

PAT 498/598 (Fall 2024)

Special Topics: Generative AI for Music and Audio Creation

Lecture 10: Diffusion Models

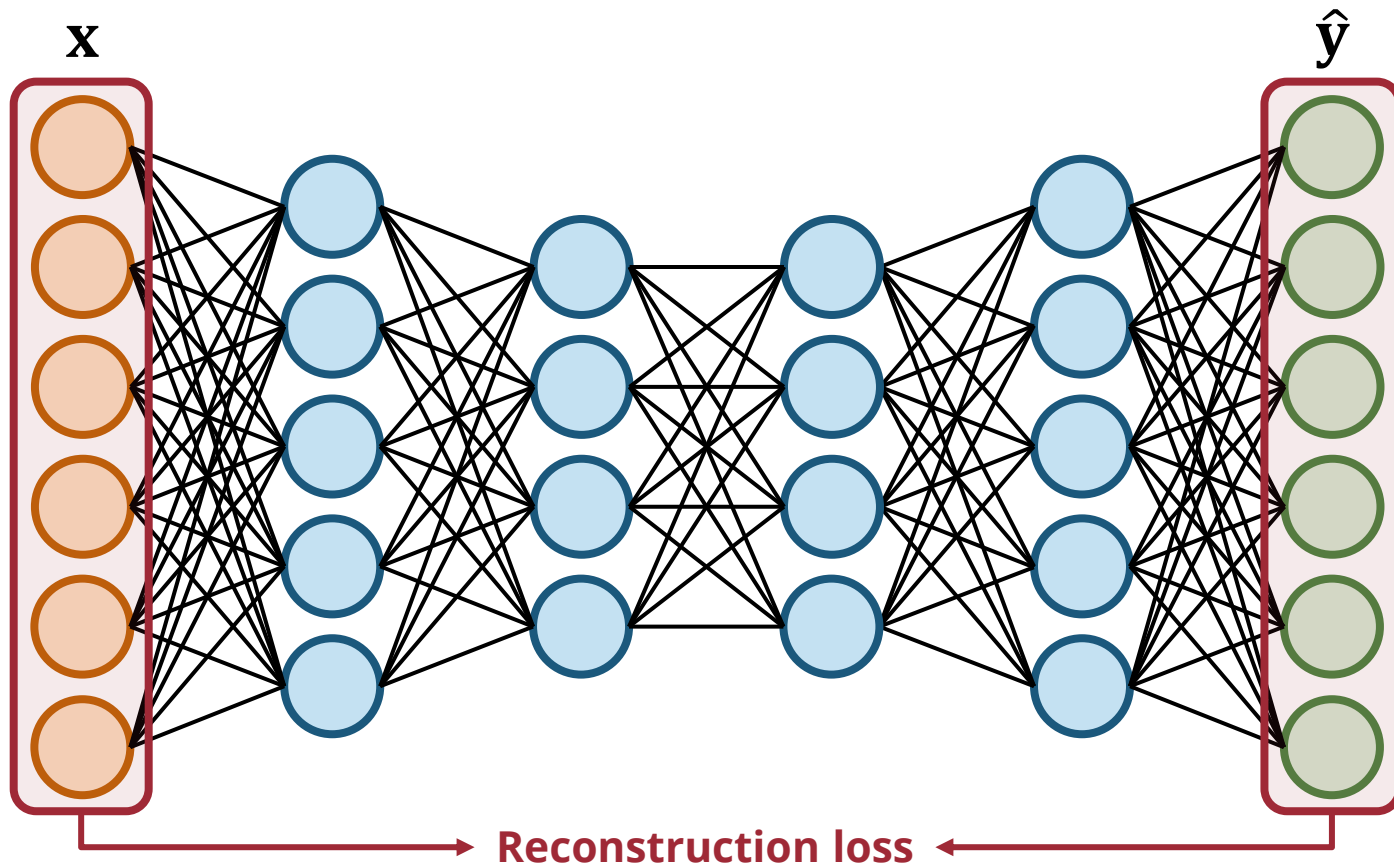
Instructor: Hao-Wen Dong



SCHOOL OF MUSIC, THEATRE & DANCE
PERFORMING ARTS TECHNOLOGY
UNIVERSITY OF MICHIGAN

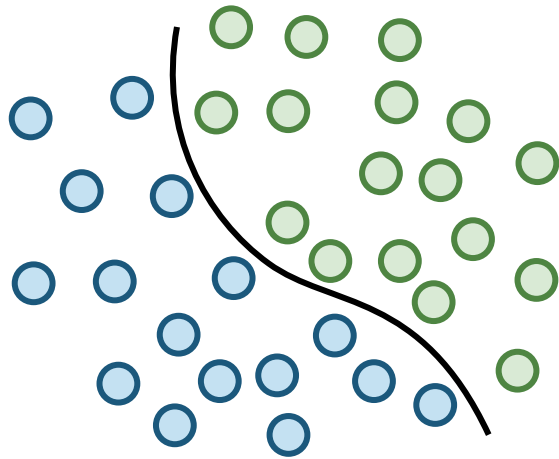
(Recap) Autoencoders

- A neural network where the **input and output are the same**



(Recap) Discriminative vs Generative Models

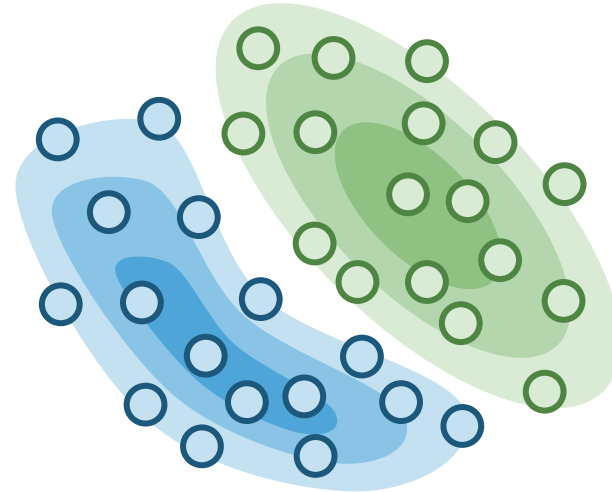
Discriminative



Discriminative models learn the decision boundary

$$P(y|x)$$

Generative

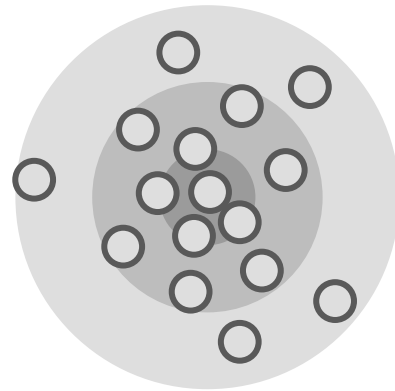


Generative models learn the underlying distribution

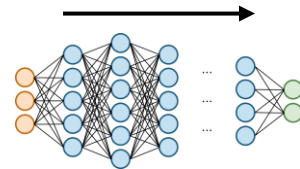
$$P(x) \text{ or } P(x|y)$$

(Recap) Generating Data from a Random Distribution

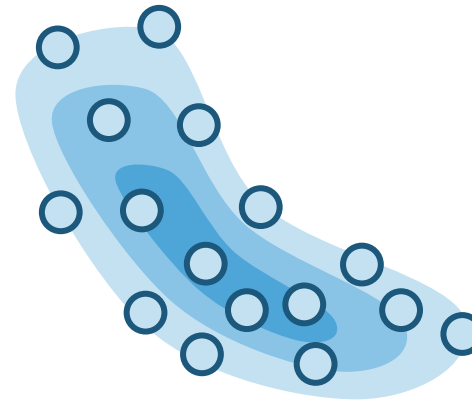
Random distribution



$P(z)$



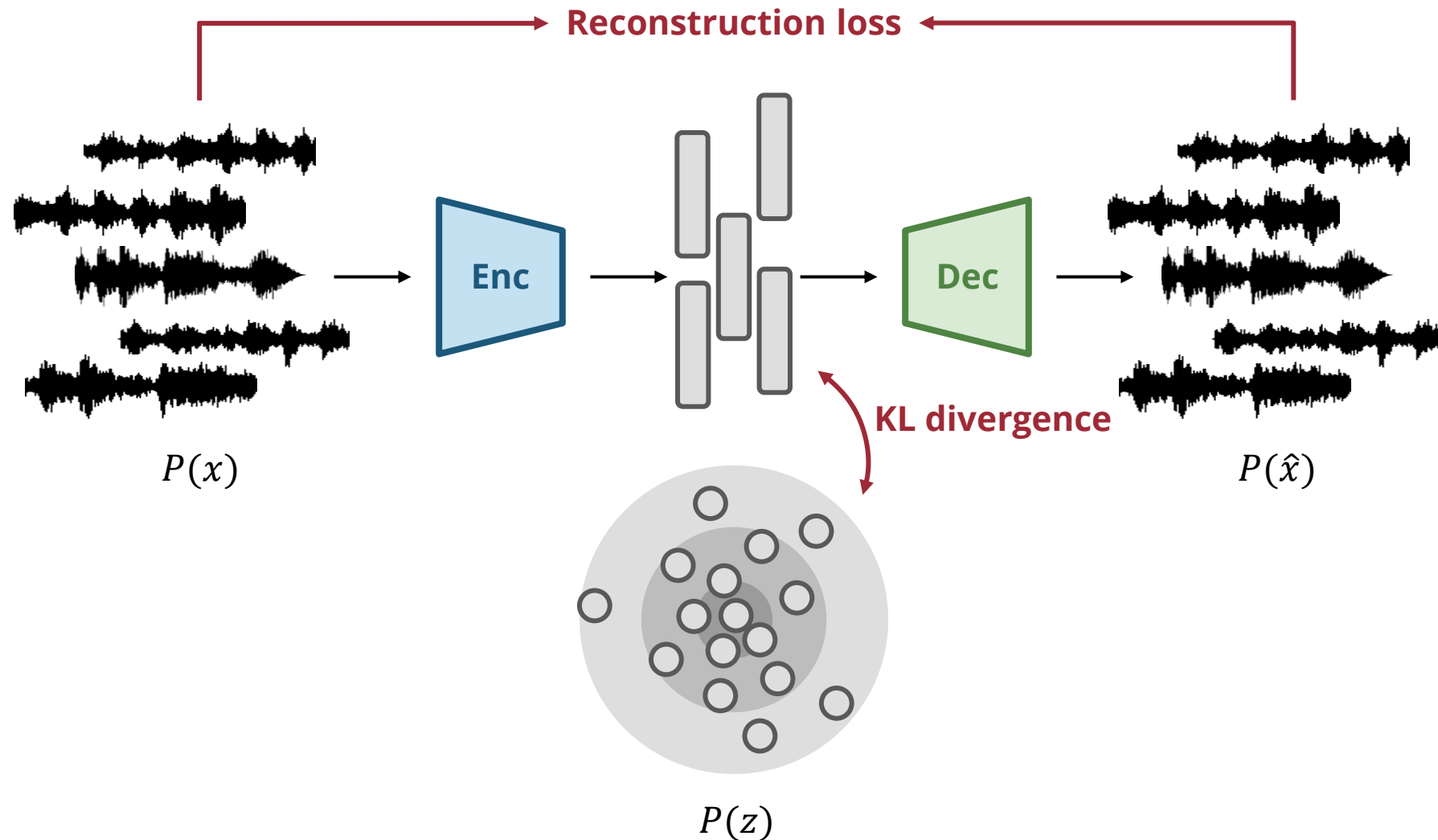
Data distribution



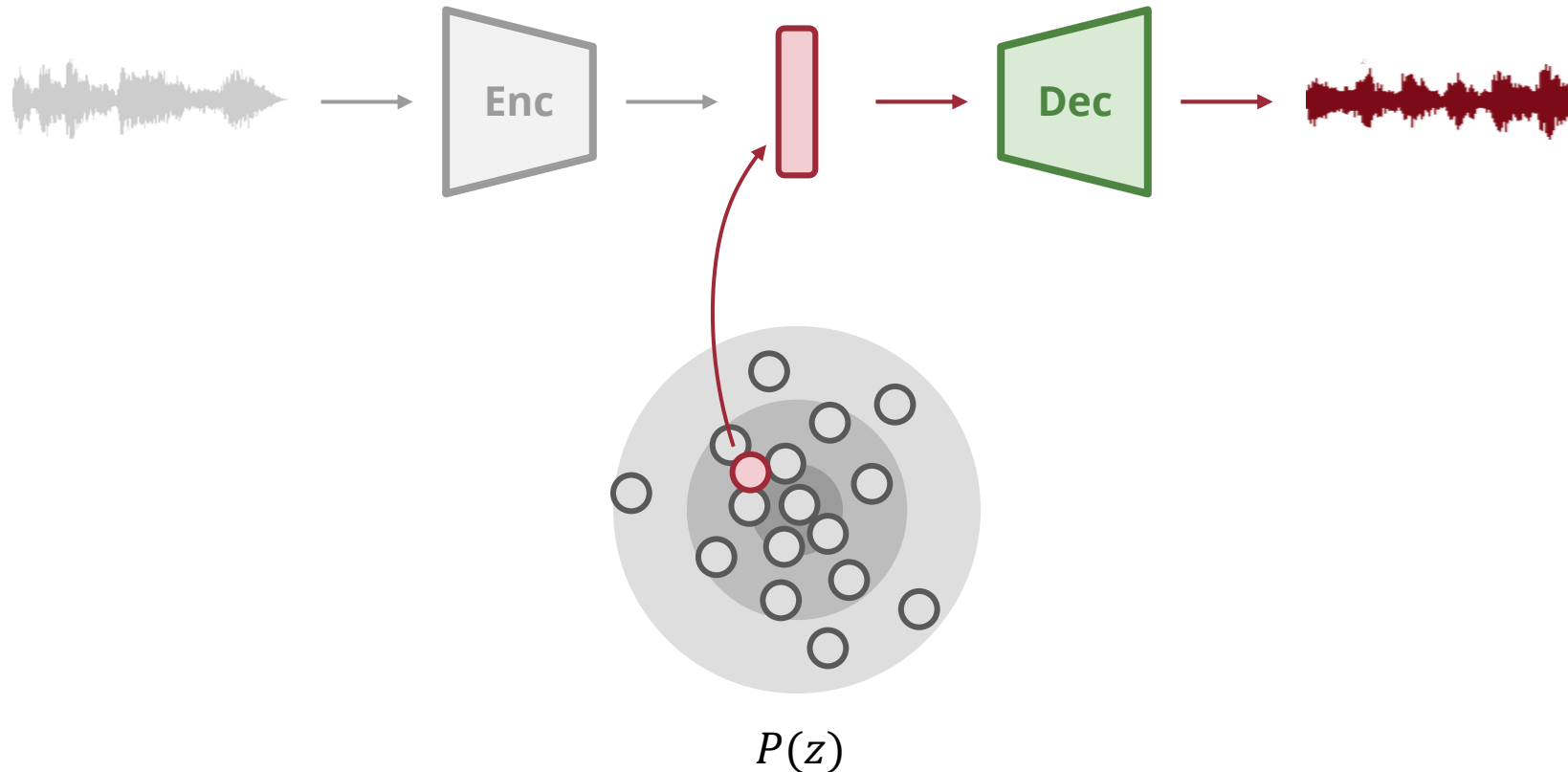
$P(x)$

If we can learn this mapping, we can easily generate new samples from the data distribution

(Recap) Variational Autoencoders (VAEs) – Training

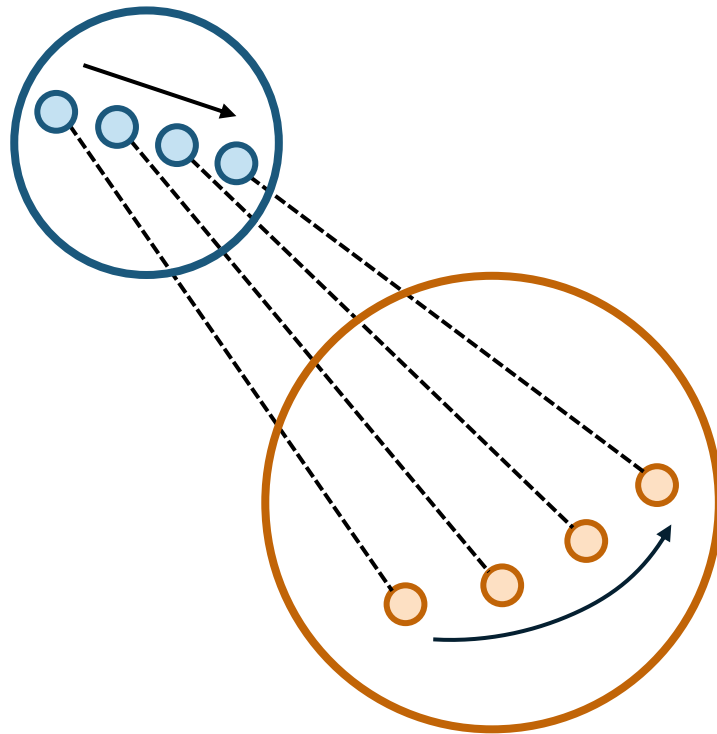


(Recap) Variational Autoencoders (VAEs) – Generation

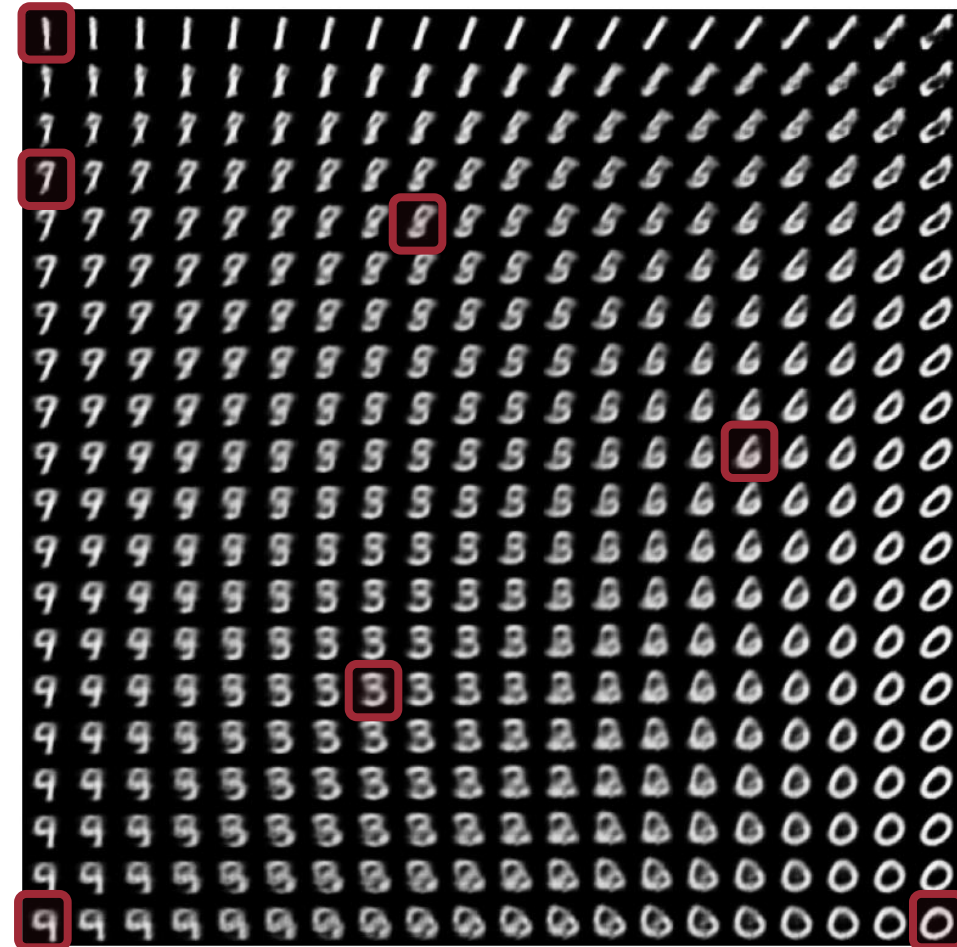


(Recap) Decoding the Latent Space of a VAE

Latent space



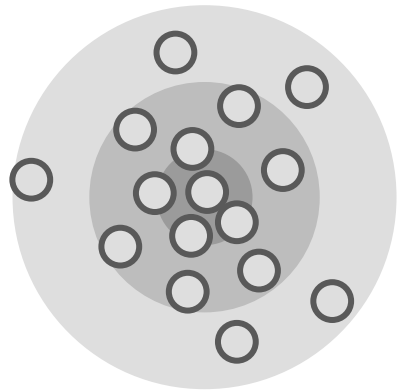
Data space



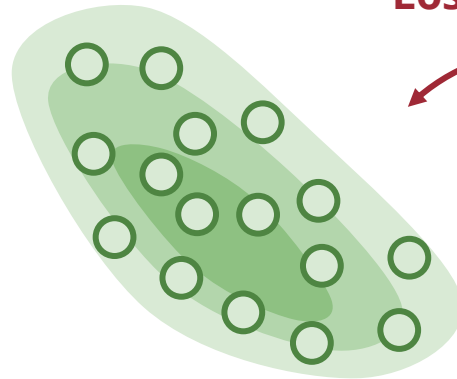
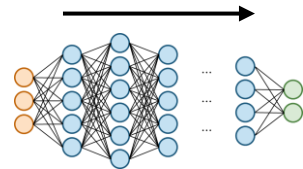
(Source: tensorflow.org)

(Recap) A Loss Function for Distributions

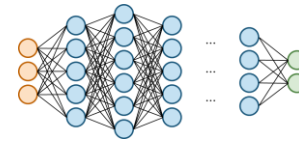
Random distribution



$P(z)$

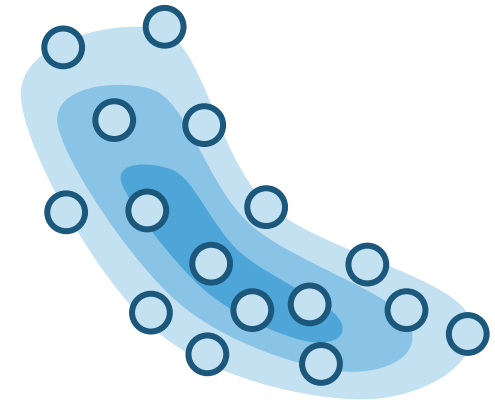


$P(\hat{x})$



Loss function?

Data distribution

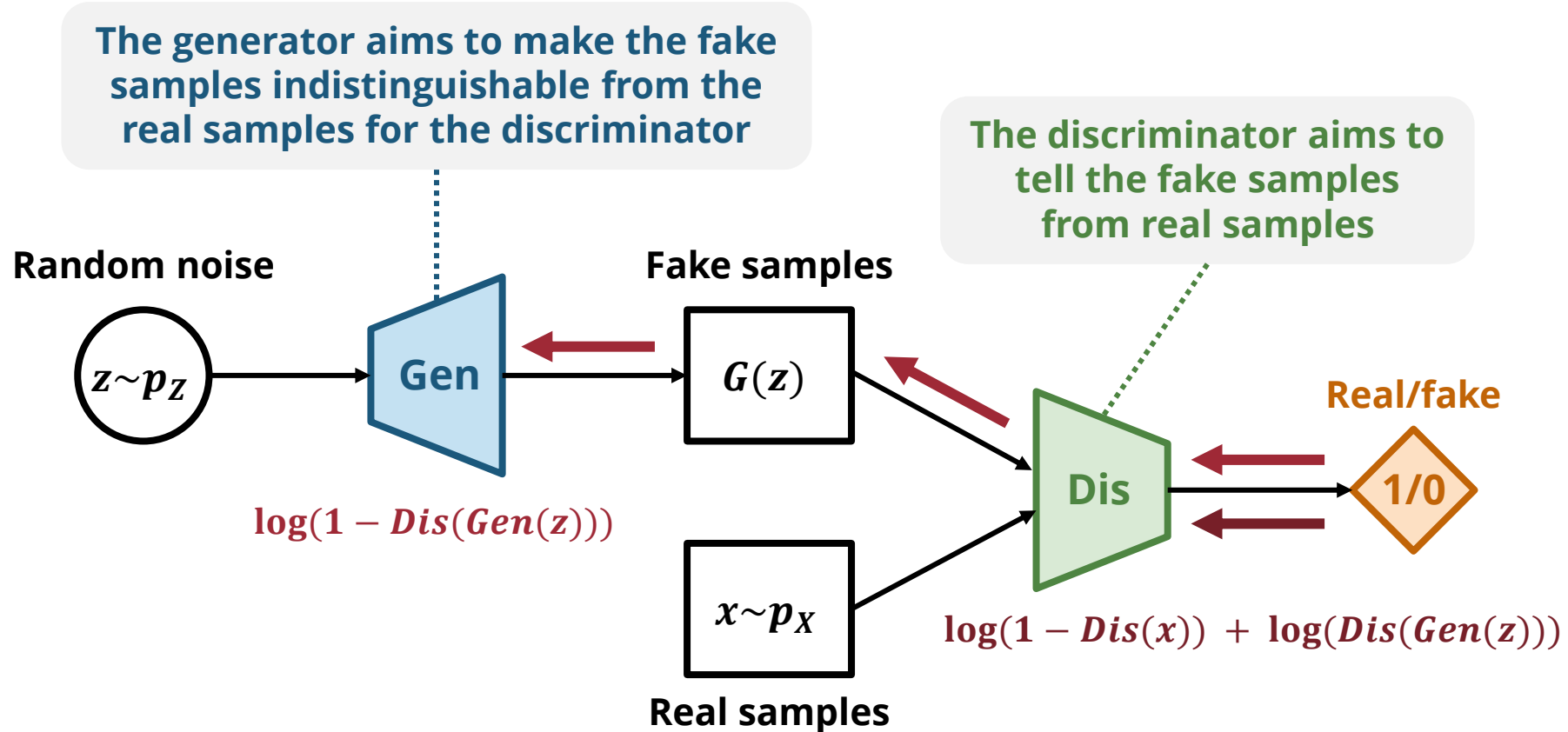


$P(x)$

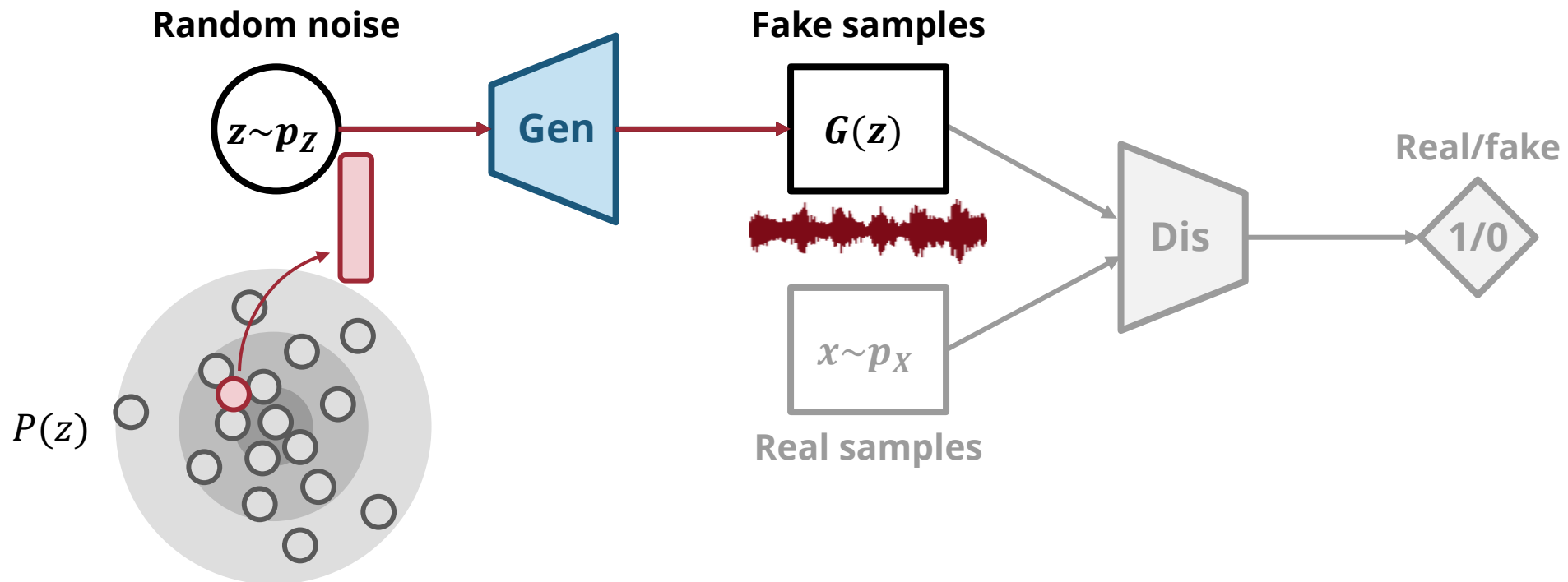
Unfortunately, no easy way to measure the difference between two distributions

But what about another neural network!?

(Recap) Generative Adversarial Nets (GANs) – Training

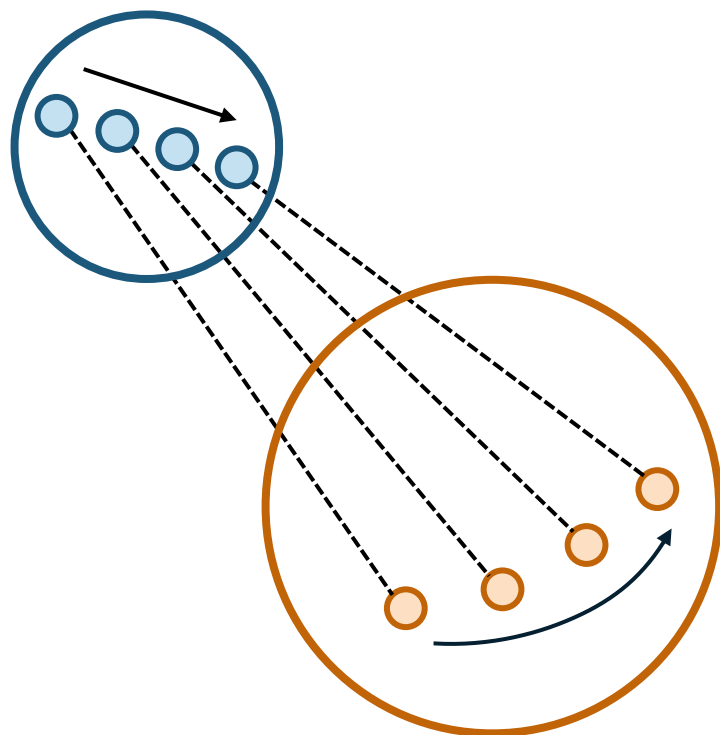


(Recap) Generative Adversarial Nets (GANs) – Generation

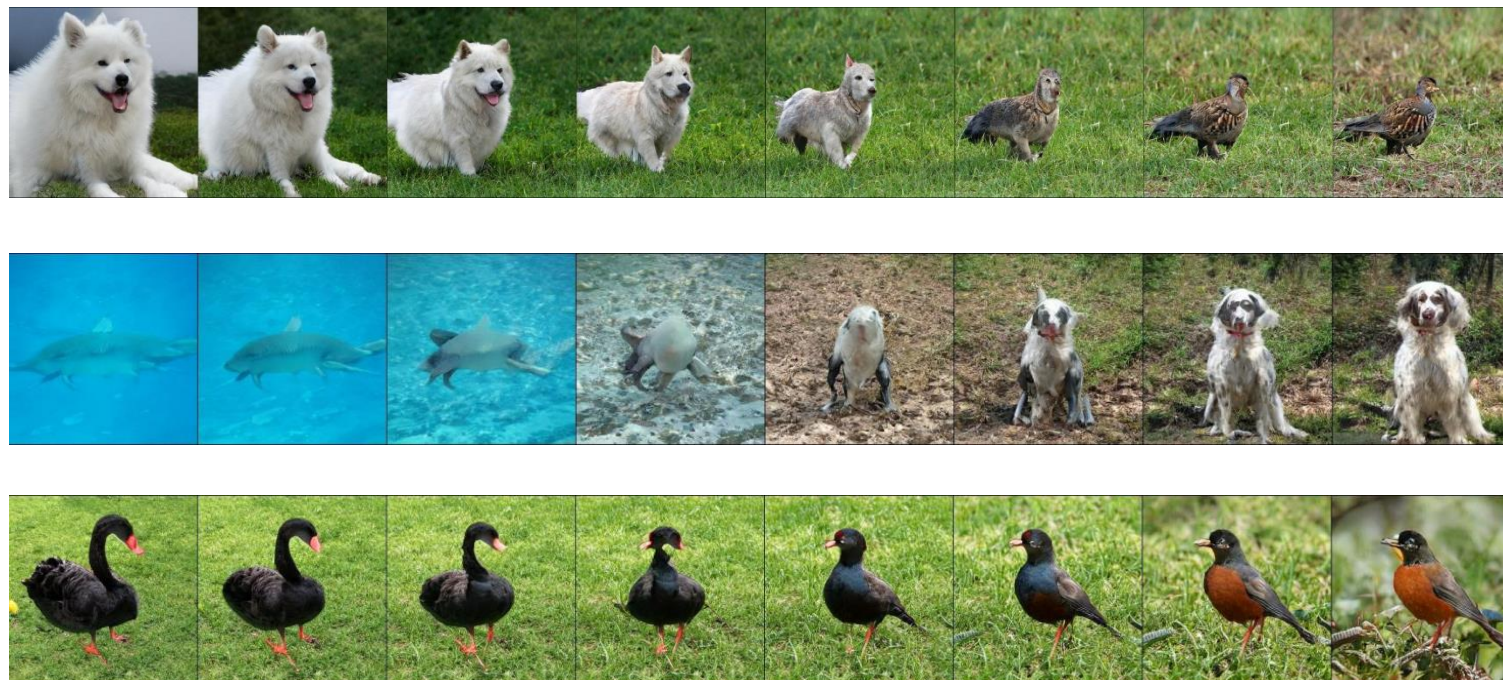


(Recap) Interpolation on the Latent Space

Latent space



Data space

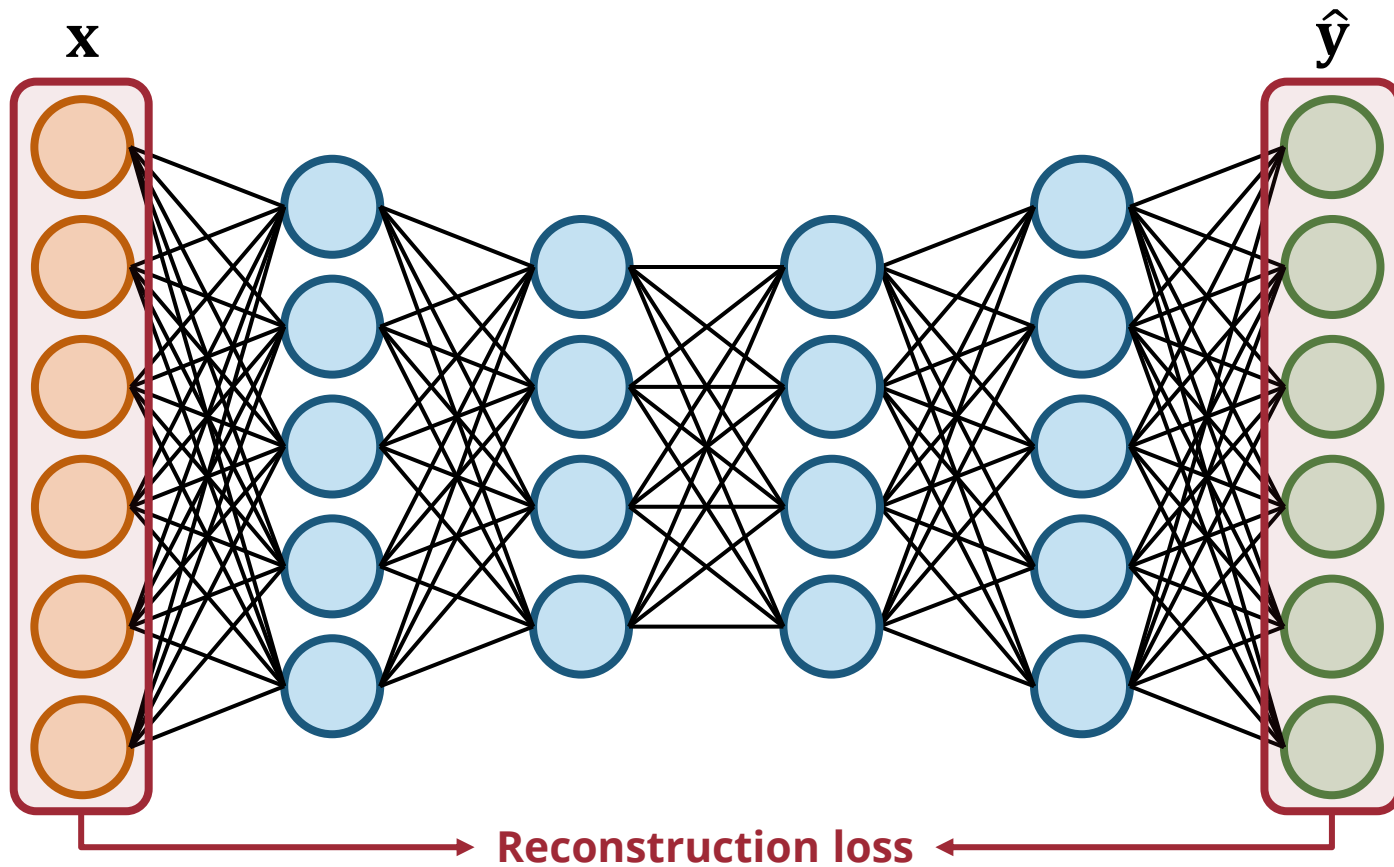


(Source: Brock et al., 2019)

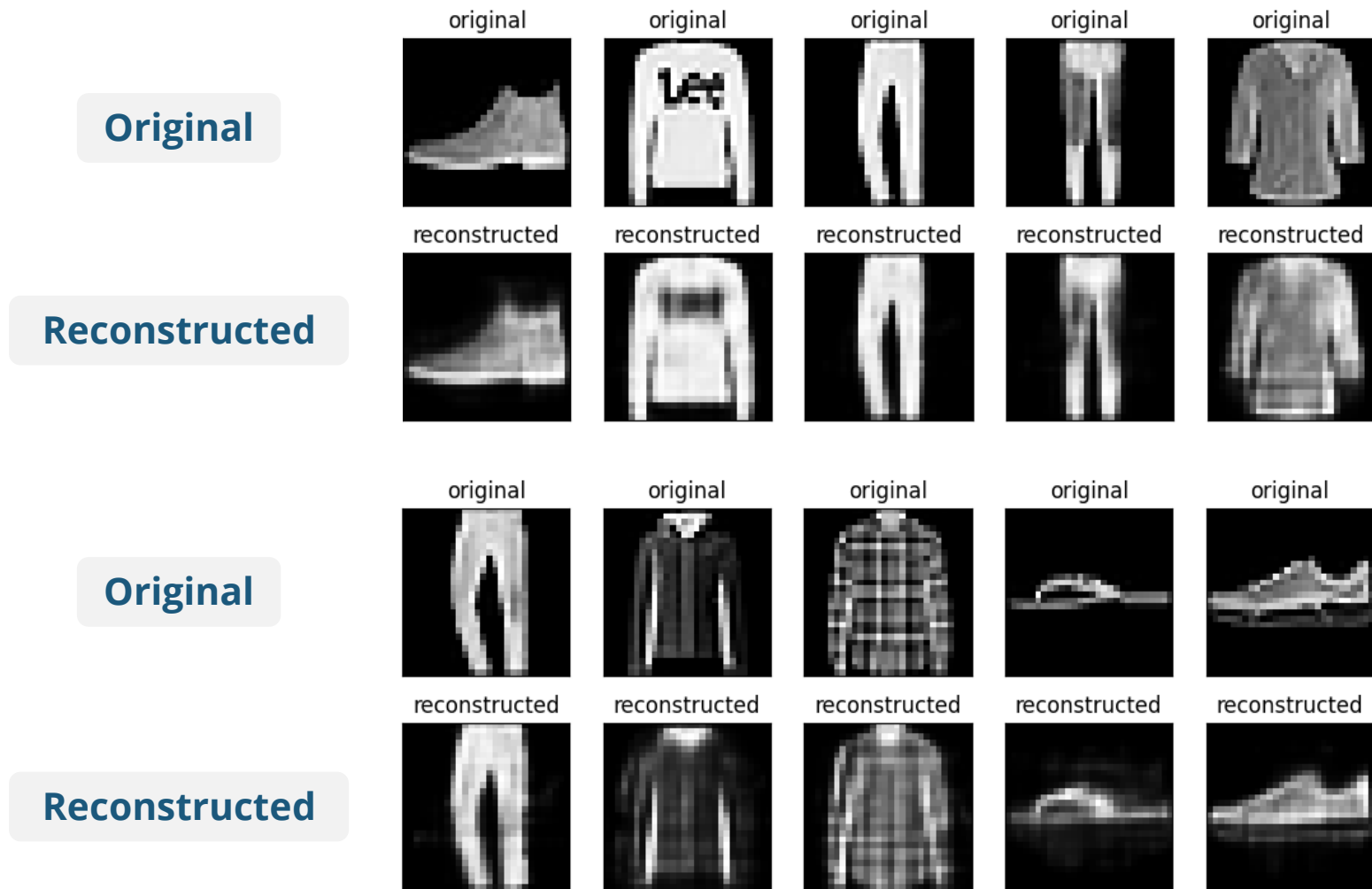
Diffusion Models

(Recap) Autoencoders

- A neural network where the **input and output are the same**

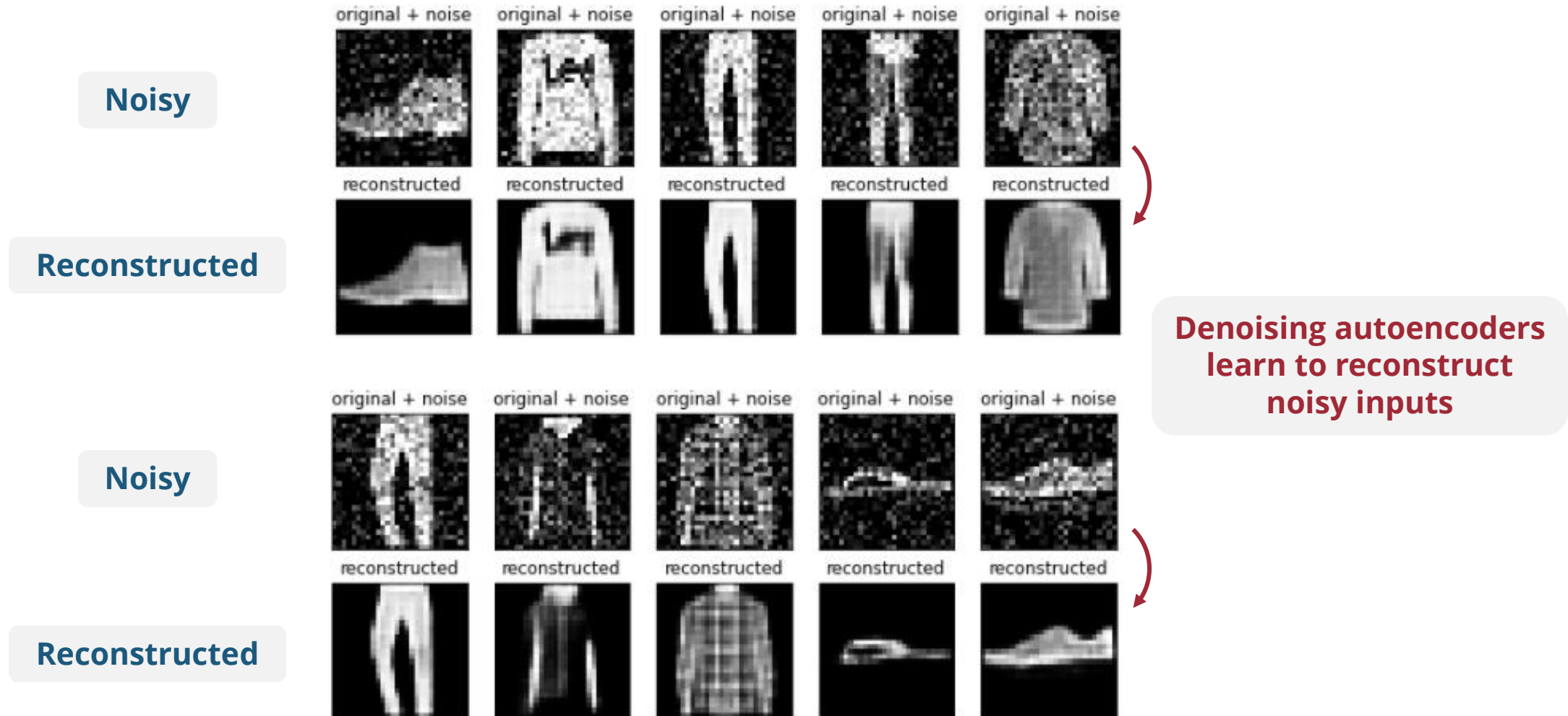


(Recap) Autoencoders – Reconstruction Examples



(Source: tensorflow.org)

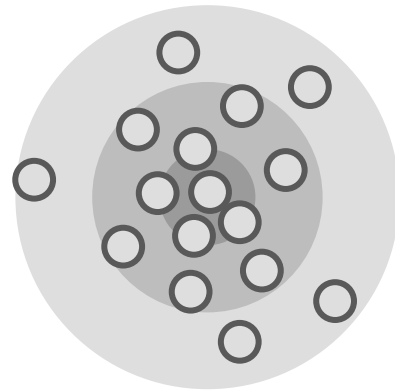
Denoising Autoencoders



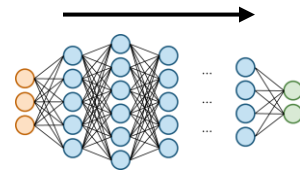
(Source: tensorflow.org)

(Recap) Generating Data from a Random Distribution

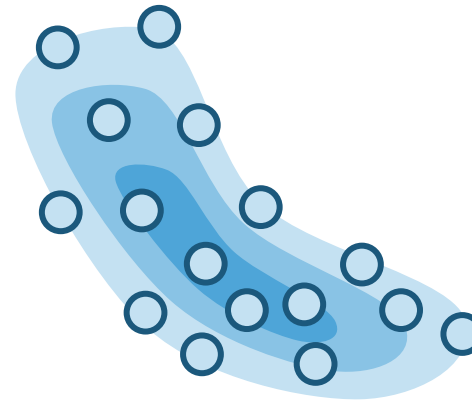
Random distribution



$P(z)$



Data distribution

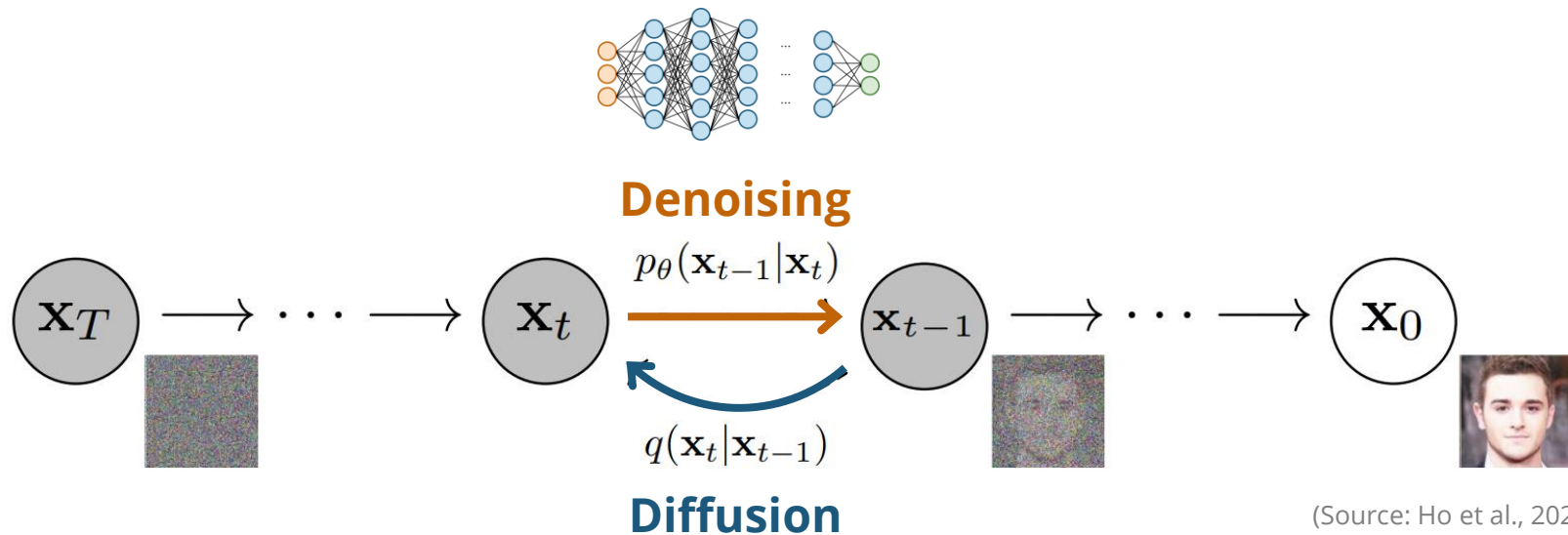


$P(x)$

If we can learn this mapping, we can easily generate new samples from the data distribution

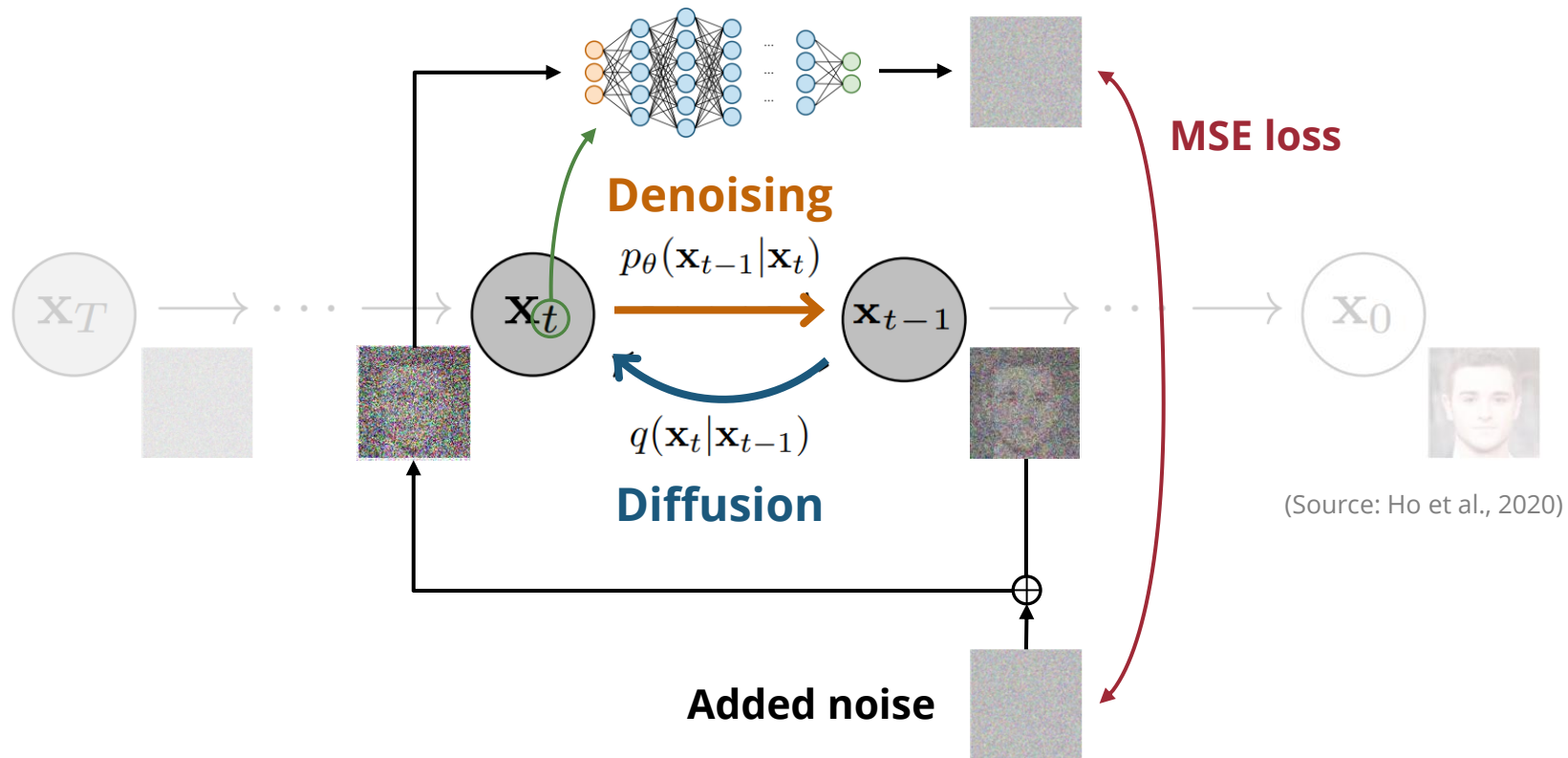
Diffusion Models

- **Intuition**: Many denoising autoencoders stacked together



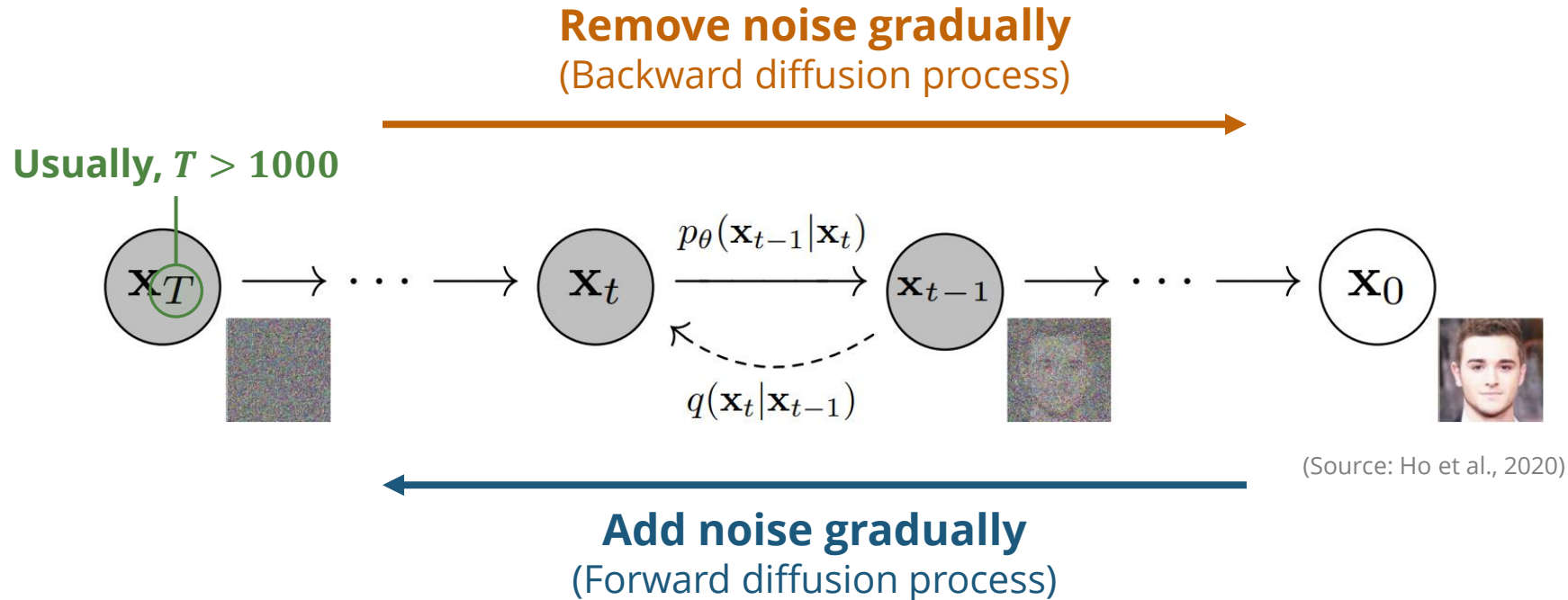
Diffusion Models – Training

- **Intuition:** Many denoising autoencoders stacked together

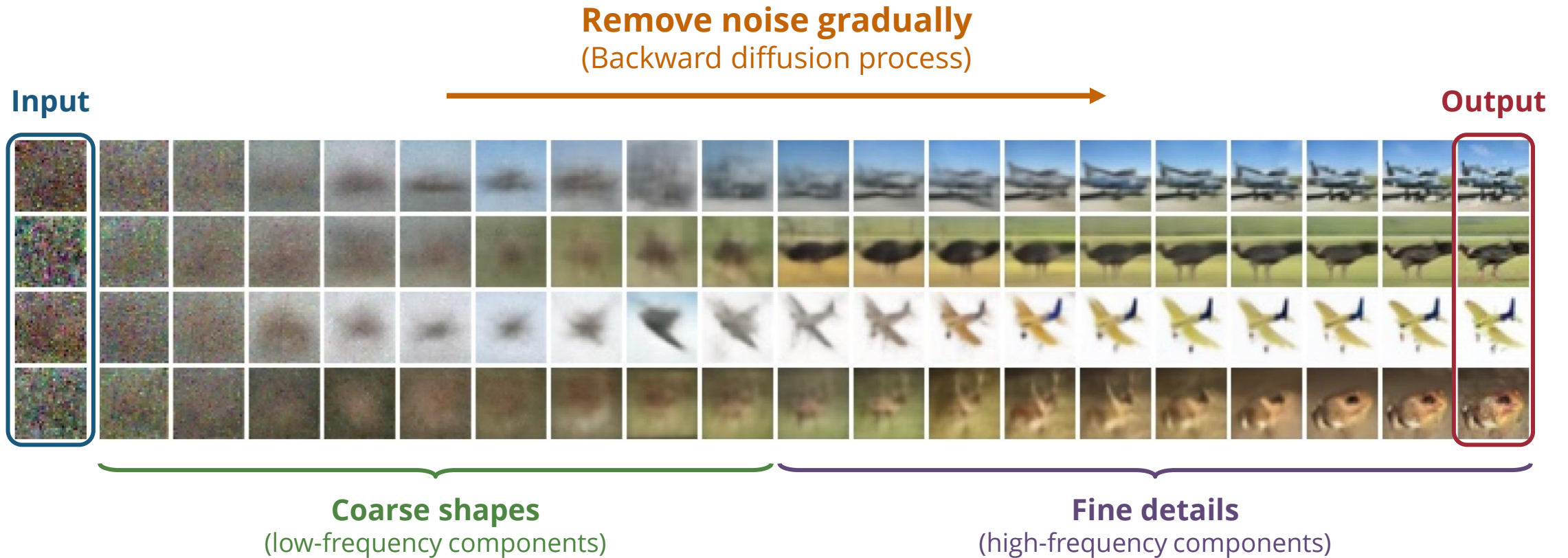


Diffusion Models

- **Intuition**: Many denoising autoencoders stacked together



Diffusion Models – Generation

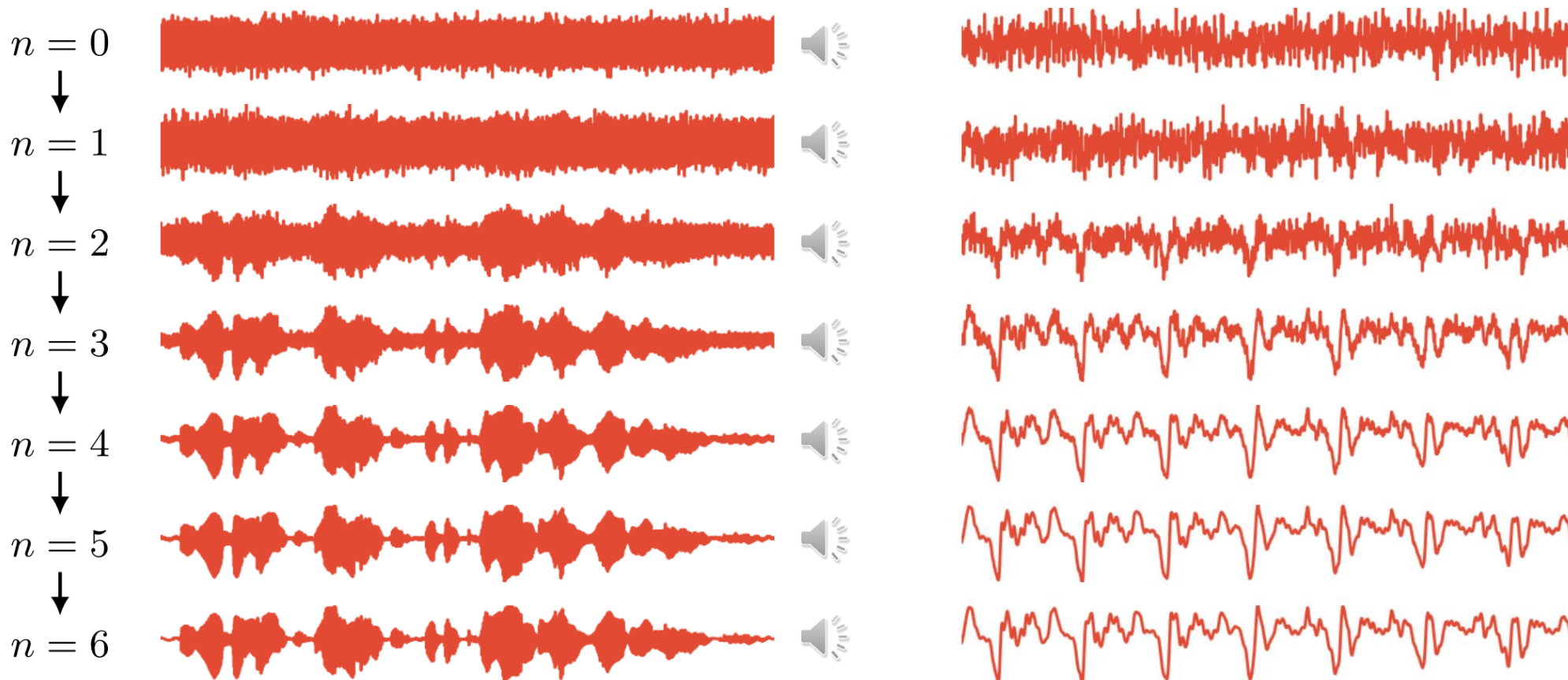


(Source: Ho et al., 2020)

WaveGrad – Diffusion Model for Waveforms

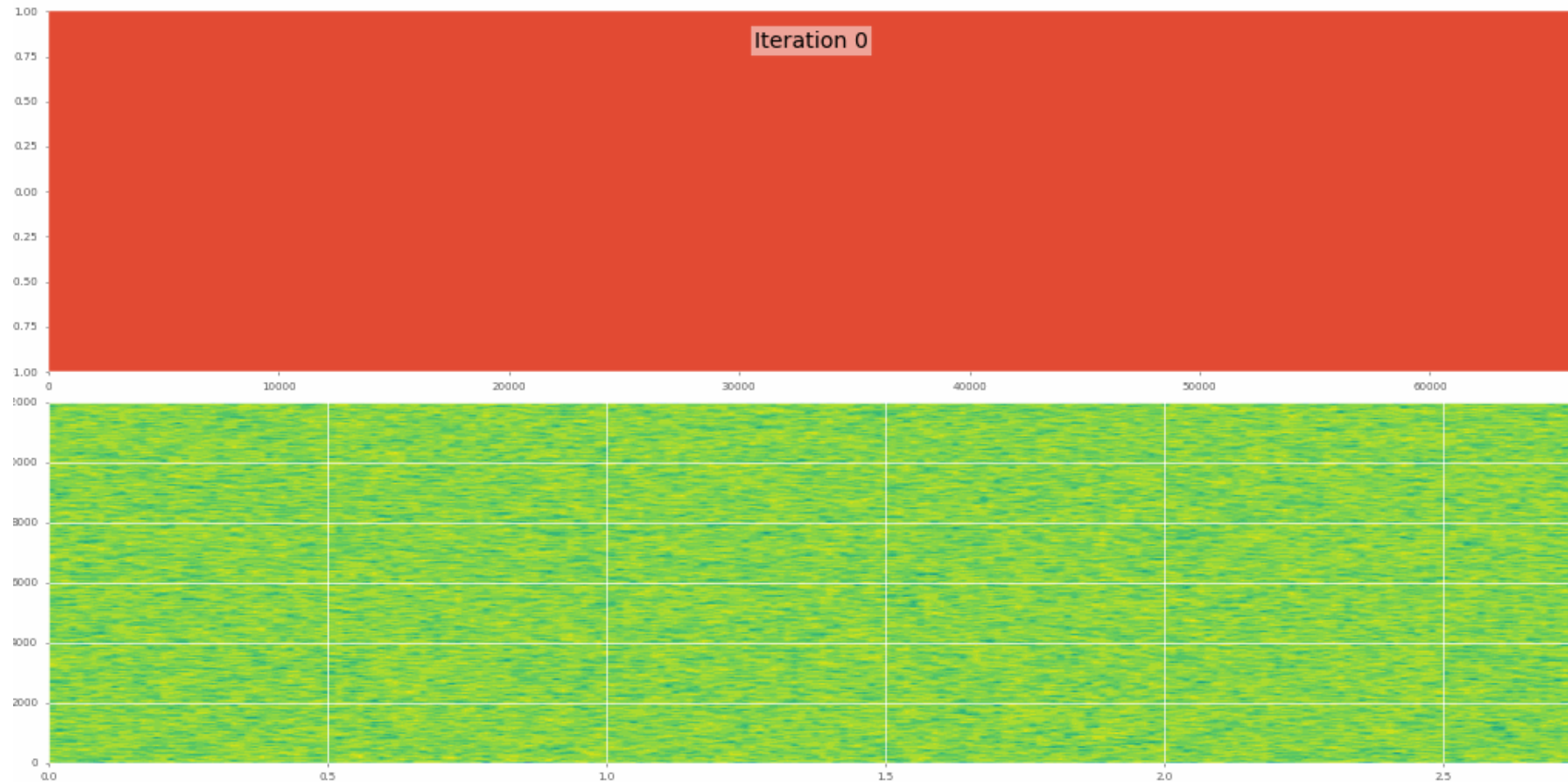
Text: Here are the match lineups for the Colombia Haiti match.

Zoom in



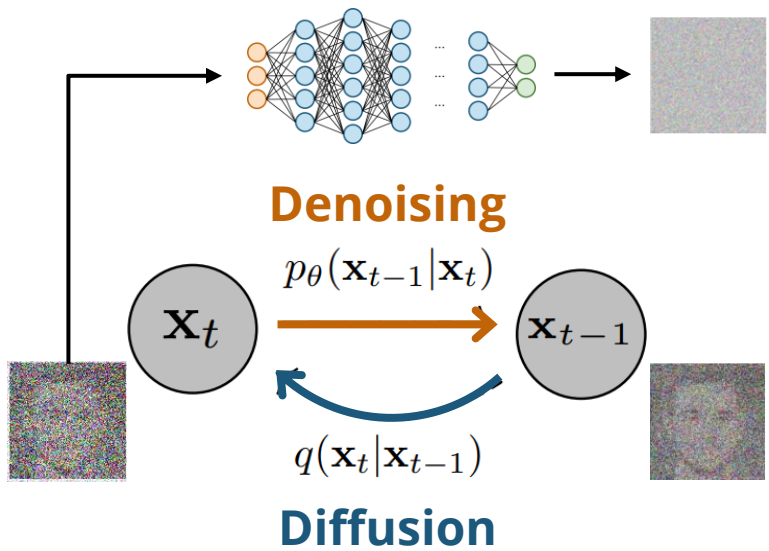
(Source: Chen et al., 2021)

WaveGrad – Diffusion for Waveforms

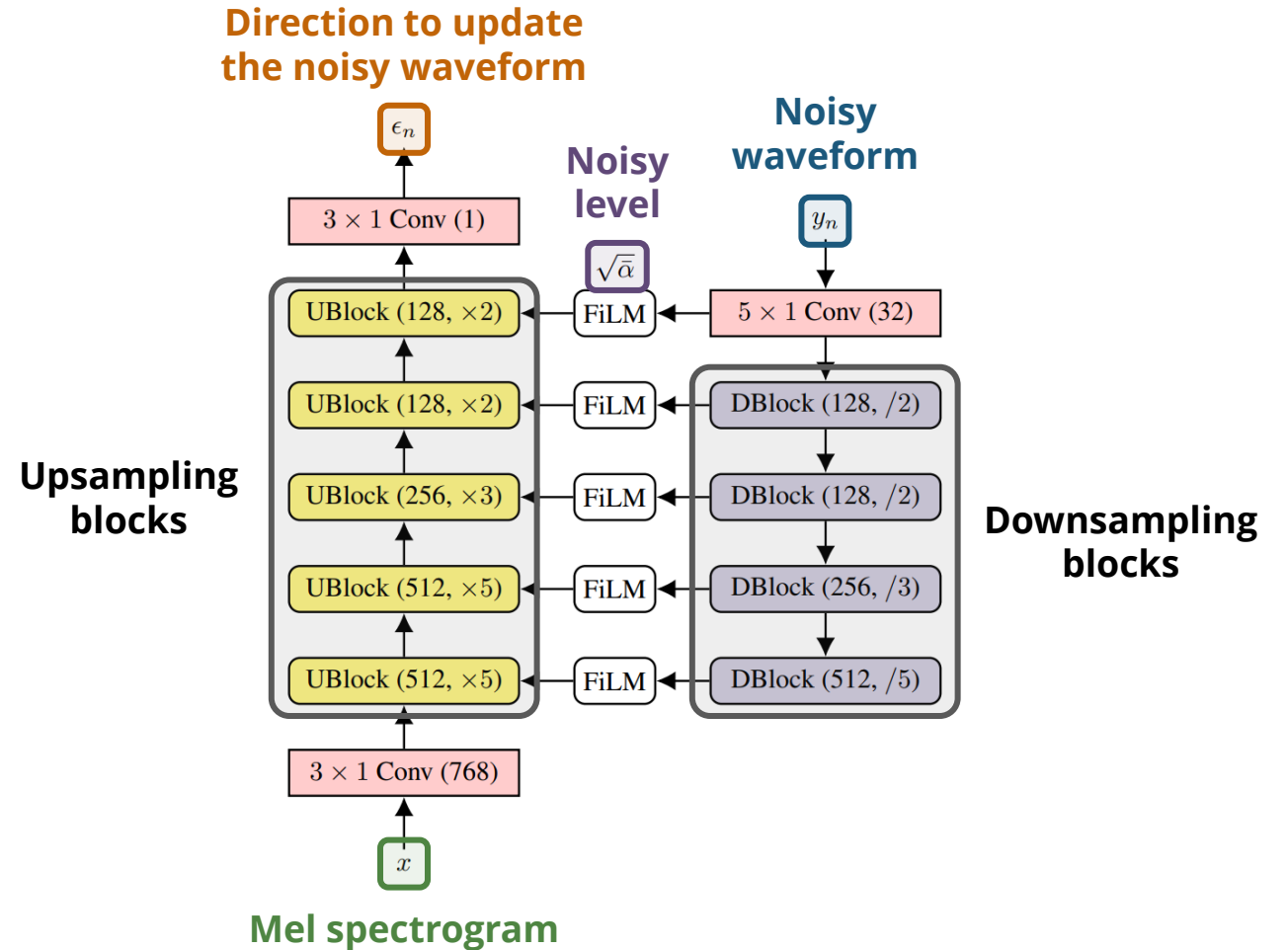


(Source: Chen et al., 2021)

WaveGrad – Diffusion for Waveforms

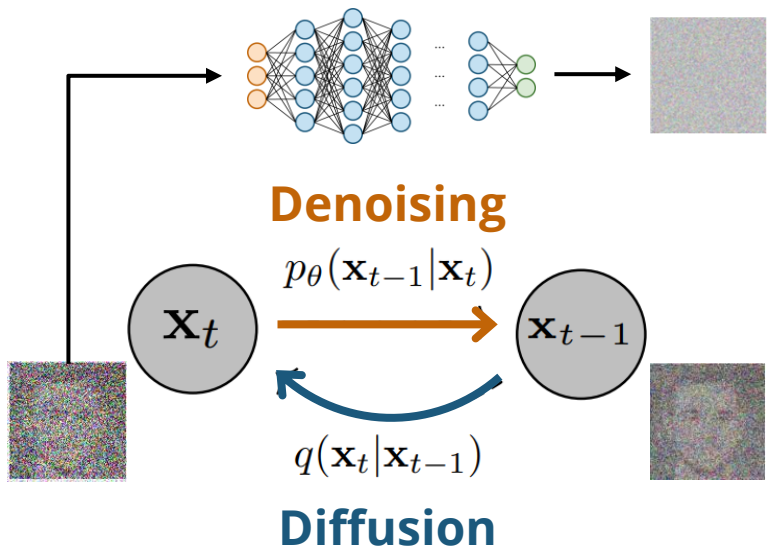


(Source: Ho et al., 2020)

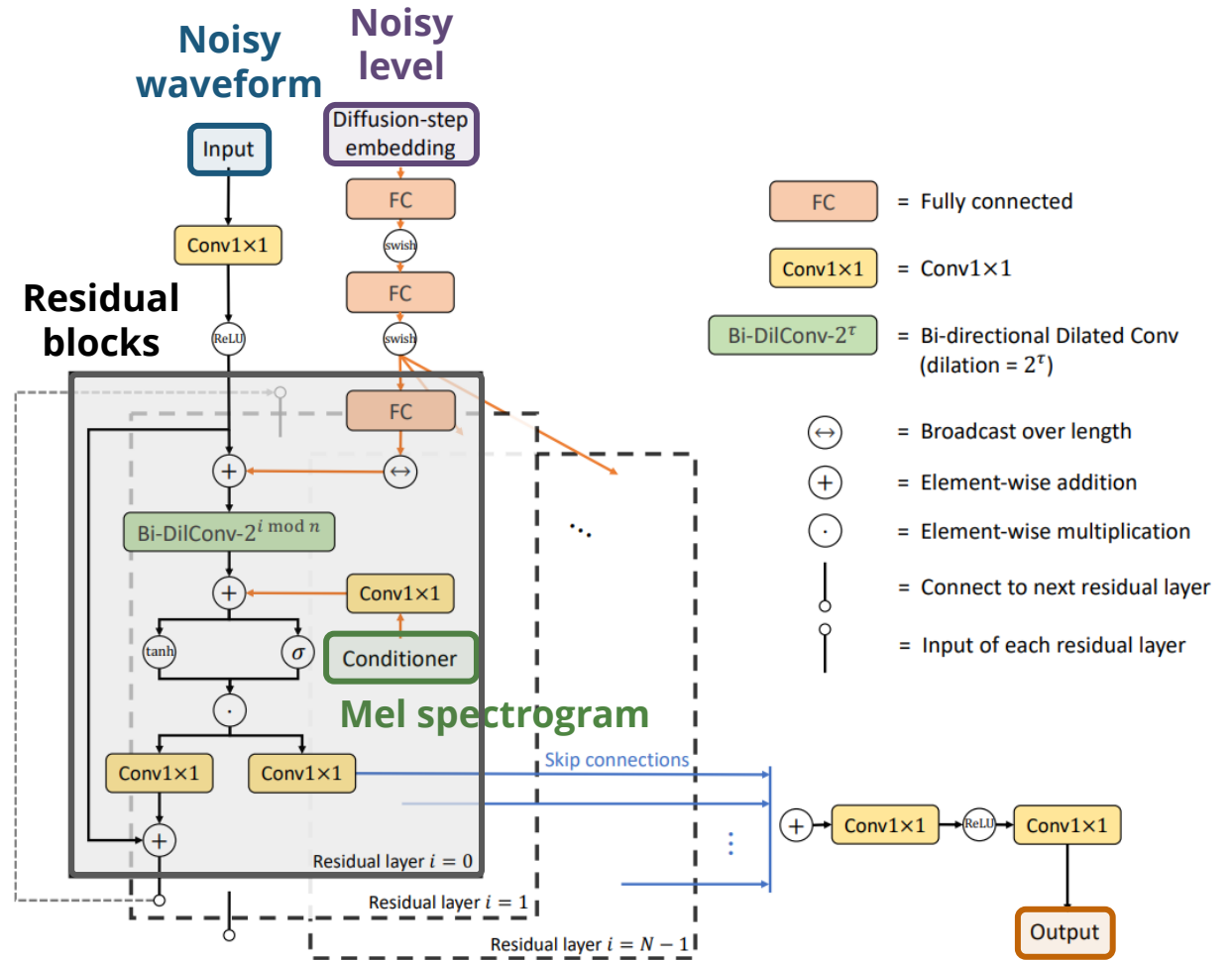


(Source: Chen et al., 2021)

DiffWave – Another Diffusion Model for Waveforms



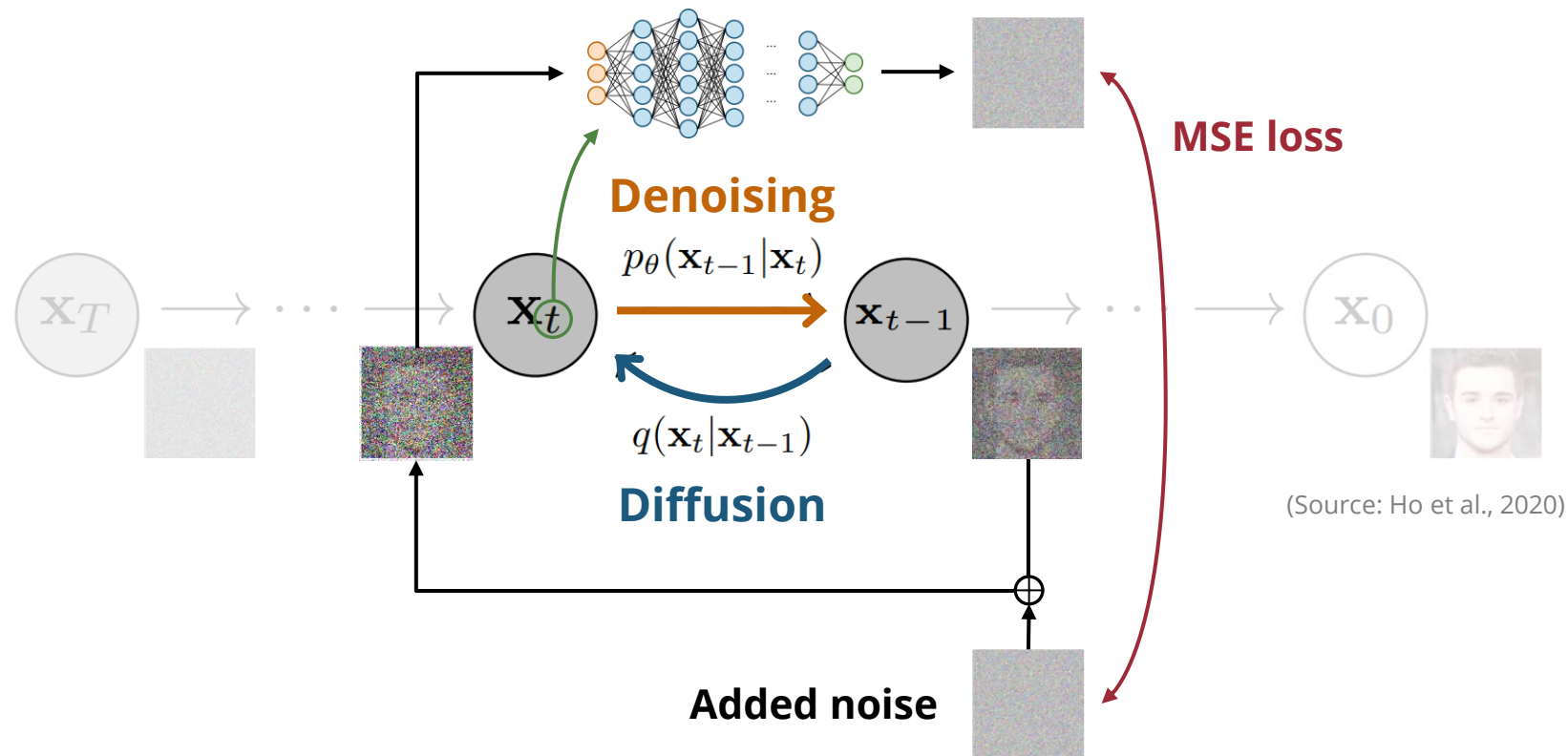
(Source: Ho et al., 2020)



(Source: Kong et al., 2021)

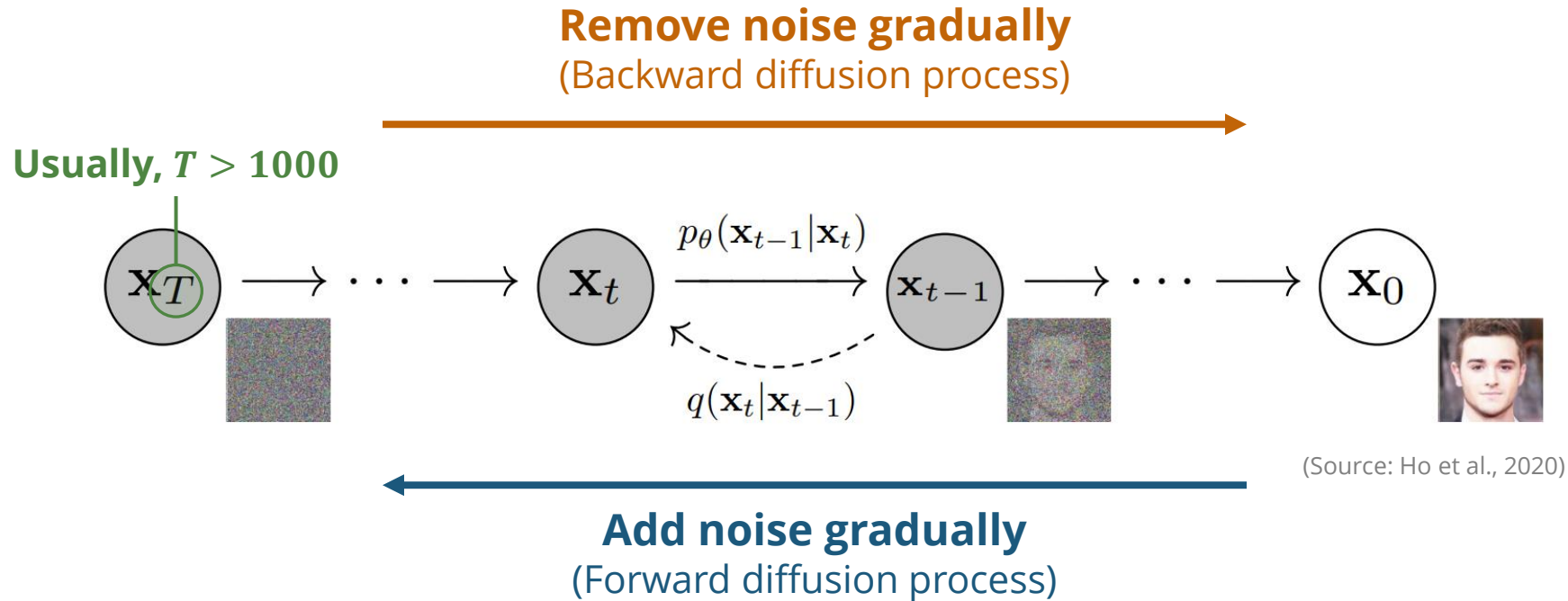
(Recap) Diffusion Models – Training

- **Intuition:** Many denoising autoencoders stacked together



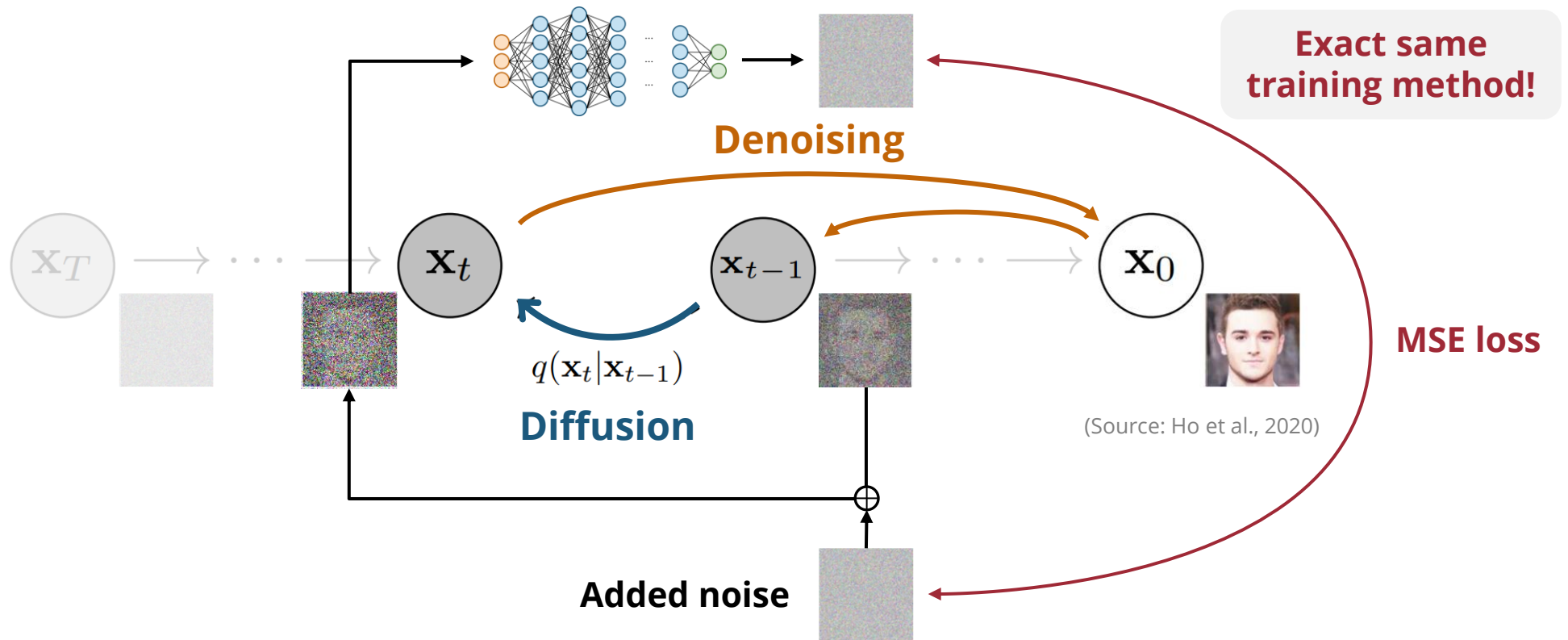
(Recap) Diffusion Models

- **Intuition**: Many denoising autoencoders stacked together



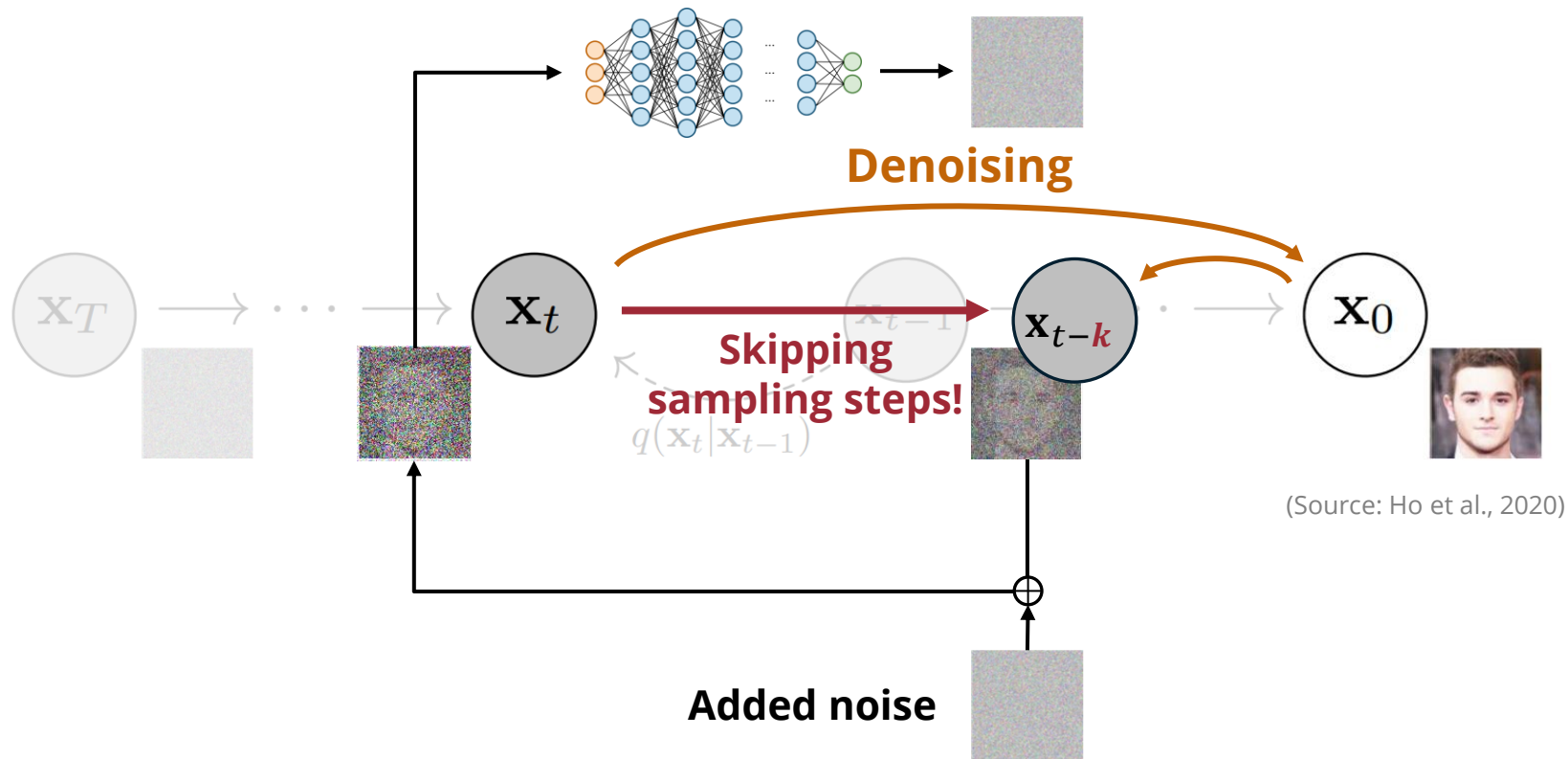
Fast Sampling for Diffusion Models

- **Intuition:** Skip some sampling steps



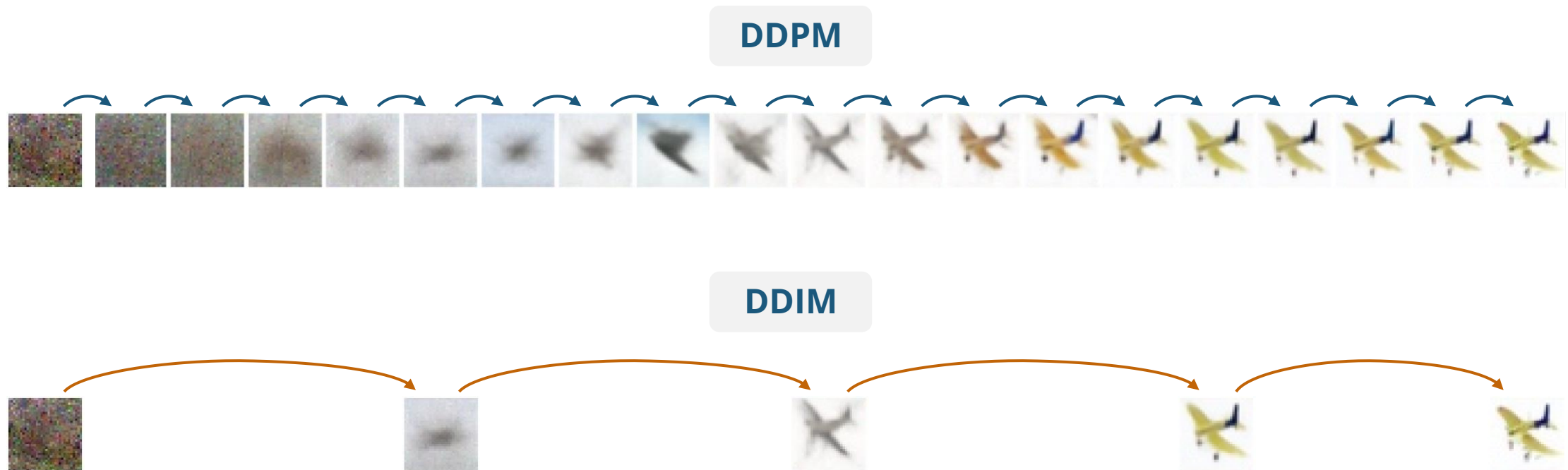
Fast Sampling for Diffusion Models

- **Intuition**: Skip some sampling steps



Fast Sampling for Diffusion Models

- **Intuition**: Skip some sampling steps

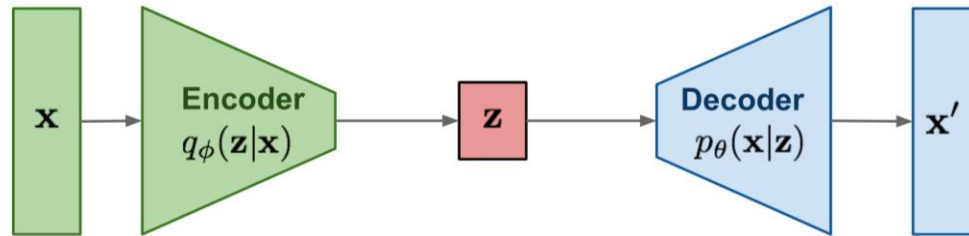


(Source: Ho et al., 2020)

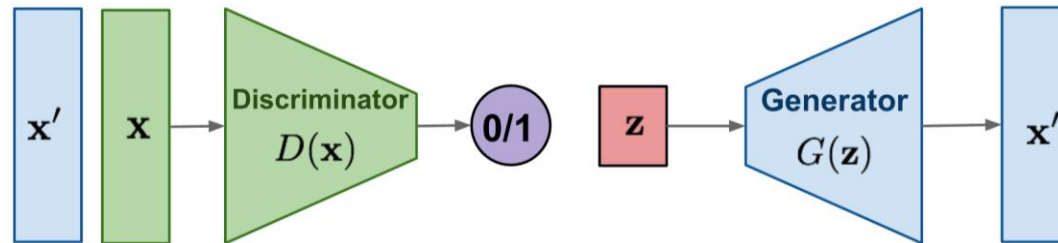
Comparison of Deep Generative Models

Comparison of Deep Generative Models

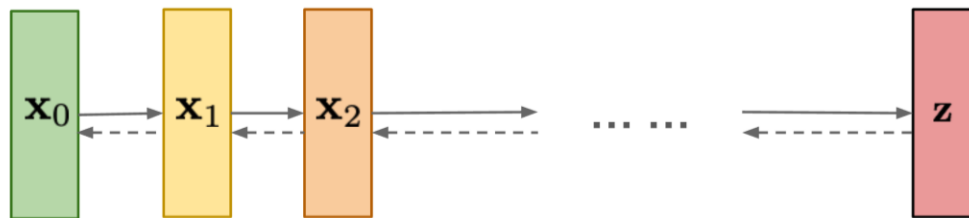
VAE: maximize variational lower bound



GAN: Adversarial training

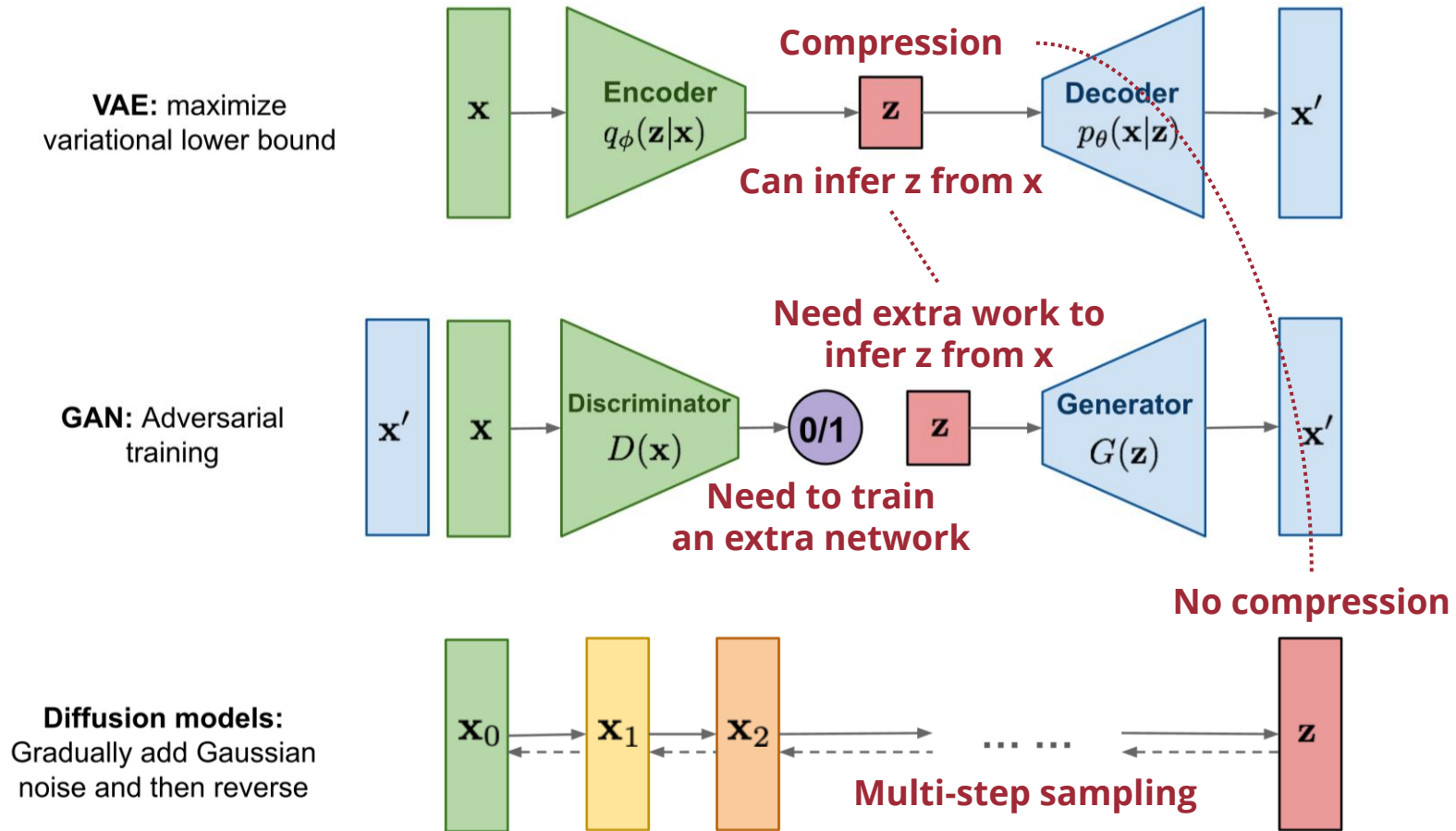


Diffusion models:
Gradually add Gaussian noise and then reverse



(Source: Weng, 2021)

Comparison of Deep Generative Models



(Source: Weng, 2021)

Network Architectures vs Training Frameworks

Network architectures

Multilayer perceptron (MLP)

Convolutional neural networks (CNNs)

Recurrent neural networks (RNNs)

Transformers

ResNets

U-Nets

⋮

Training frameworks

Autoregressive

Autoencoders

Variational autoencoders (VAEs)

Generative adversarial networks (GANs)

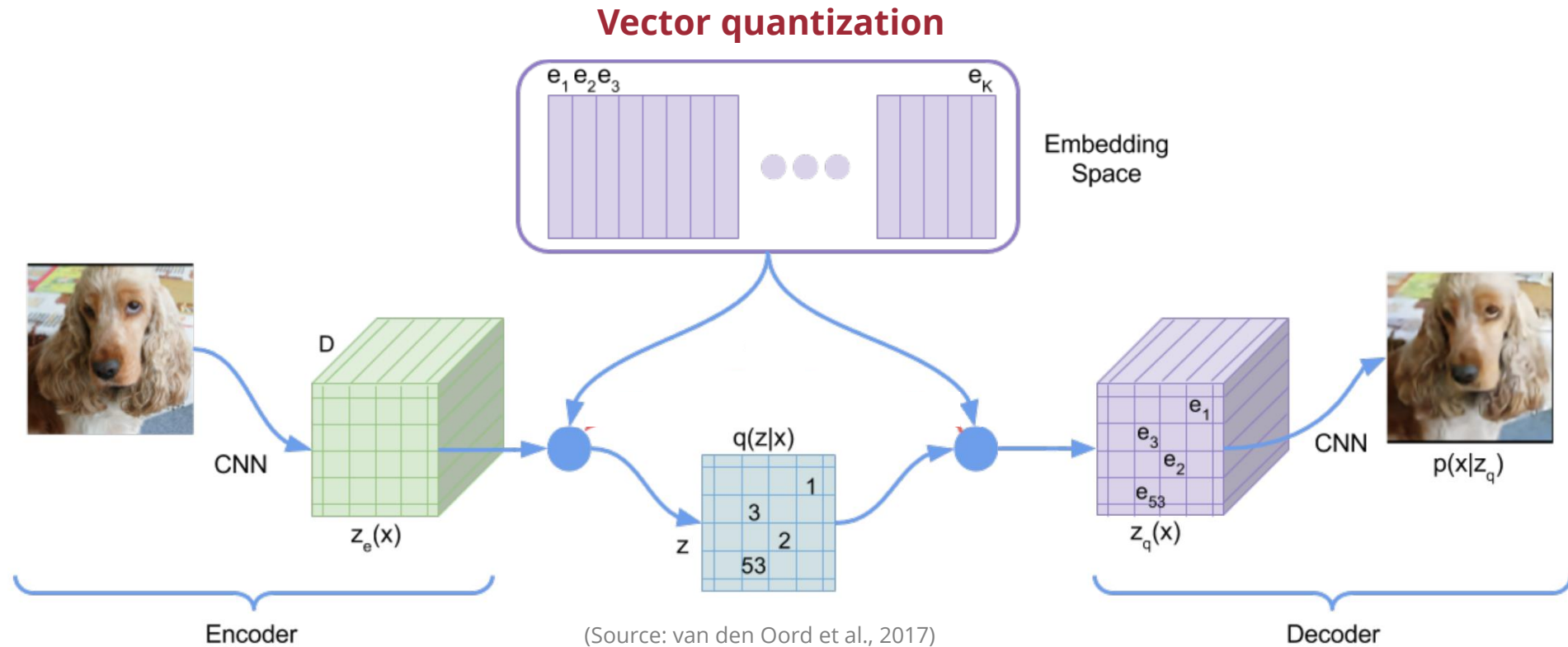
Diffusion models

Consistency models

⋮

Latent Diffusion Models

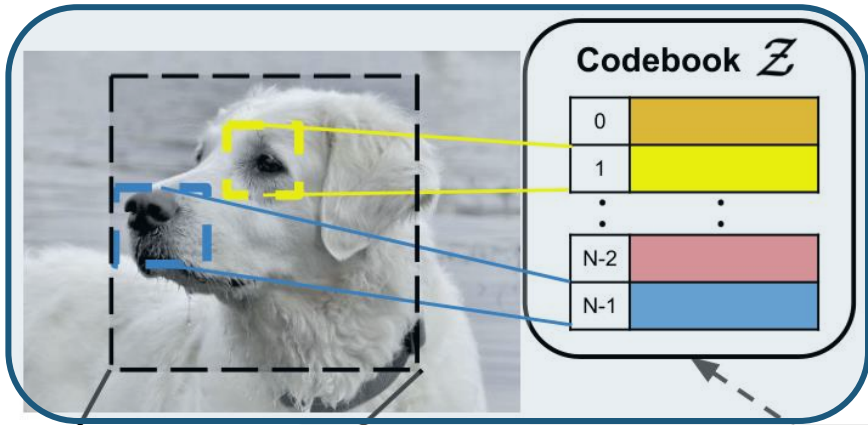
(Recap) Vector-Quantized VAEs (VQVAEs)



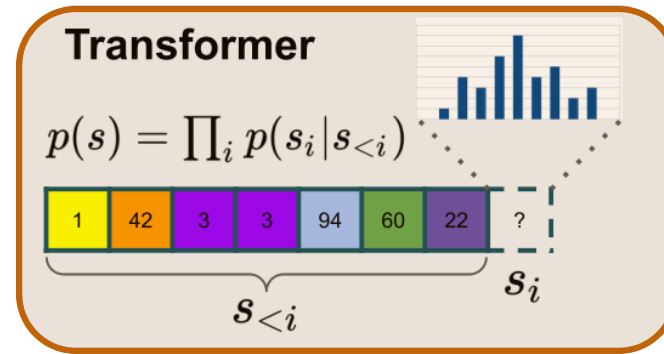
Allow only a fixed number of vectors to be used in the bottleneck layer

VQGAN

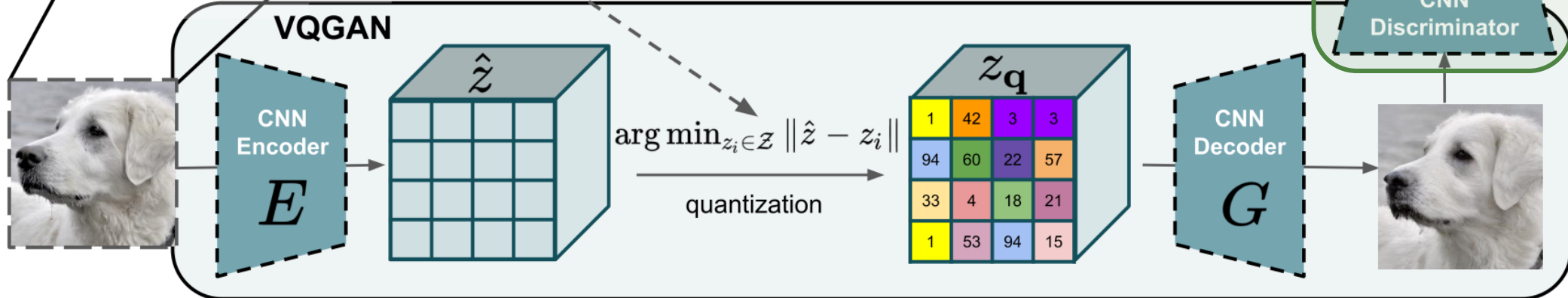
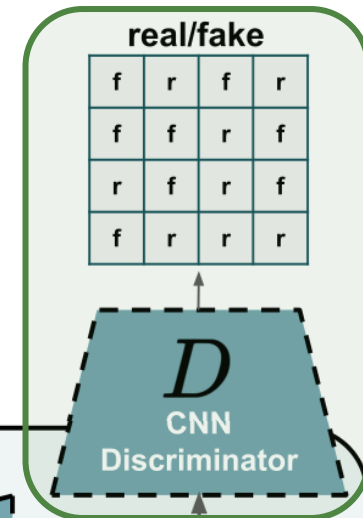
Each path is encoded into a latent code



A transformer-based language model trained with the latent codes



Patch discriminator



(Source: Esser et al., 2021)

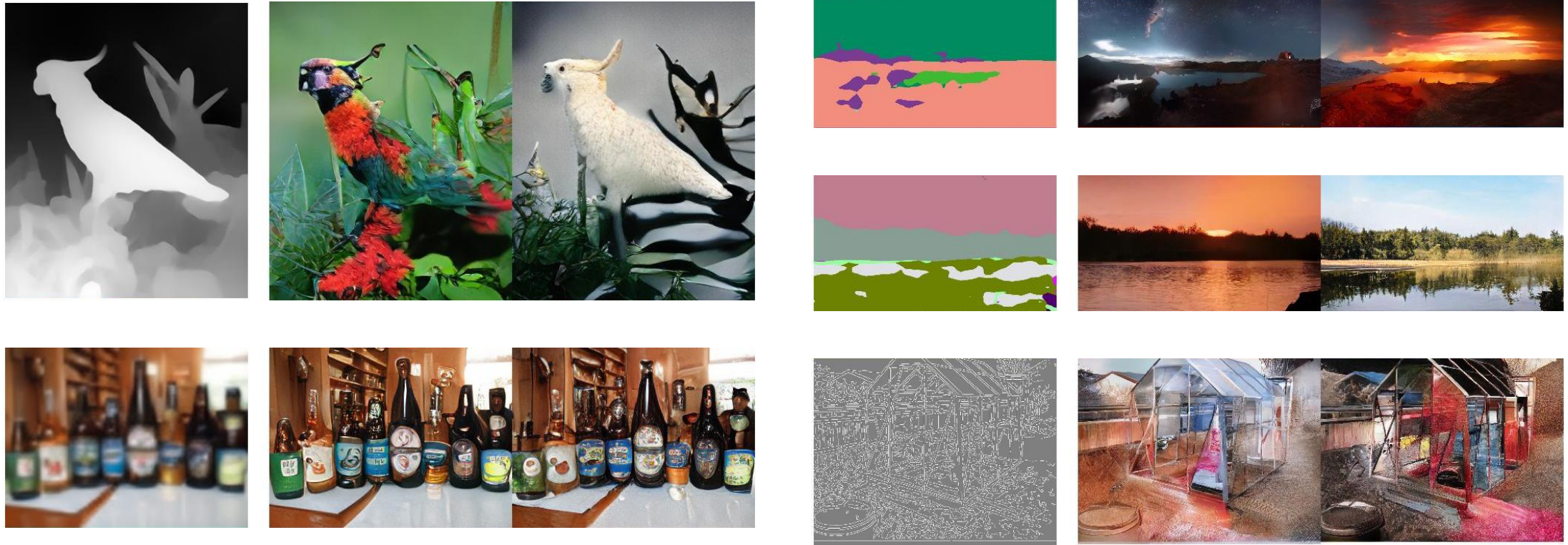
A VQGAN is a VQVAE equipped with adversarial loss

VQGAN – Conditional Generation



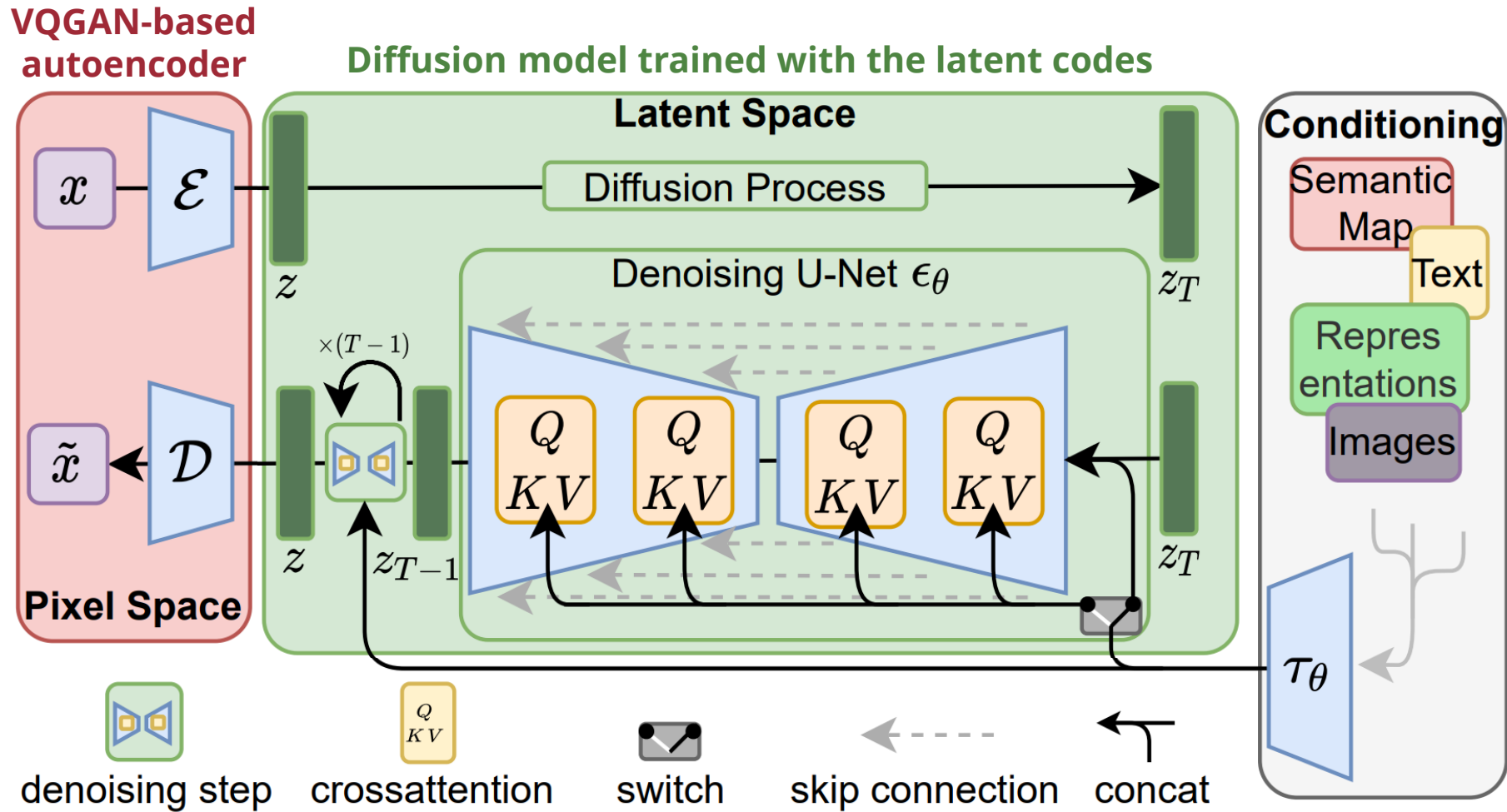
(Source: Esser et al., 2021)

VQGAN – Conditional Generation



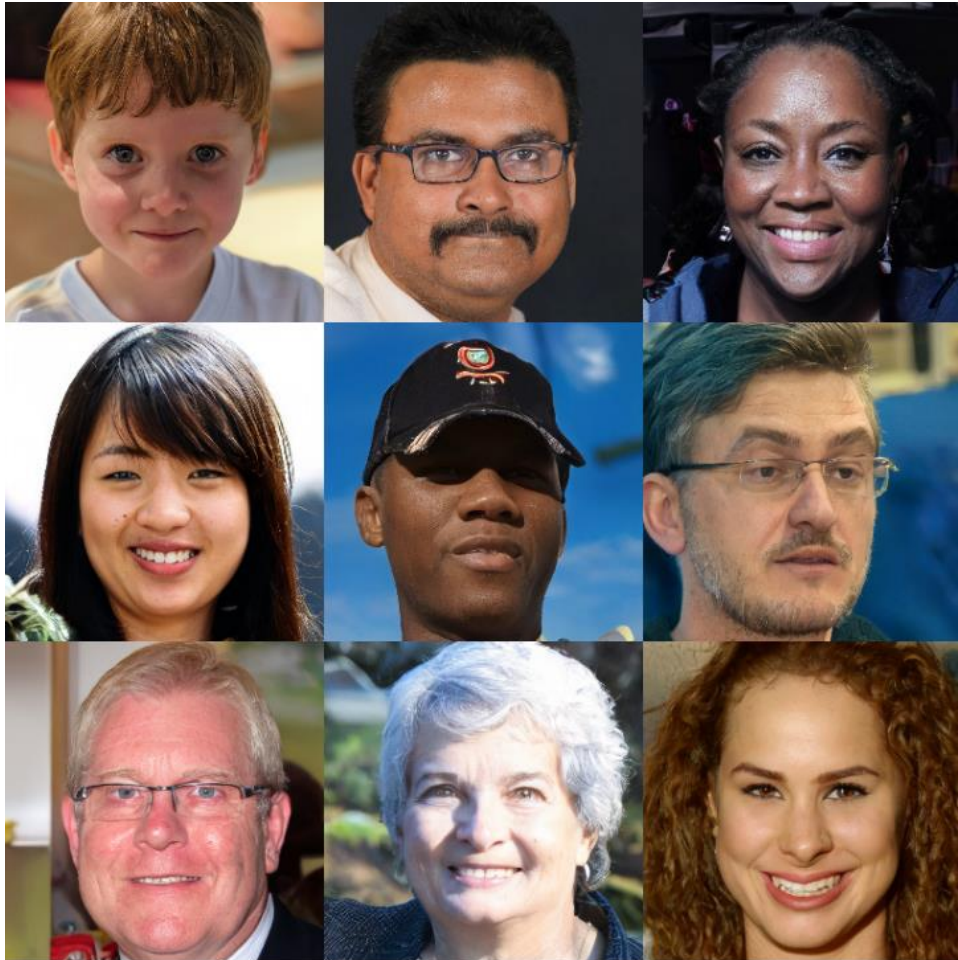
(Source: Esser et al., 2021)

Latent Diffusion Models (LDMs)



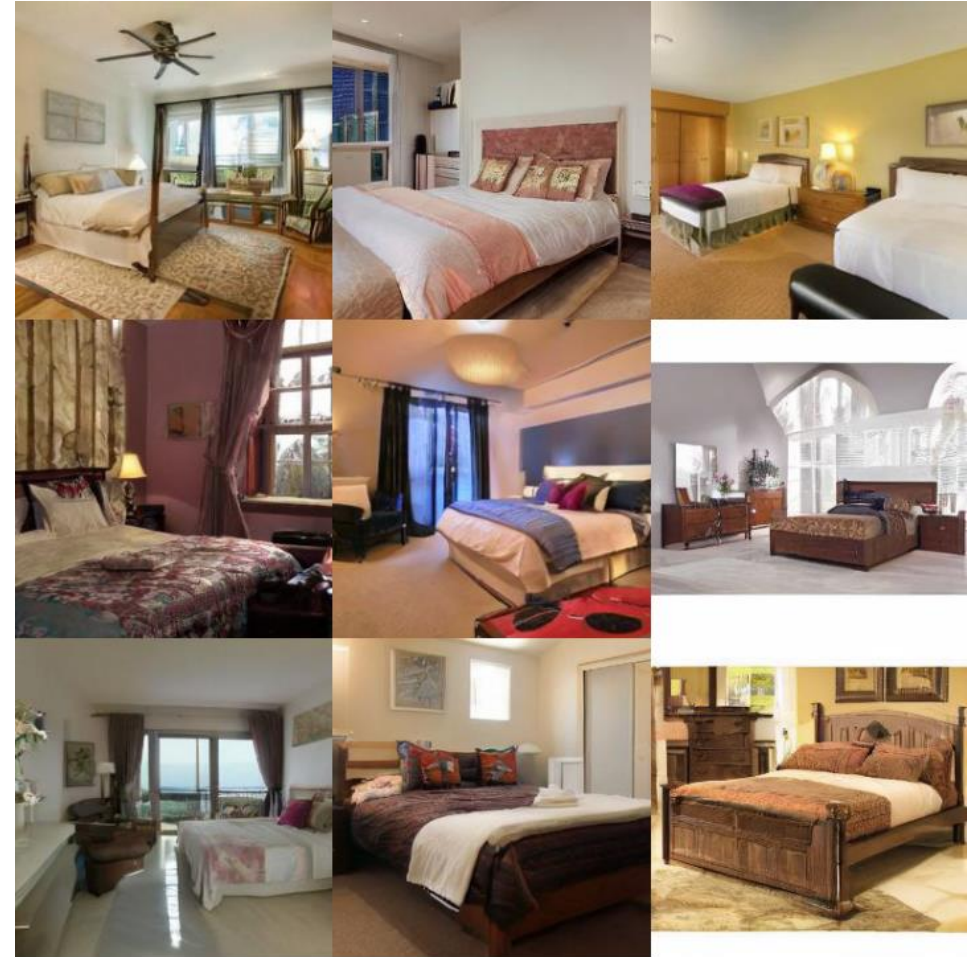
(Source: Rombach et al., 2022)

LDMs – Examples



(Source: Rombach et al., 2022)

LDMs – Examples



(Source: Rombach et al., 2022)

LDMs – Semantic Synthesis



(Source: Rombach et al., 2022)



(Recap) Network Architectures vs Training Frameworks

Network architectures

Multilayer perceptron (MLP)

Convolutional neural networks (CNNs)

Recurrent neural networks (RNNs)

Transformers

ResNets

U-Nets

⋮

Training frameworks

Autoregressive

Autoencoders

Variational autoencoders (VAEs)

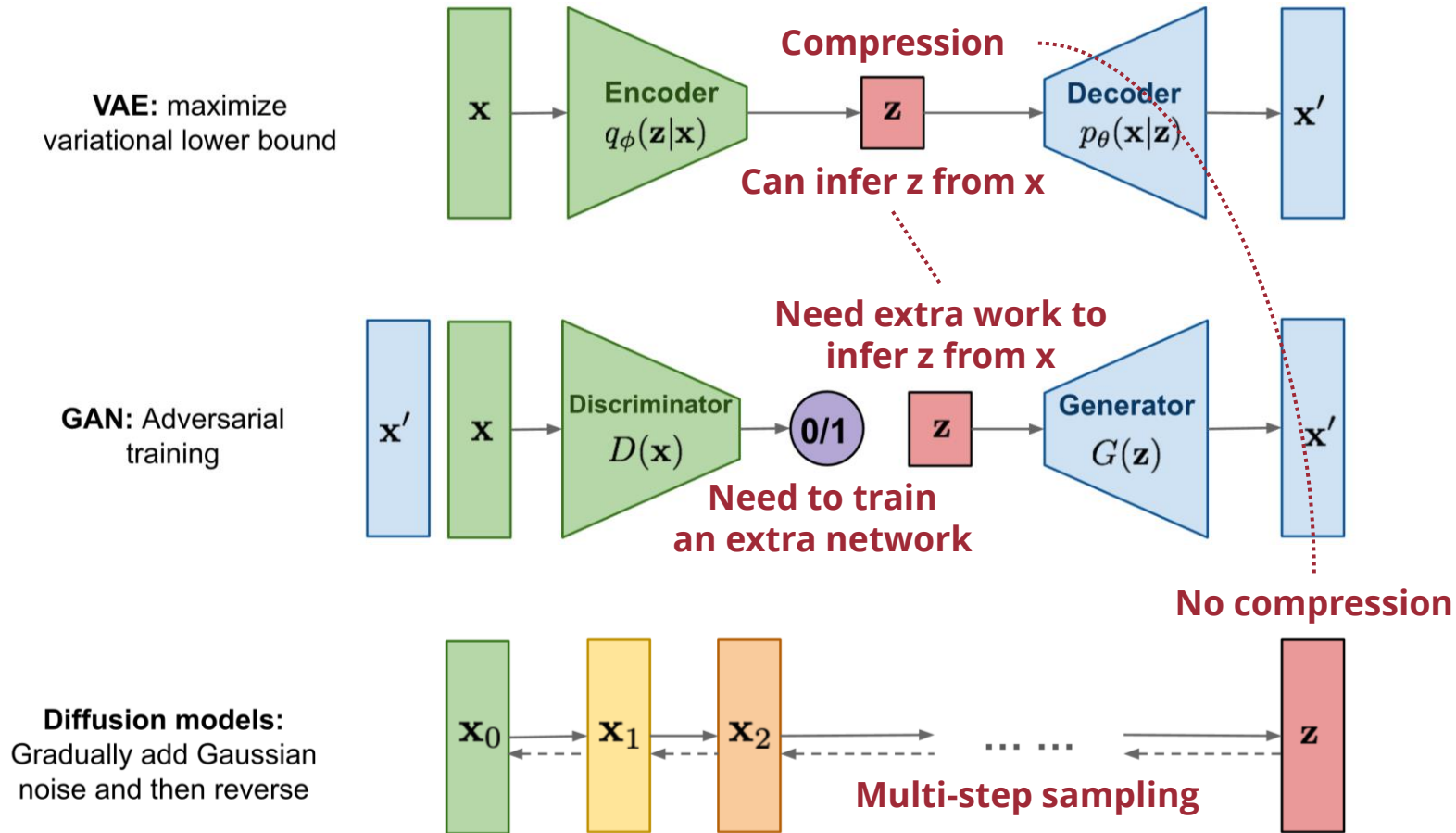
Generative adversarial networks (GANs)

Diffusion models

Consistency models

⋮

(Recap) Comparison of Deep Generative Models



(Source: Weng, 2021)