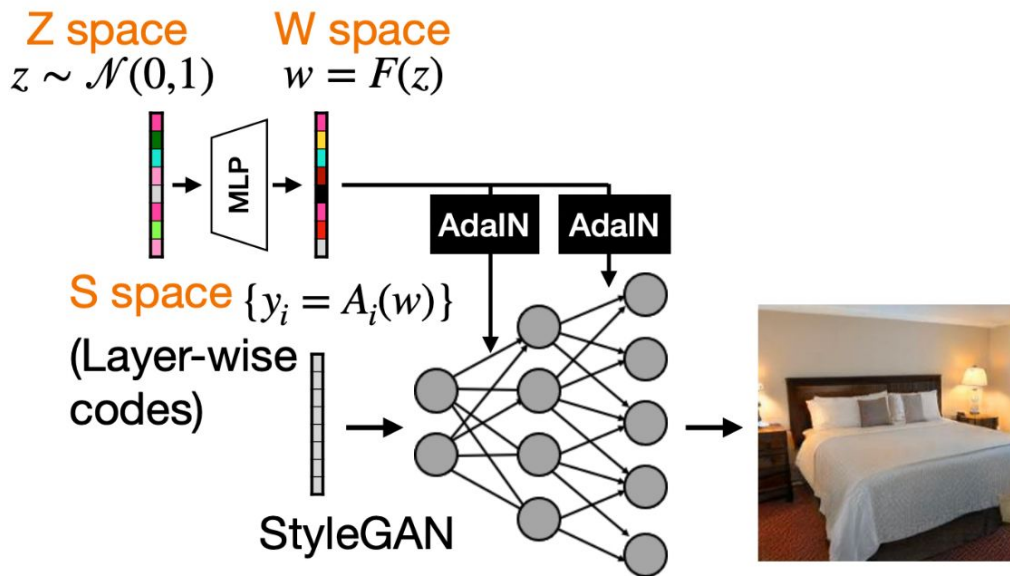


StyleGANs - The Disentangled Latent Space

- There are various definitions for disentanglement.
- A common goal is a latent space that consists of linear subspaces, each of which controls one factor of variation.



Which latent space is more disentangled?



Reconstruction Error (Xu et al. CVPR'21)

Space	MSE	FID
W space	0.0601	22.24
S space	0.0464	18.48

Disentanglement (Wu et al. CVPR'21)

Space	Disentanglement
Z space	0.31
W space	0.54
S space	0.75

GAN Inversion

Applying the pretrained GAN model to image processing tasks

GAN inversion:

$$x^* = \operatorname{argmin}_x ||G(x) - I||$$

Colorization:

$$x^* = \operatorname{argmin}_x ||\operatorname{rgb2gray}(G(x)) - I_{\text{gray}}||$$

Super-resolution:

$$x^* = \operatorname{argmin}_x ||\operatorname{down}(G(x)) - I_{\text{small}}||$$



(a) Image Reconstruction

(b) Image Colorization



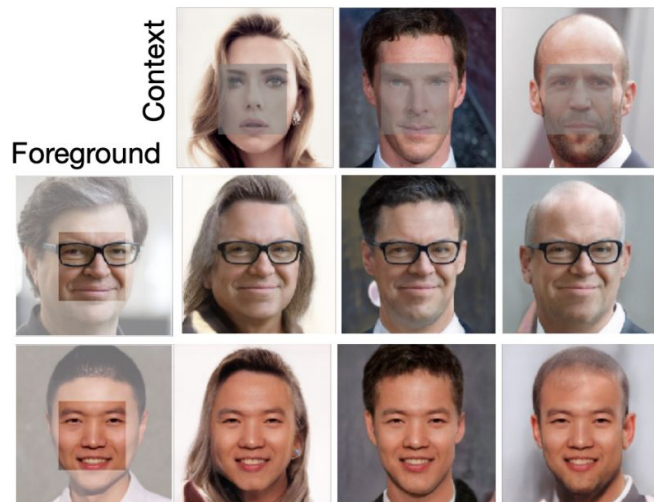
(d) Image Denoising

(e) Image Inpainting

Gu, Shen, Zhou. Image Processing Using Multi-Code GAN Prior. CVPR'20

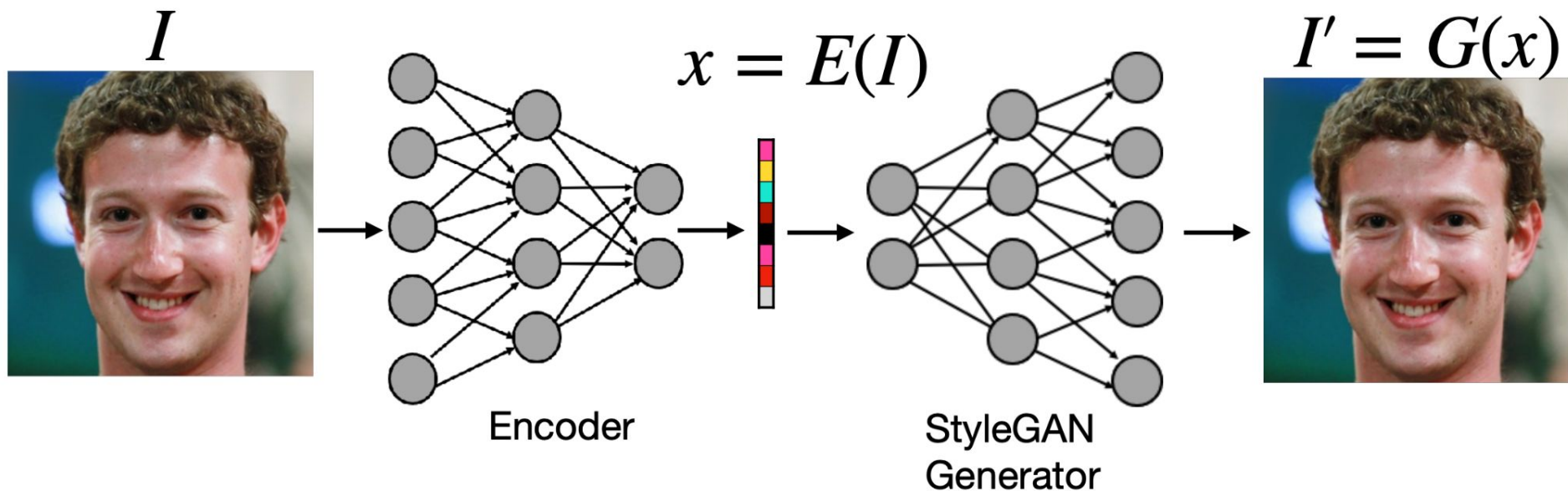
Masked optimization

$$x^* = \operatorname{argmin}_x ||m \cdot G(x) - m \cdot I_{\text{context}}||$$

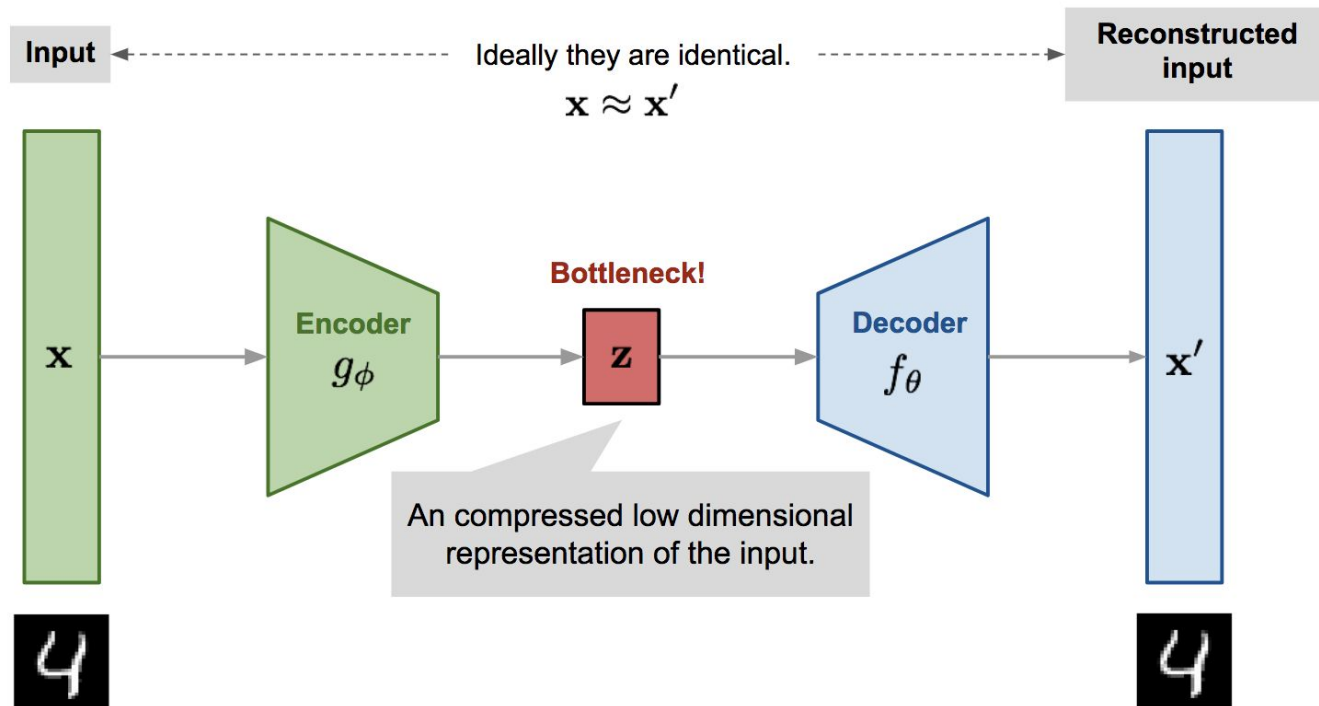


Zhu, Shen, Zhao, Zhou. In-domain GAN Inversion. ECCV'20

Encoding Real Image into StyleGAN space



AutoEncoders



AutoEncoders

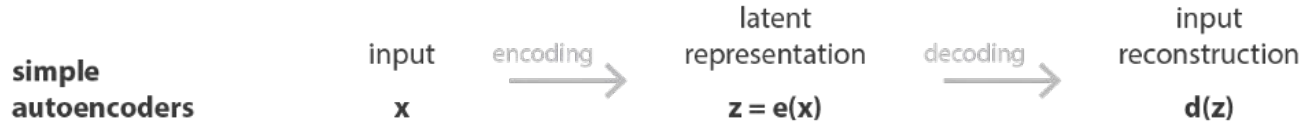
- Encoder
 - Transforms the original high-dimension input (eg., images) into the low-dimensional latent.
 - Hopefully lossless
- Decoder
 - Recovers the high-dimensional data from the encoded low-dimensional latents
- Dimensionality Reduction
 - Links to PCA

$$L_{\text{AE}}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - f_{\theta}(g_{\phi}(\mathbf{x}^{(i)})))^2$$

Variational AutoEncoders

- Vanilla autoencoder's latent space is NOT well-organized/structured to be sampled from
 - Because there is no force for the latent space to do so
- Variational AutoEncoders
 - Autoencoders whose latent space is regularized to a structured distribution (eg., Gaussian distribution)
 - The latent is now a *distribution*

Variational AutoEncoders



$$\begin{aligned} L_{\text{VAE}}(\theta, \phi) &= -\log p_{\theta}(\mathbf{x}) + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})\|p_{\theta}(\mathbf{z}|\mathbf{x})) \\ &= -\mathbb{E}_{\mathbf{z}\sim q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})\|p_{\theta}(\mathbf{z})) \end{aligned}$$

Variational AutoEncoders



$$L_{\text{VAE}}(\theta, \phi) = -\log p_{\theta}(\mathbf{x}) + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}|\mathbf{x}))$$

the reconstruction term:
negative log-likelihood

$$= -\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}))$$

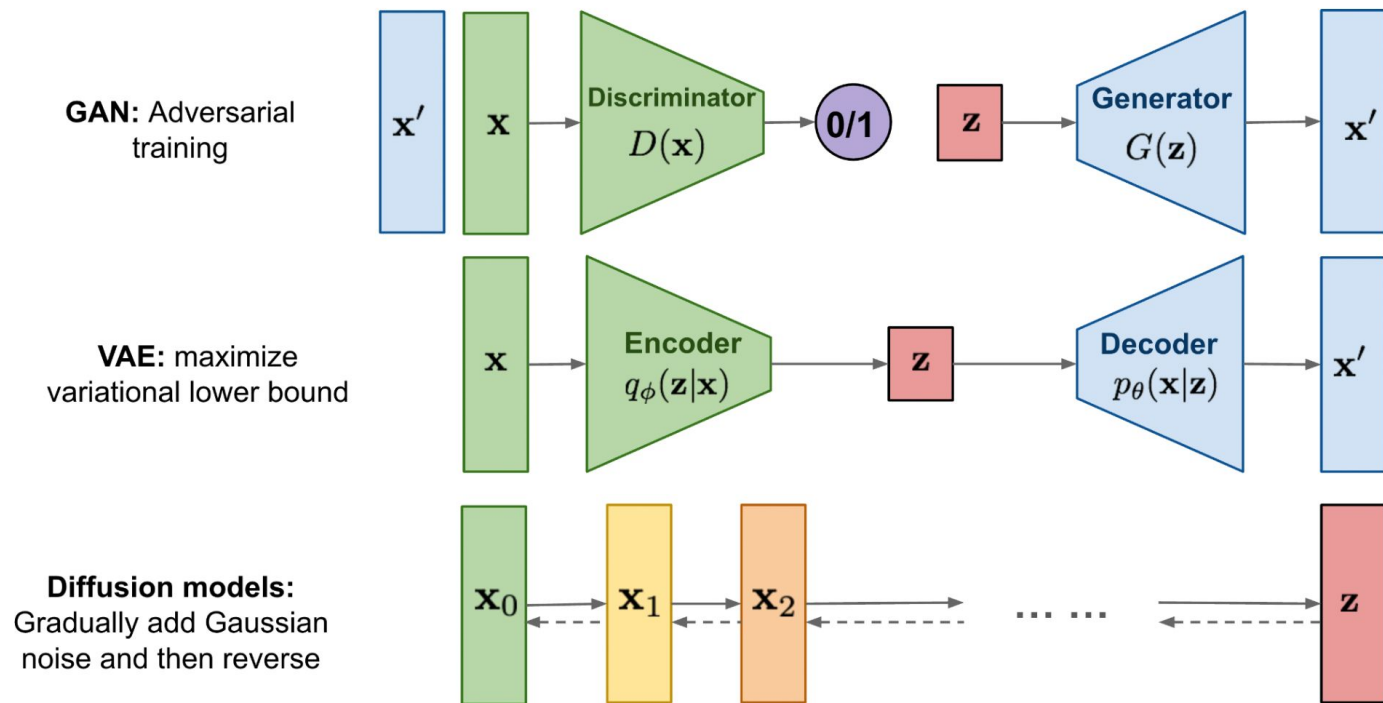
Variational AutoEncoders



$$L_{\text{VAE}}(\theta, \phi) = -\log p_{\theta}(\mathbf{x}) + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x}))$$

the regularization term: $= -\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}))$
KL divergence to the prior distribution

Summary



Recall Fidler's 3D Neural Rendering Approach

$$I \mapsto W \mapsto z \mapsto I$$

- The whole system is an autoencoder
- $I \rightarrow W$: An autoencoder approach where the decoder is a differentiable renderer $I=F(W;\alpha)$ and $P_{\phi}(W | I)$ is the encoder.
- $z \rightarrow I$: A learned StyleGAN generative model $P_{\theta}(I | z)$
- $W \rightarrow z$: Learn $f_{\psi}(z | W)$ using an autoencoder reconstruction loss
 - This can be done by the autoencoder because the latent variables of styleGANs are fairly interpretable and so the function $f_{\psi}(z | W)$ cannot be too complicated