

PAT 498/598 (Fall 2024)

# Special Topics: Generative AI for Music and Audio Creation

## Lecture 9: VAEs & GANs

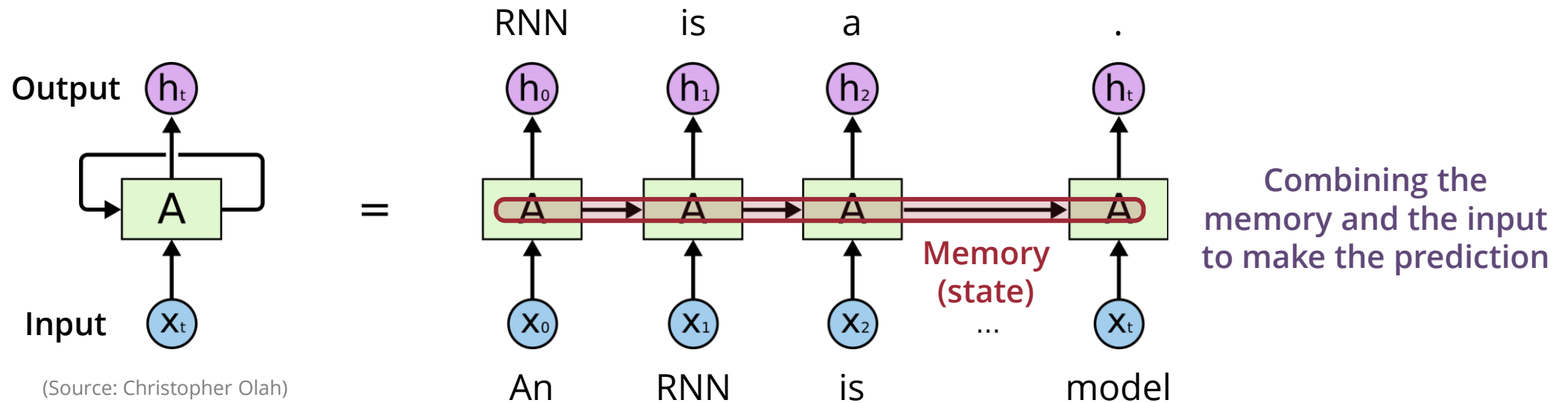
Instructor: Hao-Wen Dong



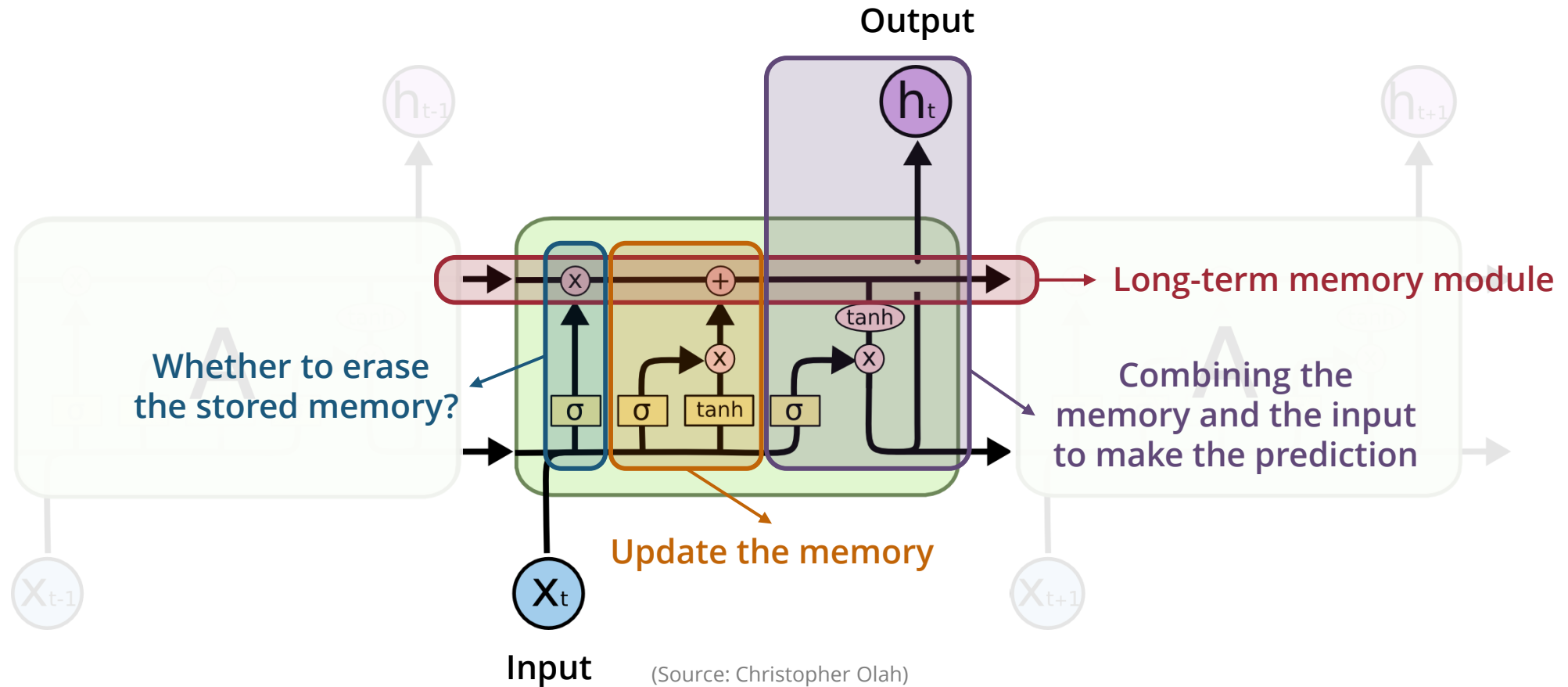
SCHOOL OF MUSIC, THEATRE & DANCE  
PERFORMING ARTS TECHNOLOGY  
UNIVERSITY OF MICHIGAN

# (Recap) What is an RNN (Recurrent Neural Network)?

- A type of neural networks that have **loops**
- Widely used for **modeling sequences** (e.g., in natural language processing)

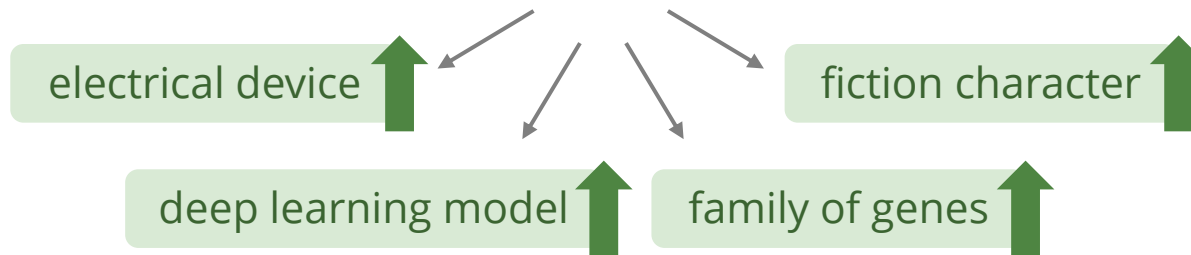


# (Recap) Demystifying LSTMs

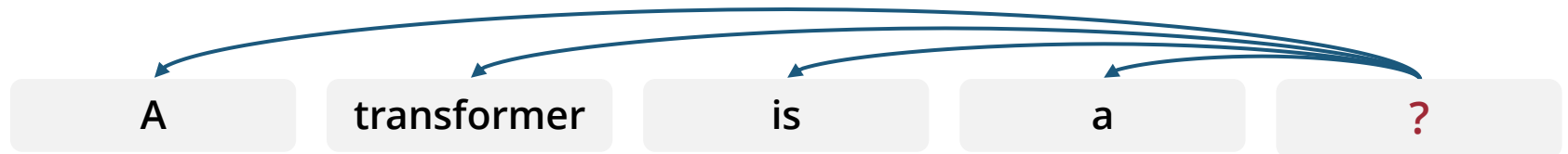


# (Recap) Demystifying Transformers

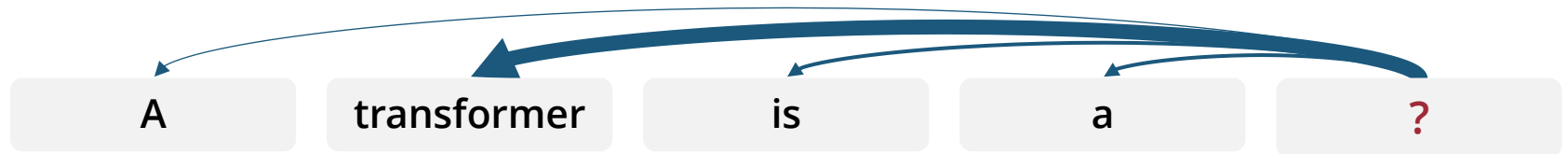
A transformer is a \_\_\_\_\_



Uniform attention



Variable attention

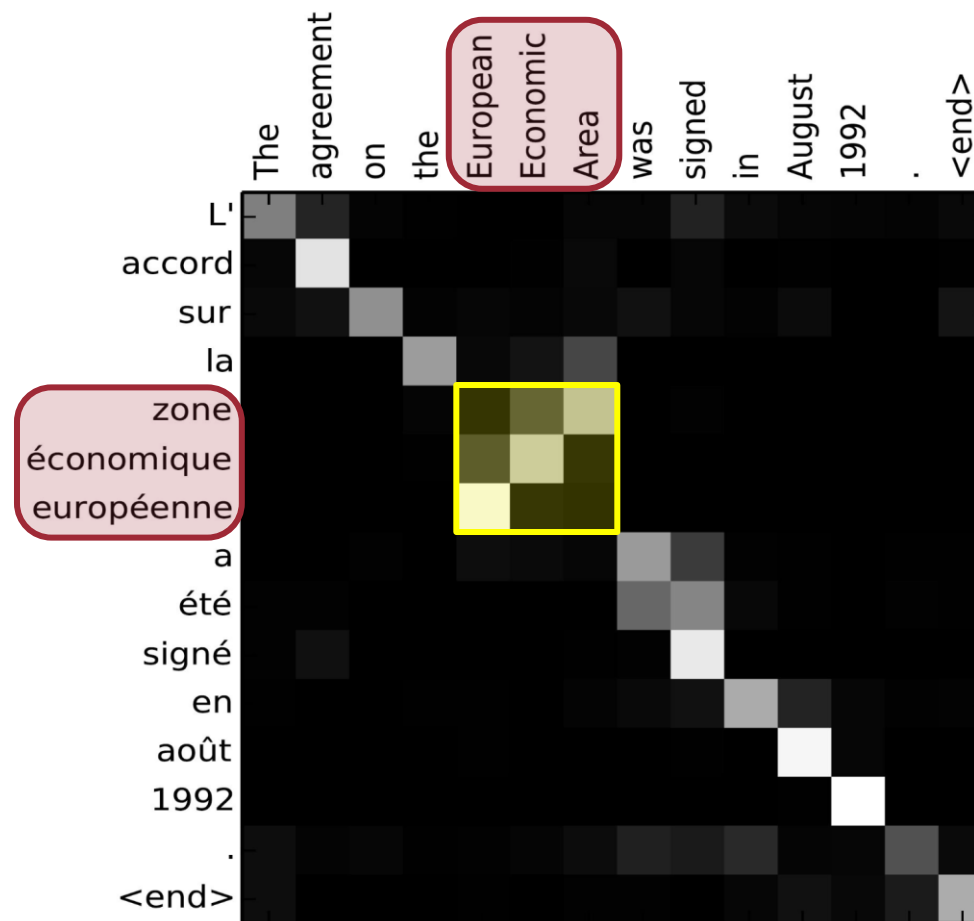


Transformers learn what to attend to from big data!

# (Recap) What does a Transformer Learn?

The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .

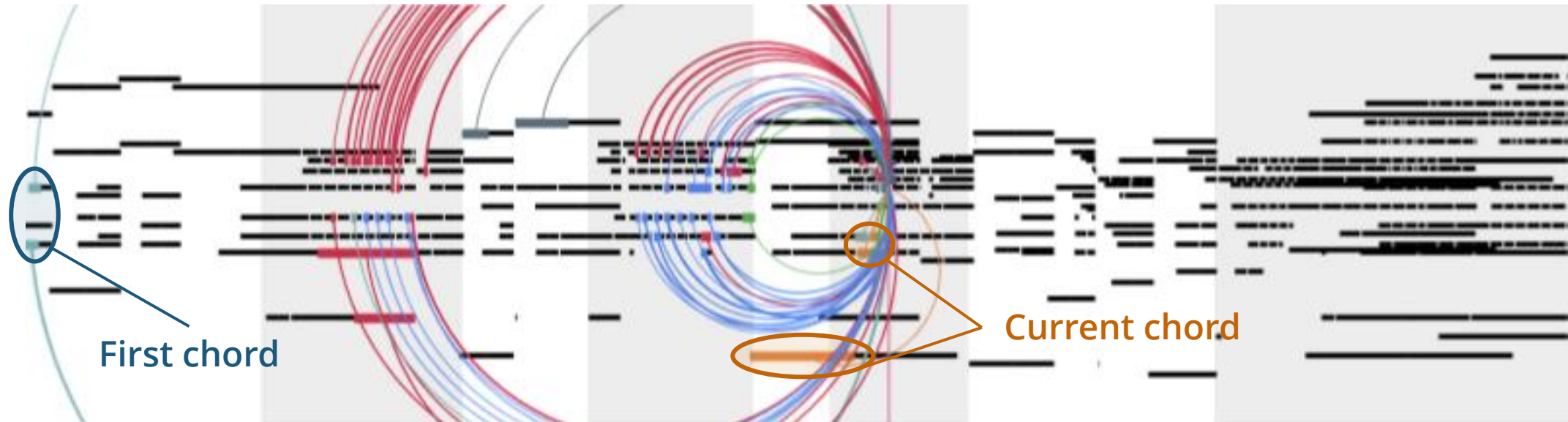
(Source: Cheng et al., 2016)



(Source: Bahdanau et al., 2015)

# (Recap) What does a Transformer Learn?

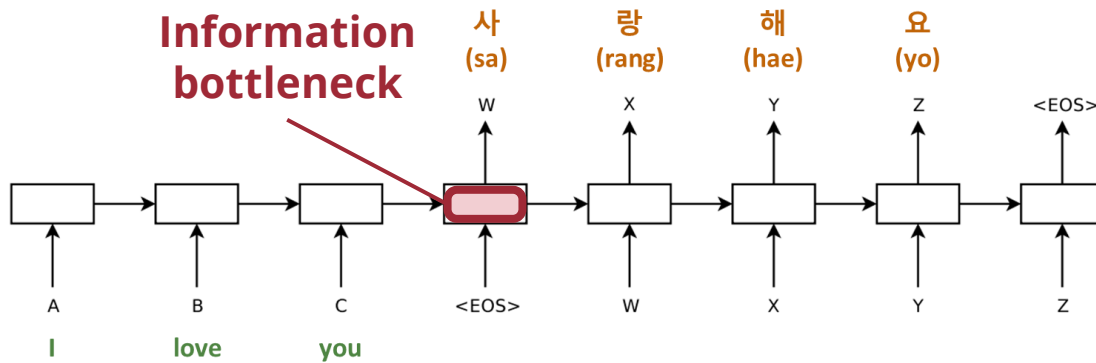
(Each color represents an attention head)



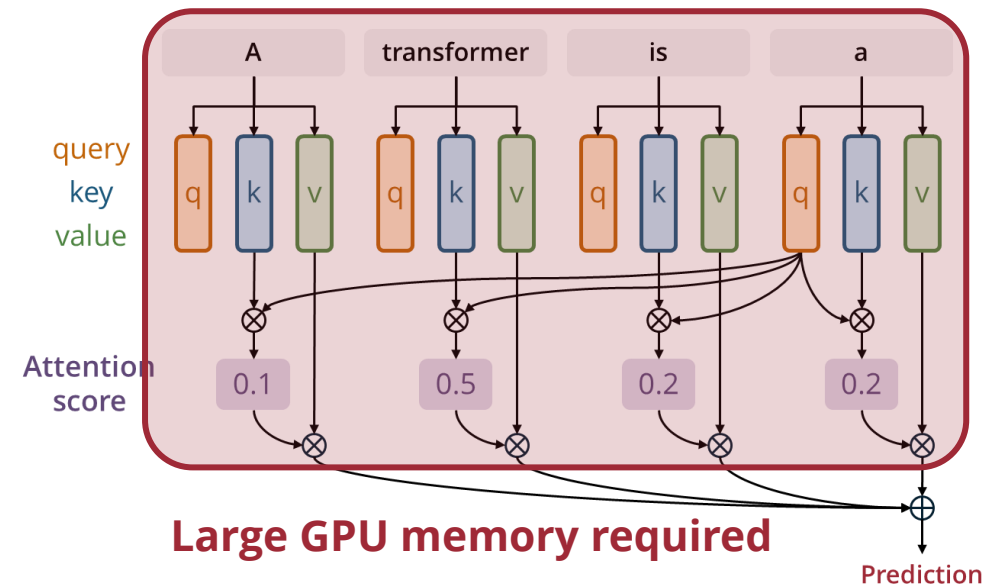
(Source: Huang et al., 2018)

# (Recap) Seq2seq vs Transformers

## Seq2seq



## Transformers

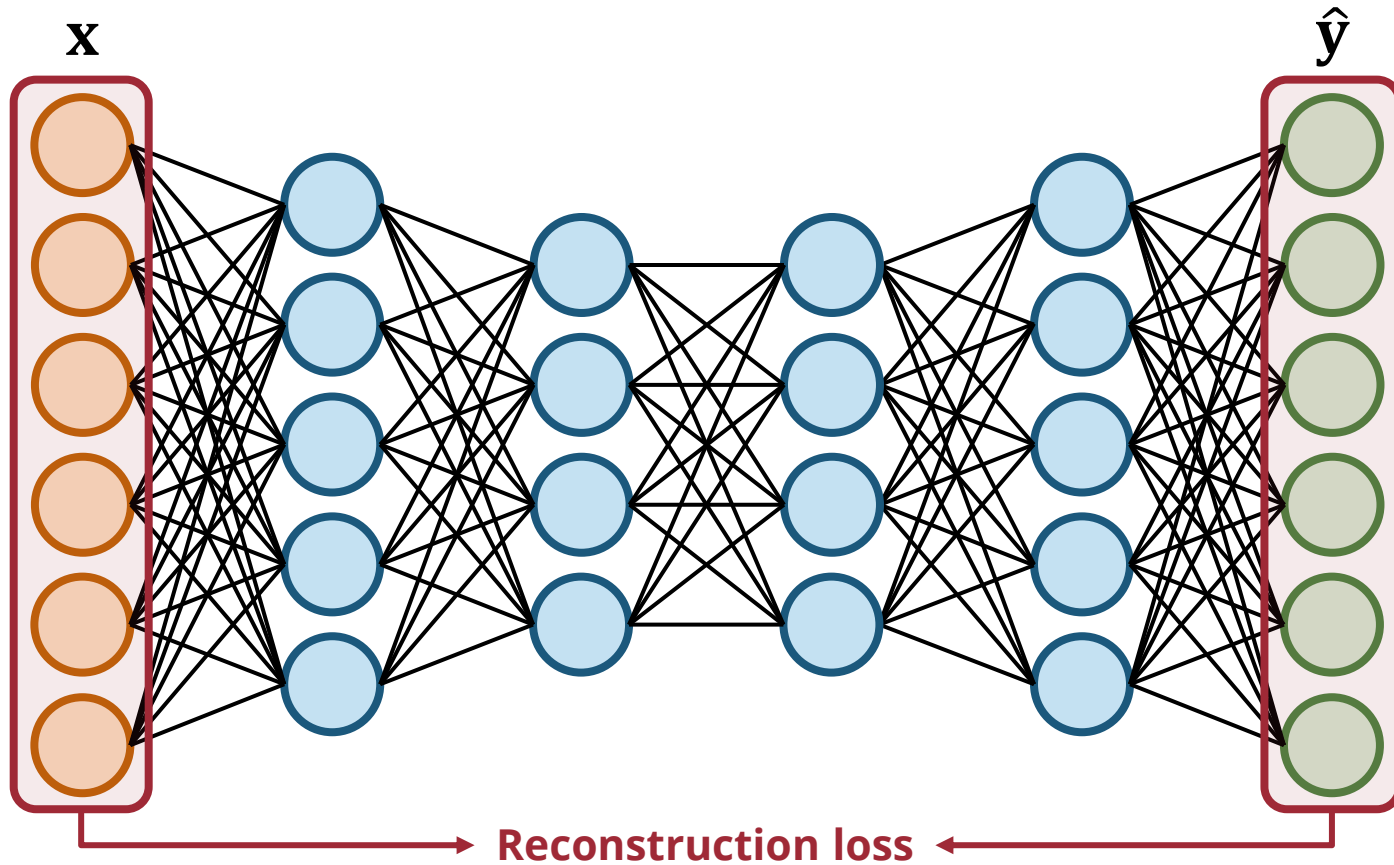


# Autoencoders



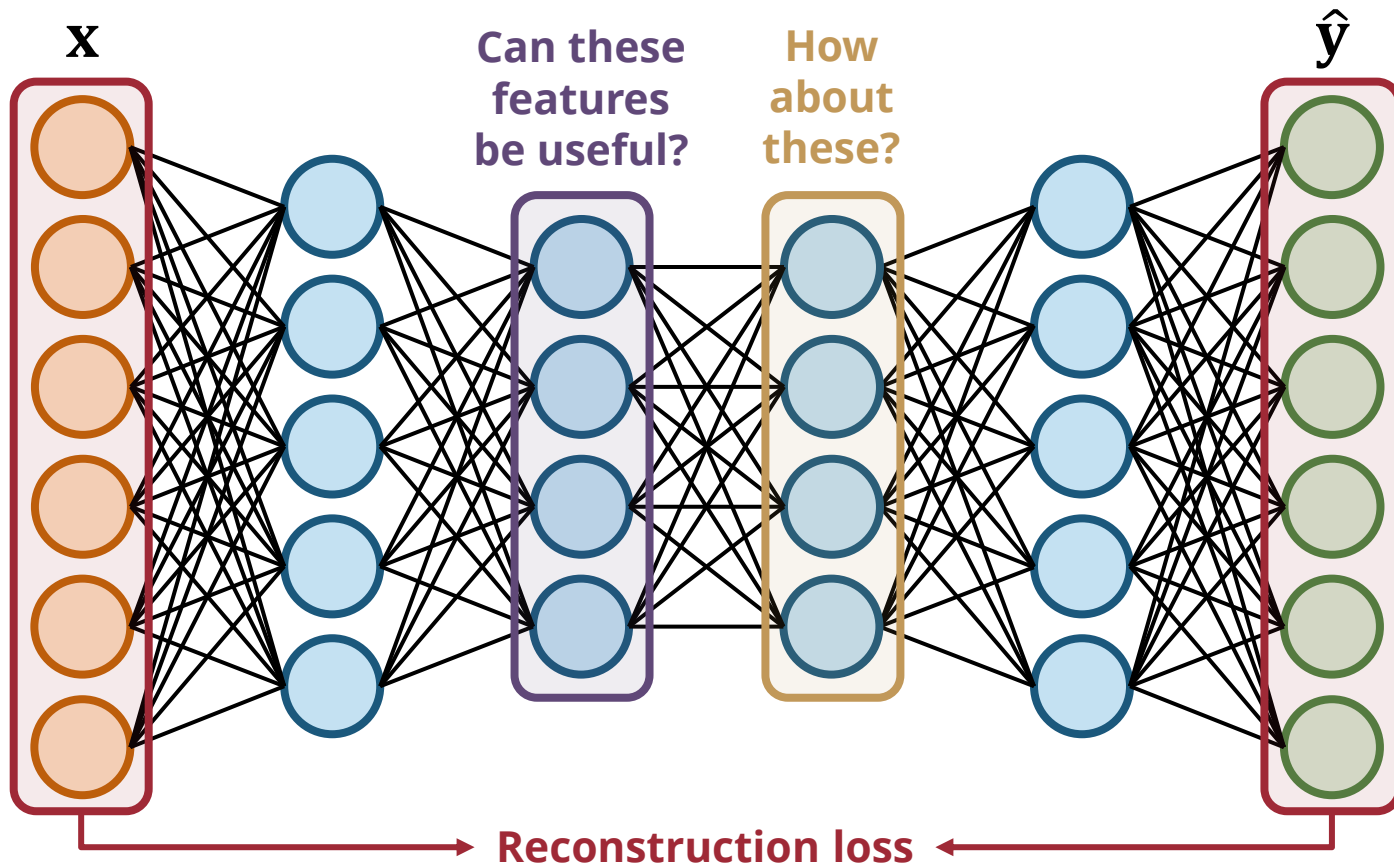
# Autoencoders

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



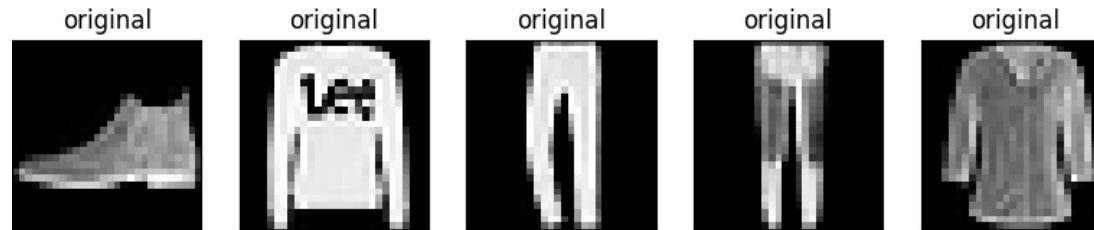
# Autoencoders

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

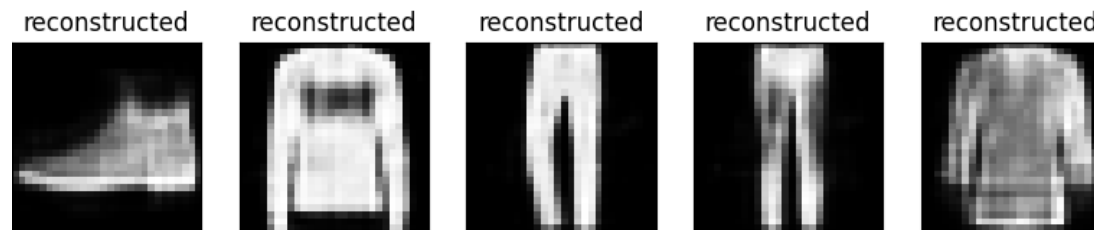


# Autoencoders – Reconstruction Examples

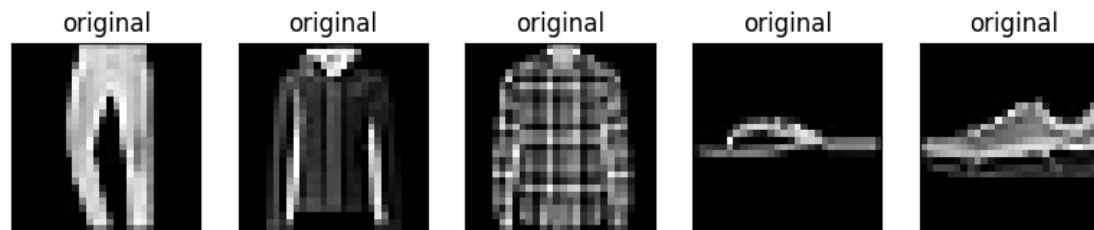
Original



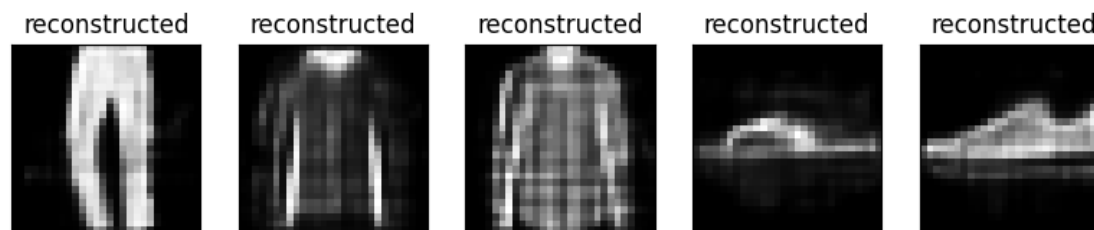
Reconstructed



Original



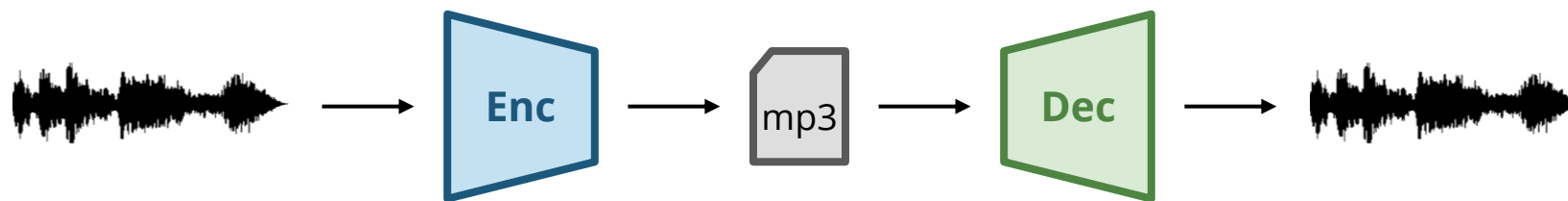
Reconstructed



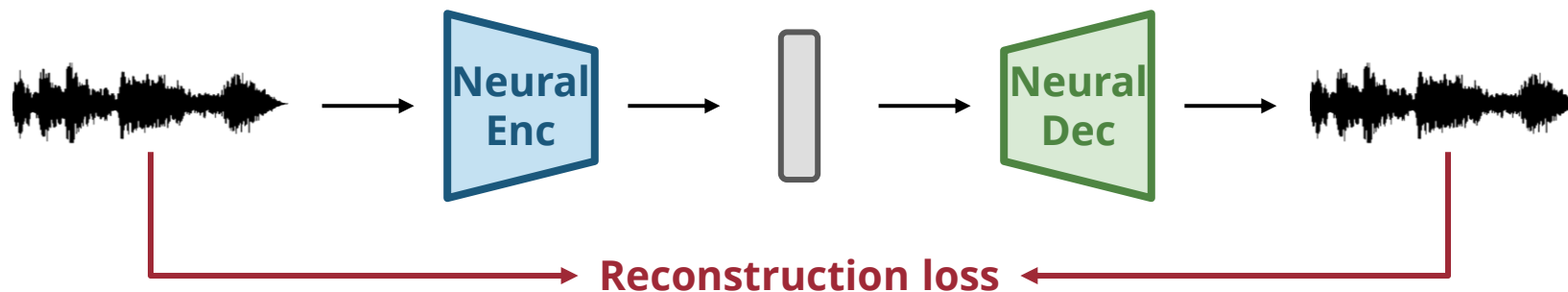
(Source: tensorflow.org)

# Codec is an Autoencoder

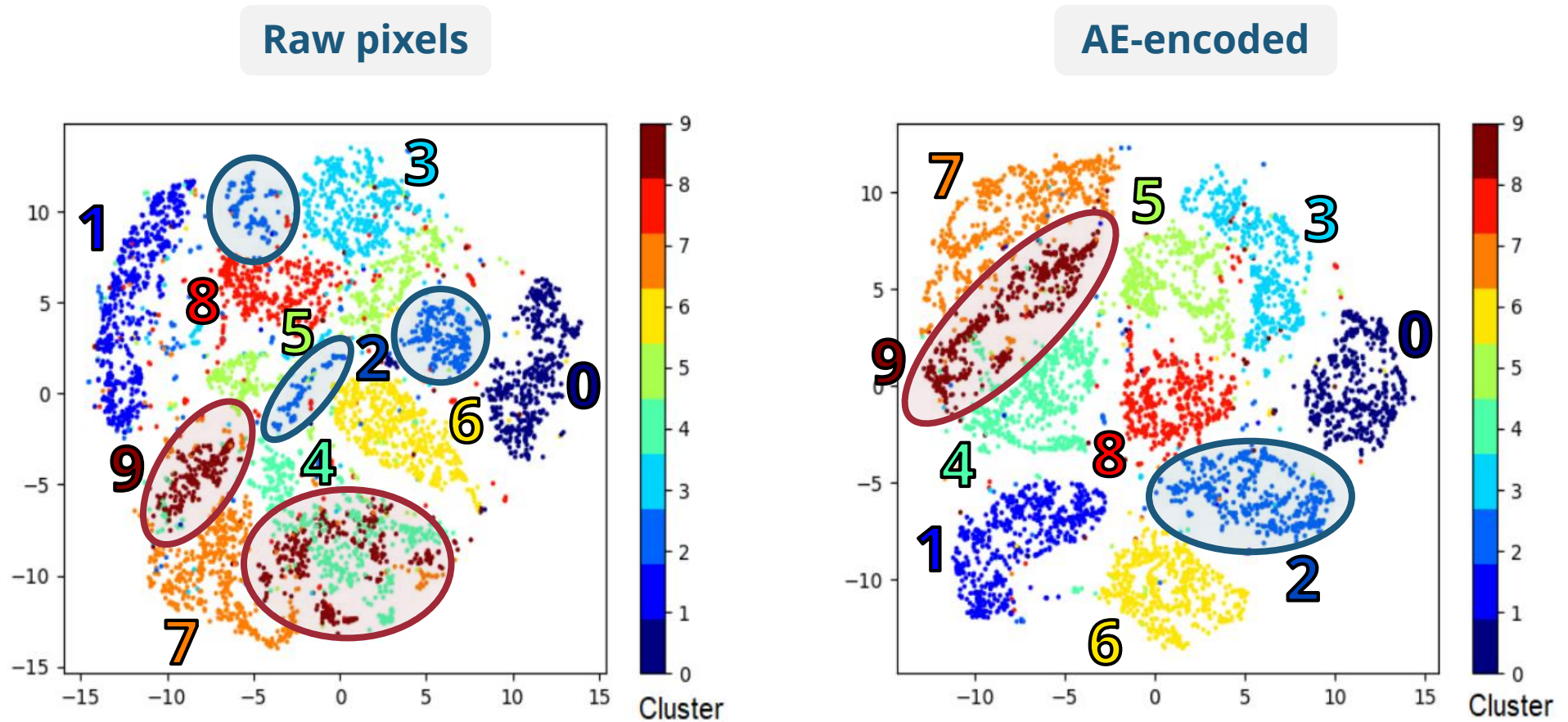
## Traditional Codec



## Neural Codec



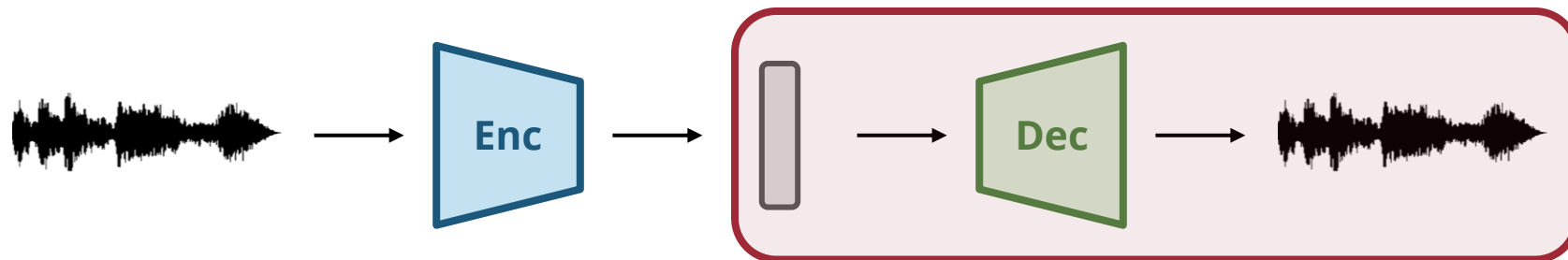
# Unsupervised Clustering with an Autoencoder



(Source: Aljalbout et al., 2020)

# Variational Autoencoders (VAEs)

# Autoencoders

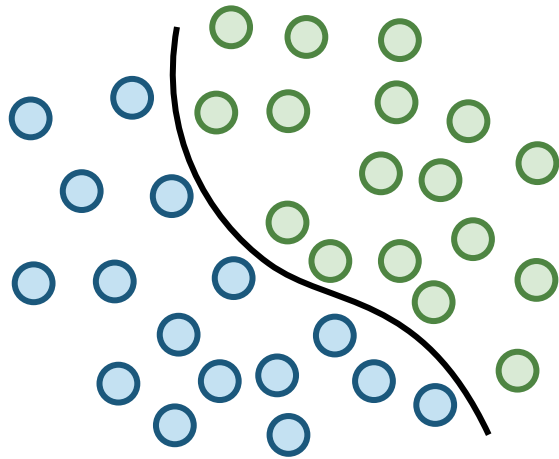


Isn't this like a generative model?

What exactly is a generative model?

# Discriminative vs Generative Models

**Discriminative**



**Discriminative models learn the decision boundary**

$$P(y|x)$$

**Generative**



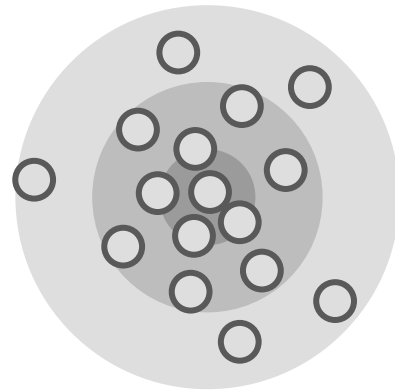
**Generative models learn the underlying distribution**

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

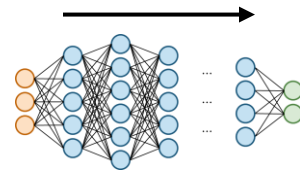


# Learning to Generate Data from Random Distributions

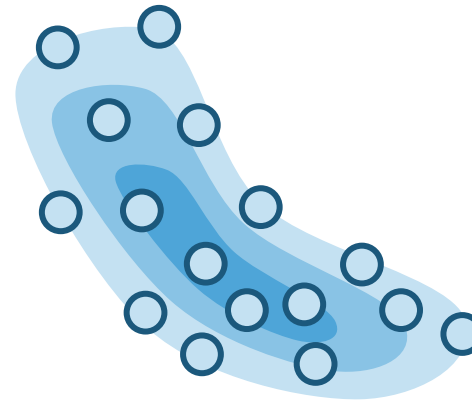
Random distribution



$P(z)$



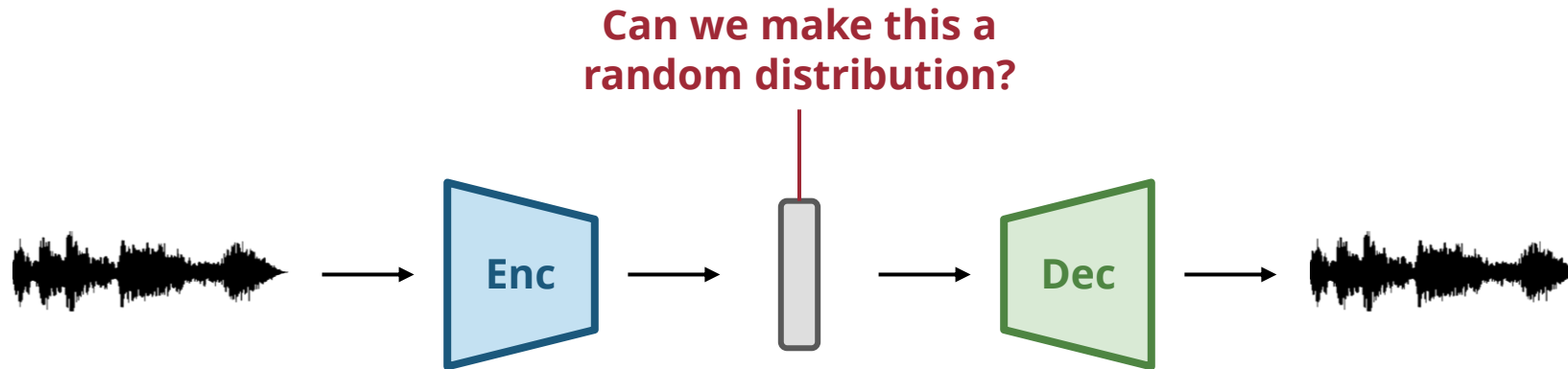
Data distribution



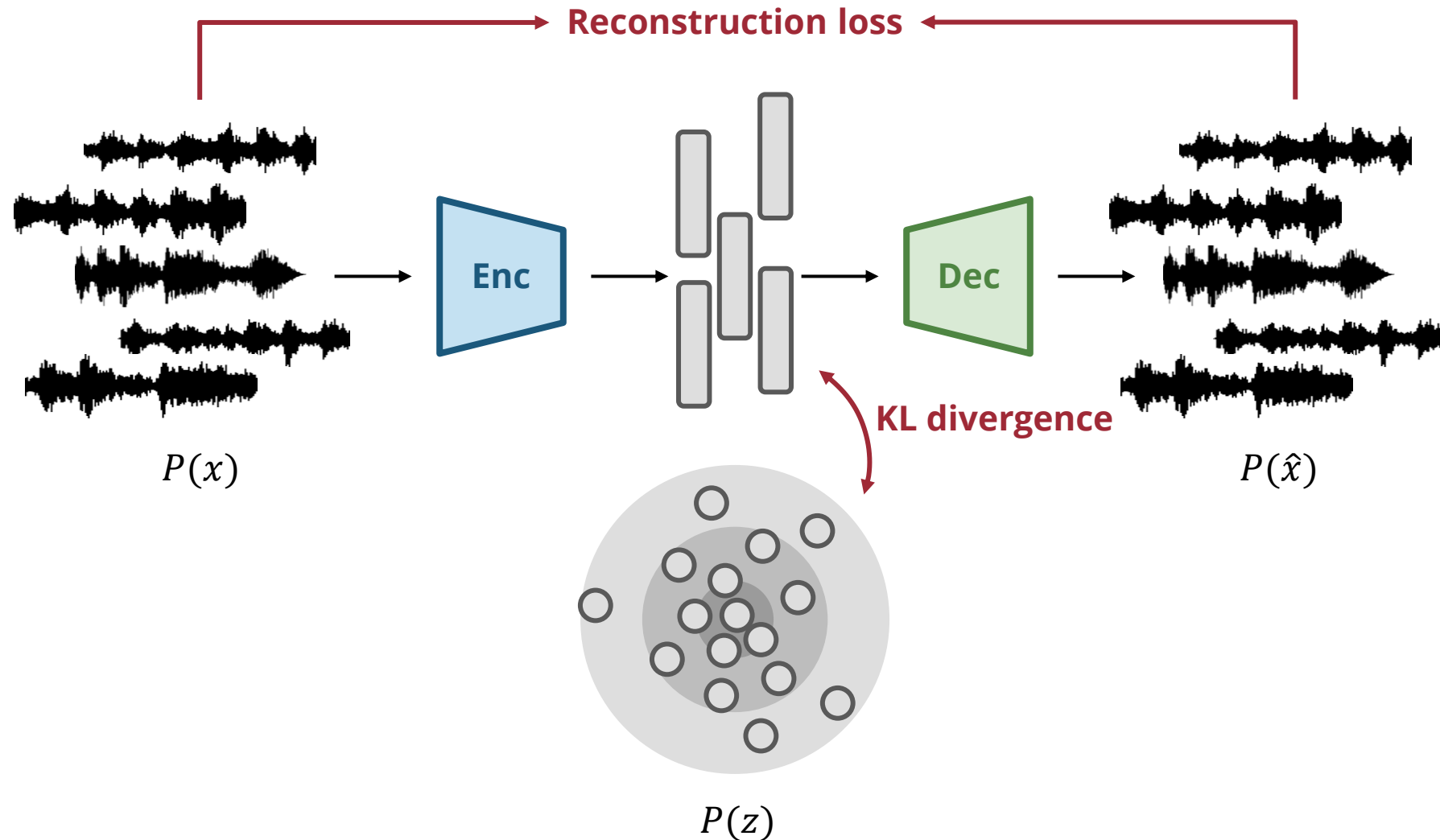
$P(x)$

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

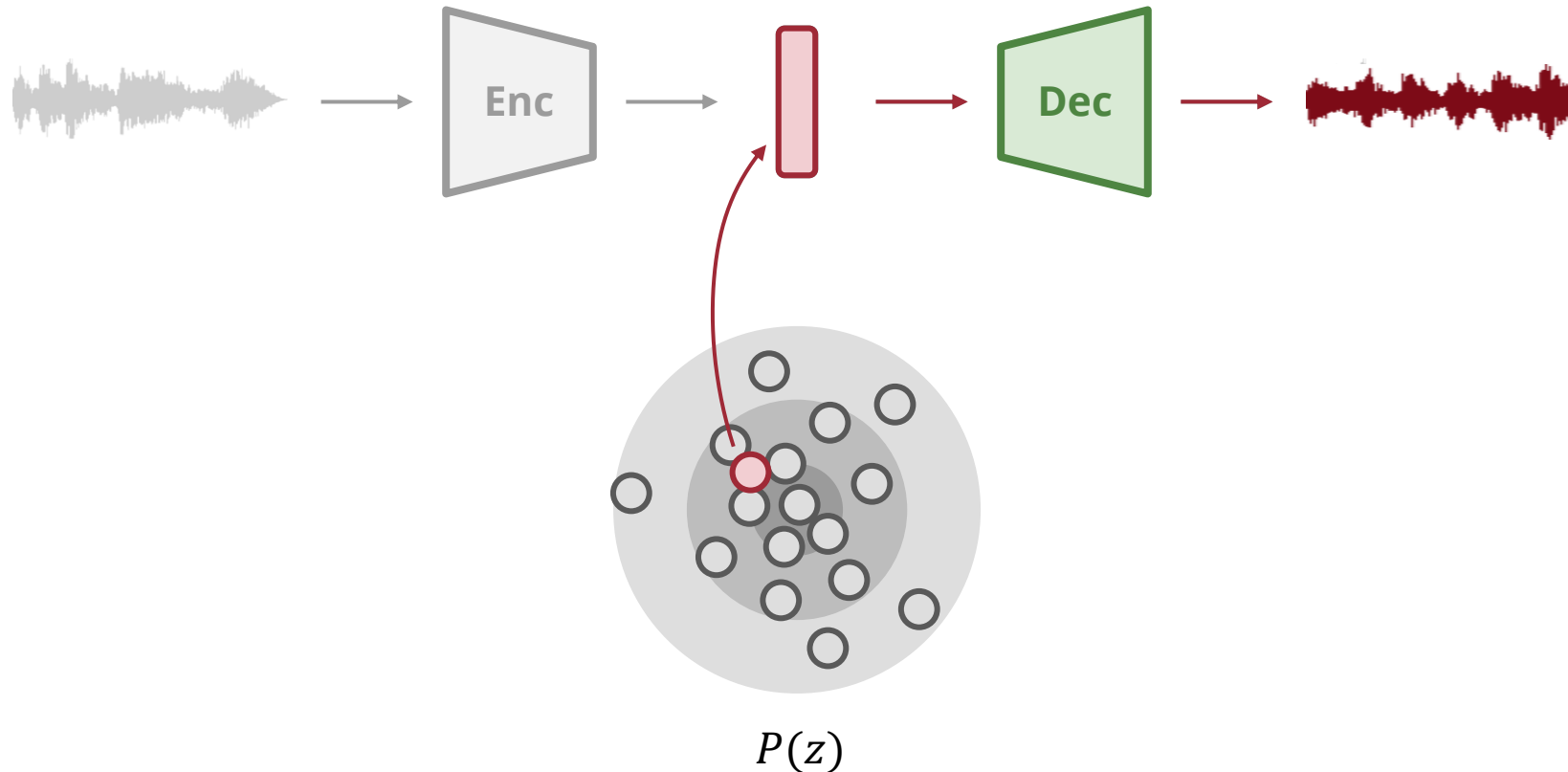
# Variational Autoencoders (VAEs)



# Variational Autoencoders (VAEs) – Training

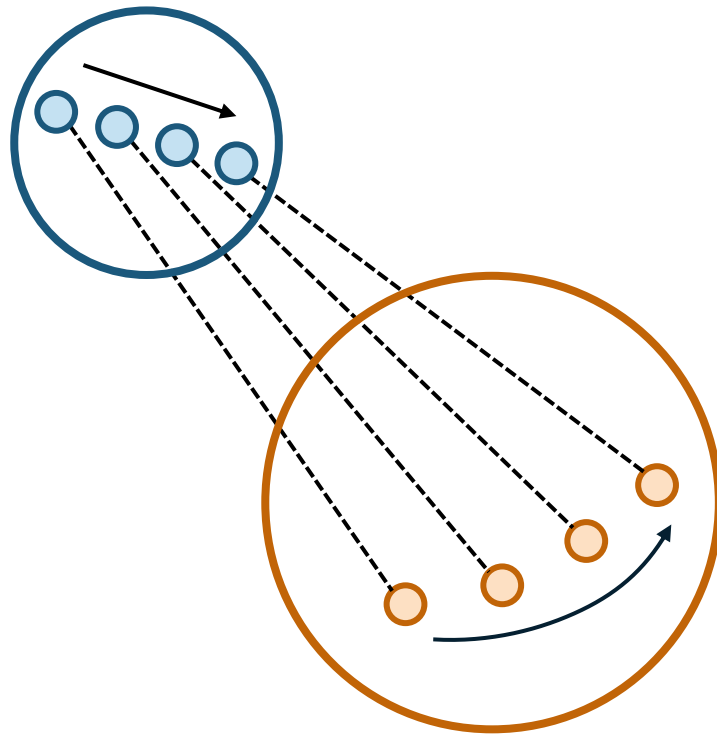


# Variational Autoencoders (VAEs) – Generation

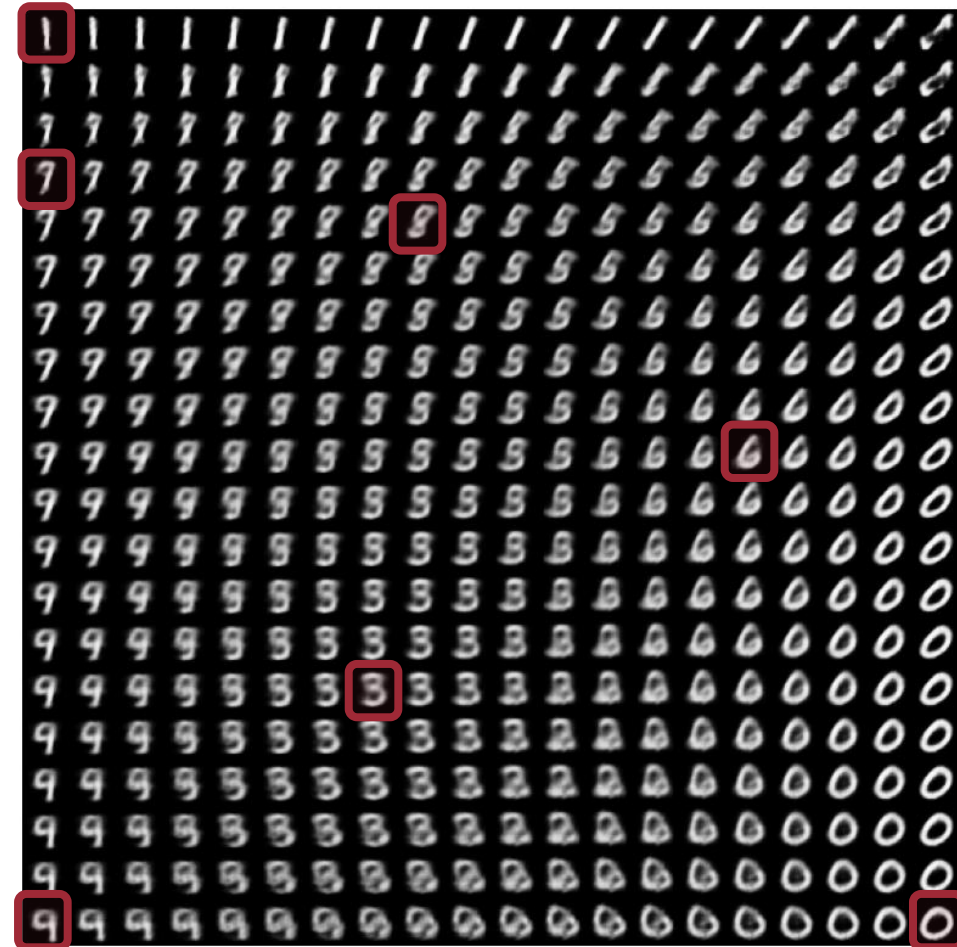


# Decoding the Latent Space of a VAE

Latent space

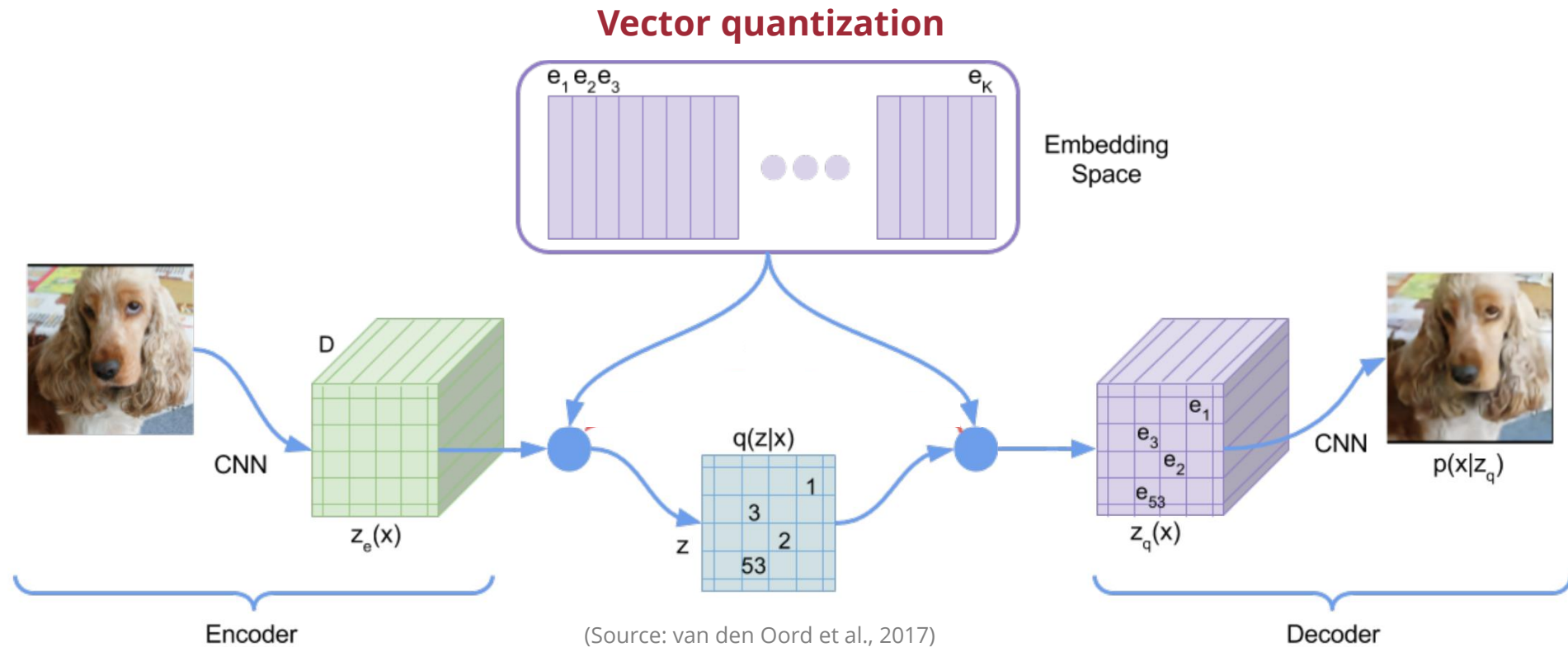


Data space



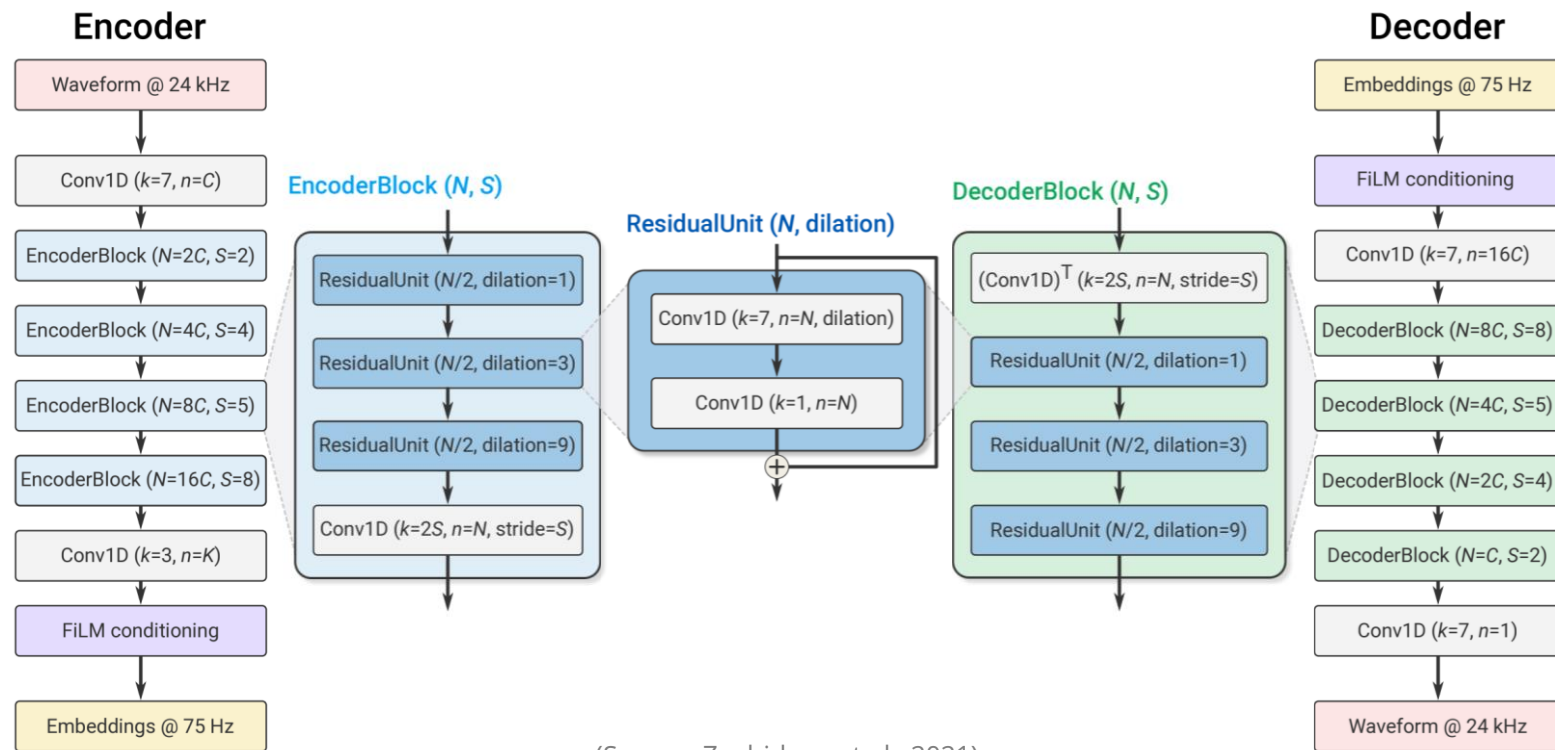
(Source: tensorflow.org)

# Vector-Quantized VAEs (VQVAEs)



# SoundStream

- **Fully-convolutional autoencoder** for audio
- Follow-up work: Encodec & Descript Audio Codec



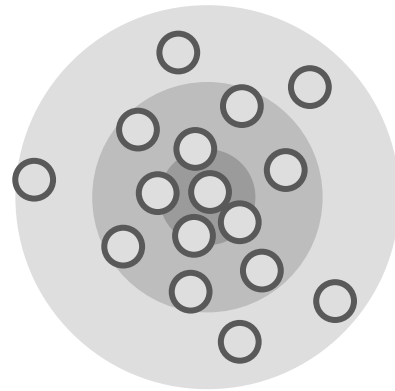
(Source: Zeghidour et al., 2021)

# Generative Adversarial Nets (GANs)

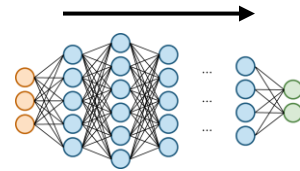


# Generating Data from a Random Distribution

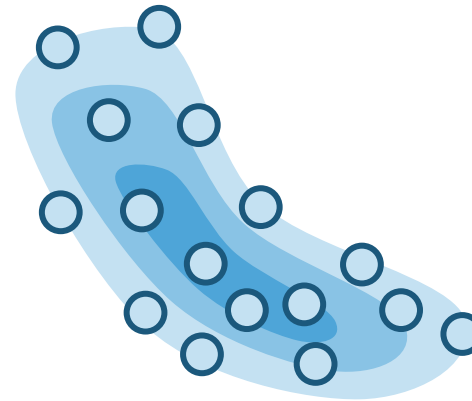
Random distribution



$P(z)$



Data distribution

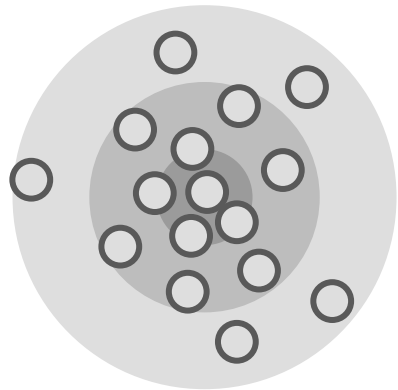


$P(x)$

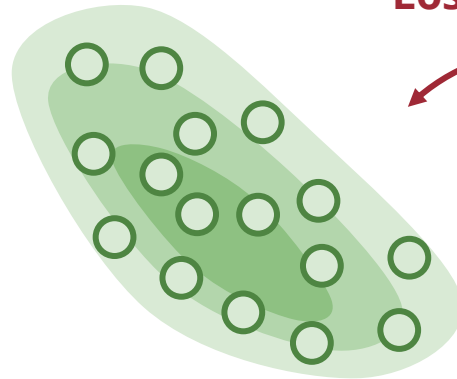
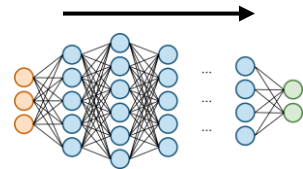
**If we can learn this mapping, we can easily generate new samples from the data distribution**

# A Loss Function for Distributions

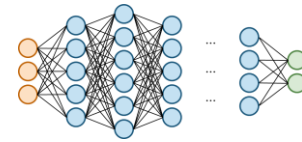
Random distribution



$P(z)$

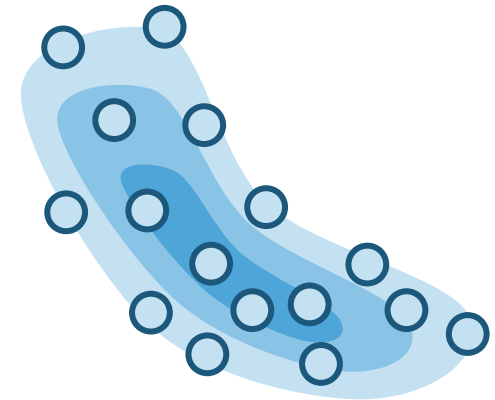


$P(\hat{x})$



Loss function?

Data distribution

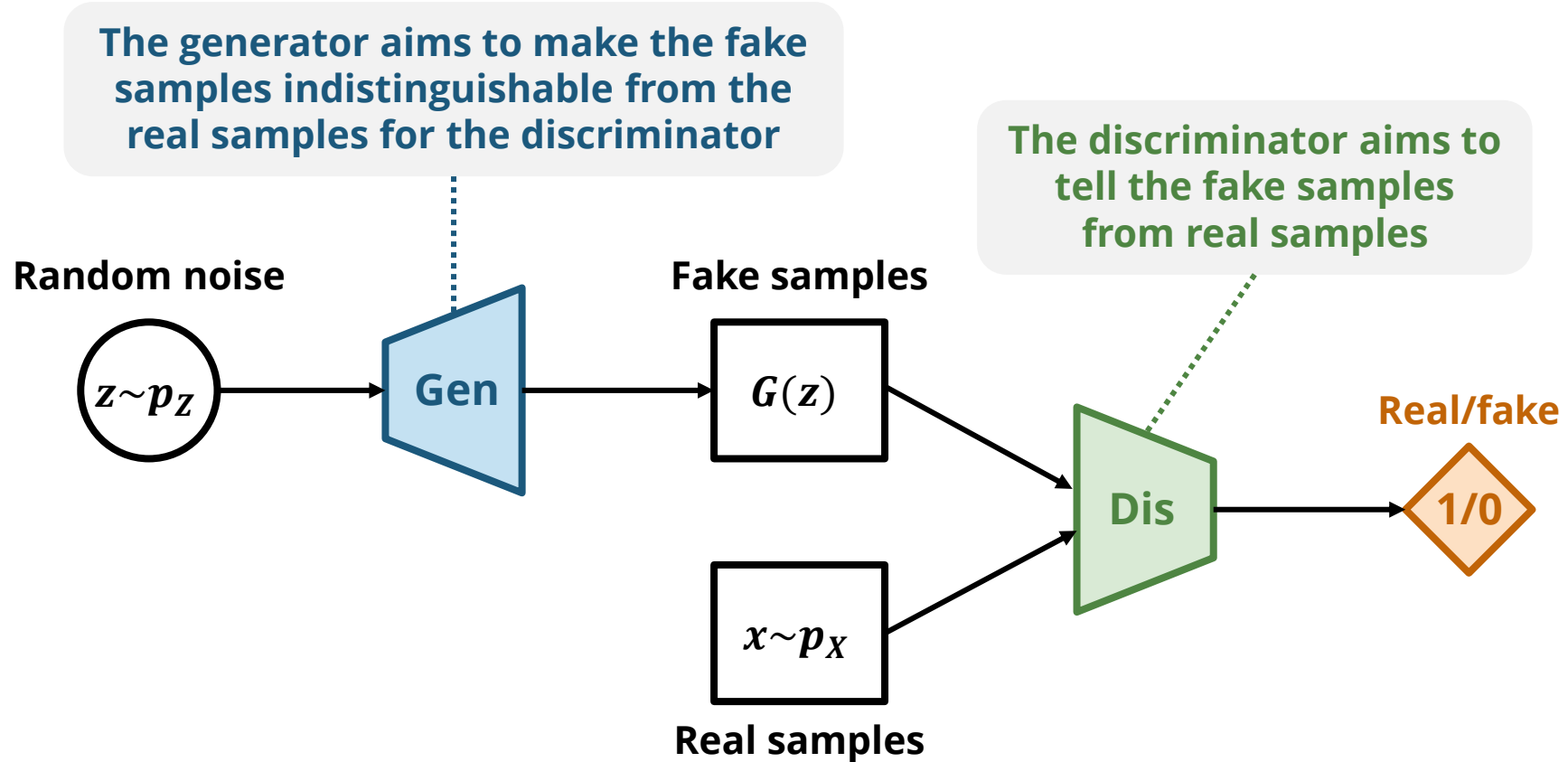


$P(x)$

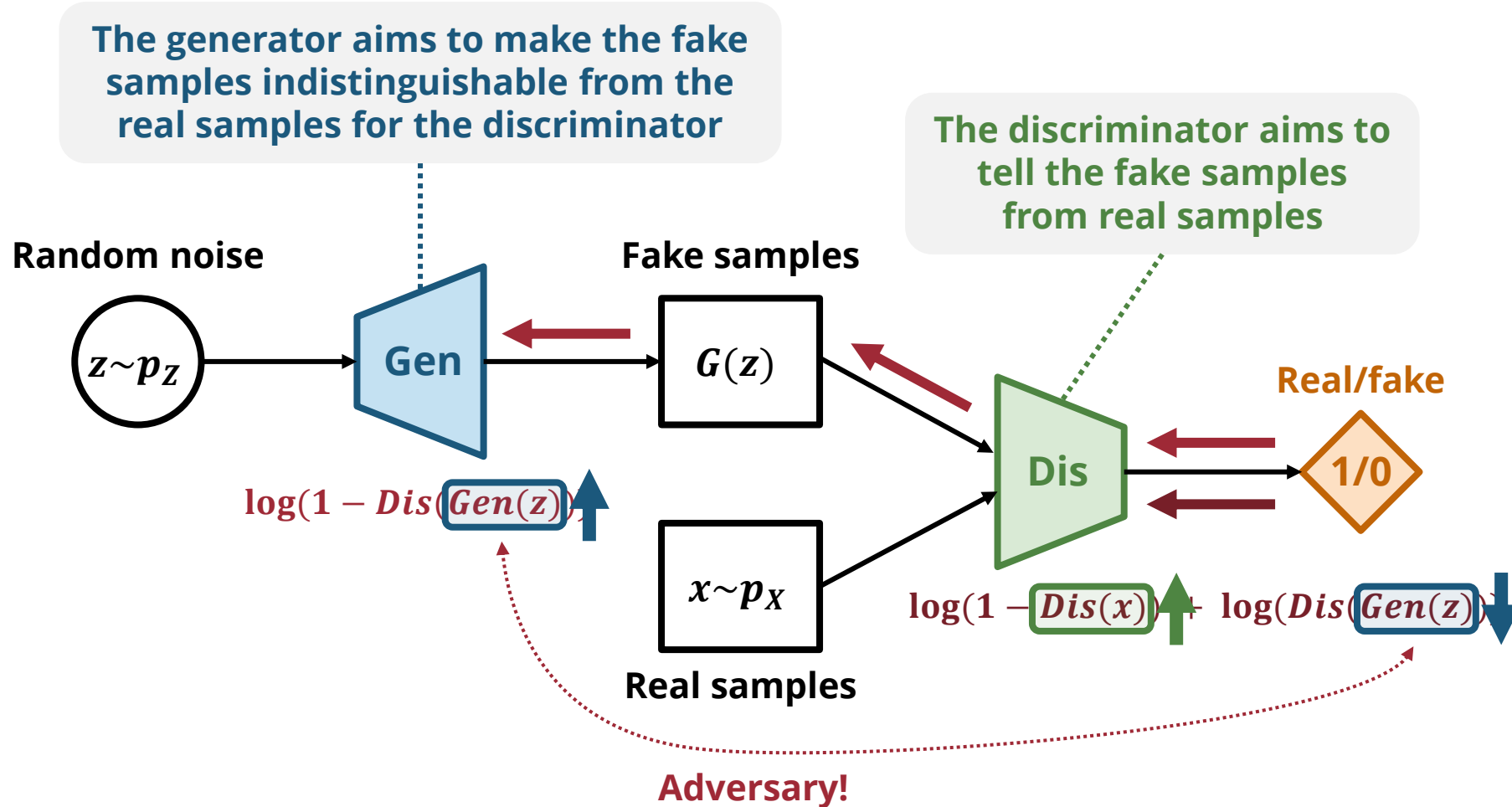
Unfortunately, no easy way to measure the difference between two distributions

But what about another neural network!?

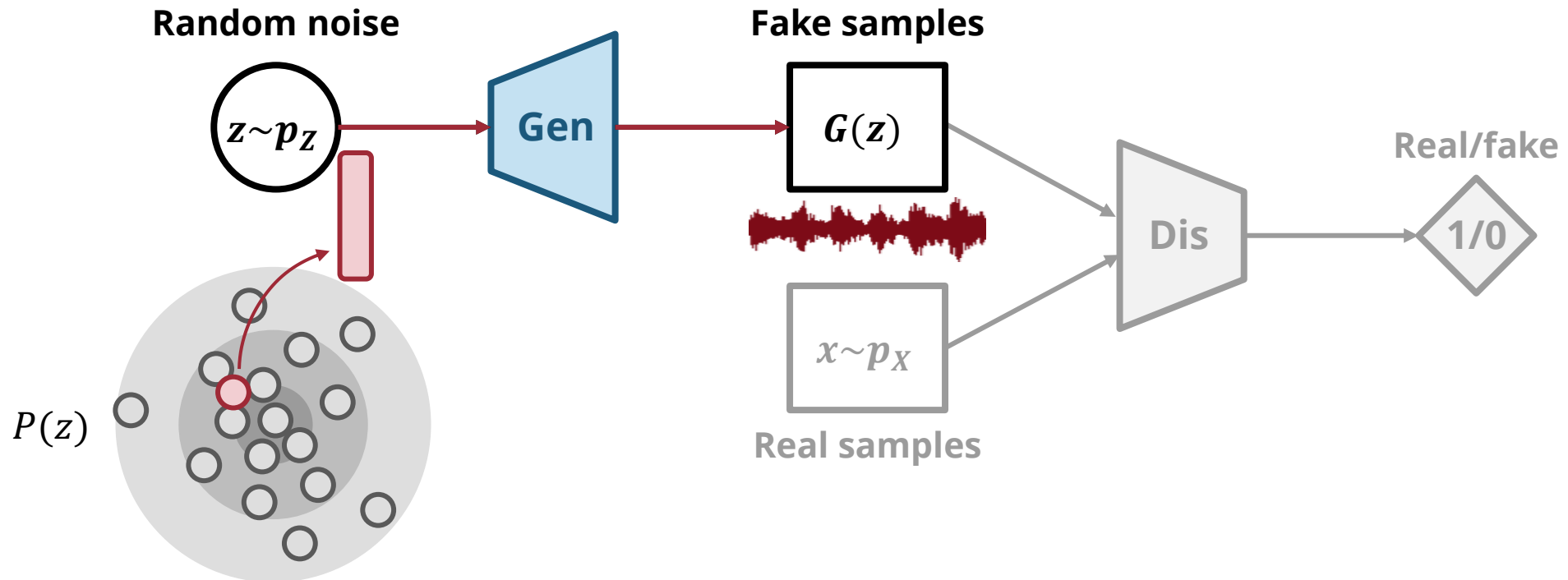
# Generative Adversarial Nets (GANs)



# Generative Adversarial Nets (GANs) – Training

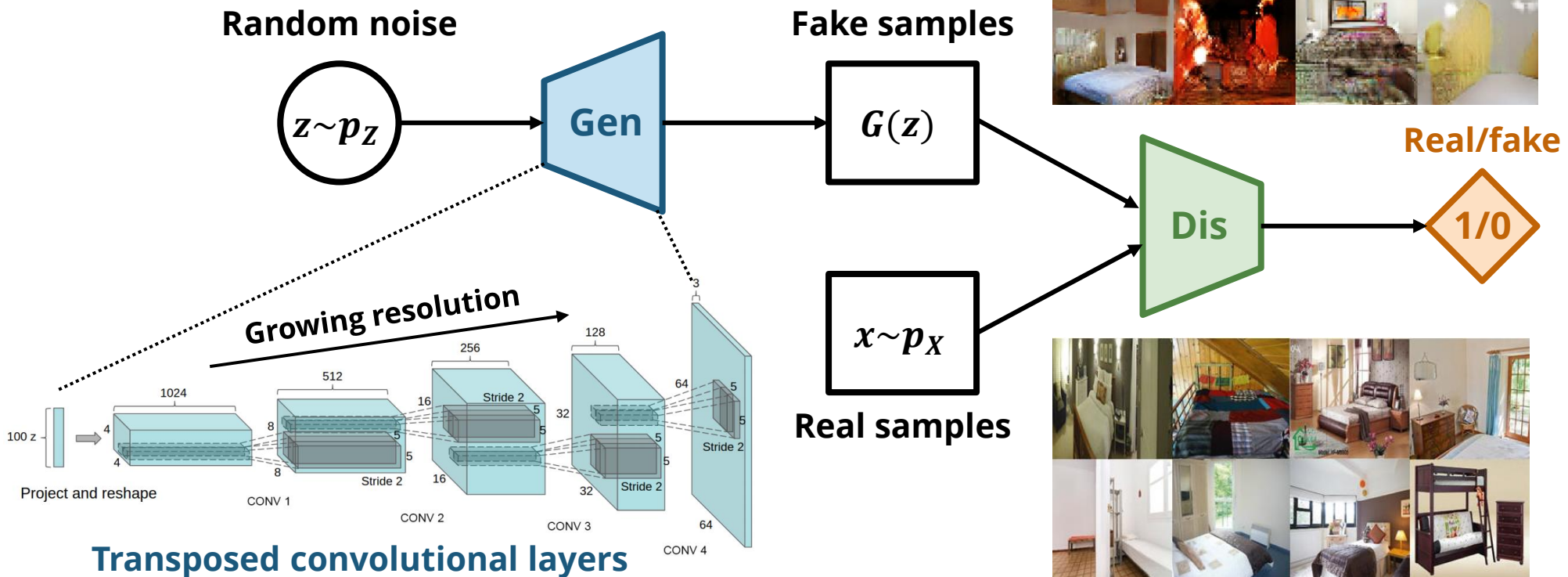


# Generative Adversarial Nets (GANs) – Generation



# Deep Convolutional GANs (DCGANs)

Use CNNs for both the generator and discriminator



# Transposed Convolution

Convolution

1	-1	-1	-1
-1	1	-1	-1
-1	-1	1	-1
-1	-1	-1	1

\*

1	-1	-1
-1	1	-1
-1	-1	1

=

9	-1
-1	9

Transposed convolution

1	-1
-1	1

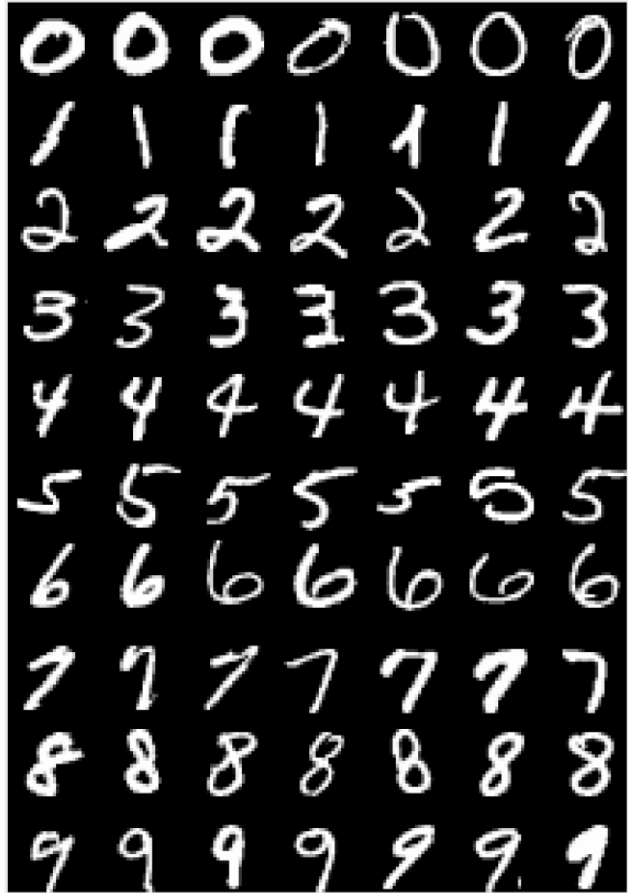
\*

1	-1	-1
-1	1	-1
-1	-1	1

=

1	0	0	1
0	4	-2	0
0	-2	4	0
1	0	0	1

# DCGAN - Examples



Groundtruth MNIST



GAN



DCGAN (ours)

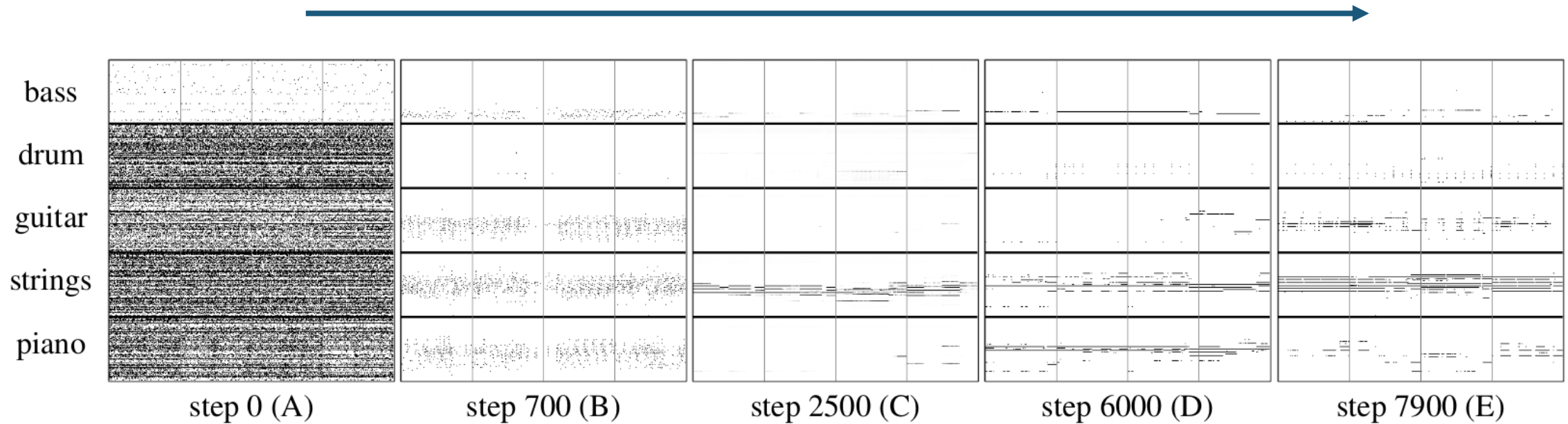
(Source: Radford et al., 2016)



# MuseGAN – A GAN for Pianorolls

The generator improves over time

So does the discriminator!



(Source: Dong et al., 2018)

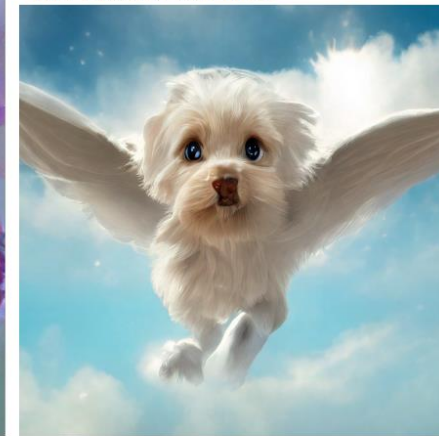
# GigaGAN: Scaling up GANs



A portrait of a human growing colorful flowers from her hair. Hyperrealistic oil painting. Intricate details.



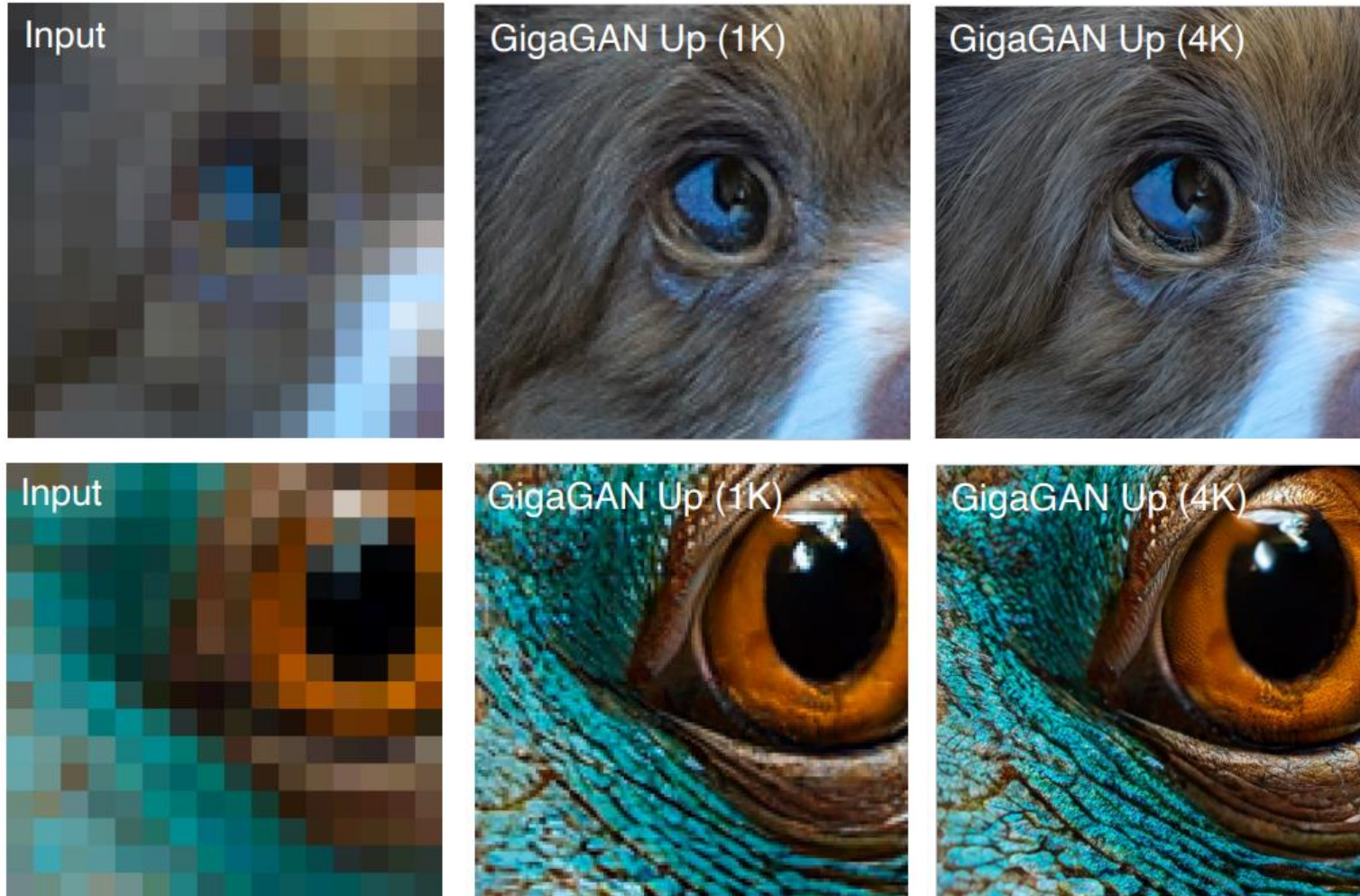
A golden luxury motorcycle parked at the King's palace. 35mm f/4.5.



a cute magical flying maltipoo at light speed, fantasy concept art, bokeh, wide sky

(Source: Kang et al., 2023)

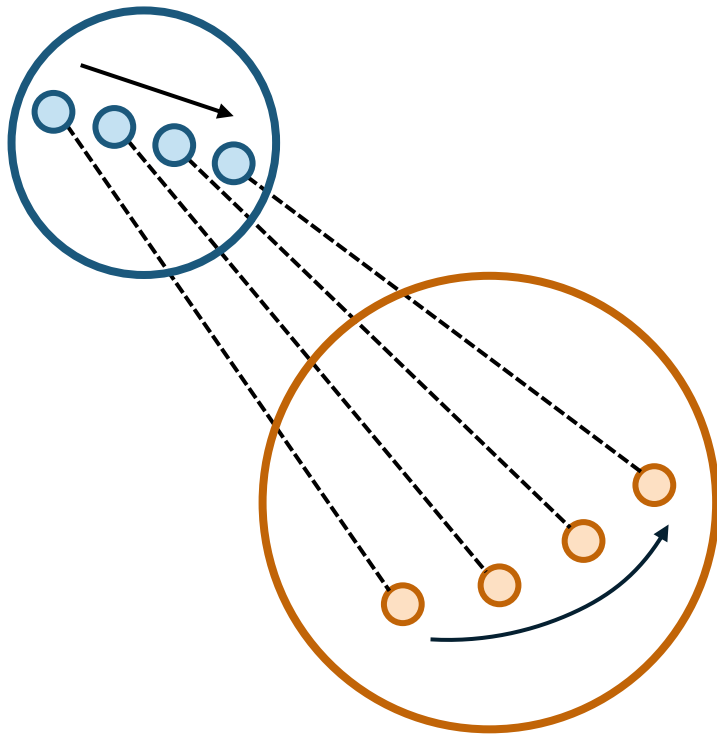
# GigaGAN for Image Super-resolution



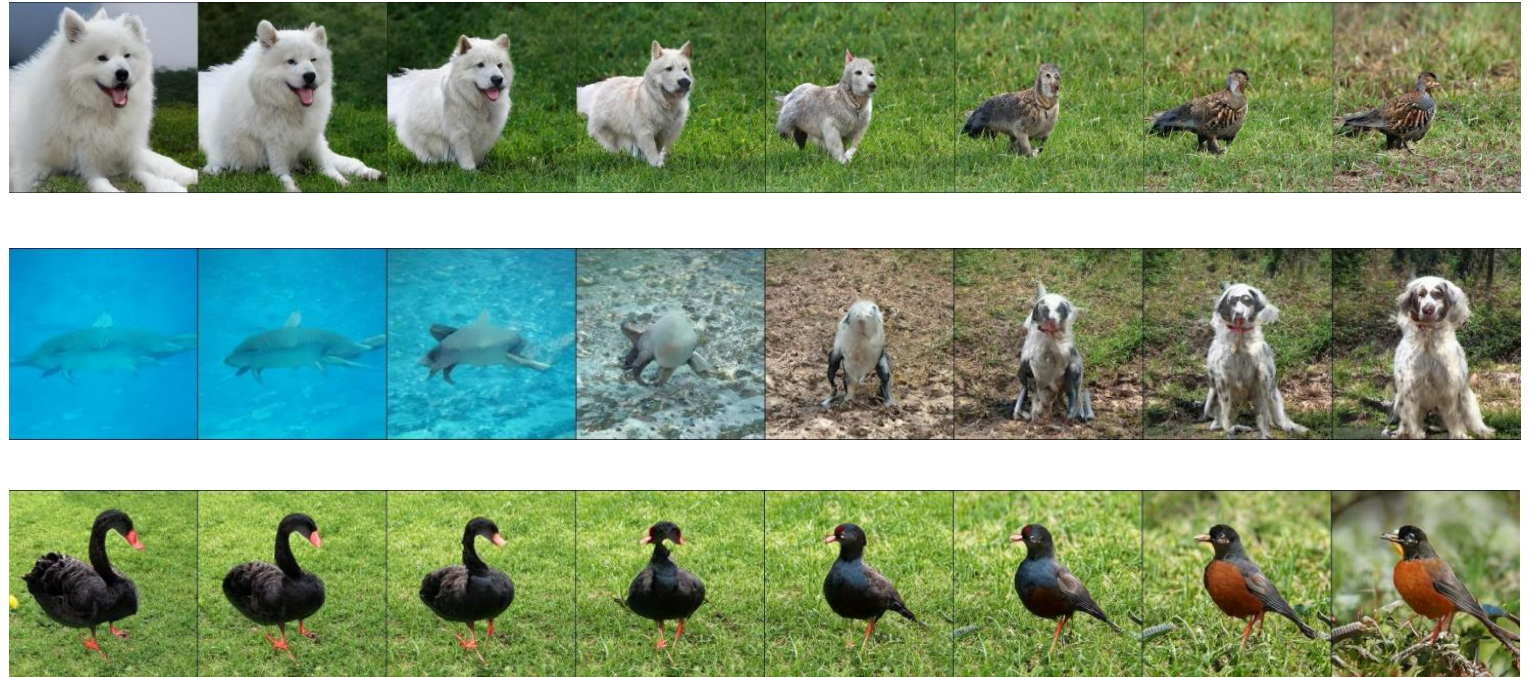
(Source: Kang et al., 2023)

# Interpolation on the Latent Space

Latent space

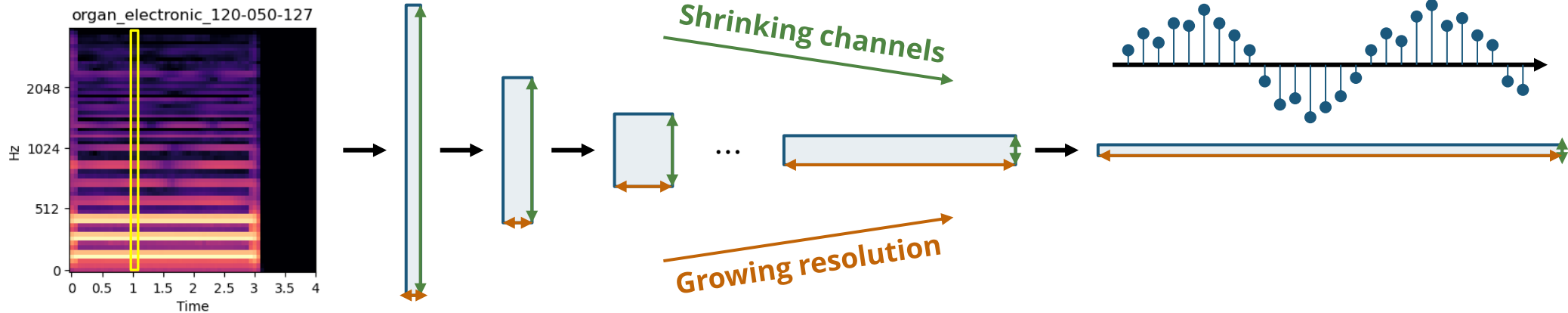
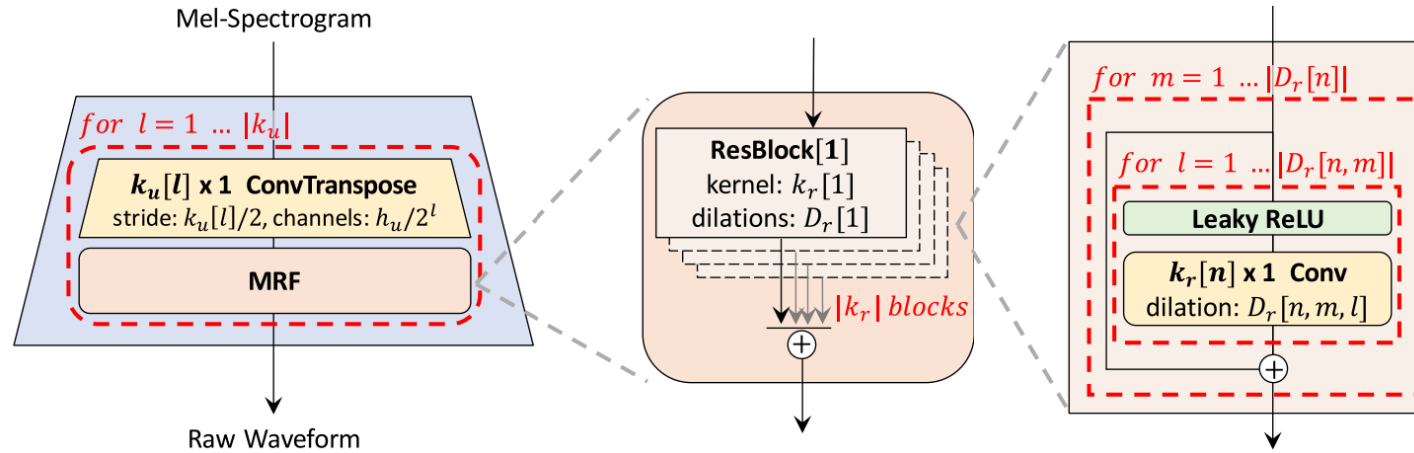


Data space



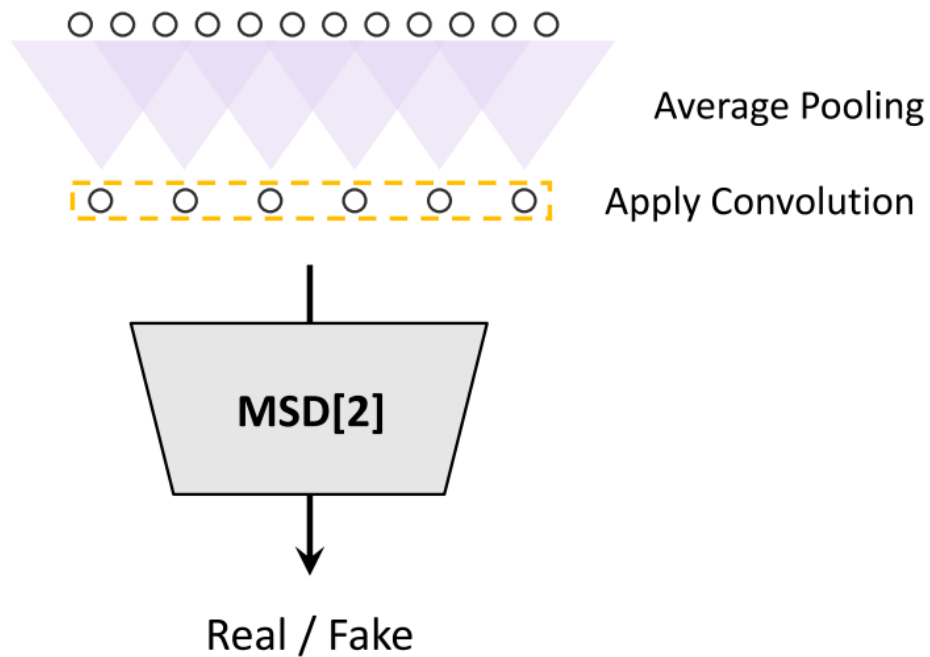
(Source: Brock et al., 2019)

# Hifi-GAN – Generator

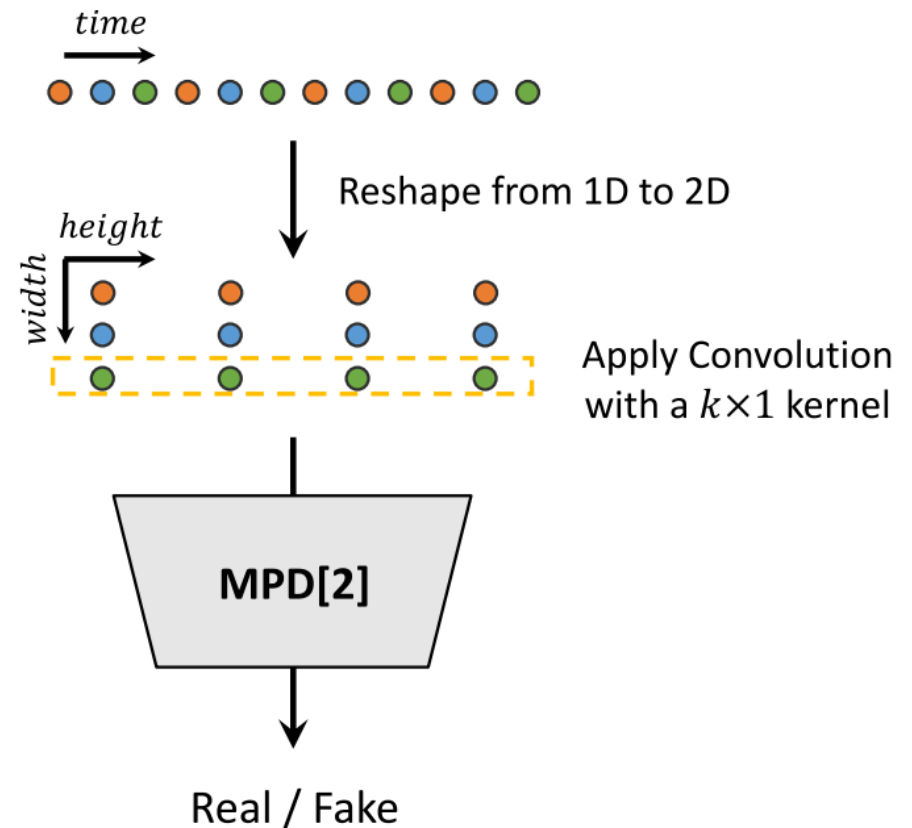


# Hifi-GAN – Discriminator

## Multi-scale discriminator

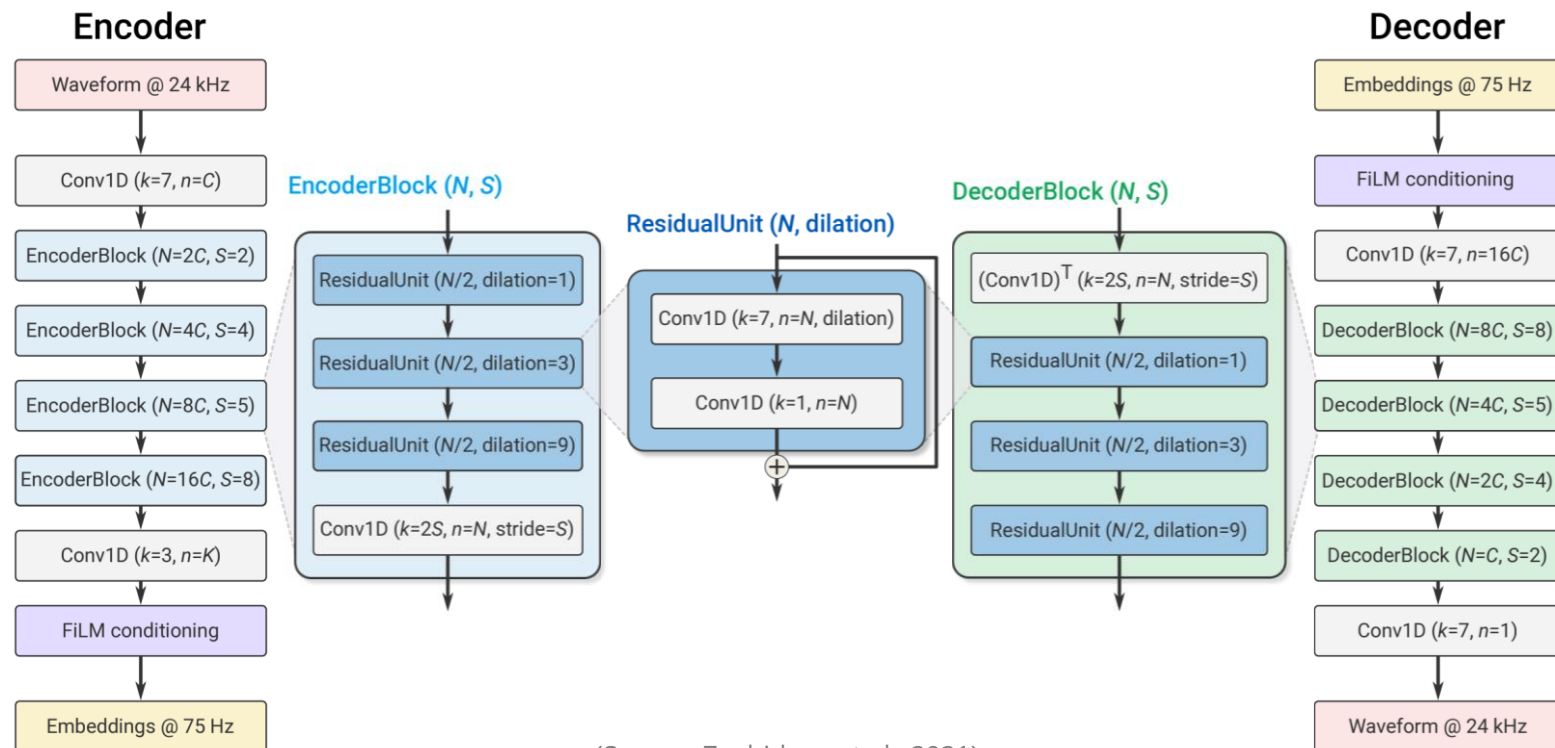


## Multi-period discriminator



# (Recap) SoundStream

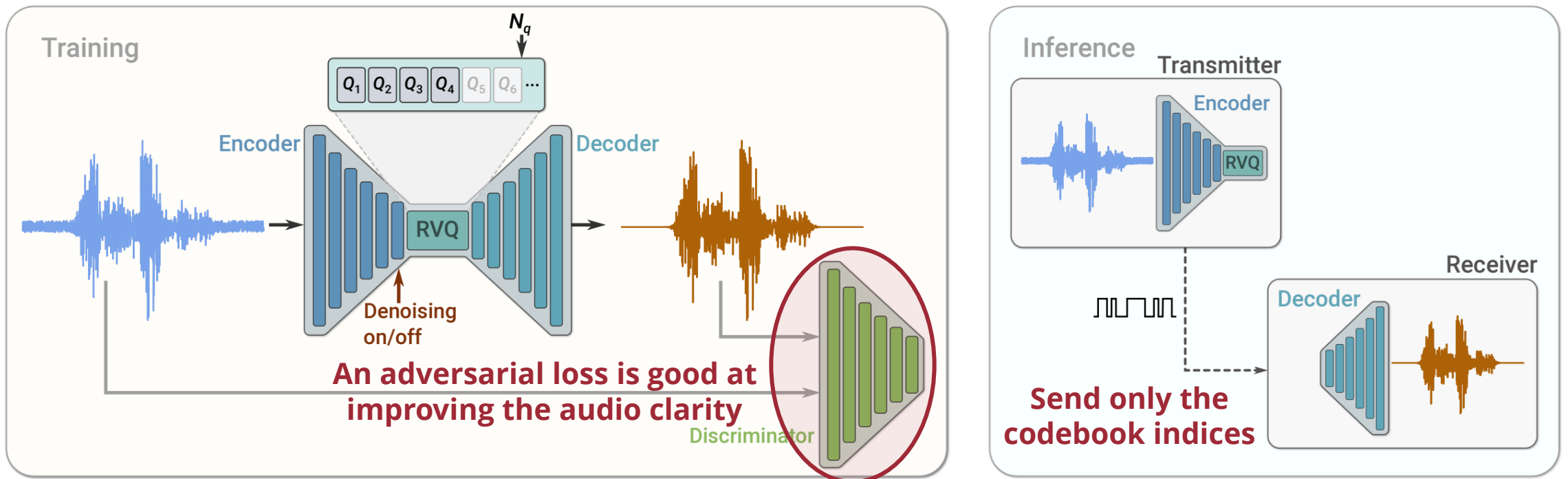
- **Fully-convolutional autoencoder** for audio
- Follow-up work: Encodec & Descript Audio Codec



(Source: Zeghidour et al., 2021)

# (Recap) SoundStream

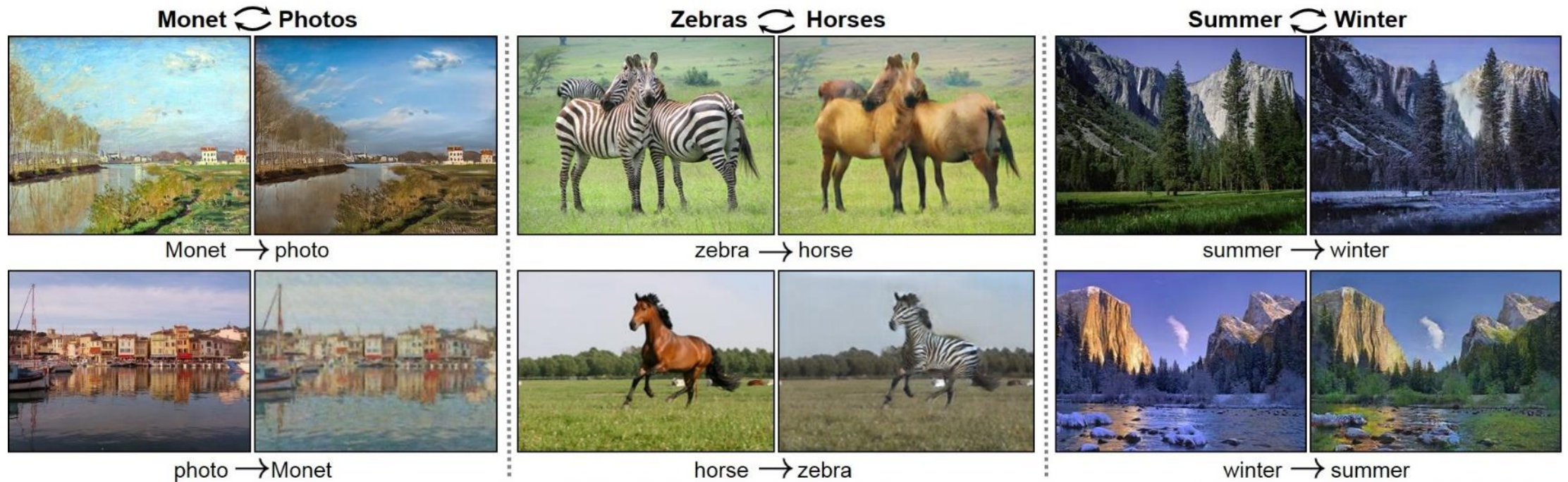
- **Fully-convolutional autoencoder** for audio
- Follow-up work: Encodec & Descript Audio Codec



(Source: Zeghidour et al., 2021)



# CycleGAN – Examples



(Source: Zhu et al., 2017)

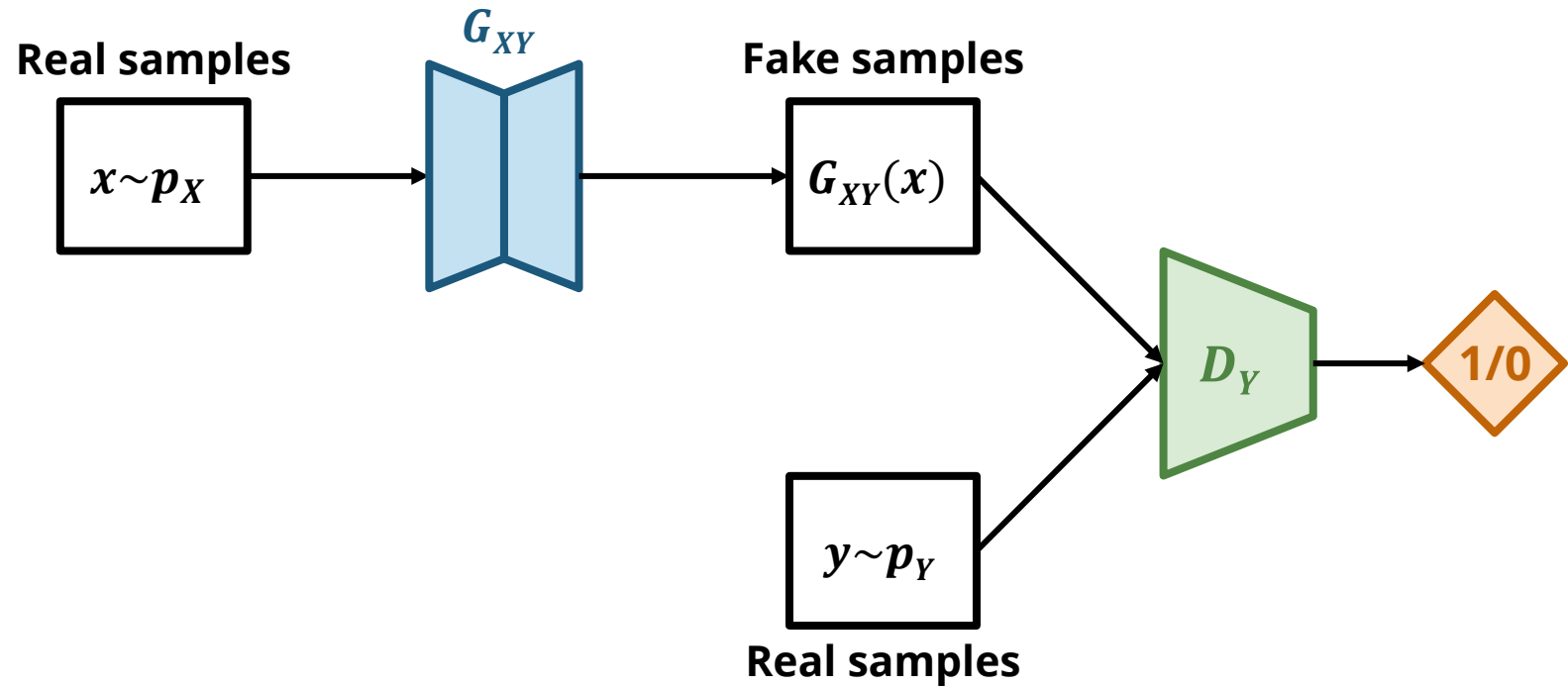
**Work without paired data!**  
For example, we only need **a collection of photos**  
and **a collection of Monet's paintings**

# CycleGAN – Goal

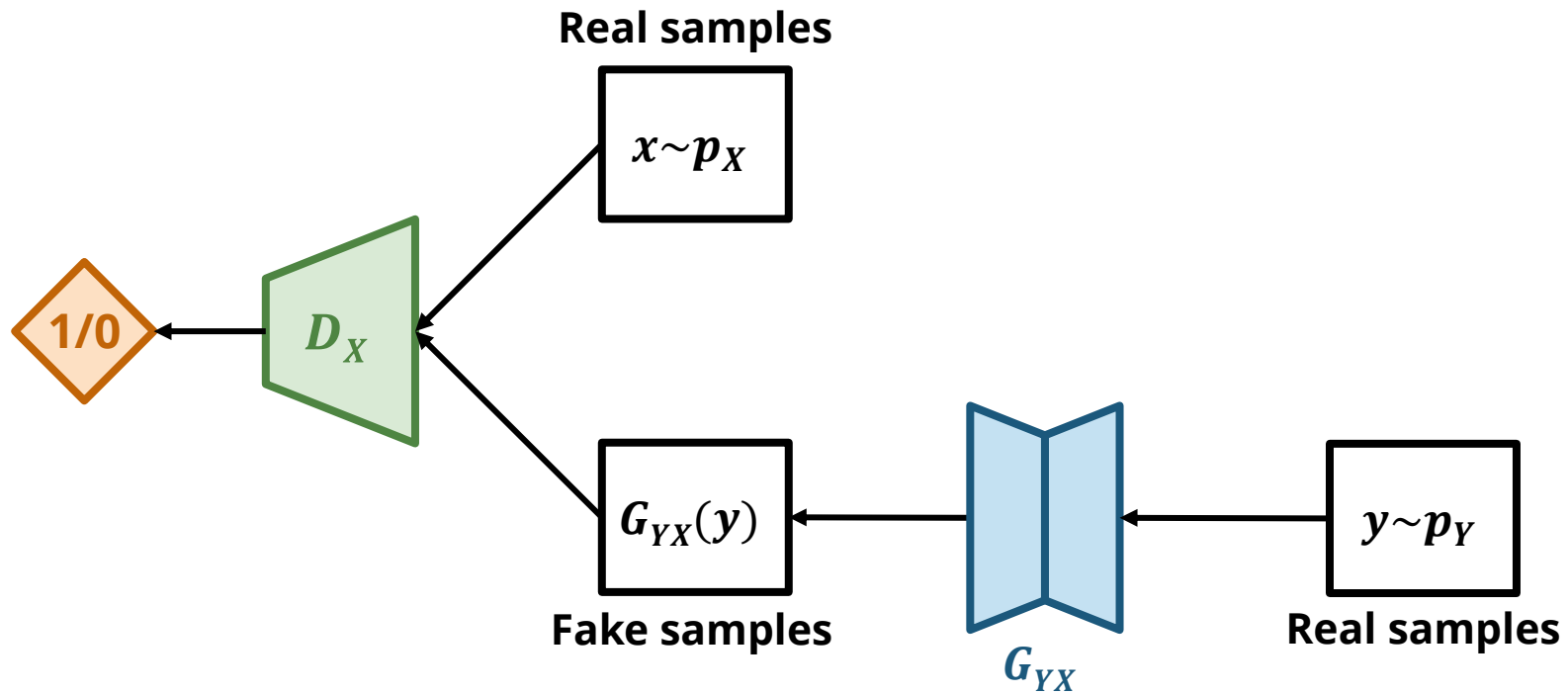


**Can we learn the mapping without paired data?**

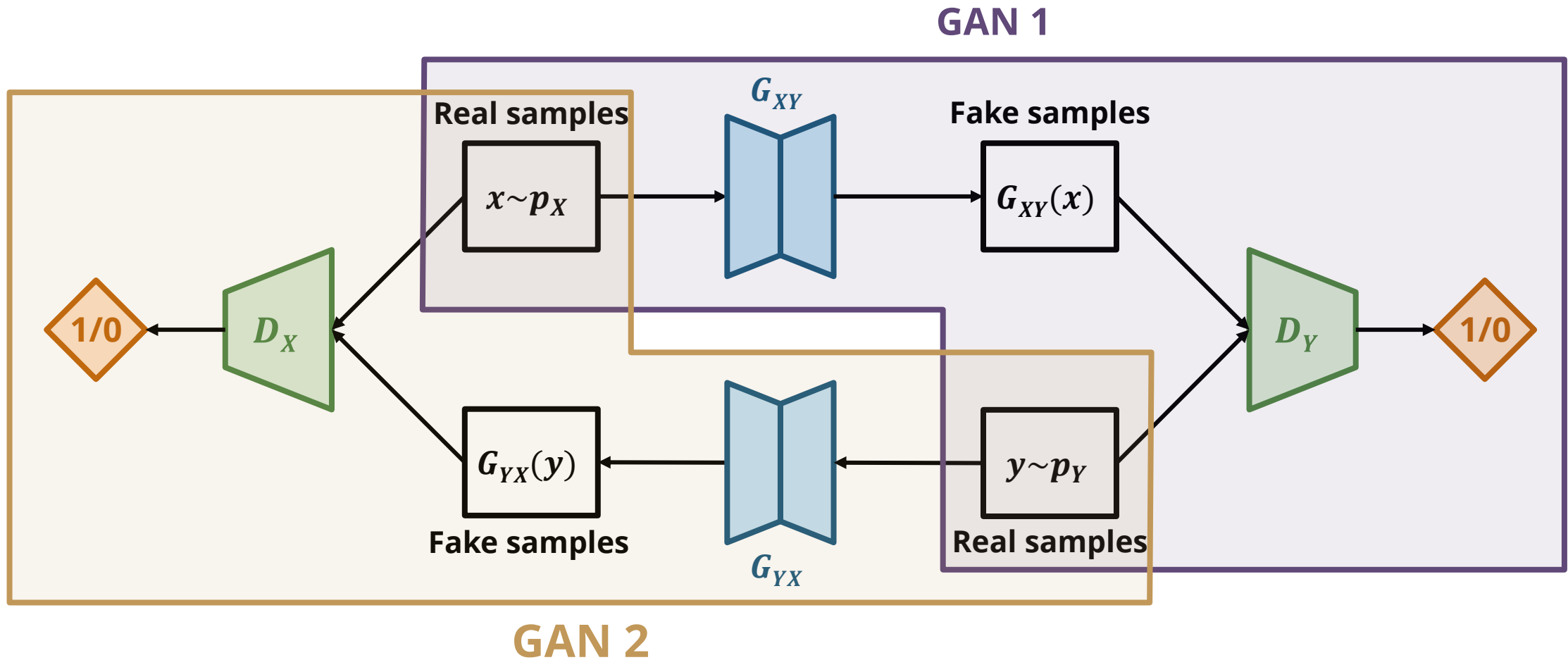
# Cycle-consistent GAN (CycleGAN)



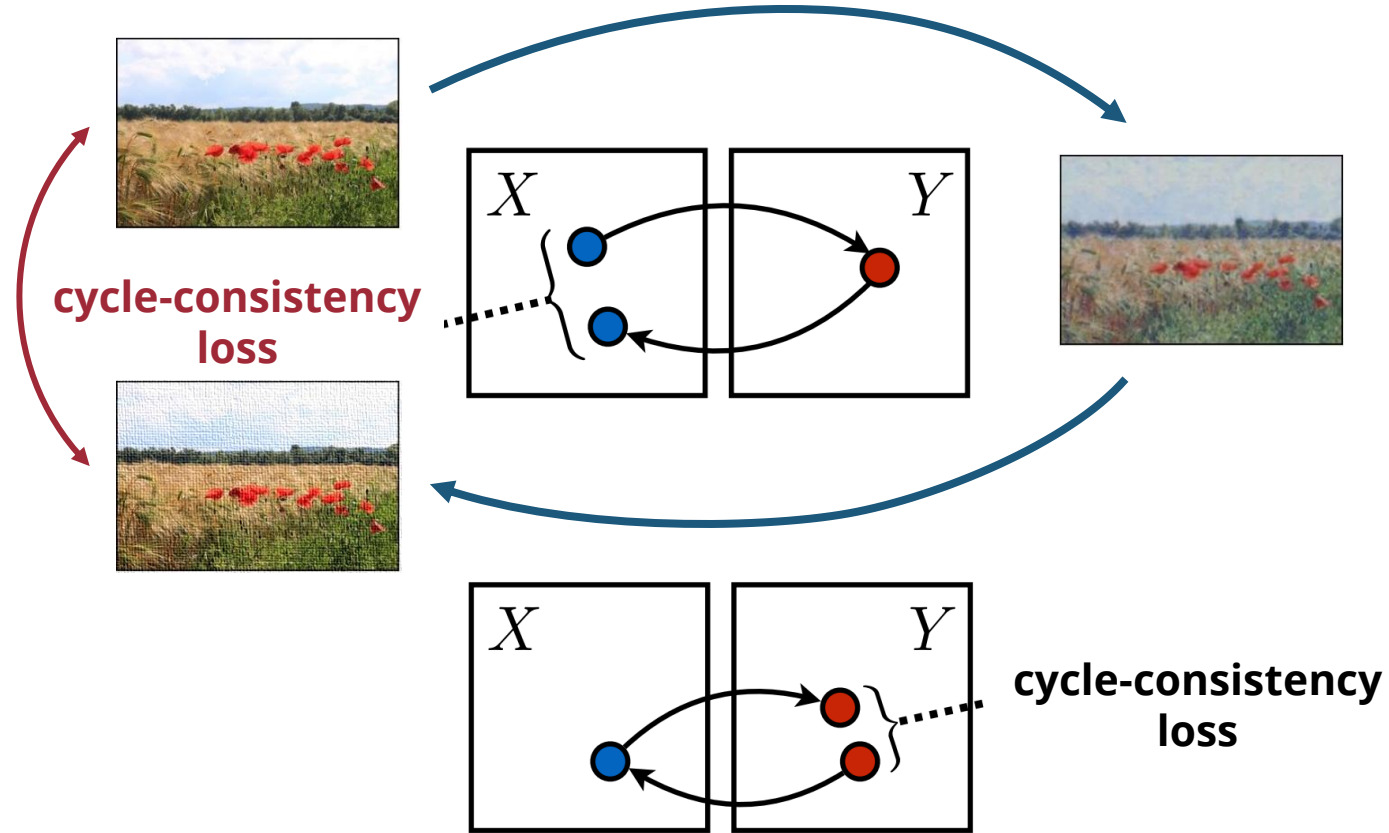
# Cycle-consistent GAN (CycleGAN)



# Cycle-consistent GAN (CycleGAN)

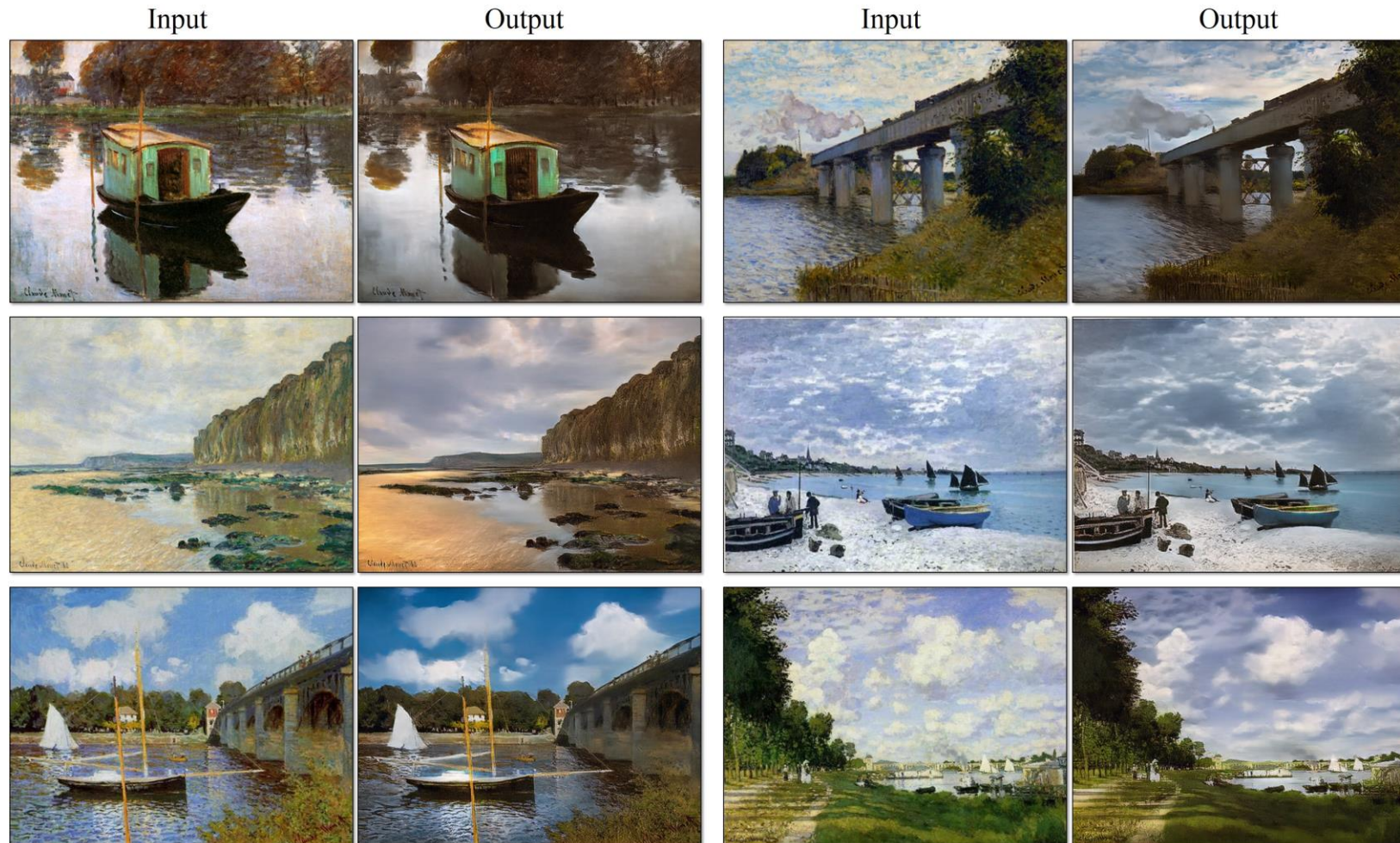


# Cycle-consistency Loss



**We only need unpaired samples in two domains**

# CycleGAN Examples – Monet to Photo



(Source: Zhu et al., 2017)

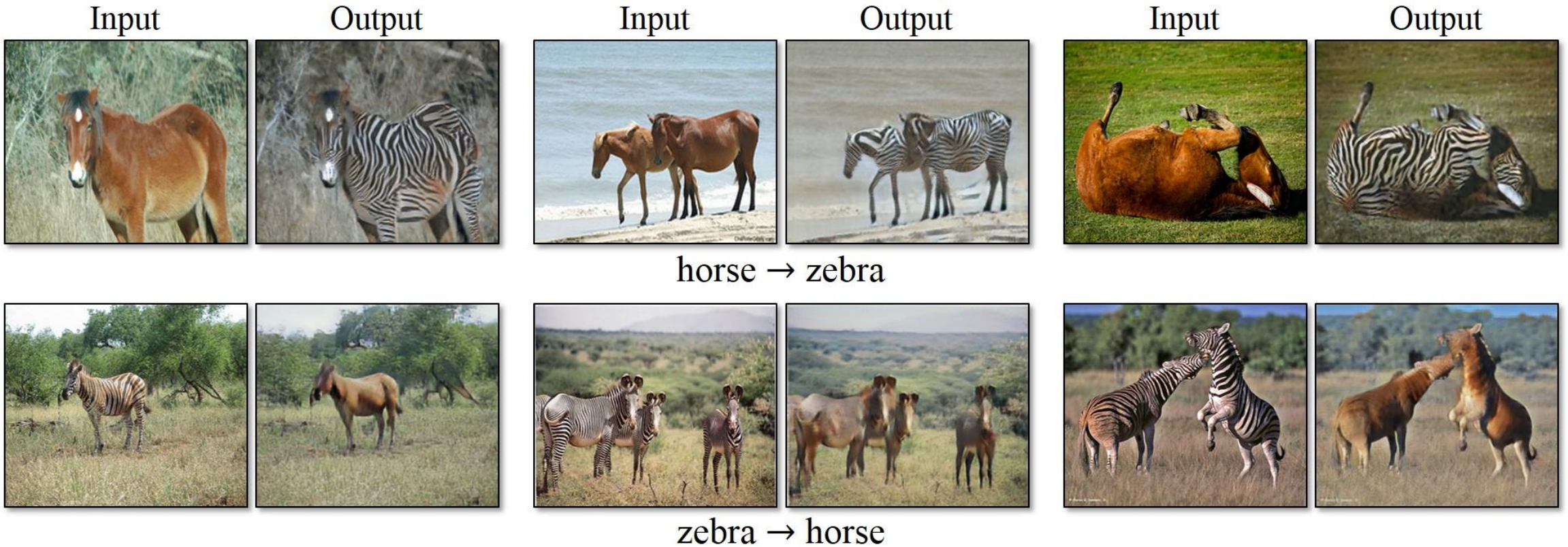
# CycleGAN Examples – Artistic Style Transfer



(Source: Zhu et al., 2017)



# CycleGAN Examples – Object Transfer



(Source: Zhu et al., 2017)

# CycleGAN Examples – Season Transfer



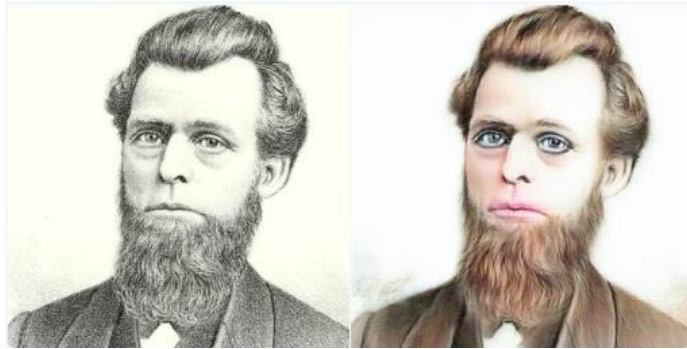
winter Yosemite → summer Yosemite



summer Yosemite → winter Yosemite

(Source: Zhu et al., 2017)

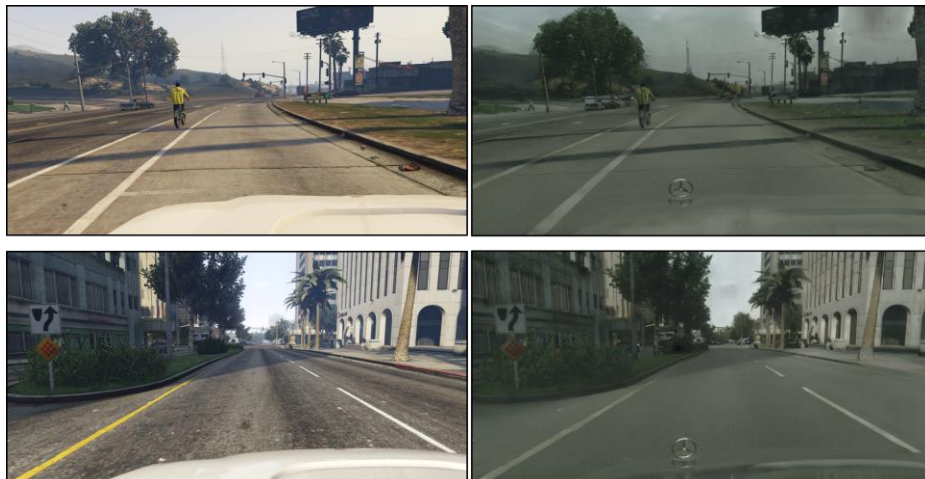
# CycleGAN Examples – And Many More!



**B&W → Color**



**Old city plan → Map → Satellite**



**GTA screenshot → Cityscape**



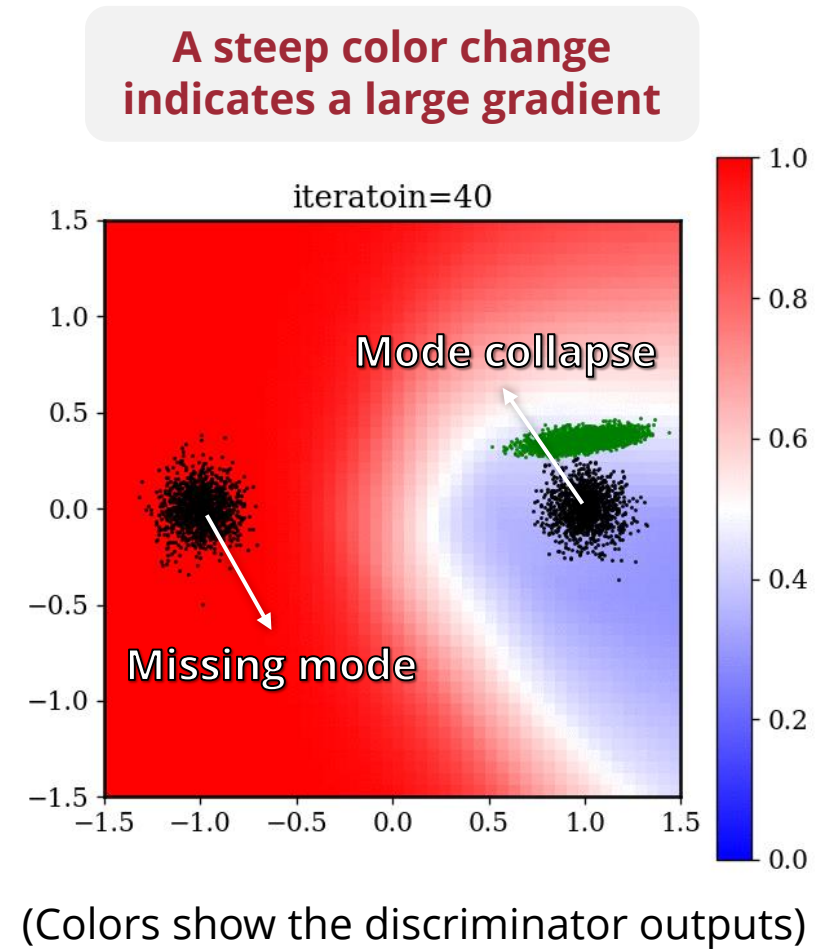
**Dog → Cat**



**Cat → Dog**

# Problems of Unregularized GANs

- **Key**—discriminator provides generator with gradients as a **guidance for improvement**
  - Discriminator has an easier job than the generator
  - Discriminator tends to provide large gradients
  - Results in unstable training of the generator
- Common failure cases
  - **Mode collapse**
  - **Missing modes**

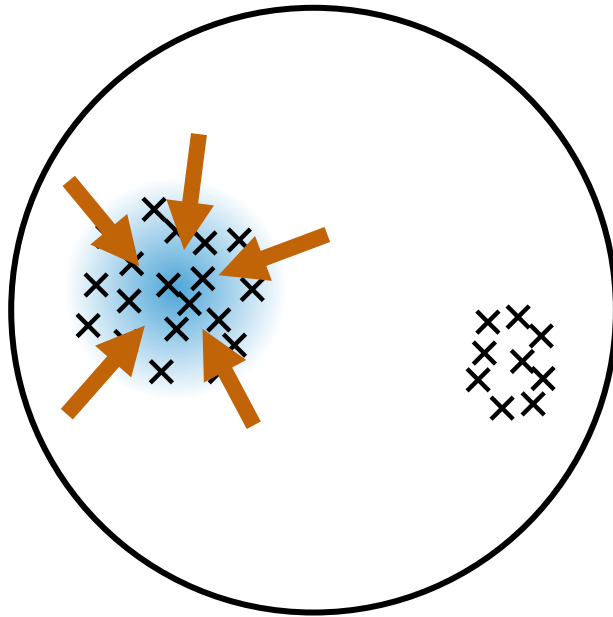


# Regularizing GANs

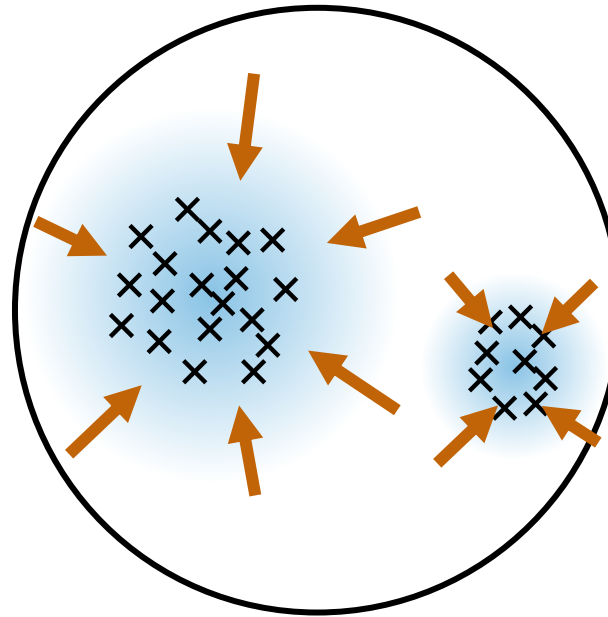
## Advantages of gradient regularization

- Provide a smoother guidance to the generator
- Alleviate mode collapse and missing modes issues

Unregularized

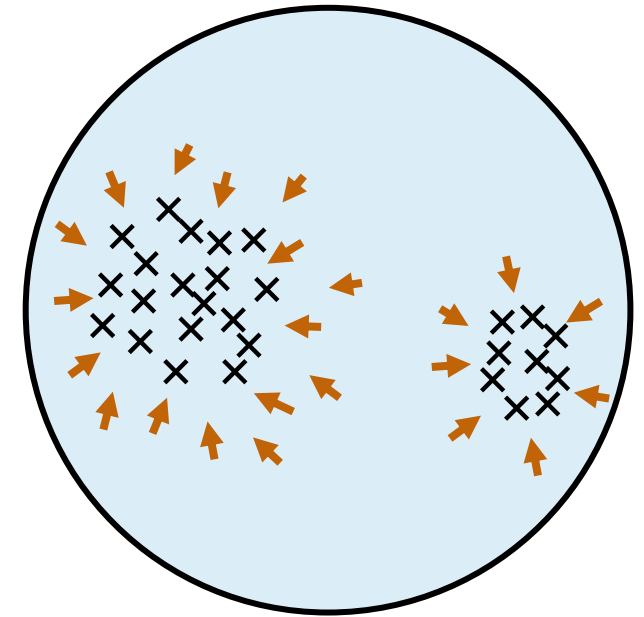


Locally regularized



Gradient clipping [1]  
Gradient penalties [2,3]

Globally regularized



Spectral normalization [4]

[1] Martin Arjovsky, Soumith Chintala, and Léon Bottou, "Wasserstein Generative Adversarial Networks," *ICML*, 2017.

[2] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville, "Improved Training of Wasserstein GANs," *NeurIPS*, 2017.

[3] Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira, "On Convergence and Stability of GANs," *arXiv preprint arXiv:1705.07215*, 2017.

[4] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida, "Spectral Normalization for Generative Adversarial Networks," *ICLR*, 2018.