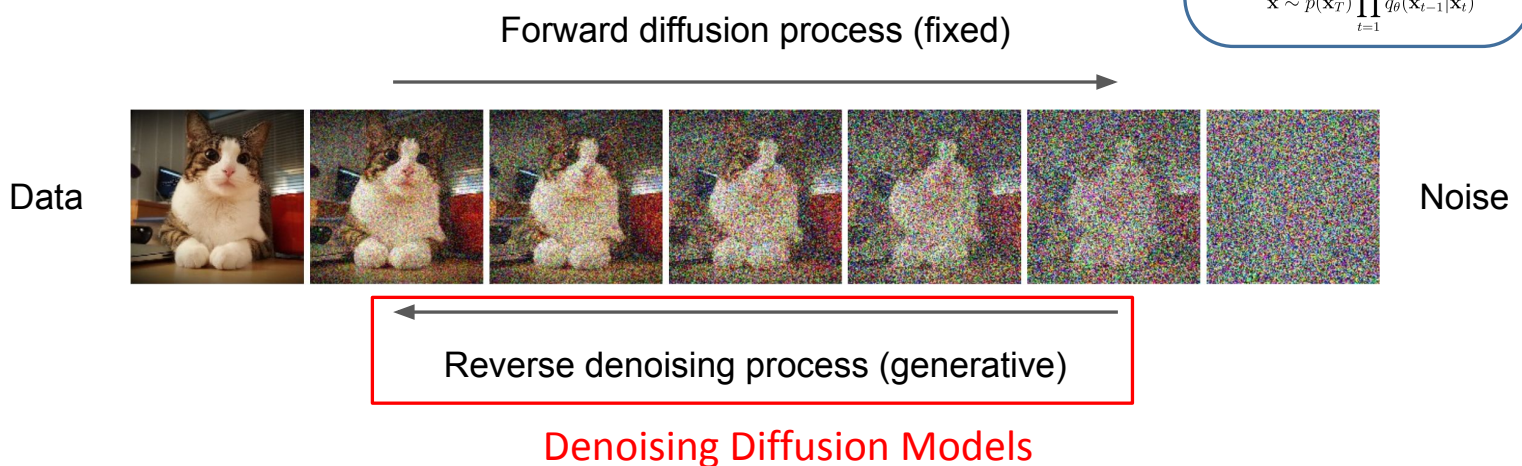


# Outline

- Generative Models Basics
- Autoregressive Models
- Autoencoder and Variational Autoencoder
- Generative Adversarial Network
- **Diffusion Models**

116

## Denoising Diffusion Models

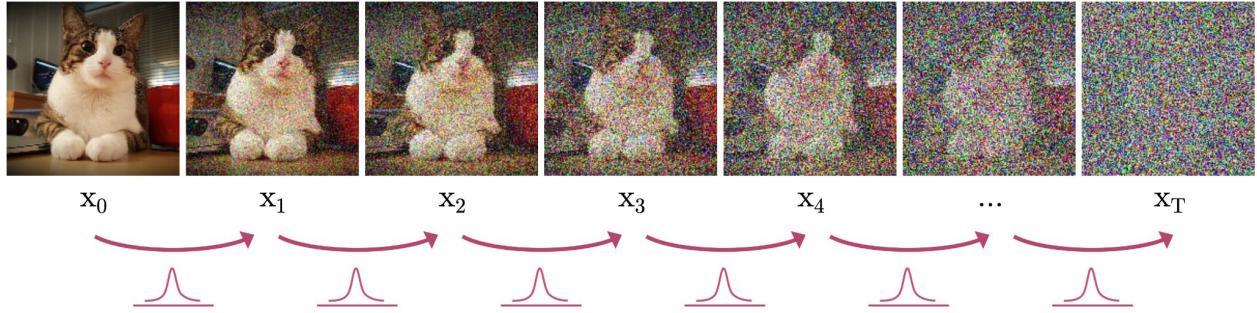


118

# Forward Diffusion Process

The formal definition of the forward process in  $T$  steps:

Forward diffusion process (fixed)

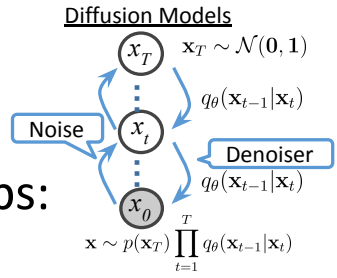


$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \rightarrow q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (\text{joint})$$

$\beta_s$  controls the *noise schedule* such that  $\mathbf{x}_T$  is complete noise.

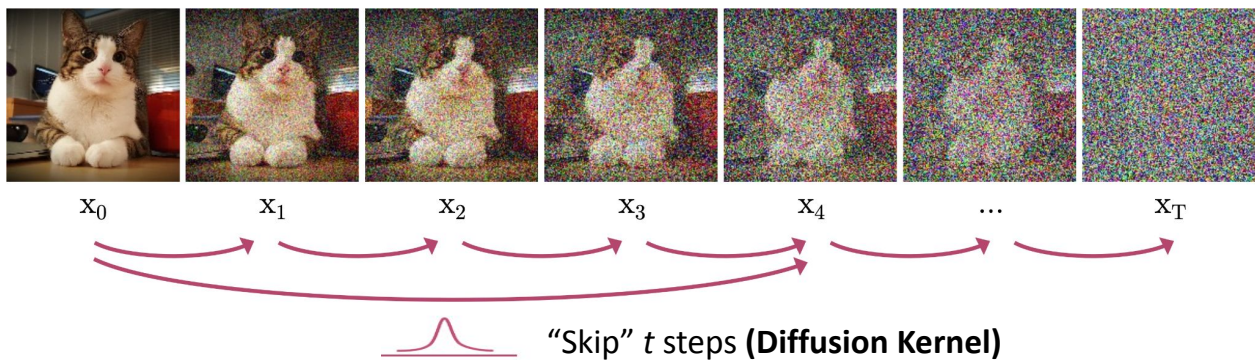
119

Slide credit: Kreis, Gao, Vahdat, CVPR 2022 Tutorial



# Diffusion Kernel

Forward diffusion process (fixed)



$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

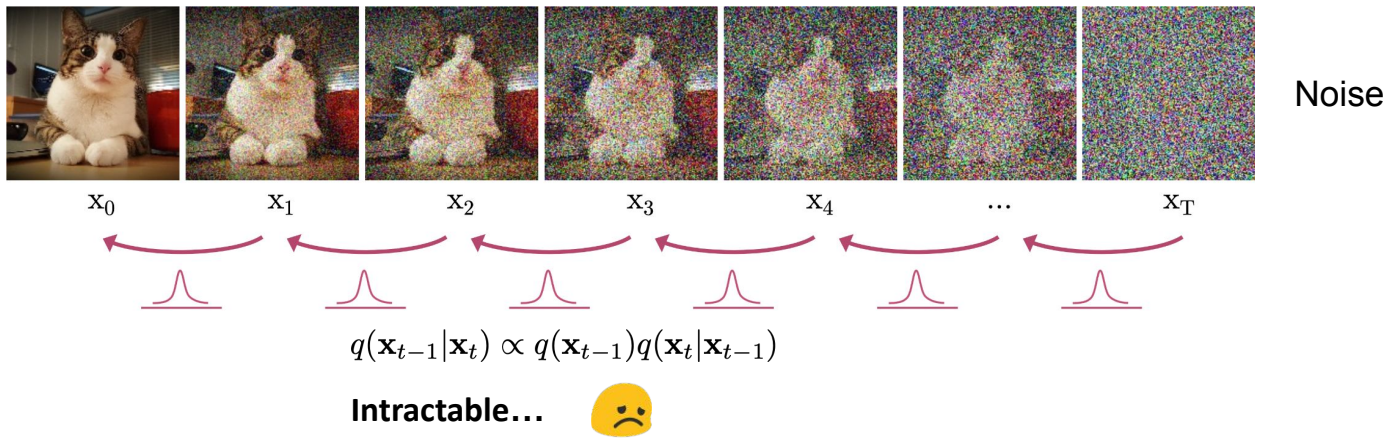
$$\text{Define } \bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$$

120

Slide credit: Kreis, Gao, Vahdat, CVPR 2022 Tutorial

# Reverse Denoising Process

Reverse denoising process (generative)



Idea: Approximate  $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ . Use a **Normal distribution** if  $\beta_t$  is small in each forward diffusion step.

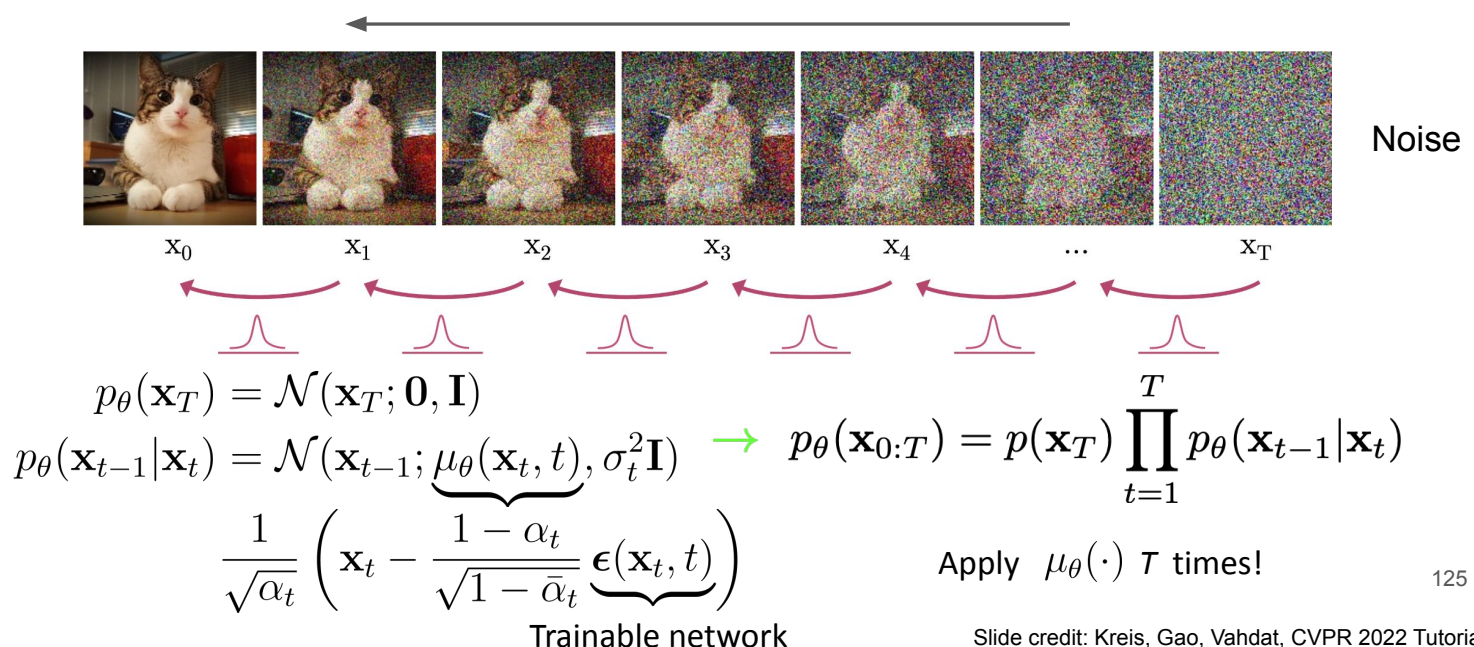
124

Slide credit: Kreis, Gao, Vahdat, CVPR 2022 Tutorial

# Reverse Denoising Process

Formal definition of forward and reverse processes in T steps:

Reverse denoising process (generative)



125

Slide credit: Kreis, Gao, Vahdat, CVPR 2022 Tutorial



# U-Net: Image-to-Image CNN network

The trainable noise prediction network is usually a U-Net

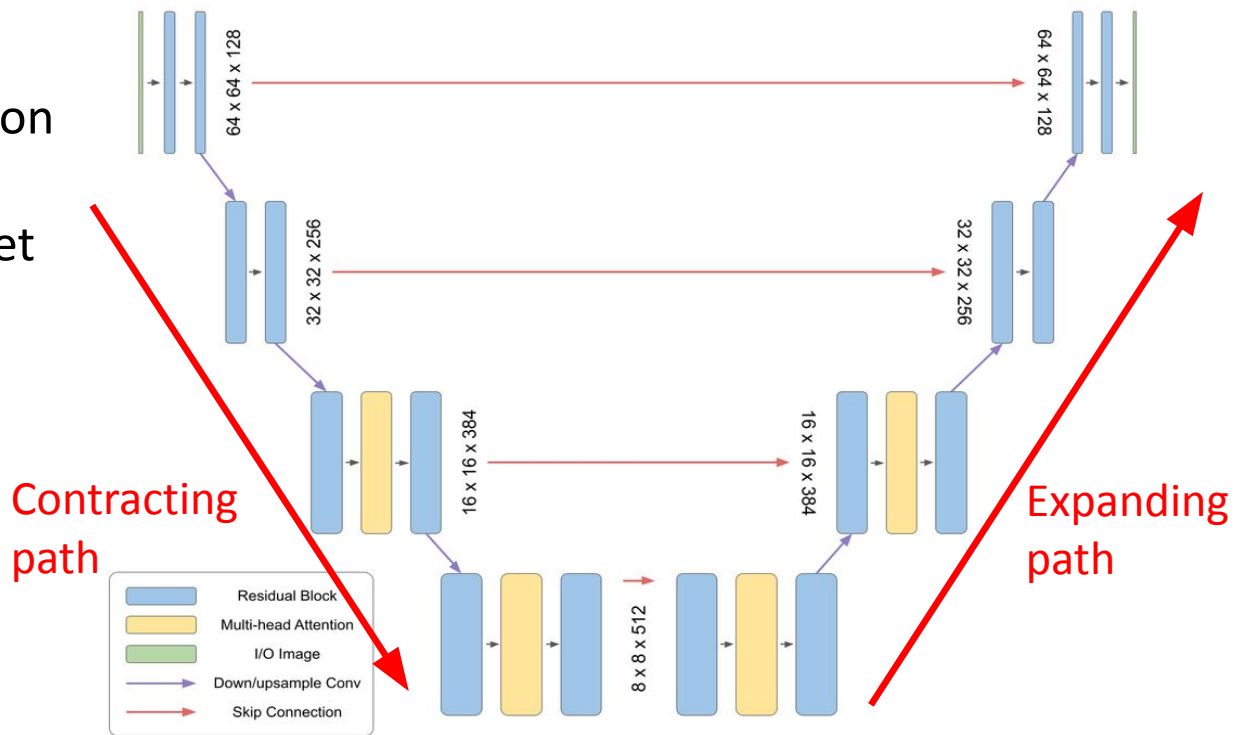


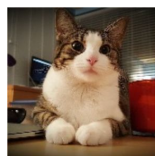
Image Source: <https://www.assemblyai.com/blog/minimagen-build-your-own-image-text-to-image-model/>

127

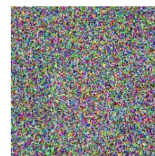
Repeat:

## Training

Sample:



$\mathbf{x}_0$



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

$t \sim \text{Uniform}(\{1, \dots, T\})$

Diffusion kernel:



$$\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$$

Add noise to clean image (with variance equivalent to  $t$  steps of noise)

Gradient update:

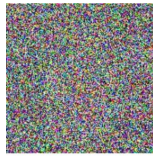
$$\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta} \left( \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2$$

$$\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta} \left( \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2$$

Give model noisy image, ask it to predict the added noise. Train with L2 loss between predicted noise and original noise

128

# Sampling



Sample:

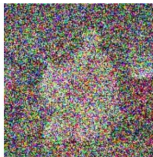
$$\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

**for**  $t = 1 \dots T$  **do**:



Sample noise:

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



Predict next step:

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \text{noise} - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\text{noise}, t) \right) + \mathbf{z}_t$$

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \mathbf{z}_t$$

129

## Summary: Training and Sample Generation

### Algorithm 1 Training

```

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

### Algorithm 2 Sampling

```

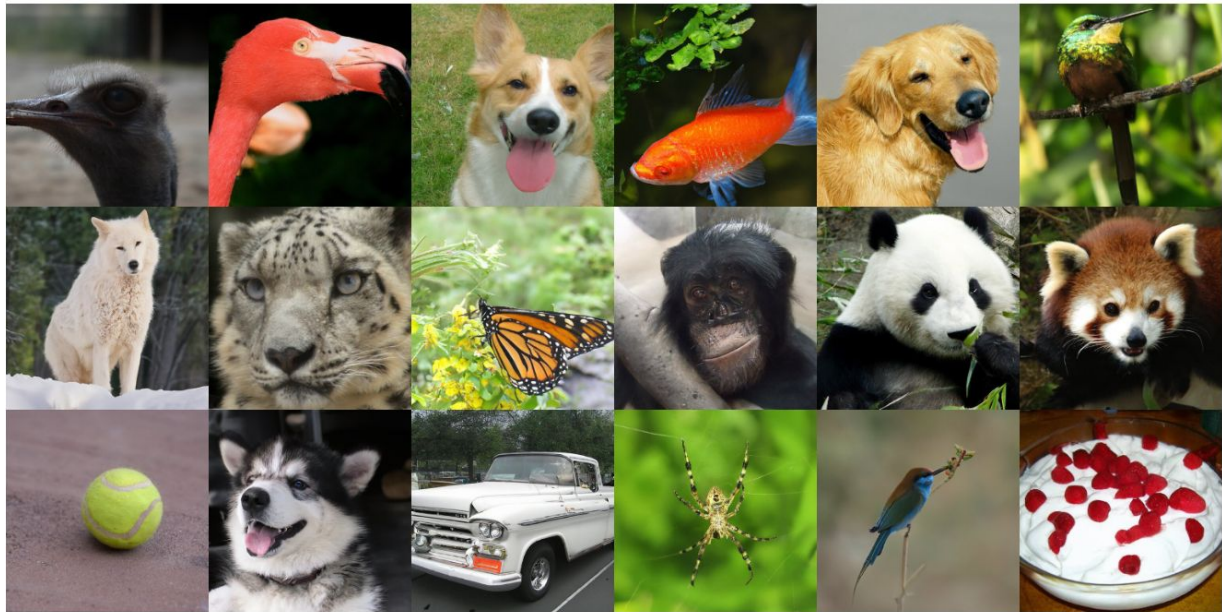
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t,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

130

# Improvements: Conditional Generation

## Guided Diffusion

512 x 512 ImageNet Conditioned Samples



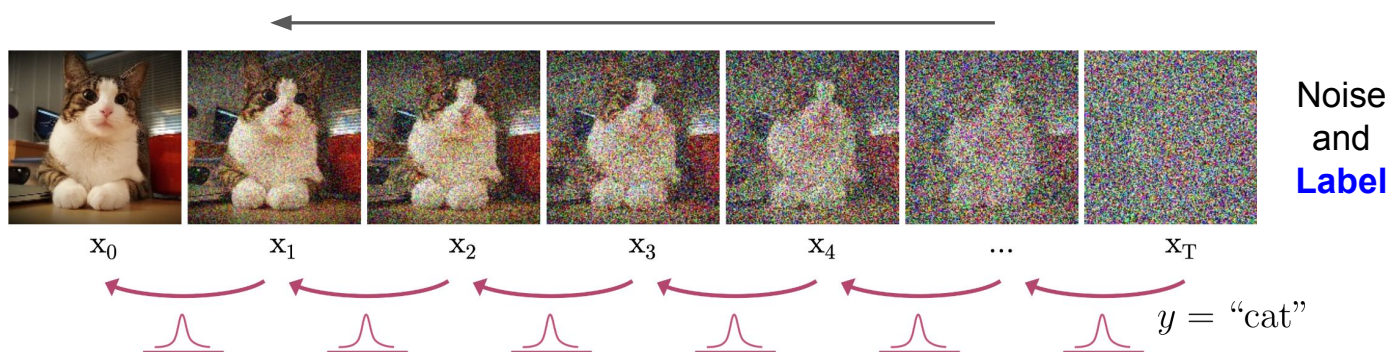
Dhariwal and Nichol "Diffusion Models Beat GANs on Image Synthesis" NeurIPS 2021

131

Data

# Improvements: Conditional Generation

**Conditional** Reverse denoising process (generative)



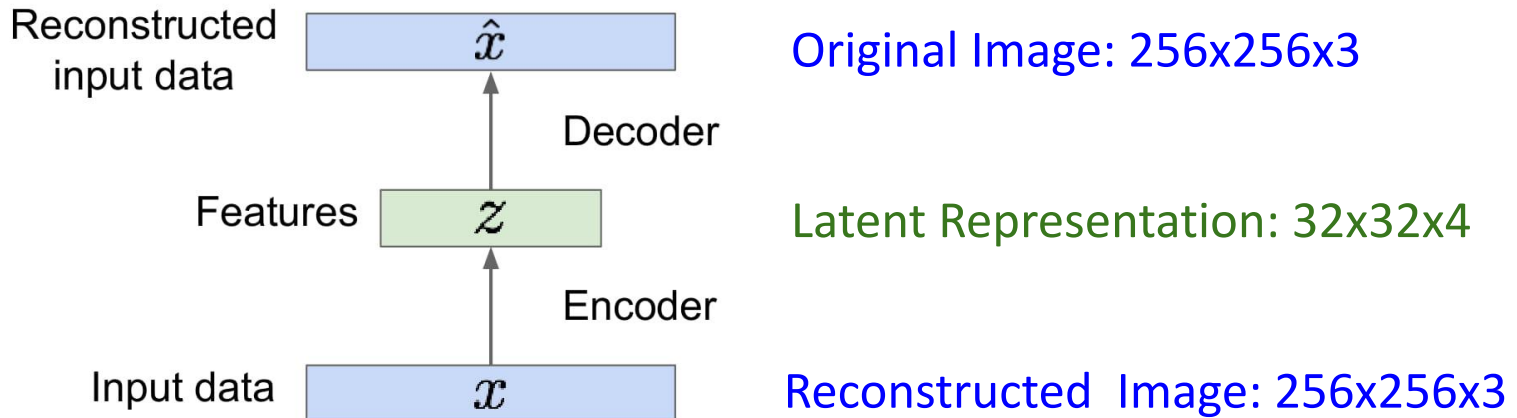
We pass the condition ( $y=\text{cat}$ ) to the noise-predicting model in addition to the noisy image and step number to “control” the generation (i.e. force the generated image to be a cat instead of a dog)

Text-to-image generation: train the diffusion model to condition on text!

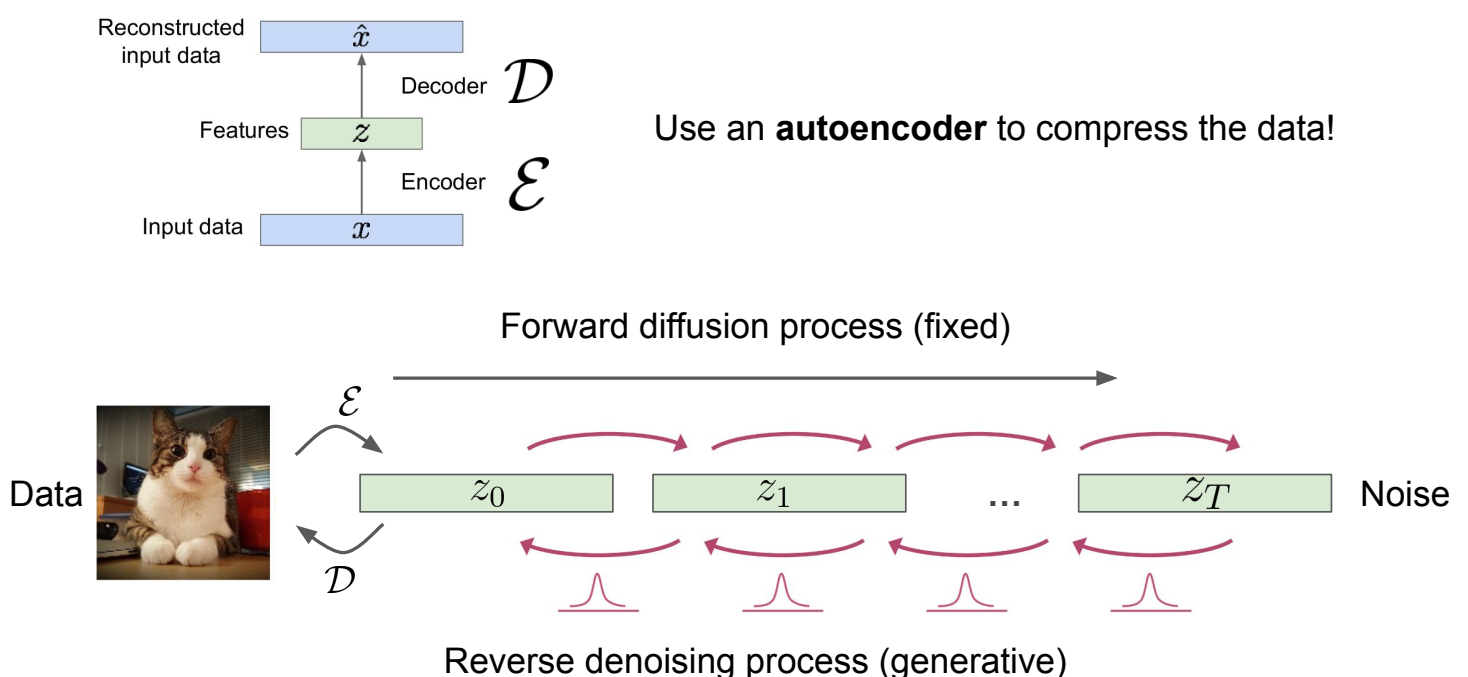
132

# Improvements: Latent Diffusion Models

Use an **autoencoder** to compress the data!



# Improvements: Latent Diffusion Models



# Most popular text-to-image generation models nowadays are Latent Diffusion Models

For example,

- OpenAI DALLÉ-2 and DALLÉ-3, and likely all images generated from text within ChatGPT
- Midjourney
- Imagen (by Google)

Diffusion models can also be used for text-to-video generation.

- OpenAI Sora is a latent diffusion model with transformer architecture

## Recap

		Pros	Cons
<b>Autoregressive (PixelRNN, PixelCNN)</b>	Explicit density model	Optimizes exact likelihood	Inefficient sequential generation, sample quality not the best
<b>Variational Autoencoders (VAE)</b>	Optimize variational lower bound on likelihood	Useful latent representation, inference queries	Sample quality not the best
<b>Generative Adversarial Networks (GANs)</b>	Game-theoretic approach	Great samples	Tricky and unstable to train, no inference queries
<b>Diffusion Models</b>	Optimize variational lower bound for denoising	Best samples! Easy to train	Inefficient sequential generation