

PAT 498/598 (Winter 2025)

# Music & AI

## **Lecture 17: Latent-based Music Generation**

Instructor: Hao-Wen Dong



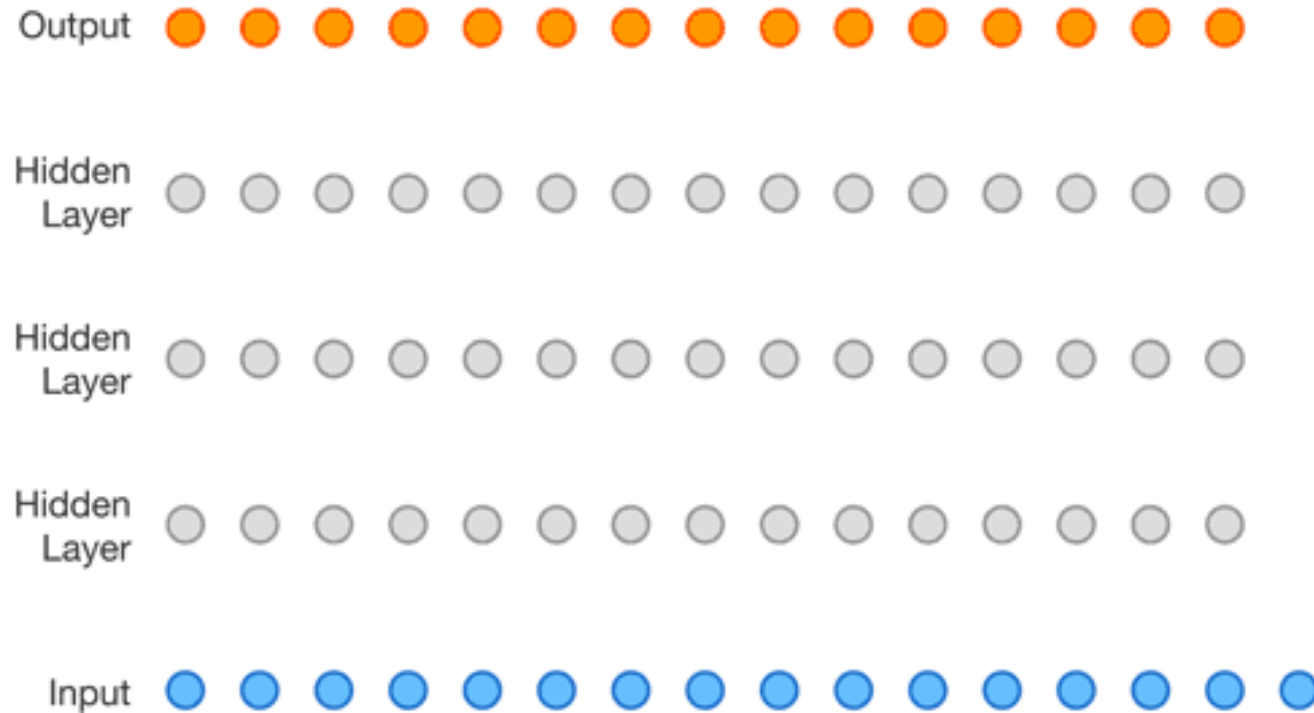
SCHOOL OF MUSIC, THEATRE & DANCE  
PERFORMING ARTS TECHNOLOGY  
UNIVERSITY OF MICHIGAN

# (Recap) Generating Waveforms using a Neural Network



(Source: van den Oord et al., 2016)

# (Recap) WaveNet (van den Oord et al., 2016)

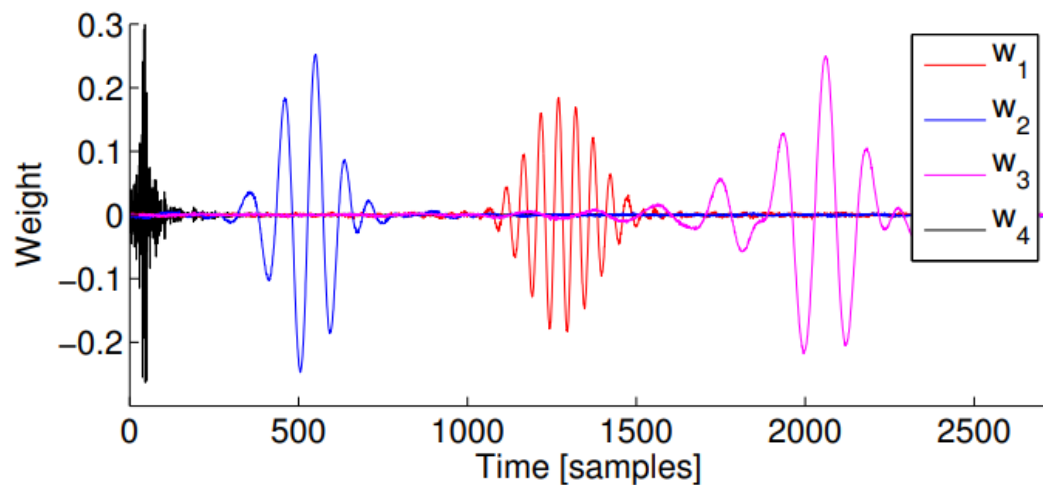


(Source: van den Oord et al., 2016)

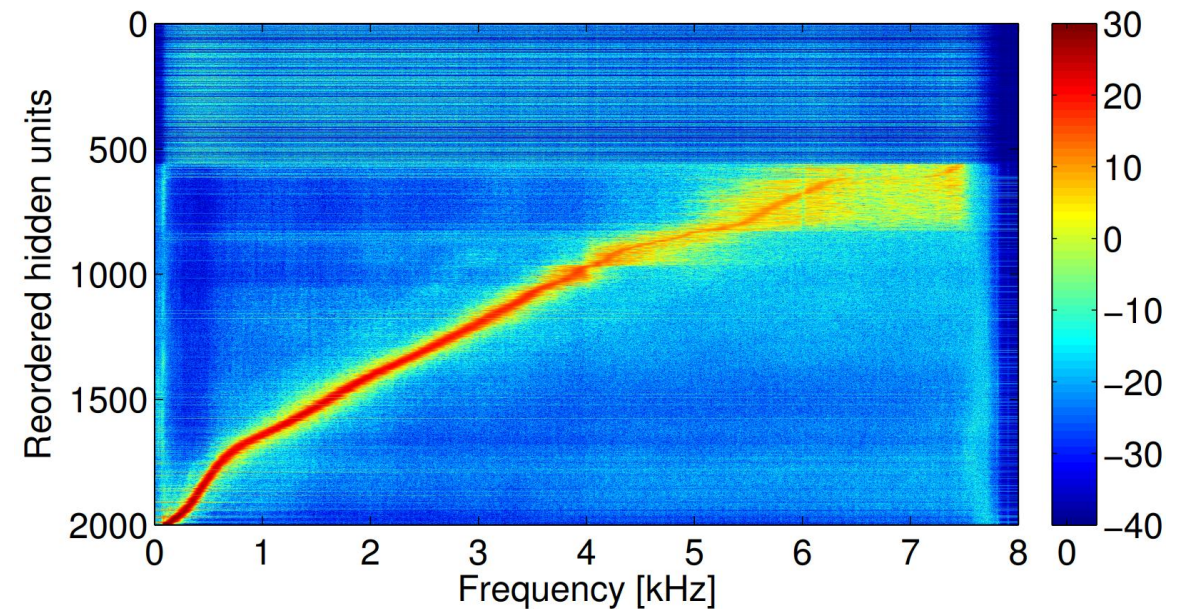
**A convolutional neural network for raw waveform generation**

# (Recap) 1D CNNs & Fourier Transform

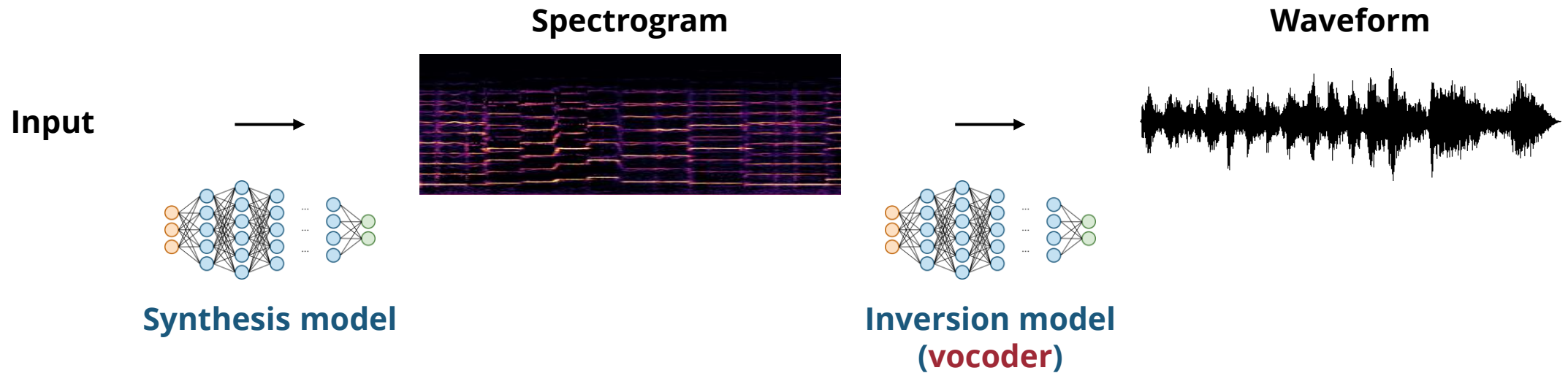
Convolution kernels learned



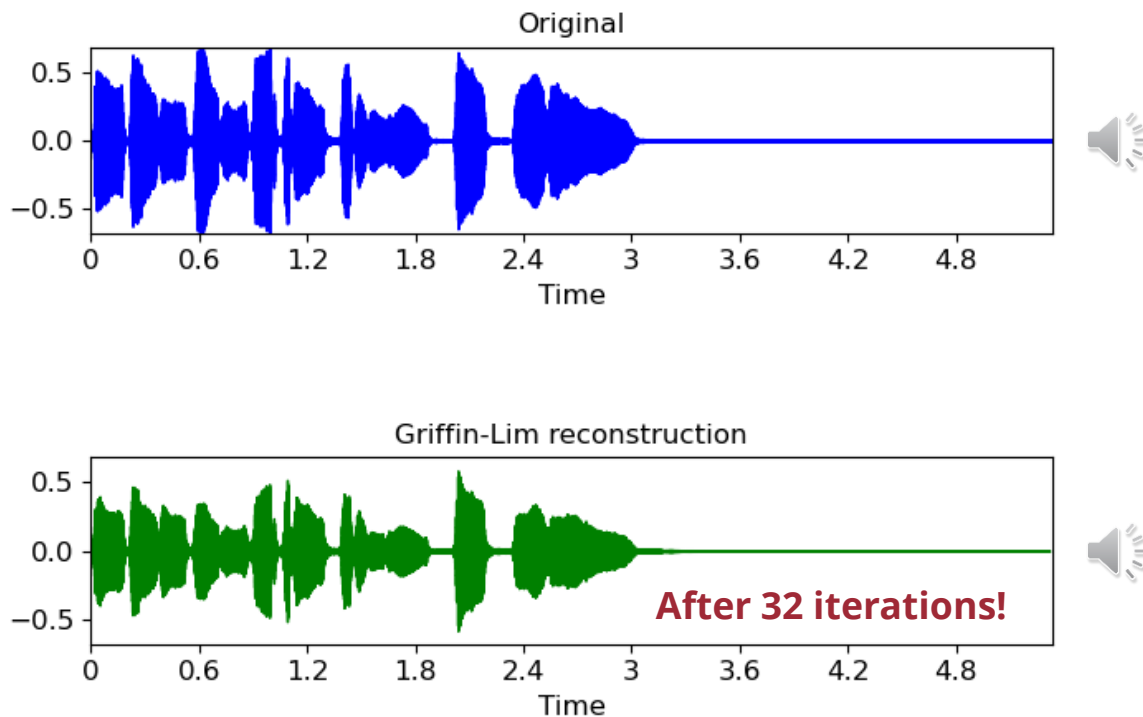
Peak frequency detected by the learned kernels



# (Recap) Frequency-domain Audio Synthesis



# (Recap) Griffin-Lim Algorithm (Griffin & Lim, 1984)



(Source: librosa documentation)

Given a magnitude-only STFT matrix



Randomly initialize the phase



$$y' = \arg \min_y (M - \text{STFT}(y))^2$$

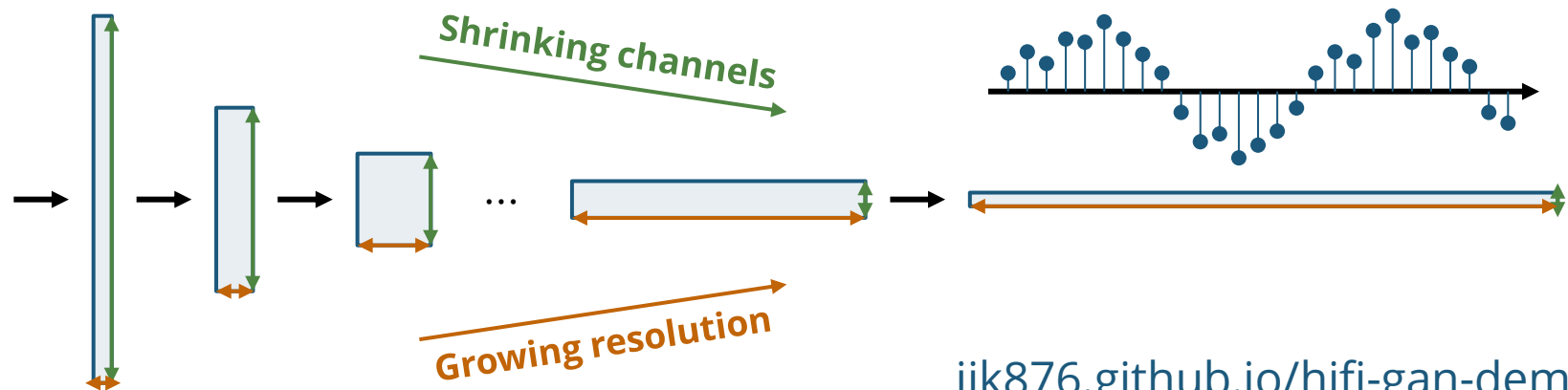
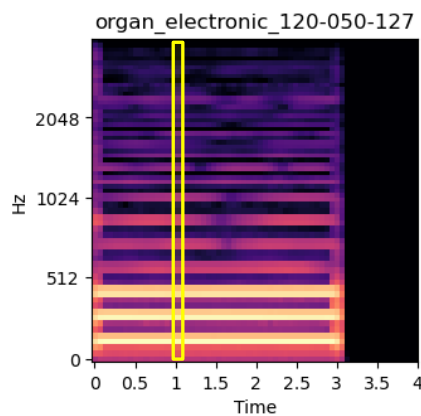
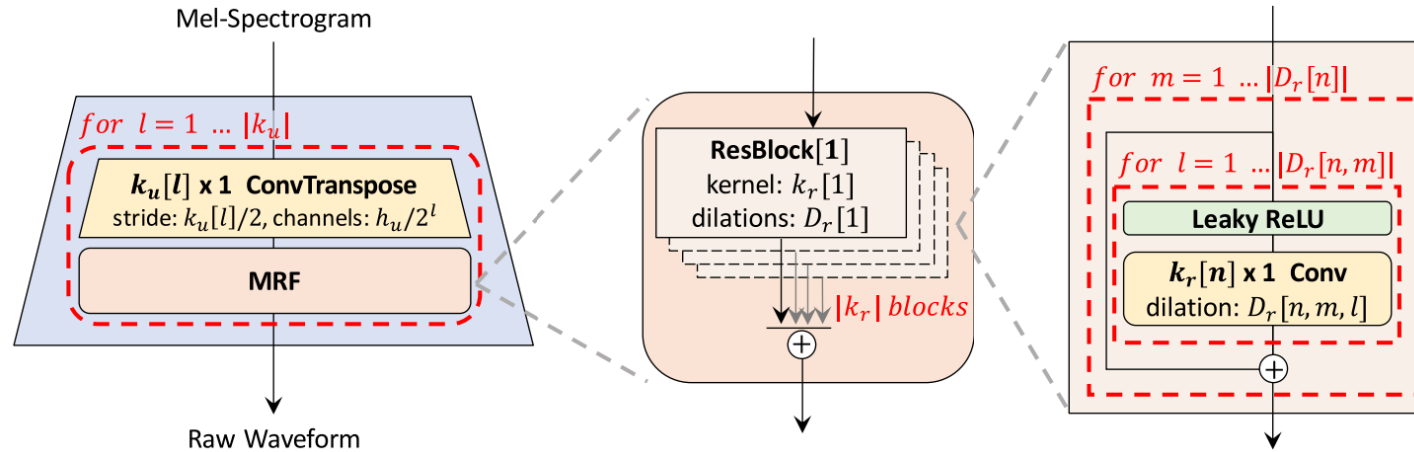
Find the signal  $y$  that minimize the MSE between the input and  $\text{STFT}(y)$



$$M' = \text{STFT}(y')$$

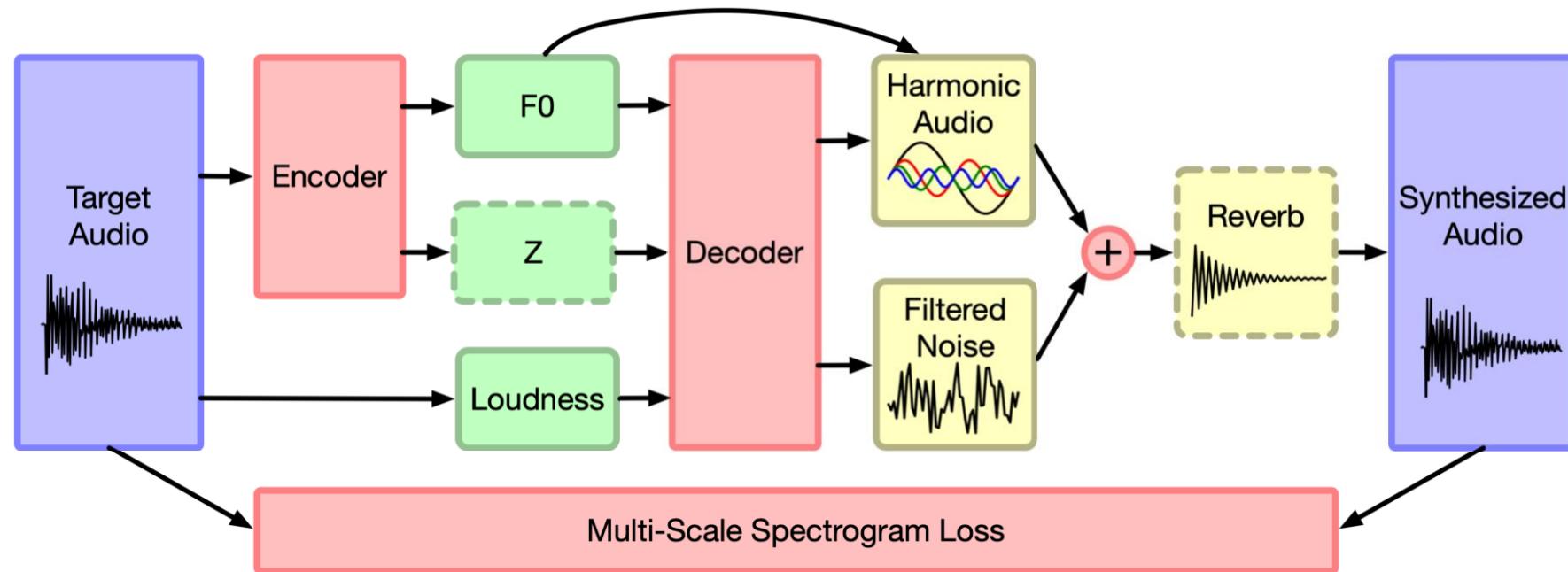
Find the STFT of the signal  $y$

# (Recap) HiFi-GAN (Kong et al., 2020)



[jik876.github.io/hifi-gan-demo](https://jik876.github.io/hifi-gan-demo)

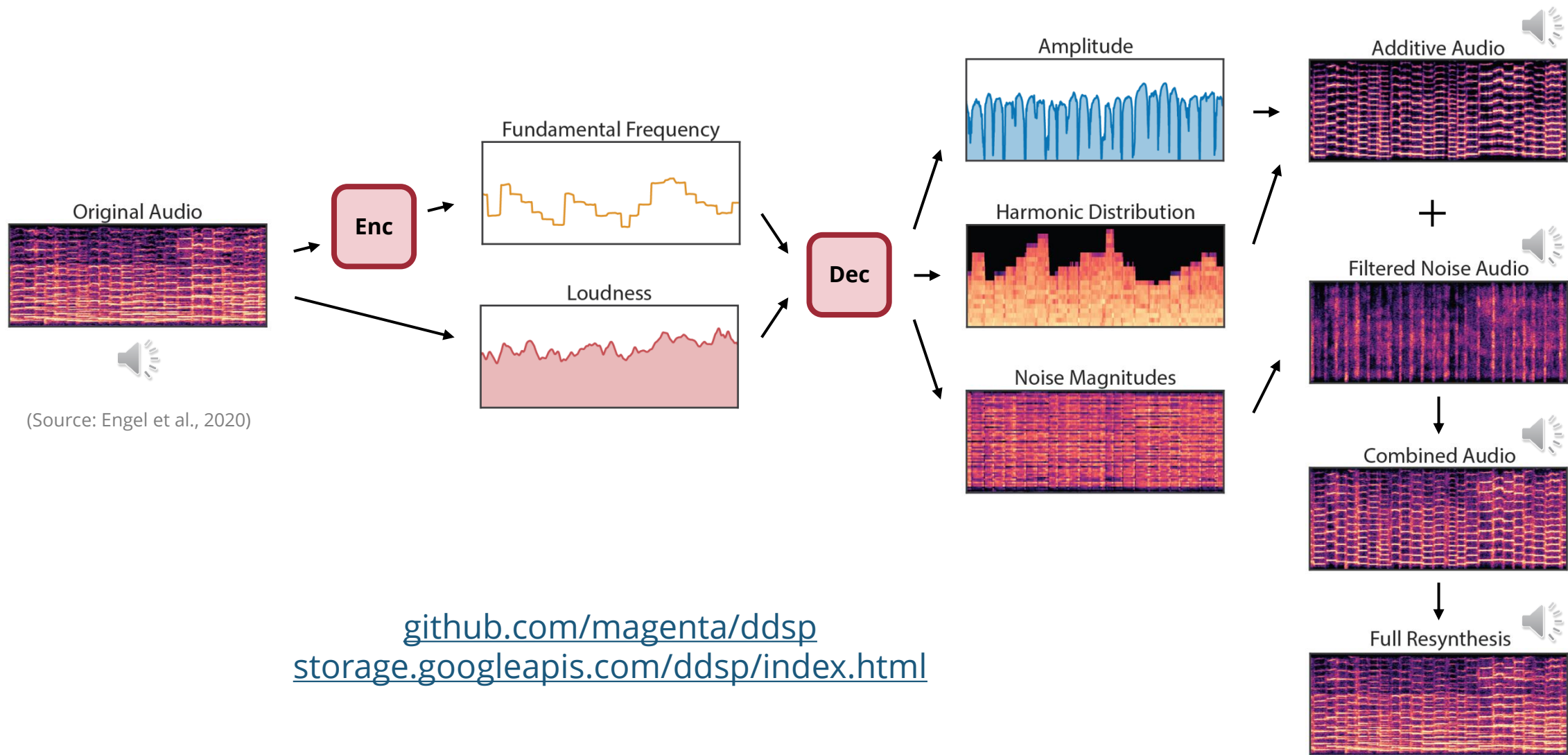
# (Recap) Differentiable DSP (DDSP) (Engel et al., 2020)



(Source: Engel et al., 2020)



# (Recap) Differentiable DSP (DDSP) (Engel et al., 2020)



# Neural Codecs

# What is a Codec?



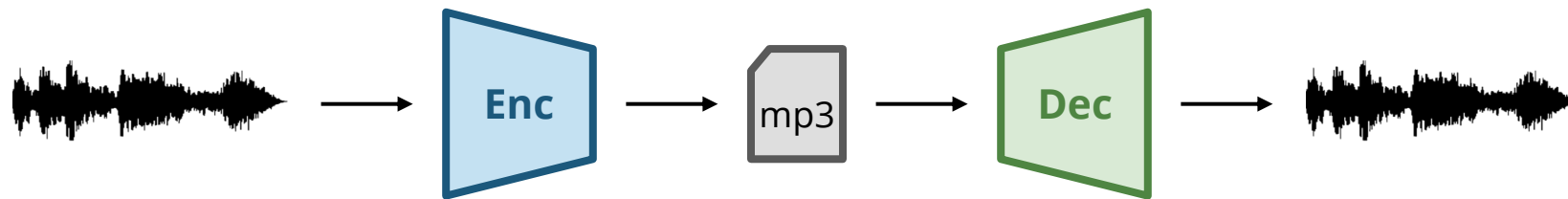
**SBC, AAC, aptX,  
aptX HD, LDAC**

# What is a Codec?

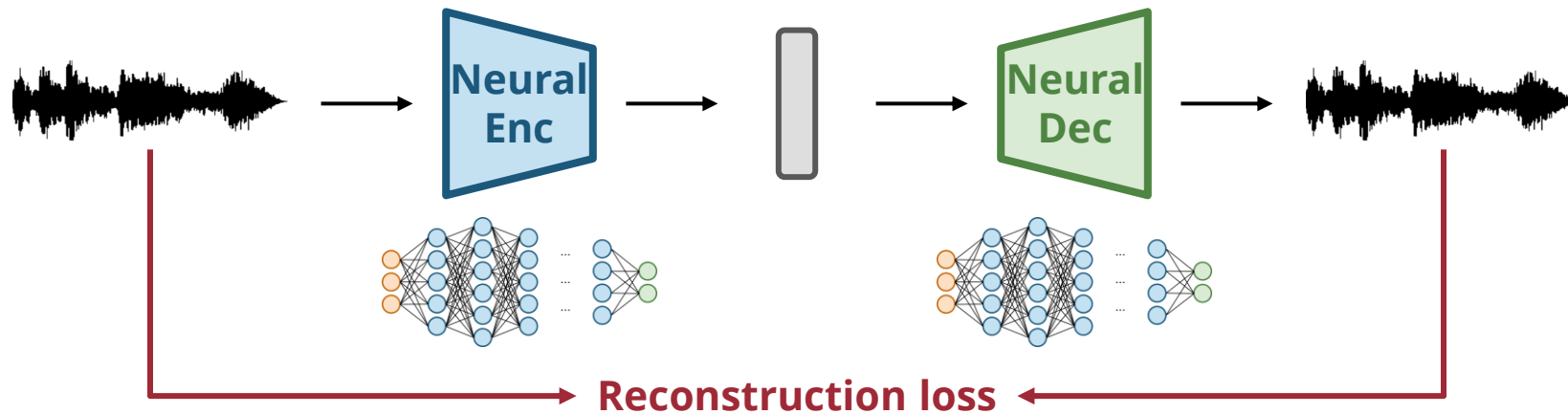


# Neural Codec

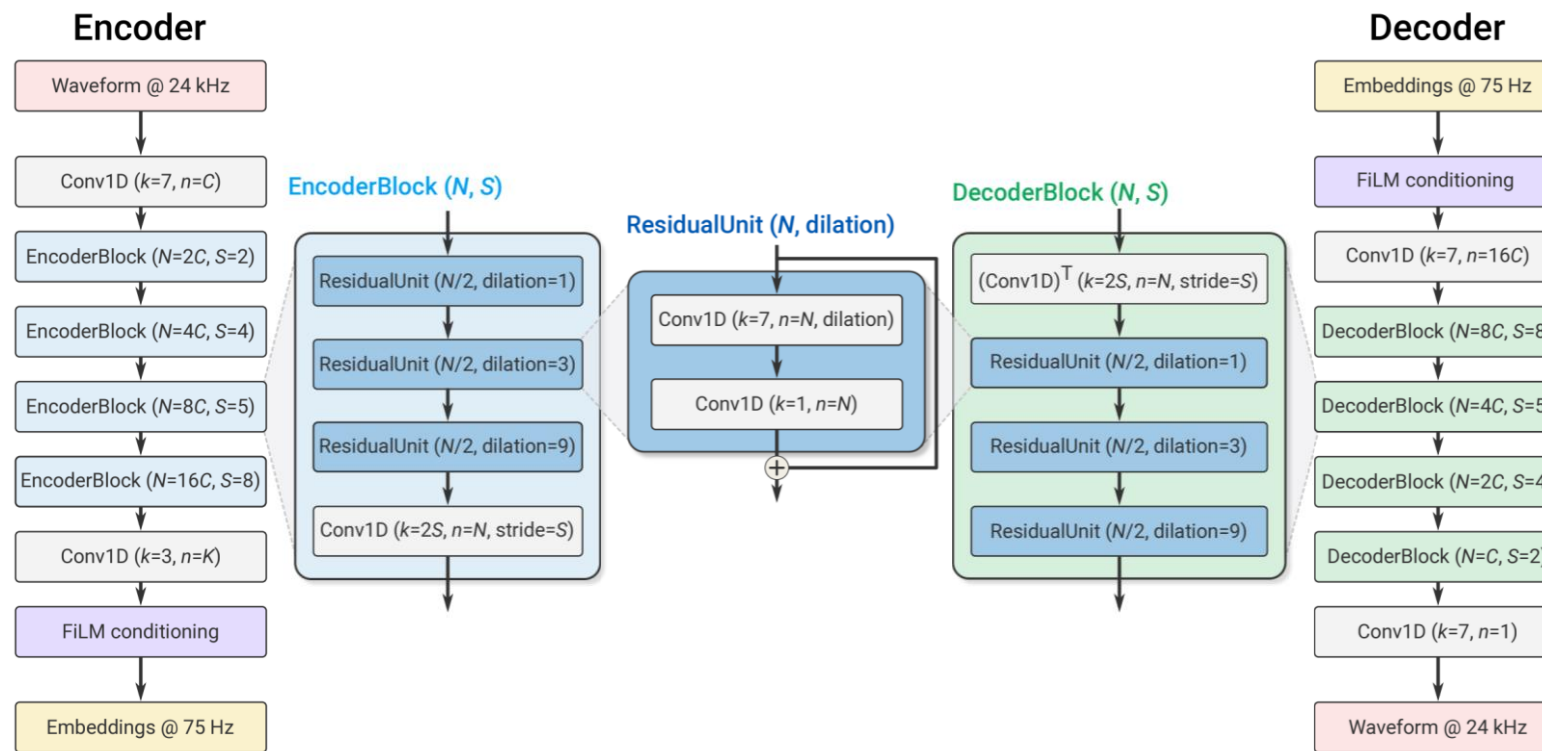
## Traditional Codec



## Neural Codec

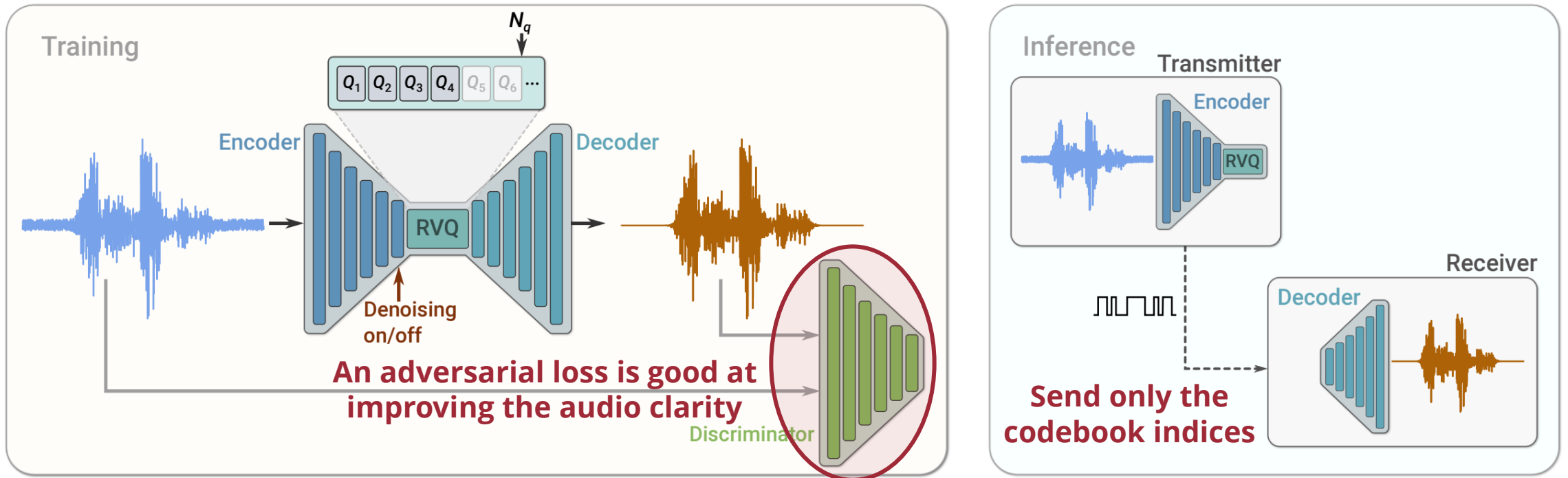


# SoundStream (Zeghidour et al., 2021)



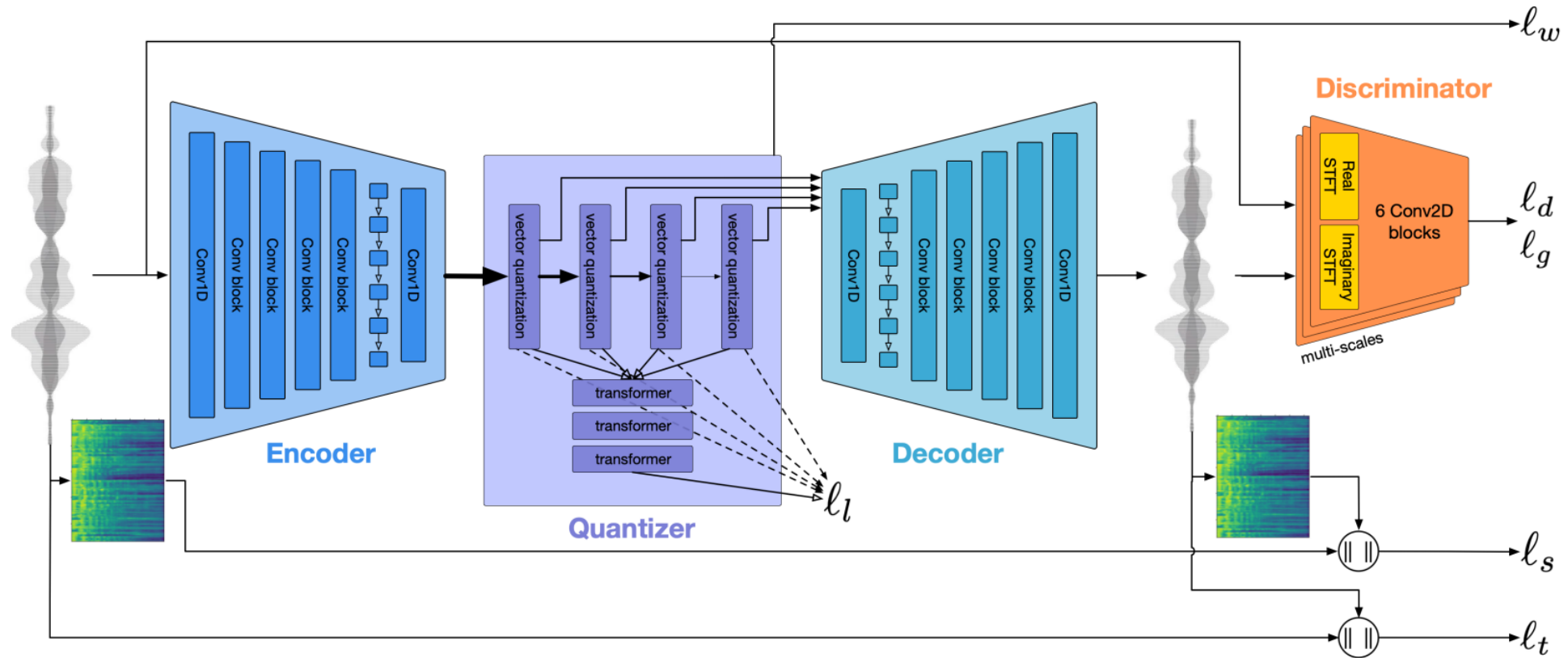
(Source: Zeghidour et al., 2021)

# SoundStream (Zeghidour et al., 2021)



(Source: Zeghidour et al., 2021)

# EnCodec (Défossez et al., 2022)



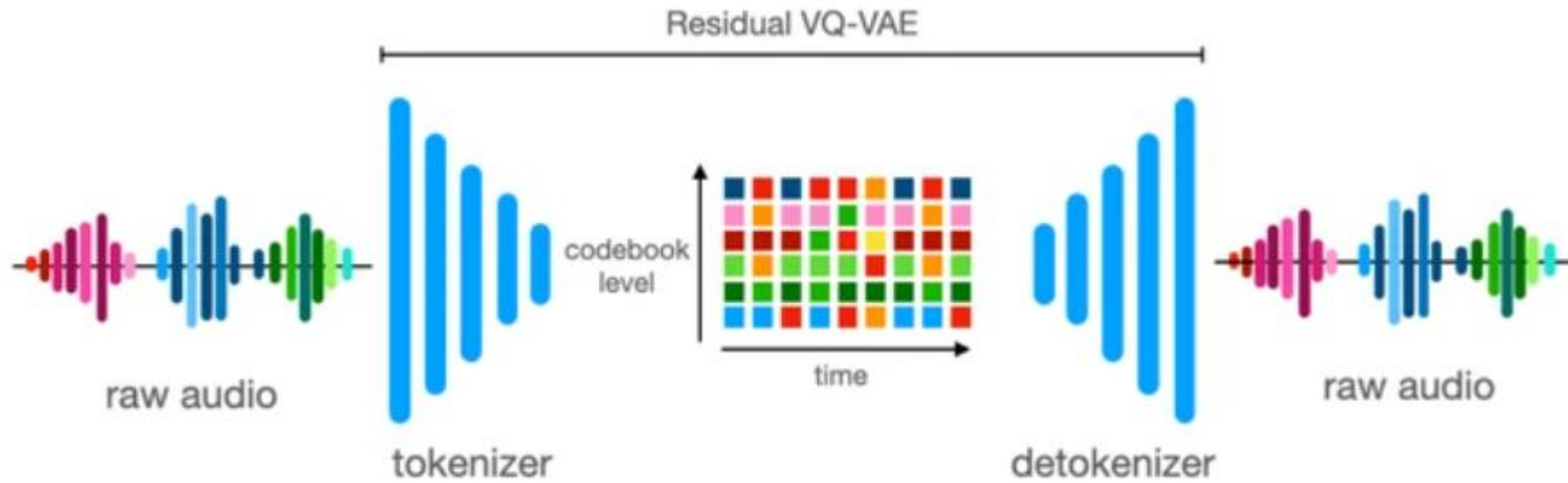
(Source: Défossez et al., 2022)

[ai.honu.io/papers/encodec/samples.html](https://ai.honu.io/papers/encodec/samples.html)

[github.com/facebookresearch/encodec](https://github.com/facebookresearch/encodec)



# Descript Audio Codec (Kumar et al., 2023)



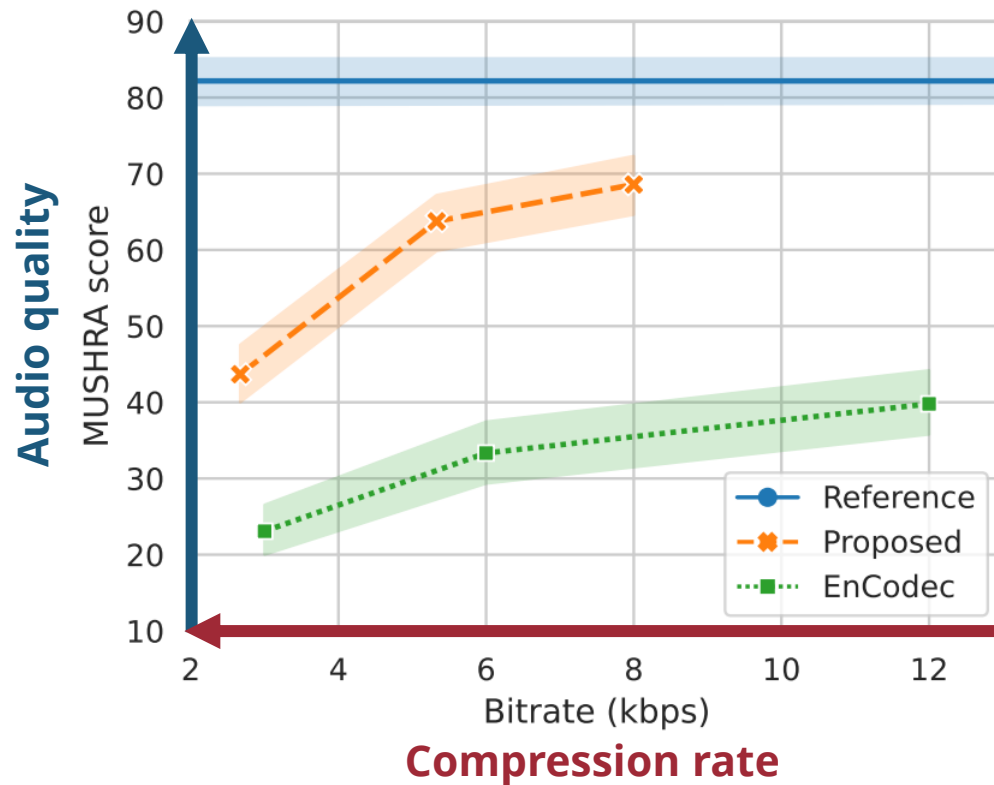
(Source: Kumar et al., 2023)

[descript.notion.site/Descript-Audio-Codec-11389fce0ce2419891d6591a68f814d5](https://descript.notion.site/Descript-Audio-Codec-11389fce0ce2419891d6591a68f814d5)

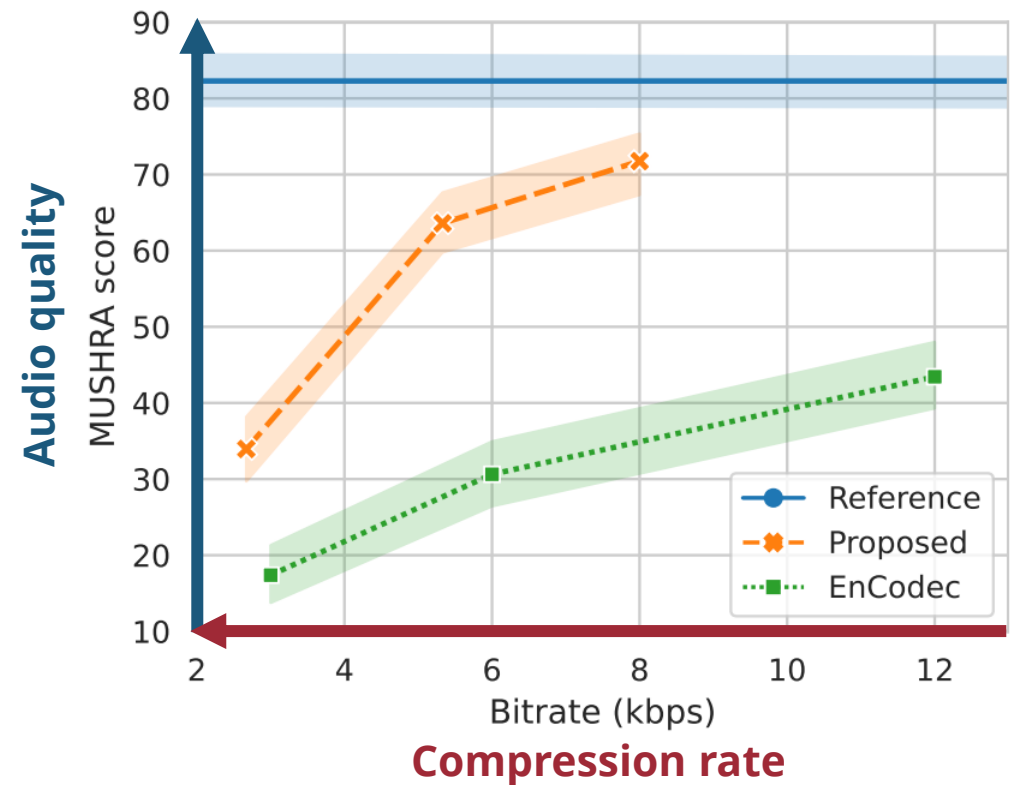
[github.com/descriptinc/descript-audio-codec](https://github.com/descriptinc/descript-audio-codec)

# Descript Audio Codec (Kumar et al., 2023)

Listening Test Results @ 44.1 kHz



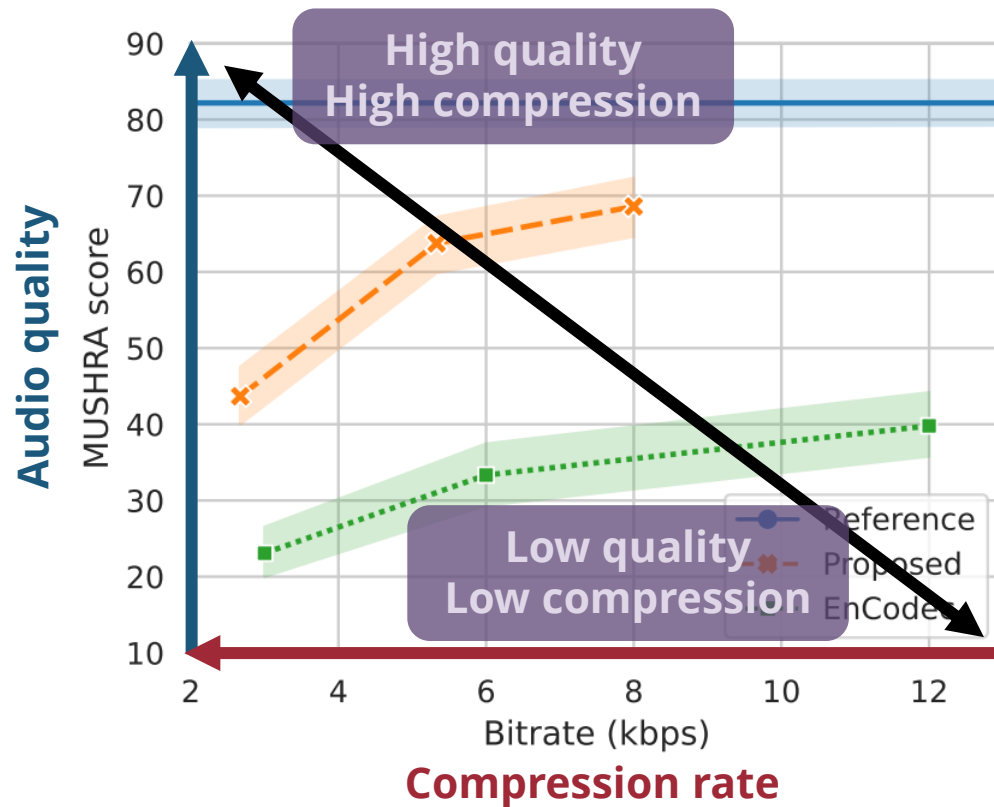
Listening Test Results @ 24 kHz



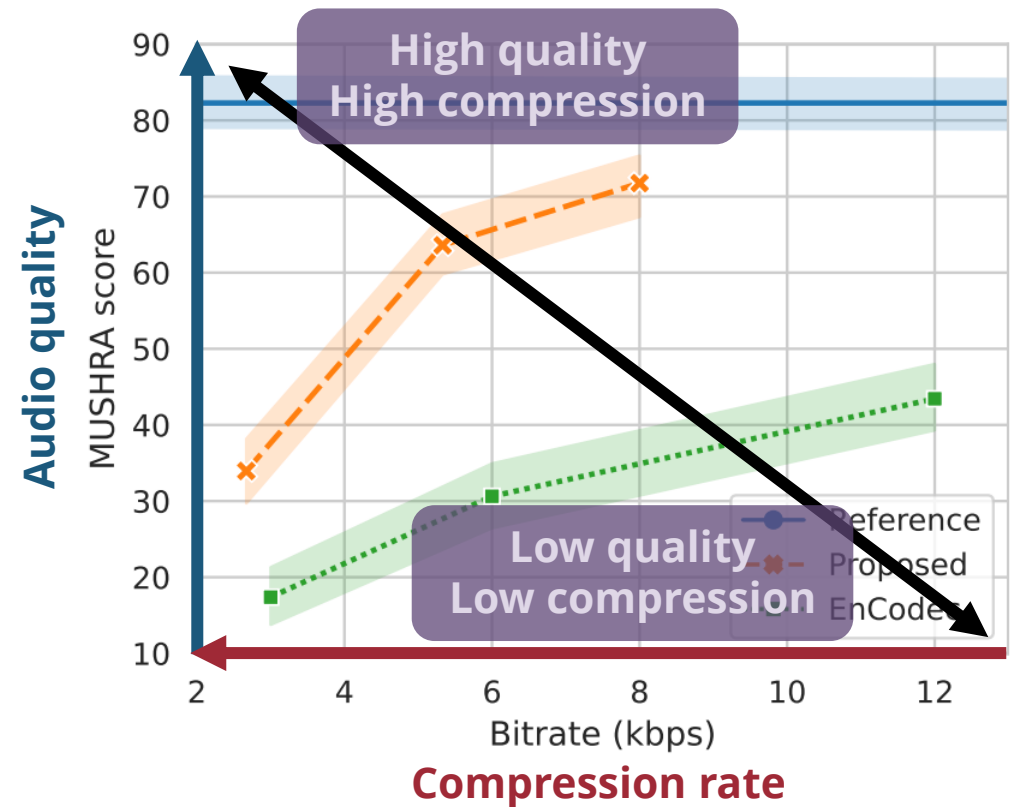
(Source: Kumar et al., 2023)

# Descript Audio Codec (Kumar et al., 2023)

Listening Test Results @ 44.1 kHz



Listening Test Results @ 24 kHz

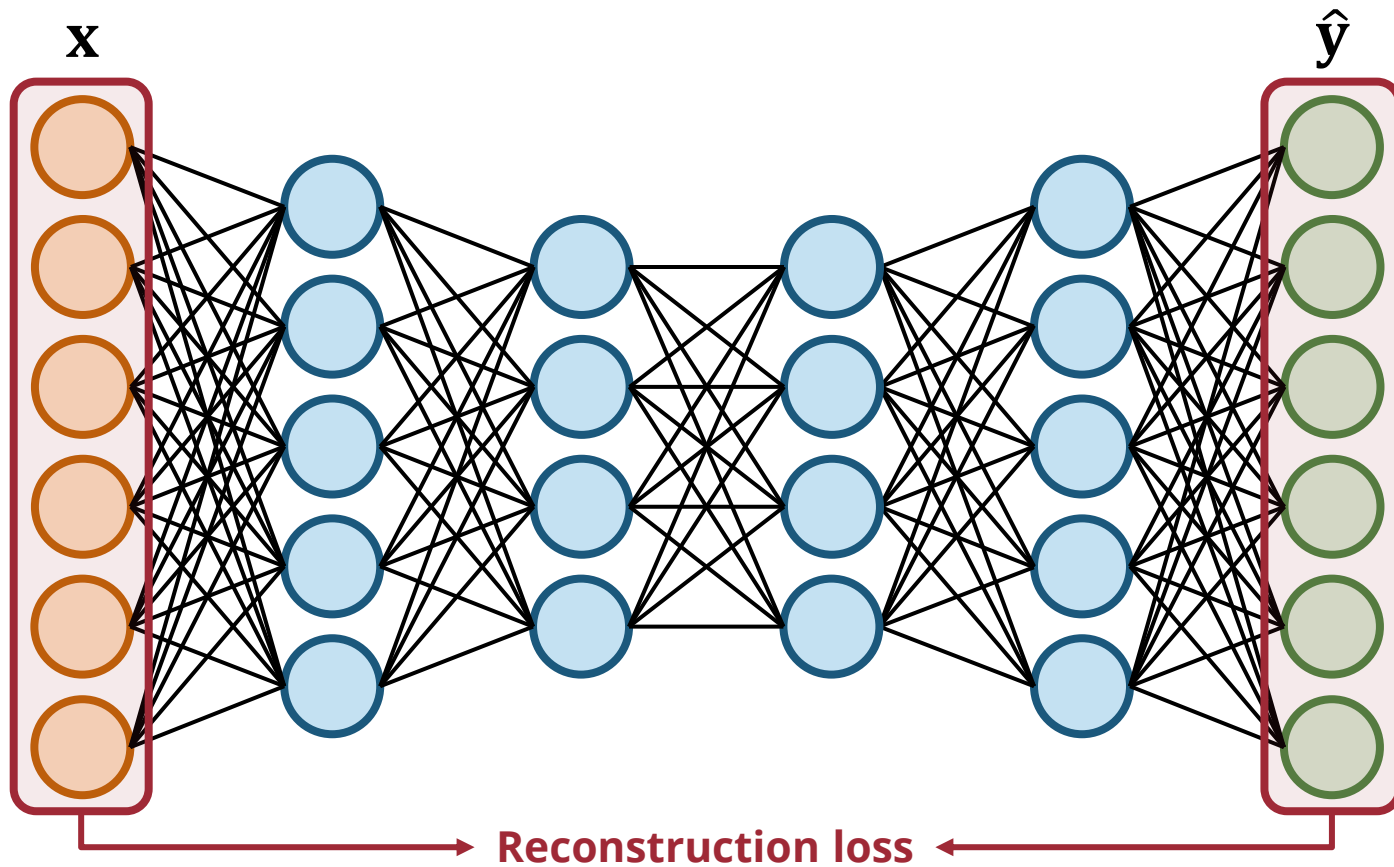


(Source: Kumar et al., 2023)

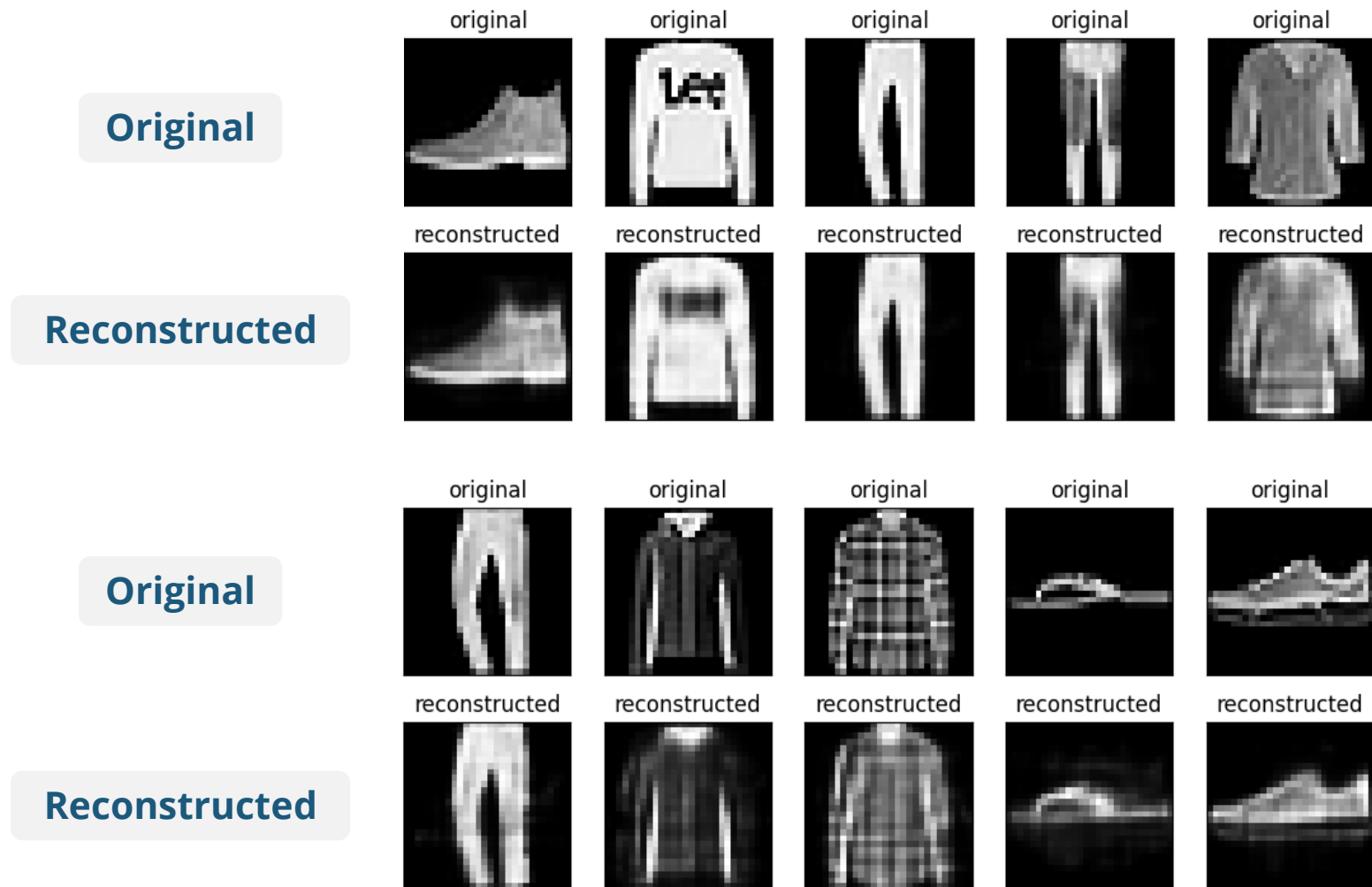
# Latent Diffusion Models

# (Recap) Autoencoders

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

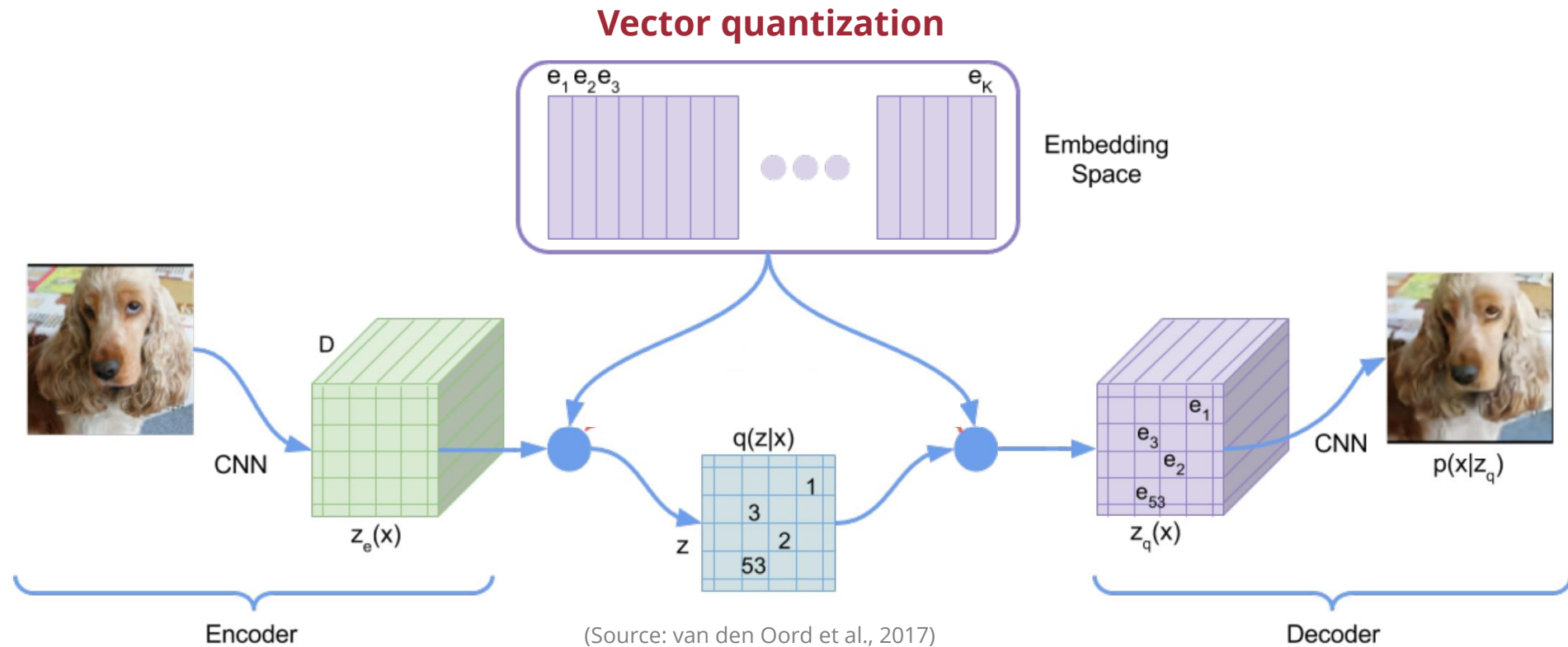


# (Recap) Autoencoders – Reconstruction Examples



(Source: tensorflow.org)

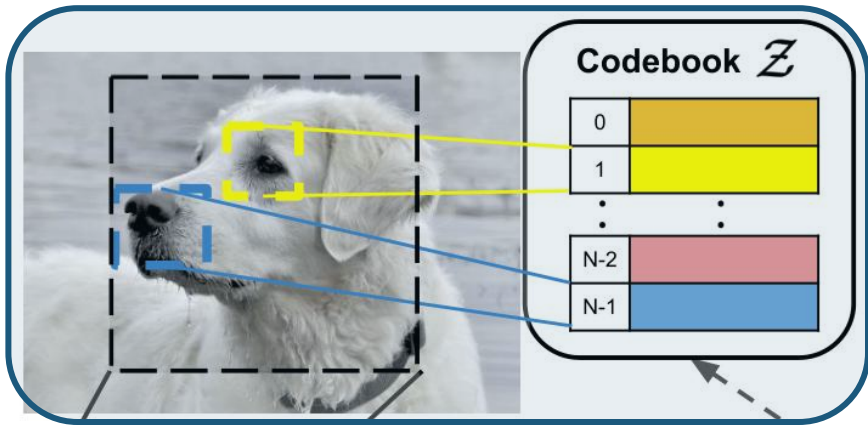
# Vector-Quantized VAEs (VQVAEs)



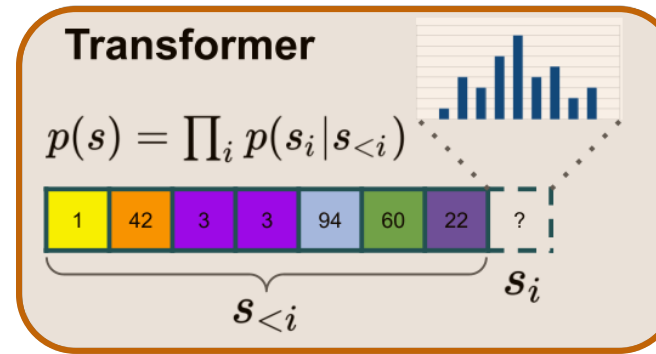
**Allow only a fixed number of vectors to be used in the bottleneck layer**

# VQGAN

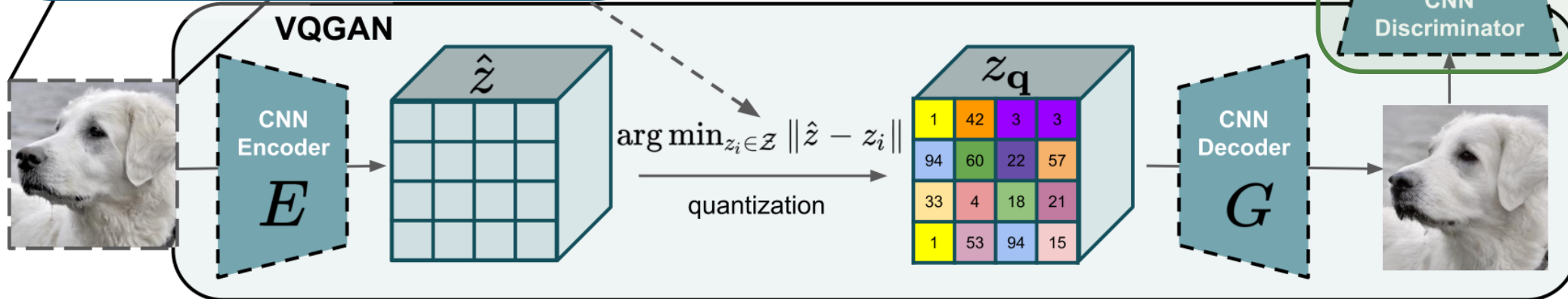
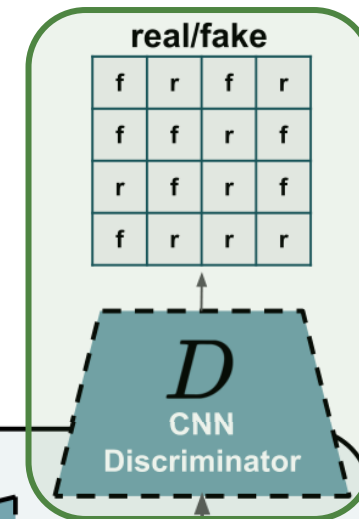
Each path is encoded into a latent code



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**

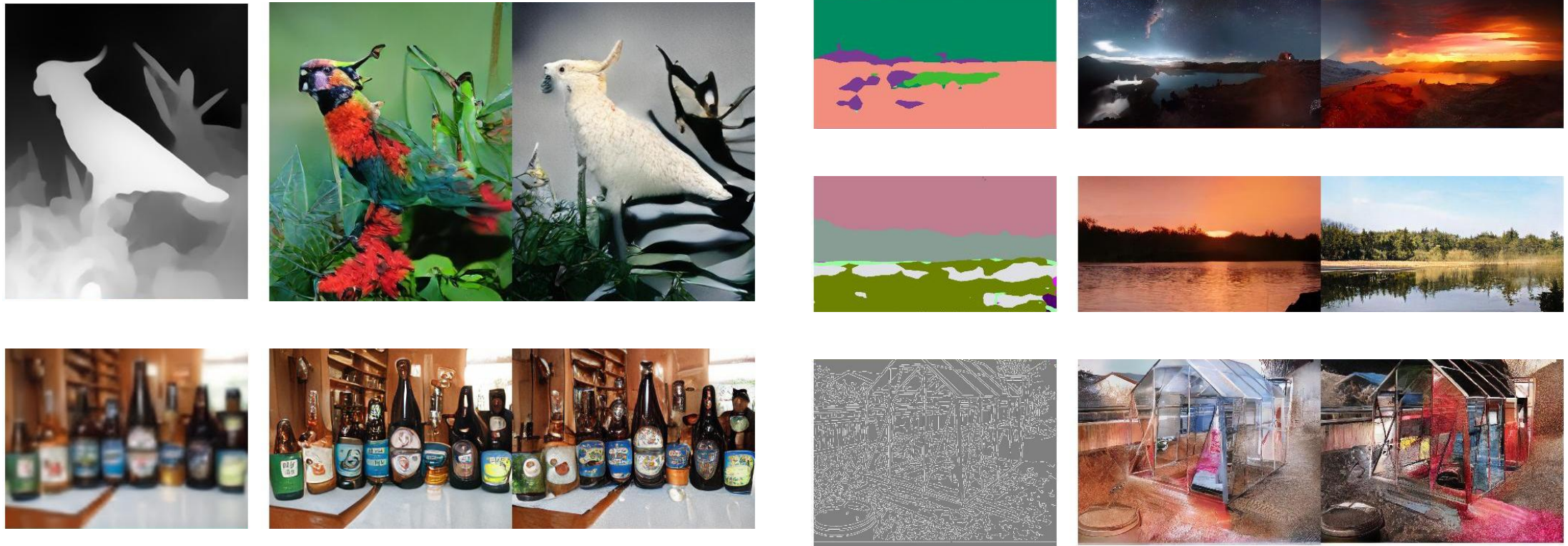


# VQGAN: Conditional Generation



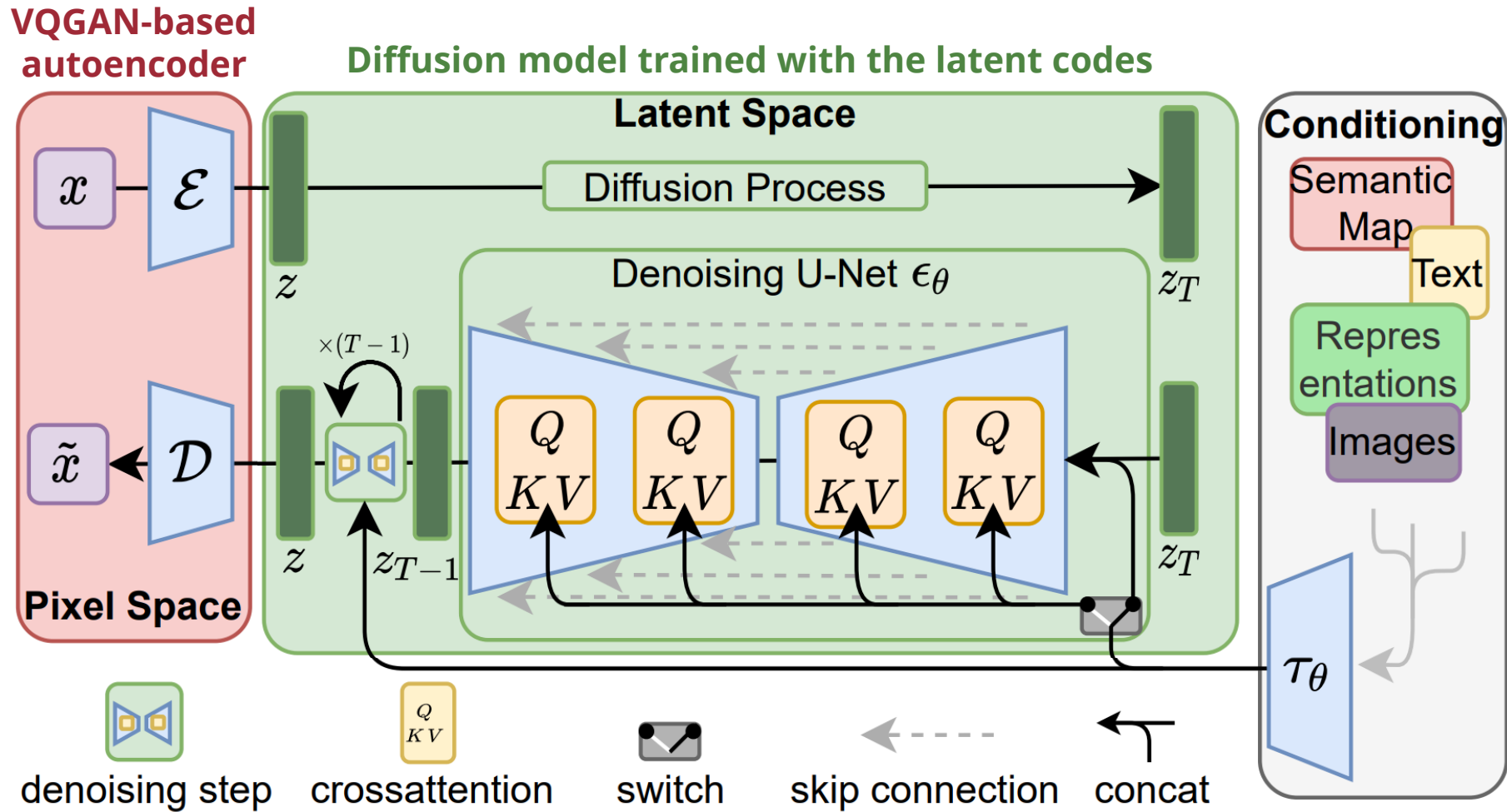
(Source: Esser et al., 2021)

# VQGAN: Conditional Generation



(Source: Esser et al., 2021)

# Latent Diffusion Models (LDMs)



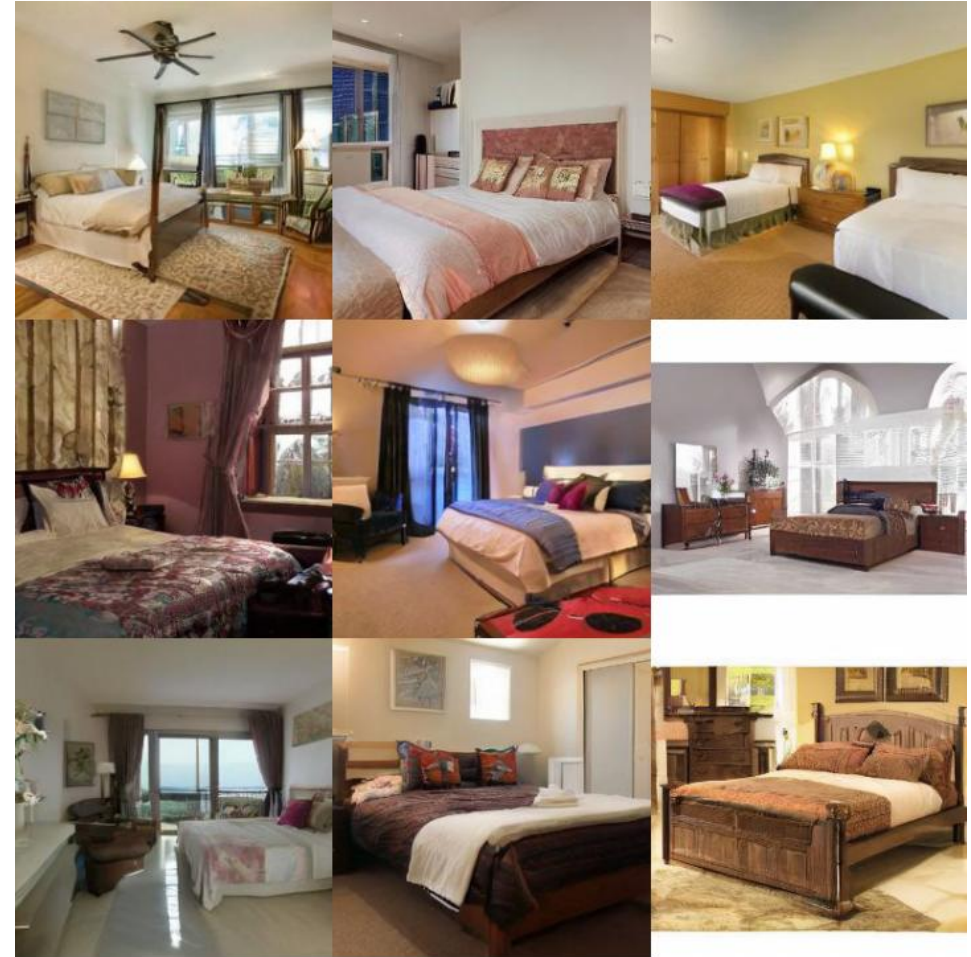
(Source: Rombach et al., 2022)

# LDMs: Examples



(Source: Rombach et al., 2022)

# LDMs: Examples



(Source: Rombach et al., 2022)

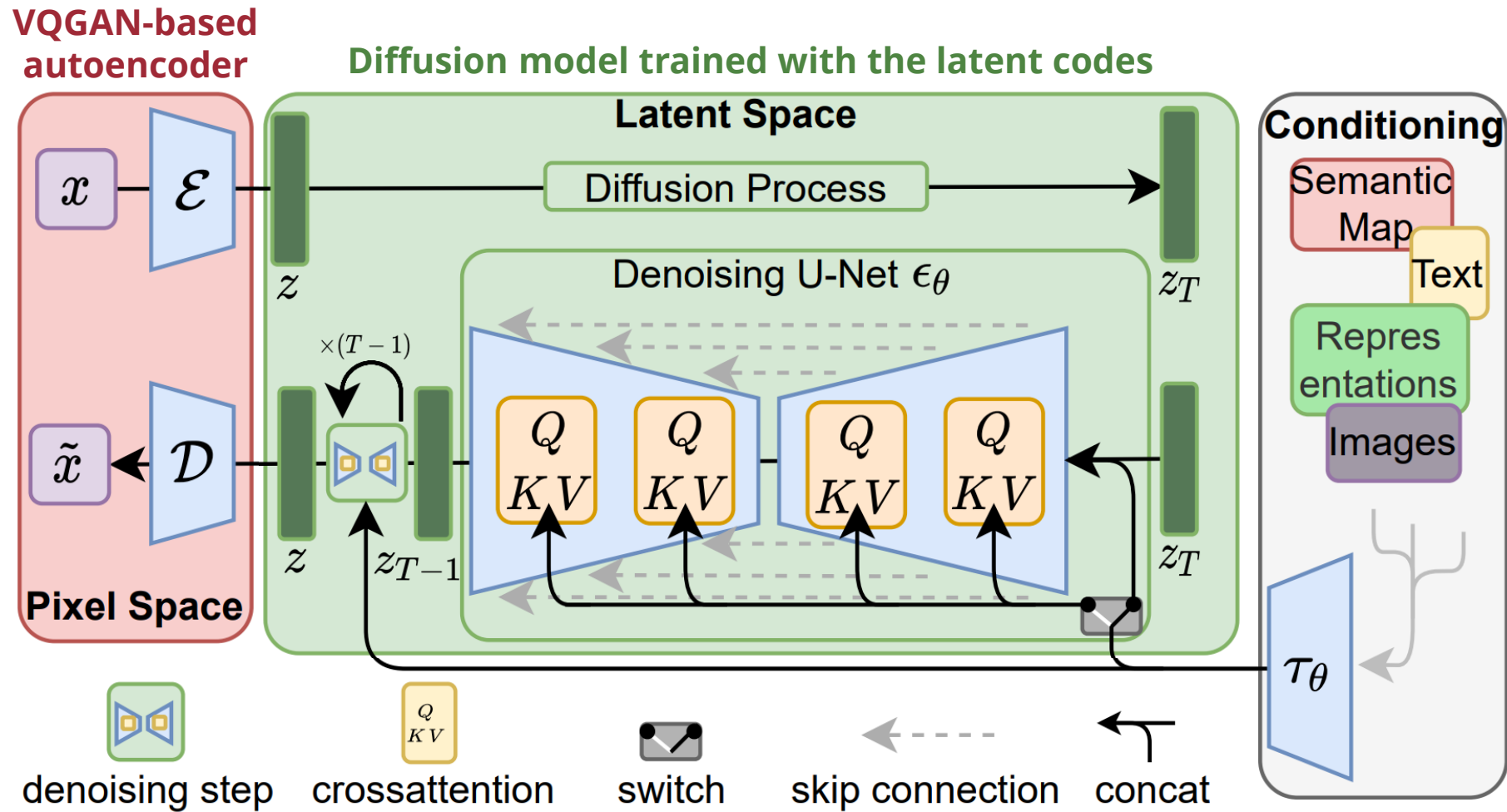
# LDMs: Semantic Synthesis



(Source: Rombach et al., 2022)



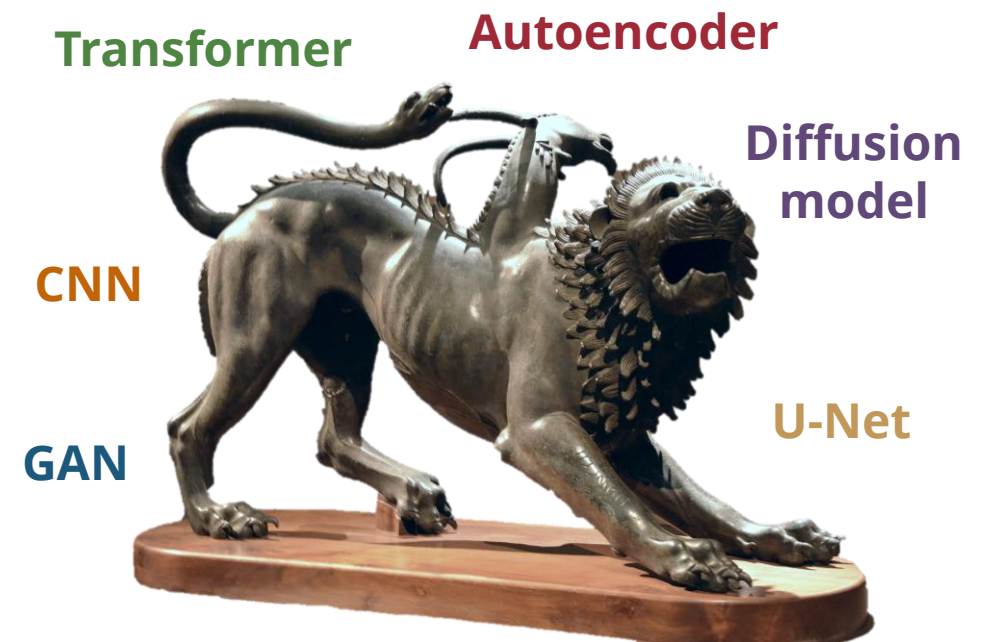
# Latent Diffusion Models (LDMs)



(Source: Rombach et al., 2022)

# Latent Diffusion Model is a Chimera

- **A neural codec**
  - An CNN-based autoencoder
  - Trained with a GAN-like adversarial loss
- **Diffusion model in the latent space**
  - A denoising U-Net
- **A conditioning module**
  - Transformer-like cross-attention mechanism

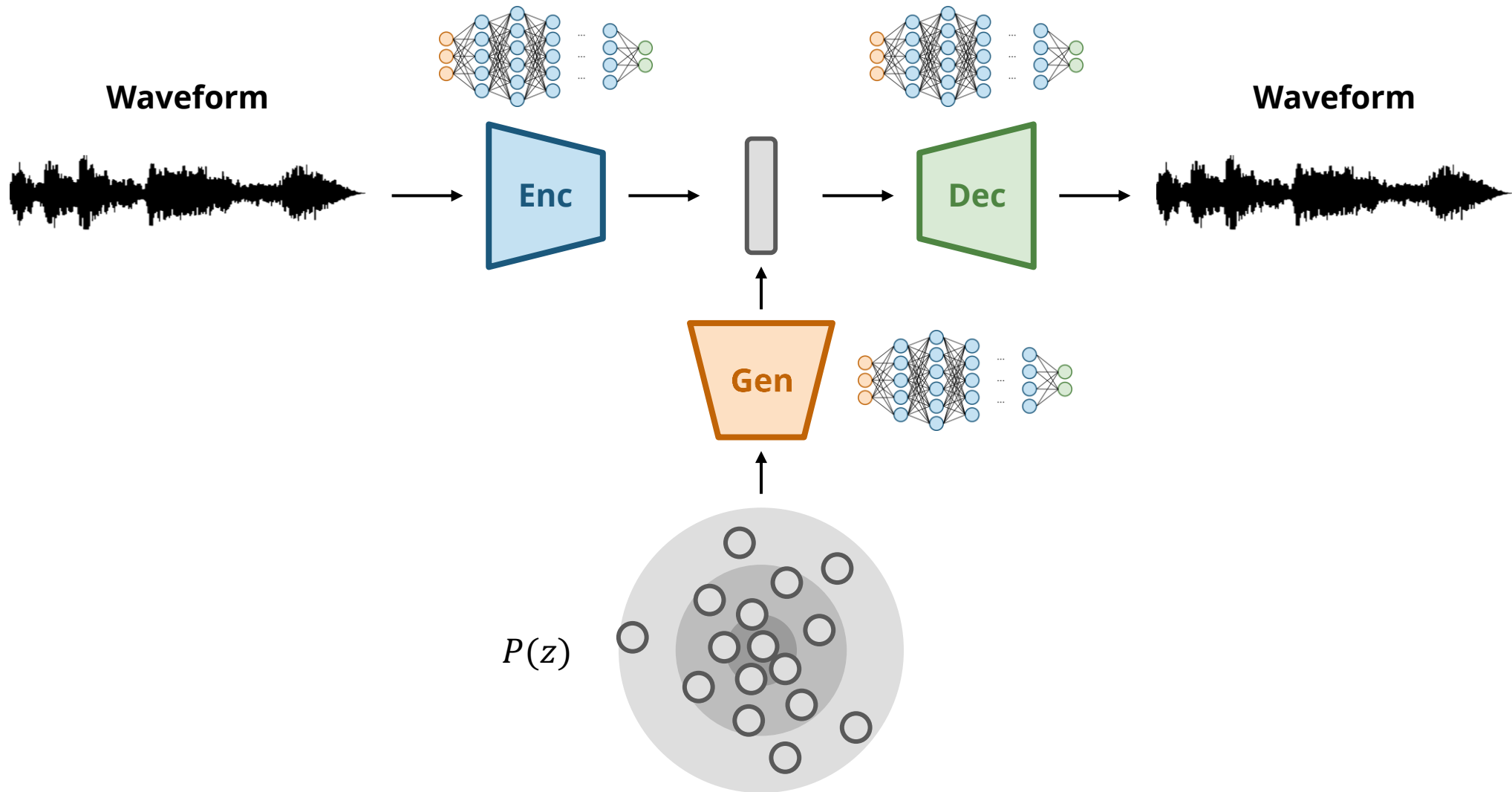


(Source: Raddato via worldhistory.org)

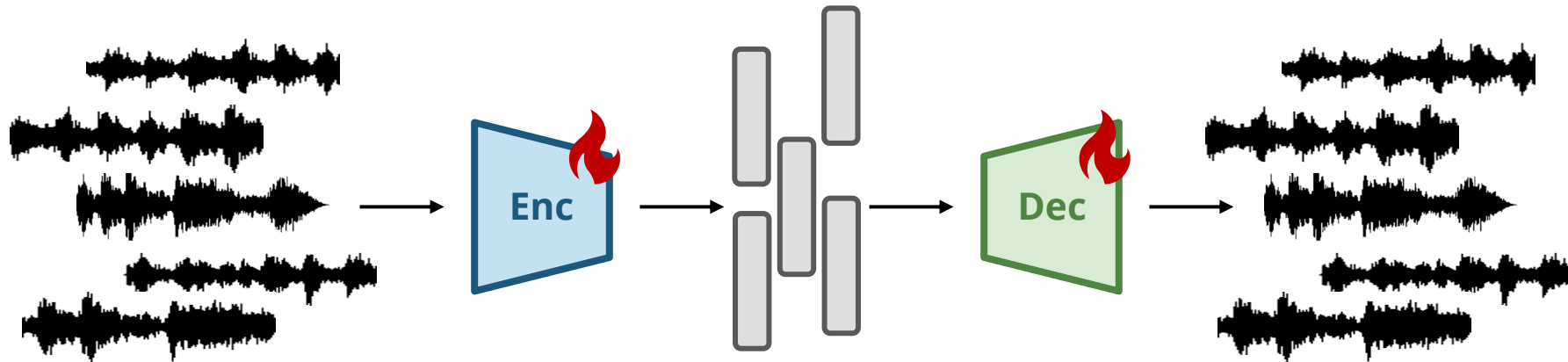


# Latent-based Audio Synthesis

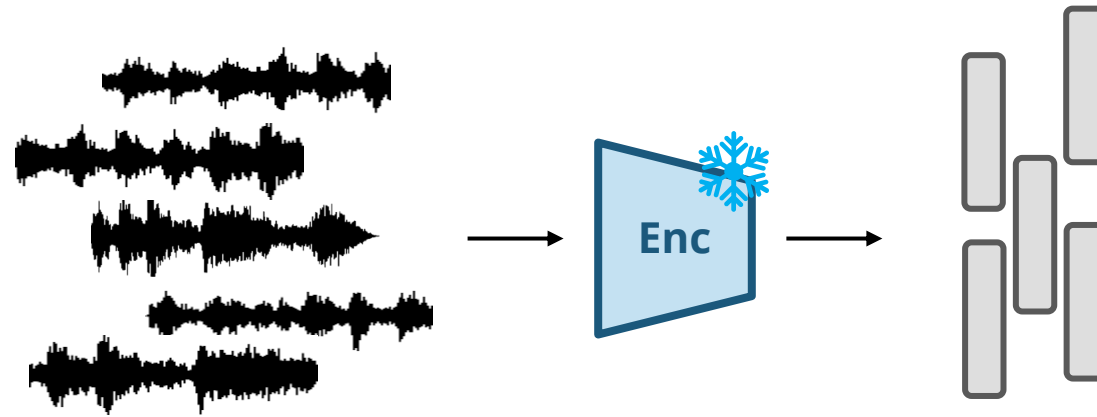
# Latent-based Audio Synthesis



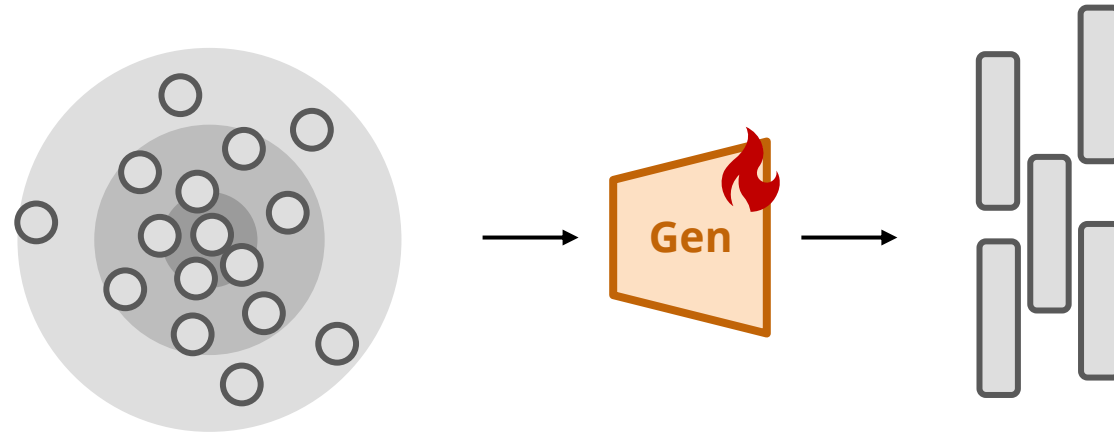
# Step 1: Train an Autoencoder



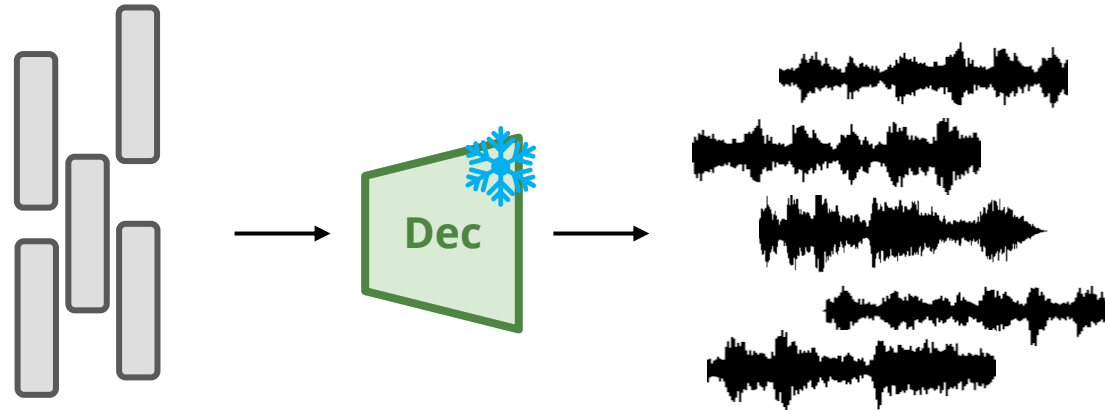
## Step 2: Compute the Latent Vectors



## Step 3: Train a Latent Generative Model

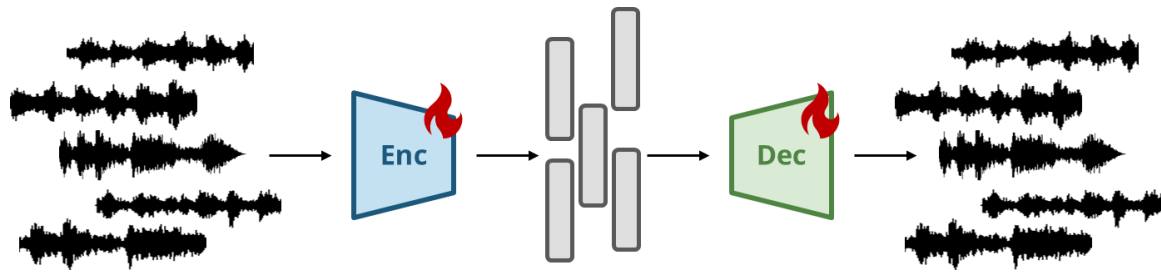


## Step 4: Decode the Latent Vectors

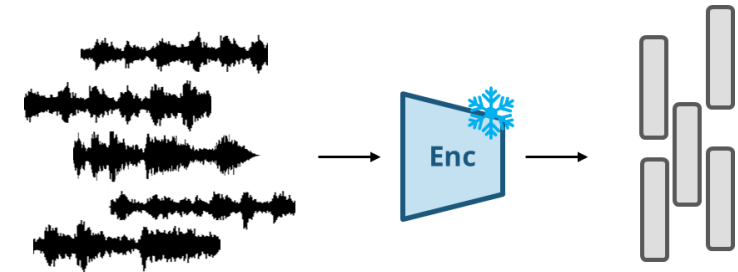


# Pipeline

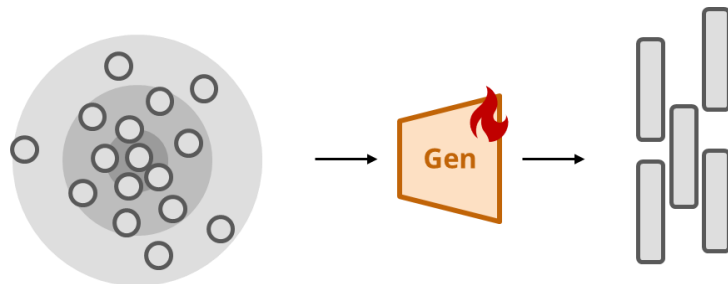
## Step 1: Train an Autoencoder



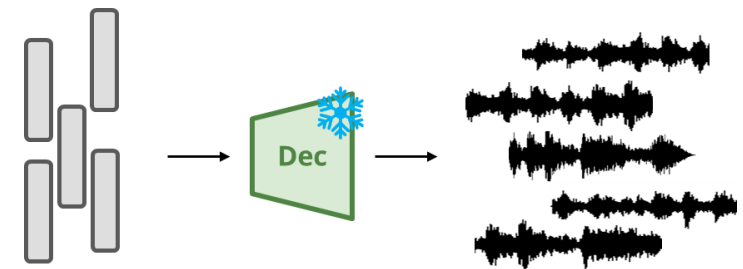
## Step 2: Compute the Latent Vectors



## Step 3: Train a Latent Generative Model

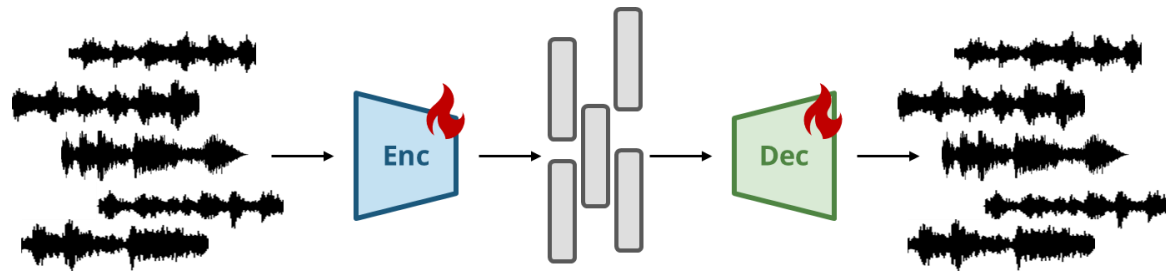


## Step 4: Decode the Latent Vectors

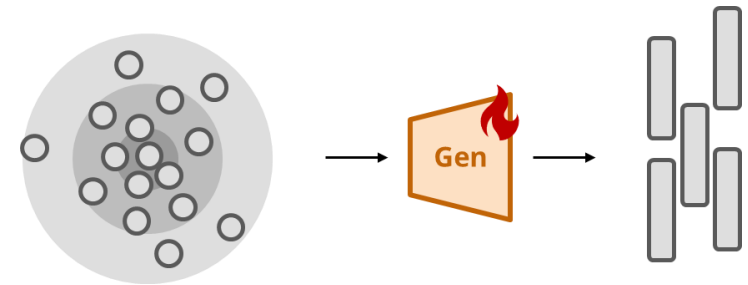


# Training

Autoencoder

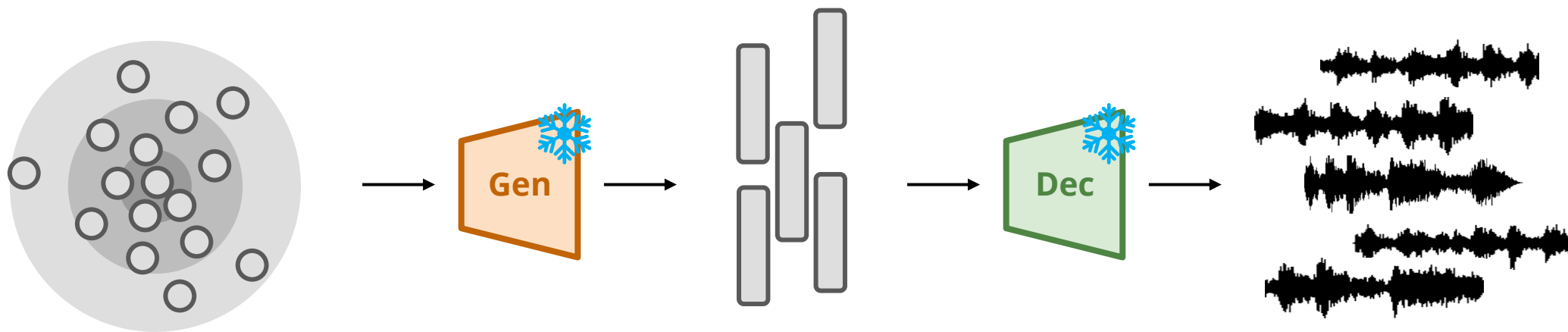


Latent Generative Model

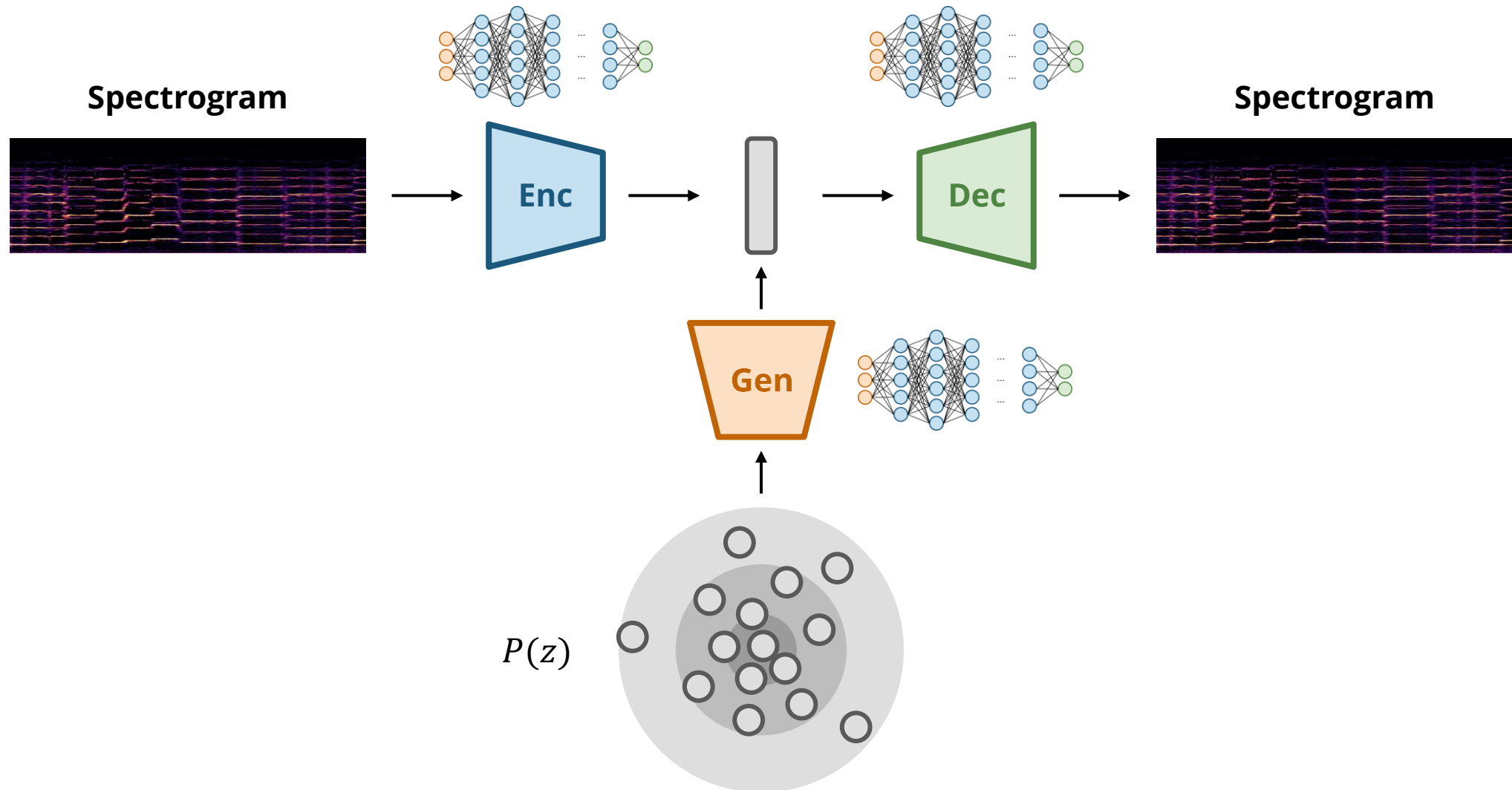




# Inference



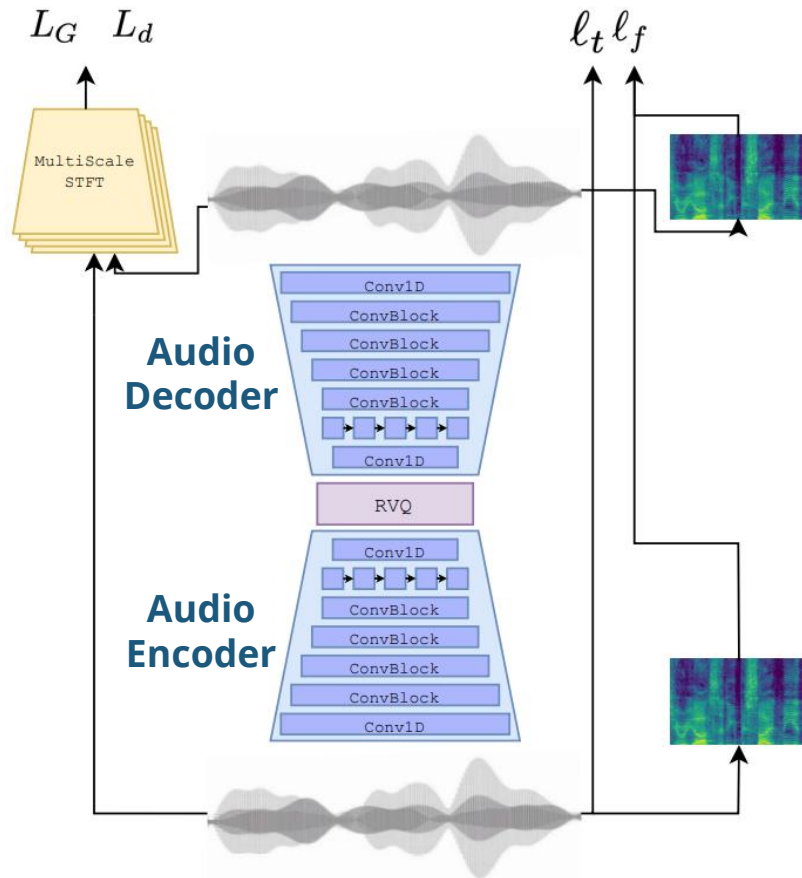
# Latent-based Audio Synthesis



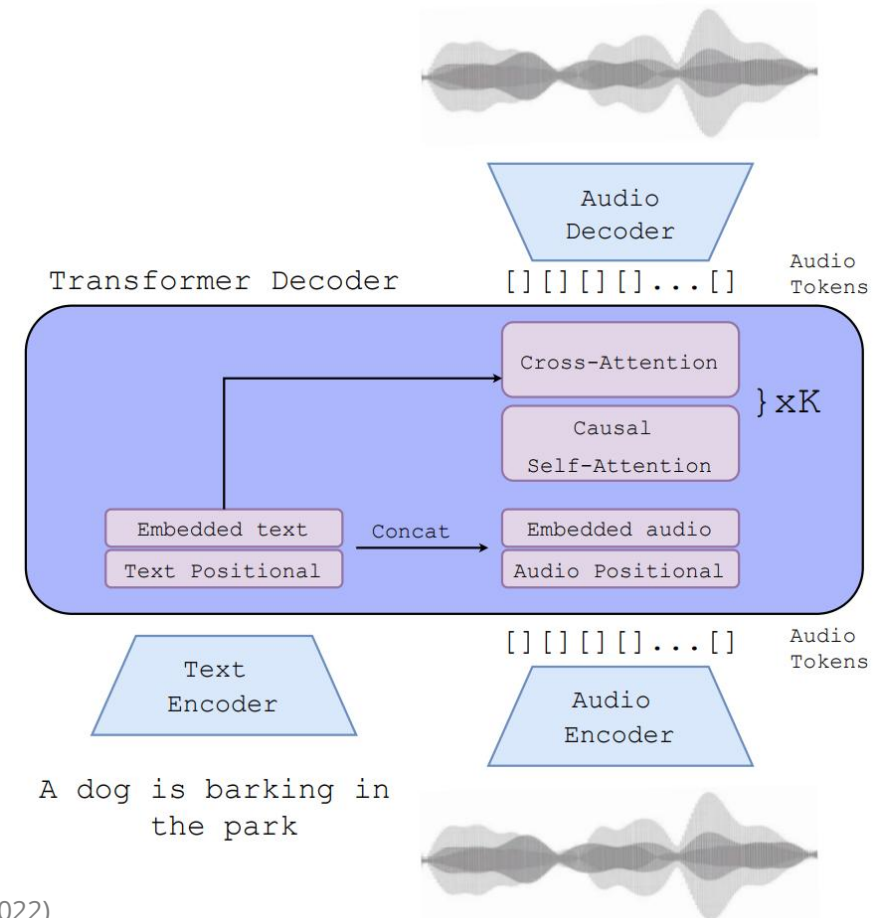
# Latent-based Audio Synthesis

# AudioGen (Kreuk et al., 2023)

## Audio Autoencoder

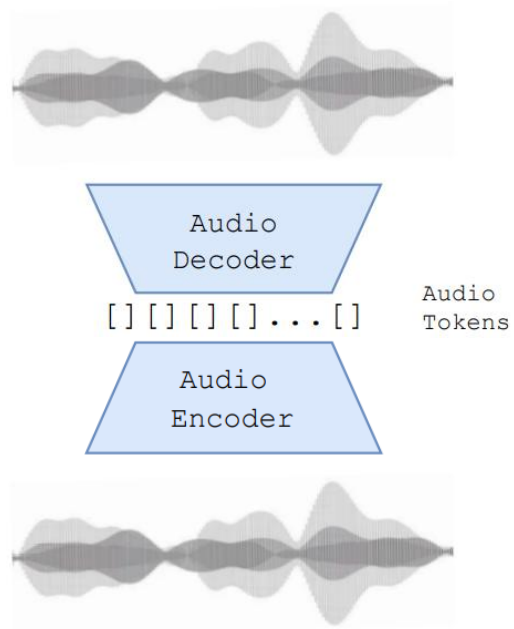


## Audio Language Model

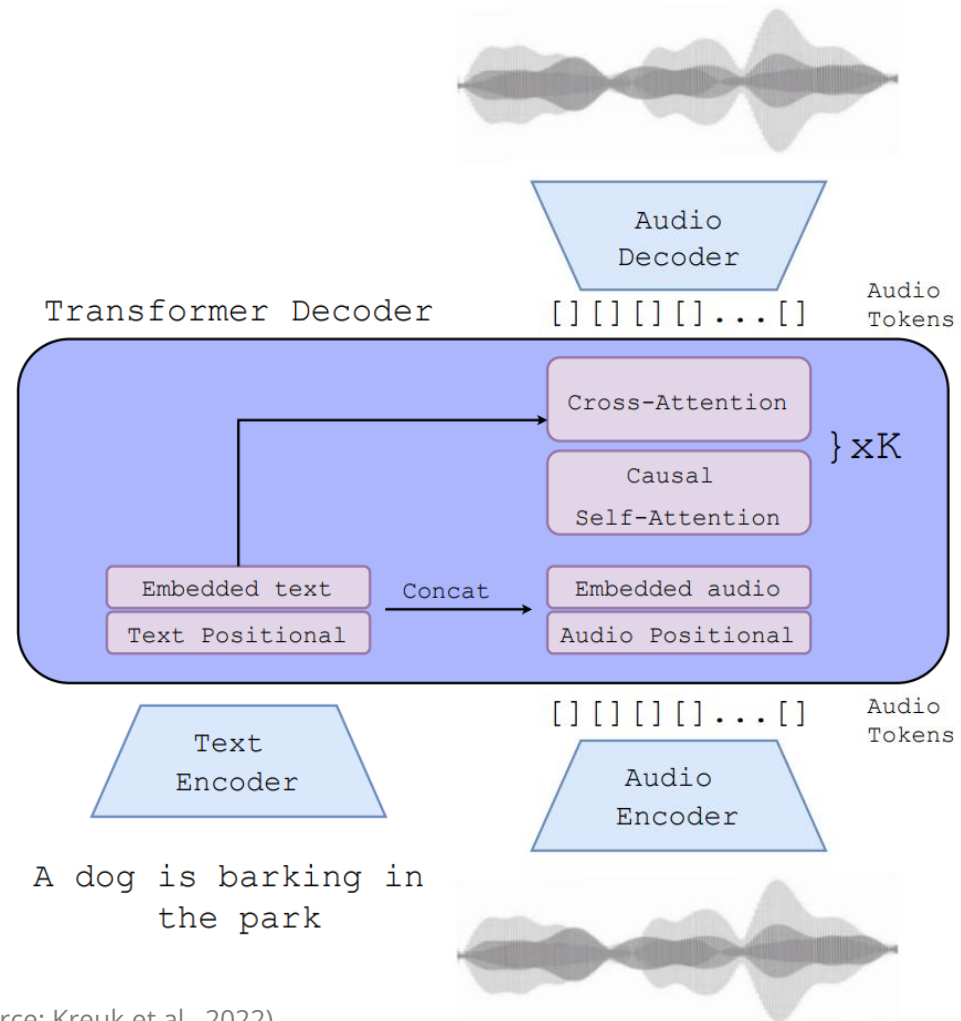


(Source: Kreuk et al., 2022)

# AudioGen (Kreuk et al., 2023)



**4k hours**  
**(speech, music, sound effects)**



(Source: Kreuk et al., 2022)

# AudioGen: Examples (Kreuk et al., 2023)



(Source: Kreuk et al., 2022)

[felixkreuk.github.io/audiogen](https://felixkreuk.github.io/audiogen)

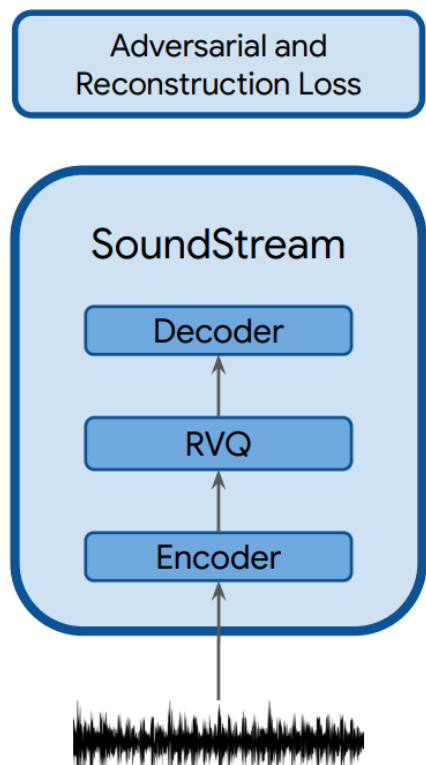
# MusicGen (Copet et al., 2023)

- AudioGen for Music
- Use EnCodec (Défossez et al., 2022) as the autoencoder
  - instead of SoundStream for AudioGen (Kreuk et al., 2023)
- **20k hours** of licensed music
  - Internal dataset      10k      High-quality (private)
  - Shutterstock        25k      Instrument-only
  - Pond5                 365k     Instrument-only

[ai.honu.io/papers/musicgen/](https://ai.honu.io/papers/musicgen/)

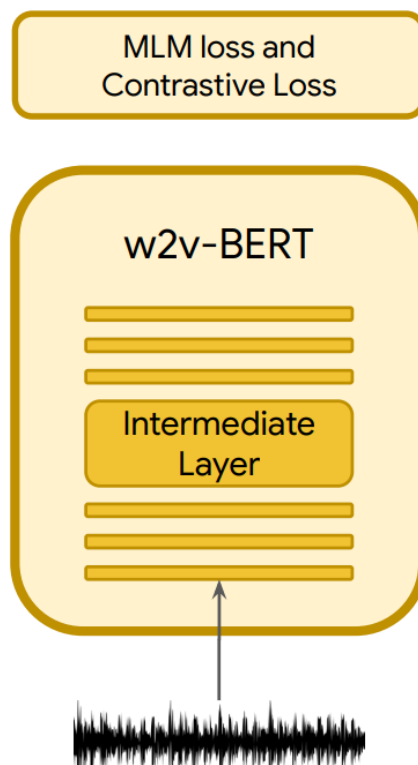
# MusicLM (Agostinelli et al., 2023)

## Audio autoencoder

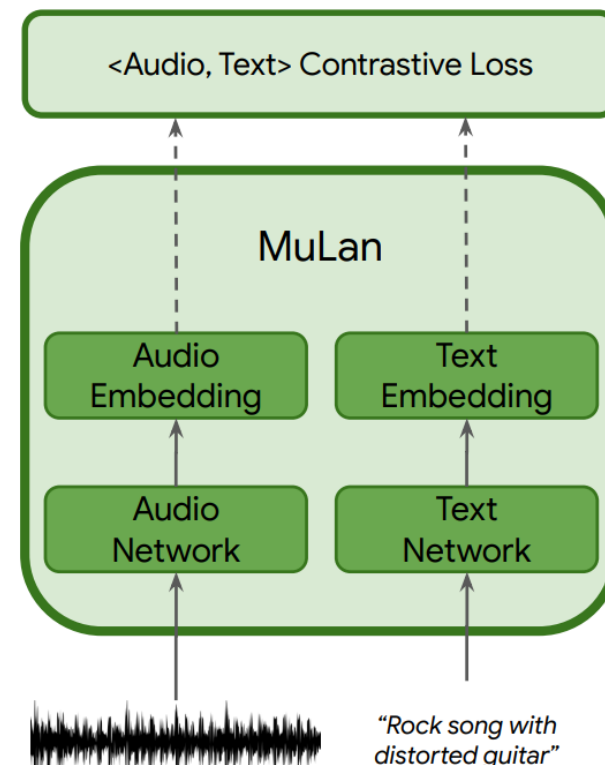


**106k songs, 8.2k hours**

## Semantic representation



## Text-music correspondence

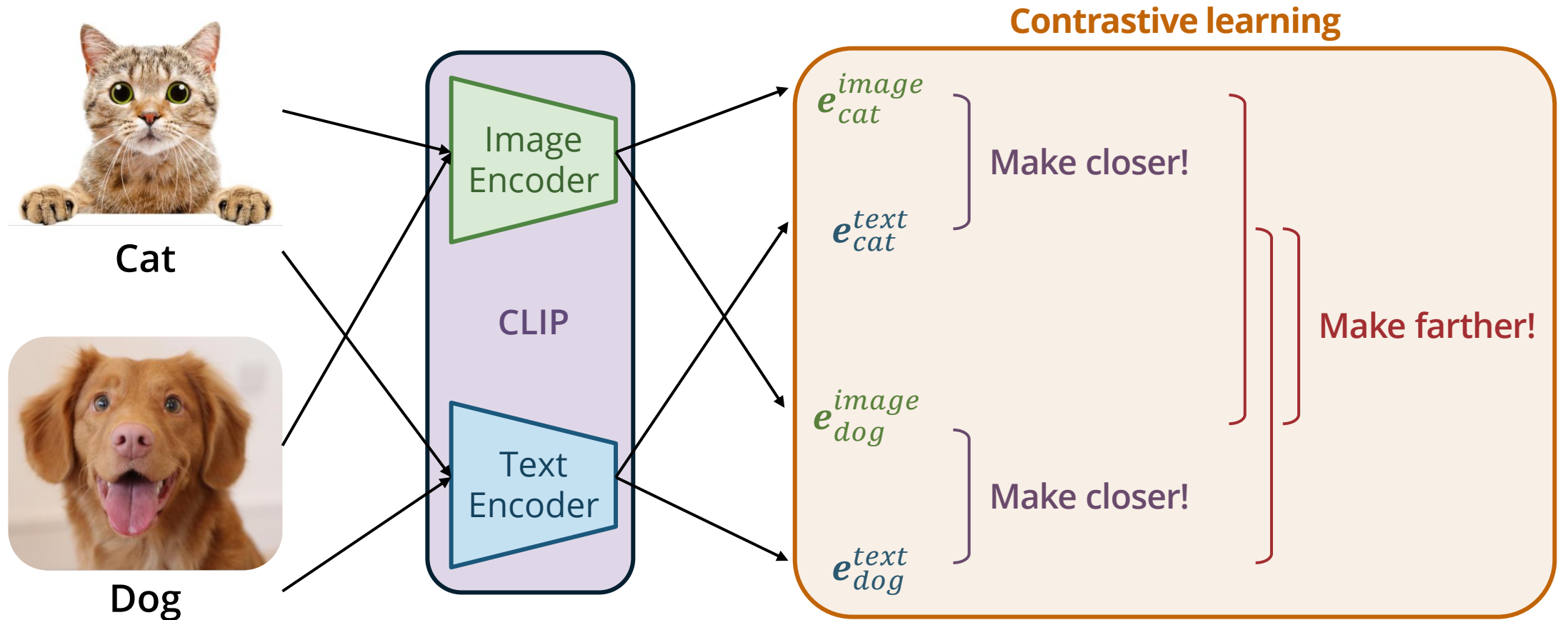


**44M 30-sec clips, 370k hours**

(Source: Agostinelli et al., 2022)

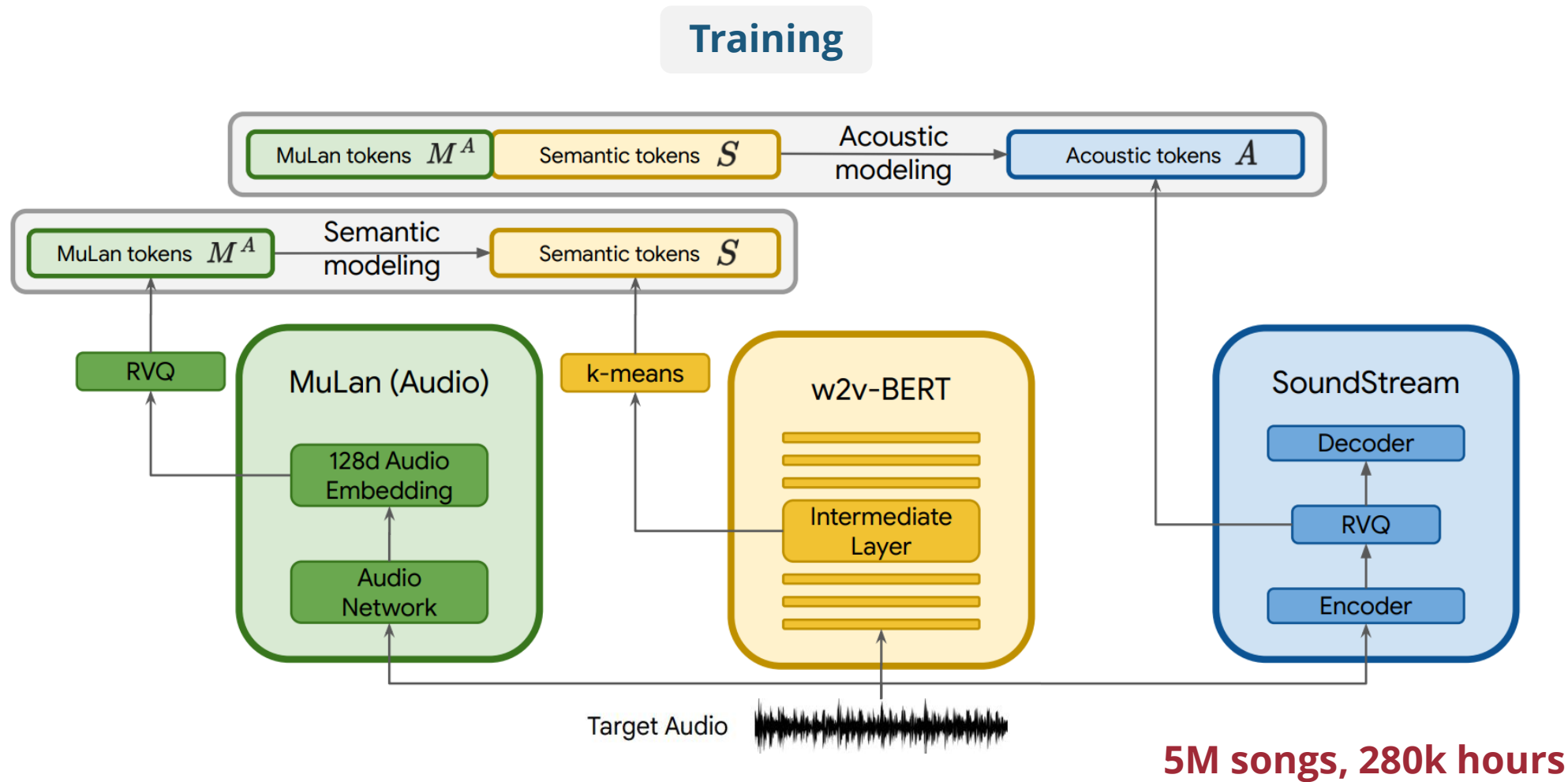


# Contrastive Language-Image Pretraining (CLIP)



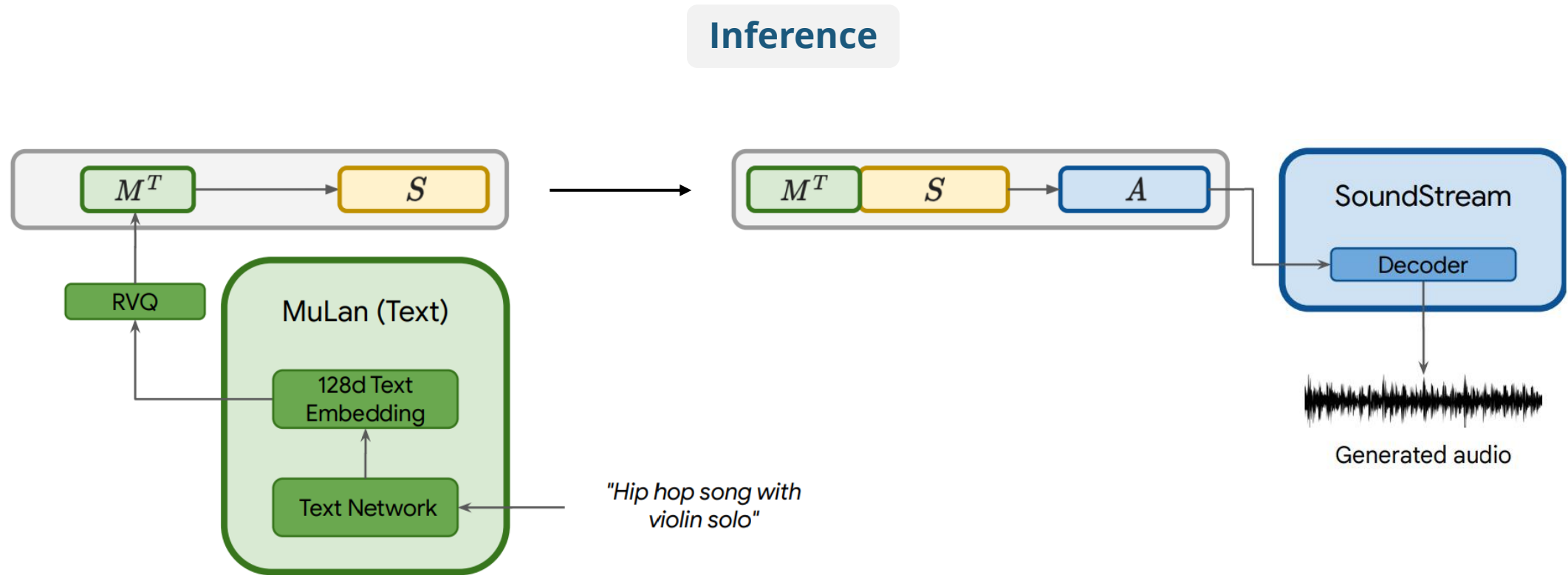
Learn a **shared embedding space** for images and texts

# MusicLM (Agostinelli et al., 2023)



(Source: Agostinelli et al., 2022)

# MusicLM (Agostinelli et al., 2023)

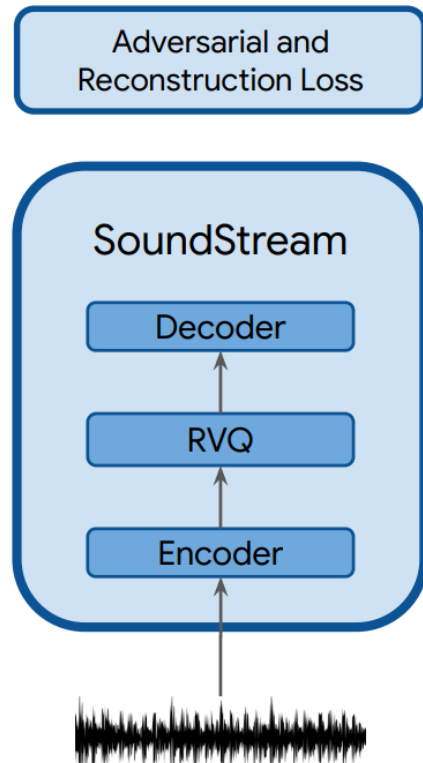


(Source: Agostinelli et al., 2022)

[google-research.github.io/seanet/musiclm/examples/](https://google-research.github.io/seanet/musiclm/examples/)

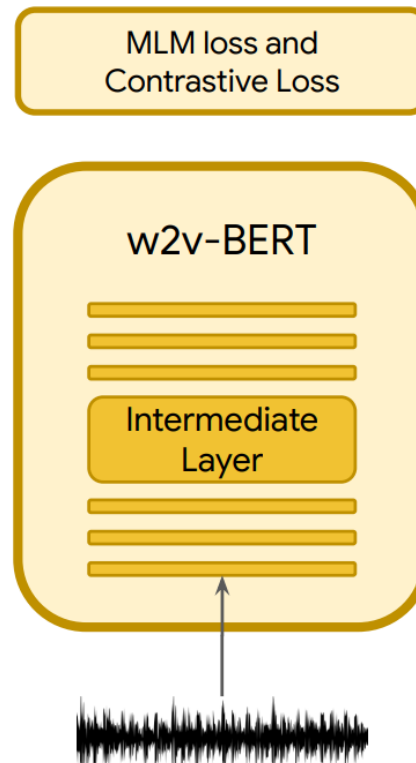
# MusicLM (Agostinelli et al., 2023)

## Audio autoencoder

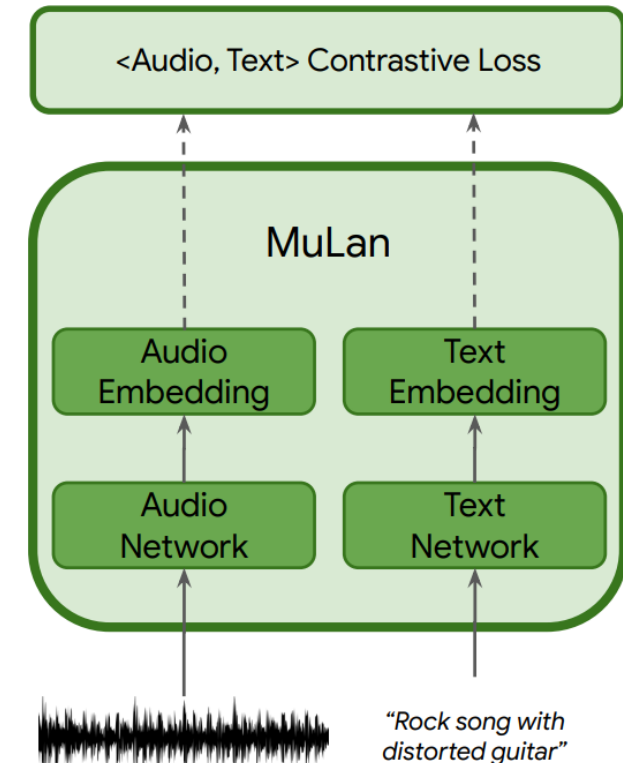


**106k songs, 8.2k hours**

## Semantic representation



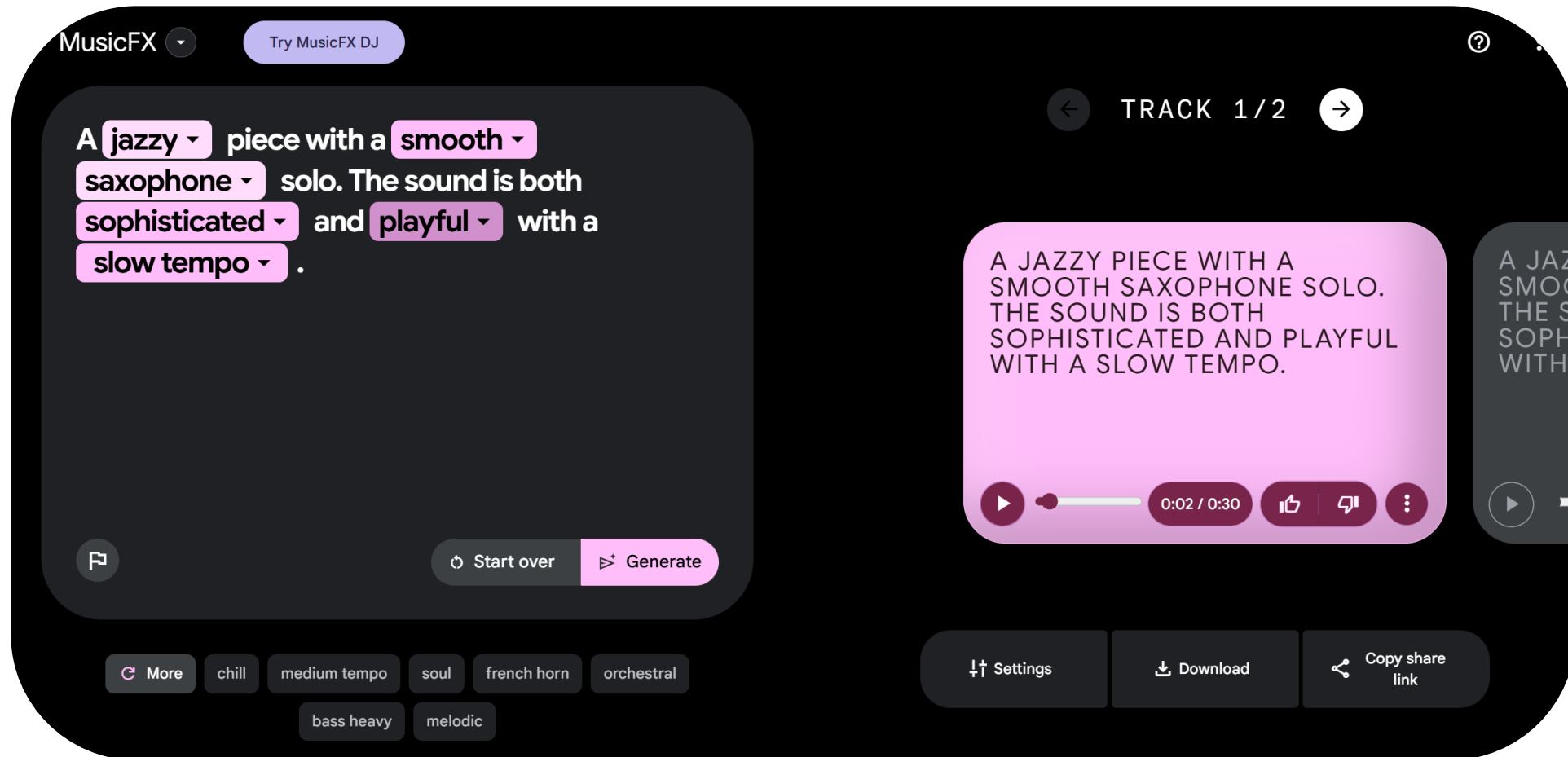
## Text-music correspondence



**44M 30-sec clips, 370k hours**

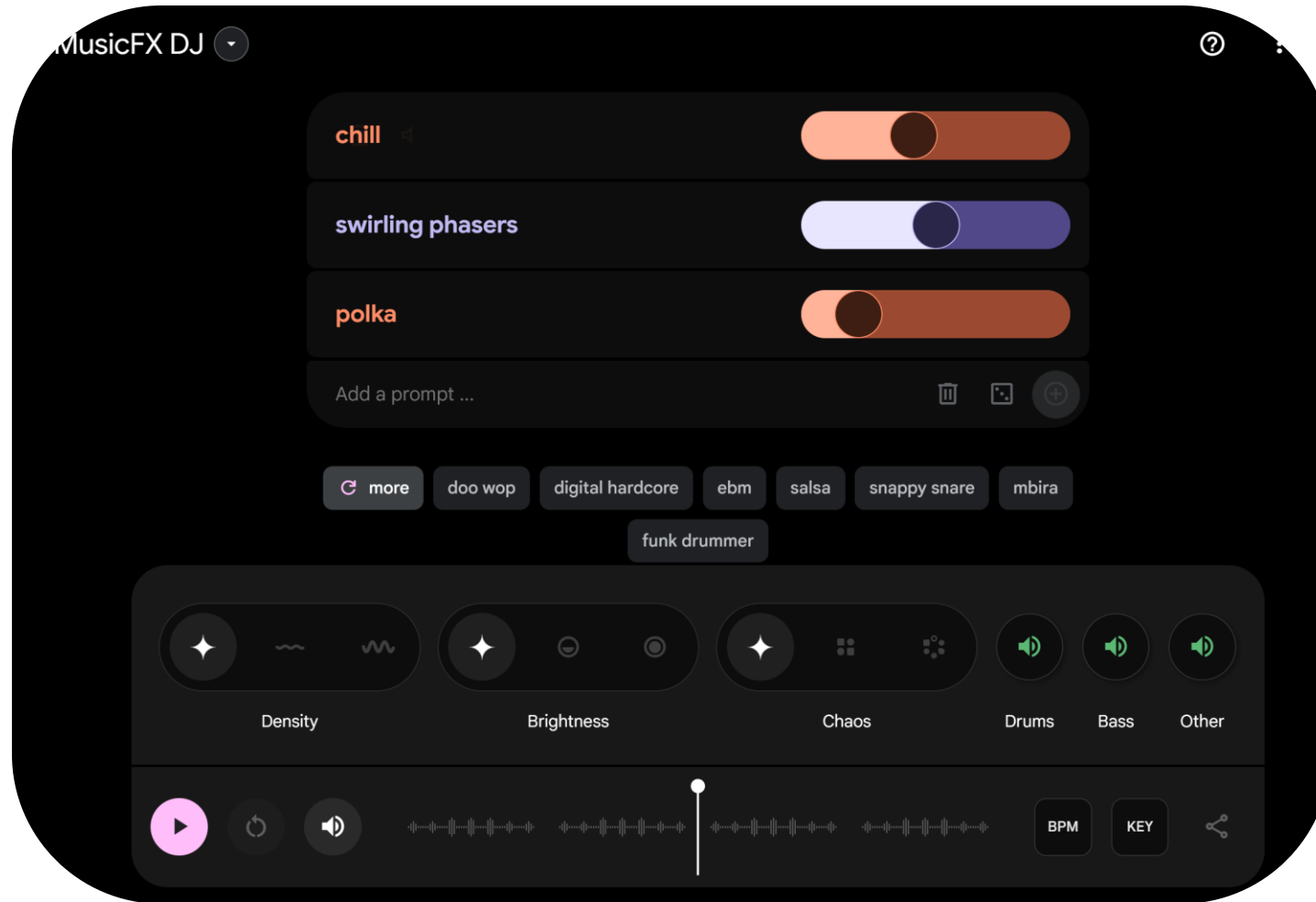
(Source: Agostinelli et al., 2022)

# Music FX (2024)



[aitestkitchen.withgoogle.com/tools/music-fx](https://aitestkitchen.withgoogle.com/tools/music-fx)

# Music FX DJ (2024)



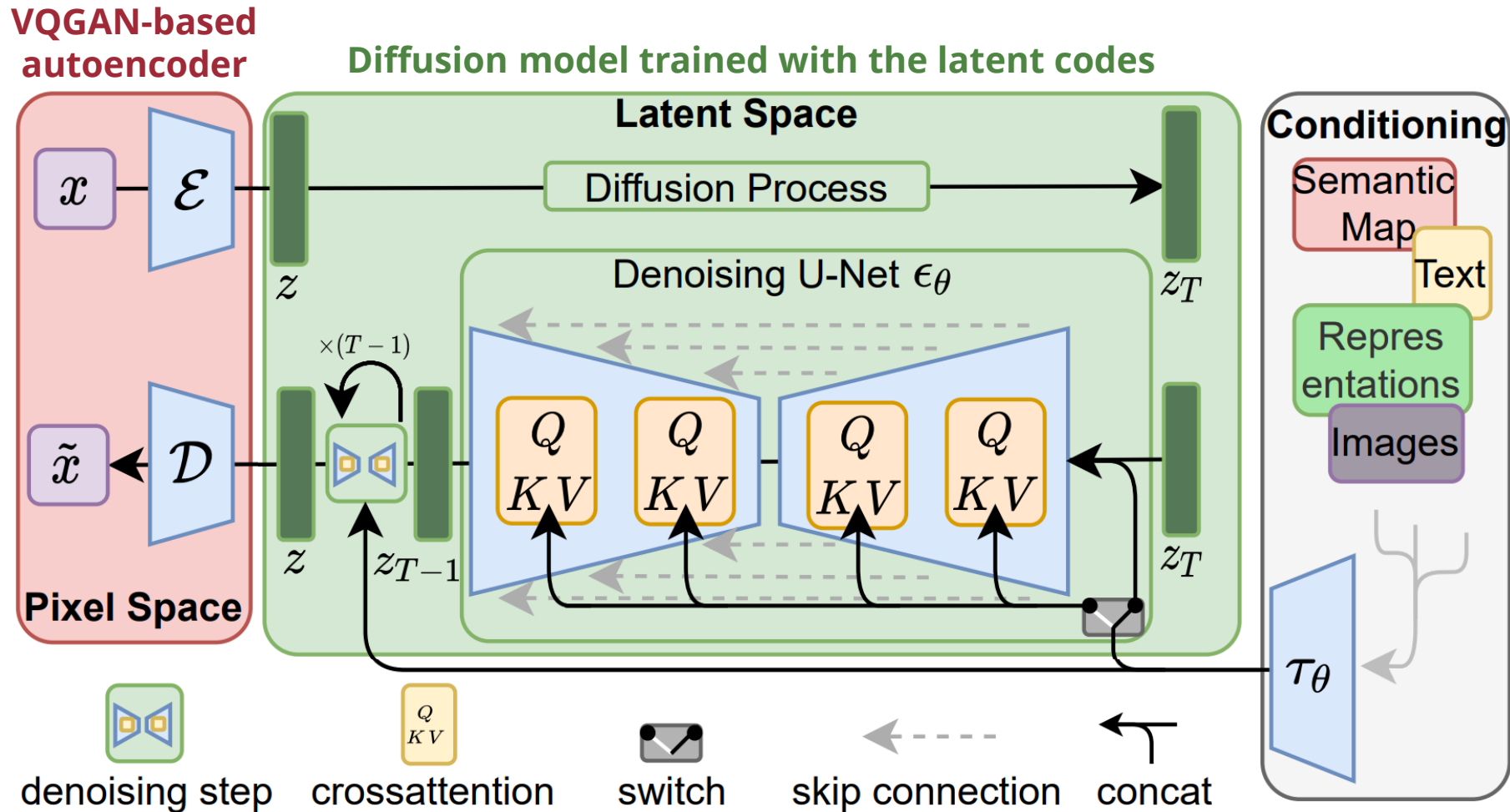
[aitestkitchen.withgoogle.com/tools/music-fx-dj](https://aitestkitchen.withgoogle.com/tools/music-fx-dj)

# Music FX DJ (2024)



[youtube.com/live/IUQW5LgBZvQ](https://youtube.com/live/IUQW5LgBZvQ)

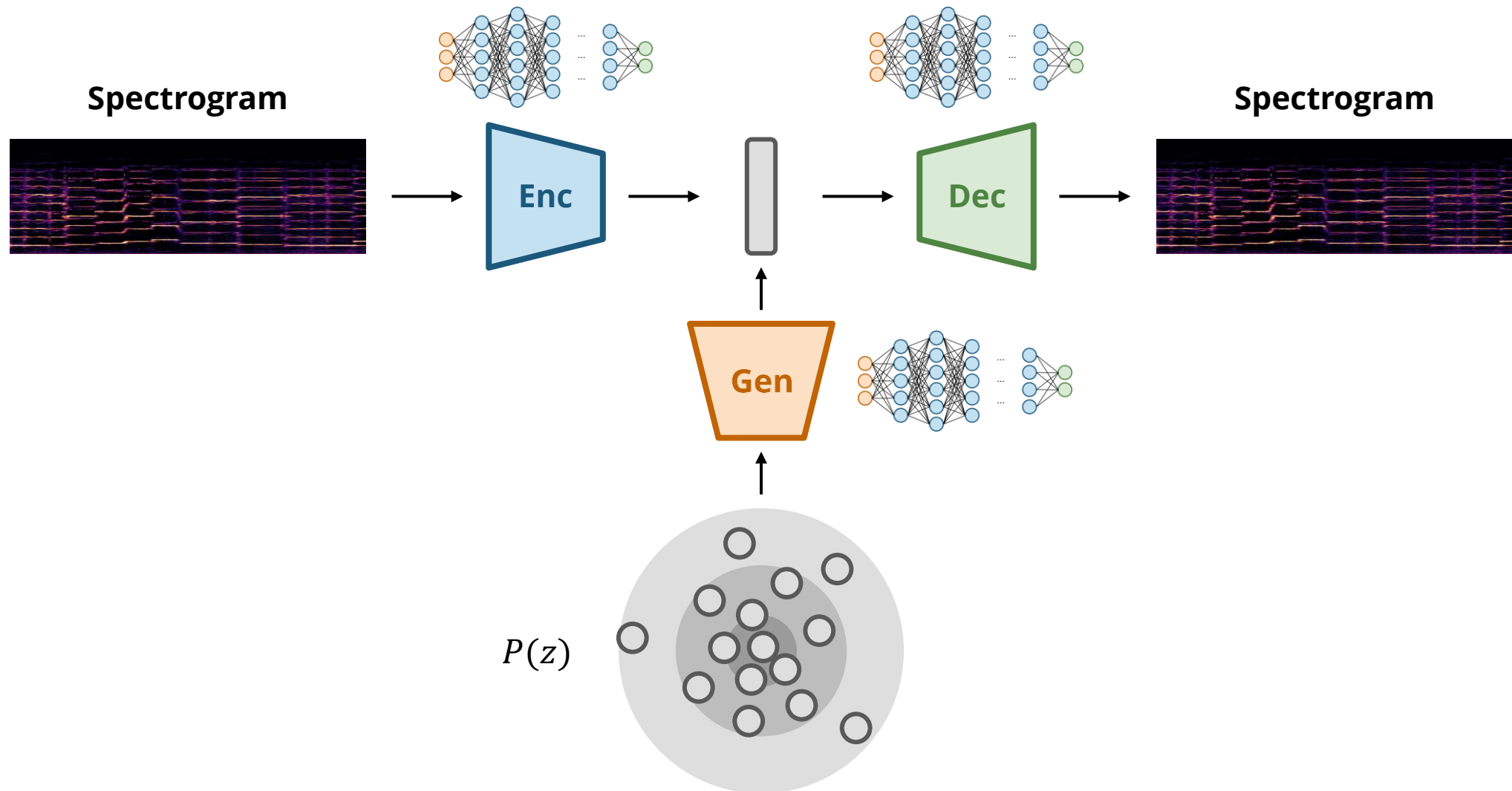
# (Recap) Latent Diffusion Models (LDMs)



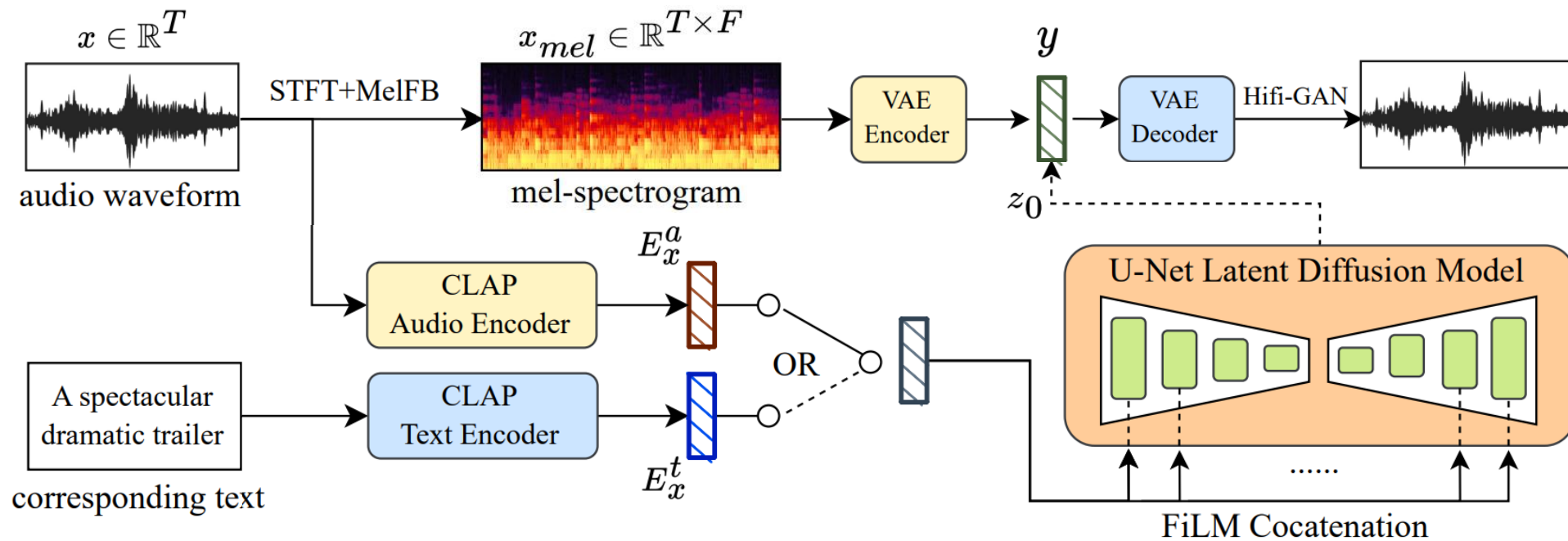
(Source: Rombach et al., 2022)



# (Recap) Latent-based Audio Synthesis



# Example: MusicLDM (Chen et al., 2023)



(Source: Ke et al., 2023)

[musicldm.github.io](https://musicldm.github.io)

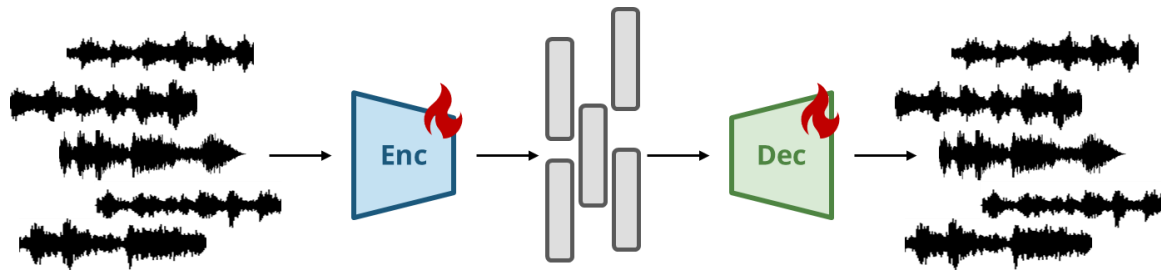
## Example: MusicLDM (Chen et al., 2023)



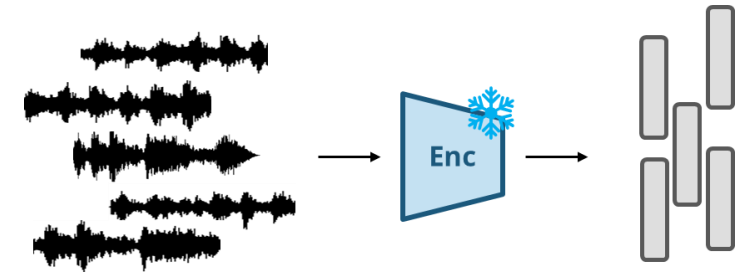
[youtu.be/DALv7ea6cv0](https://youtu.be/DALv7ea6cv0)

# Pipeline

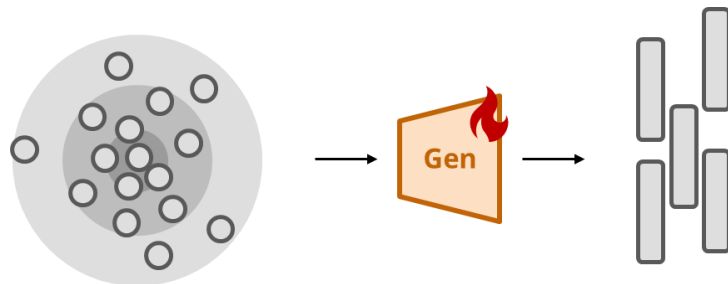
Step 1: Train an Autoencoder



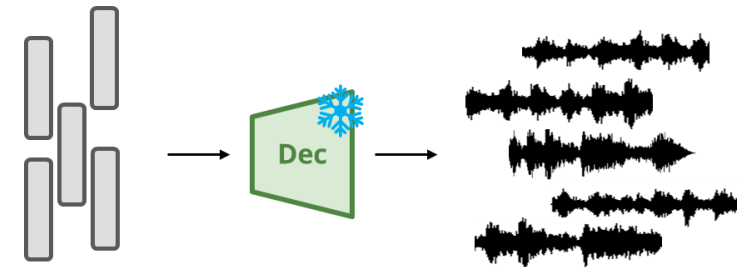
Step 2: Compute the Latent Vectors



Step 3: Train a Latent Generative Model



Step 4: Decode the Latent Vectors



# Creative Applications of Music Generation Systems

# unloop (Garcia et al., 2023)

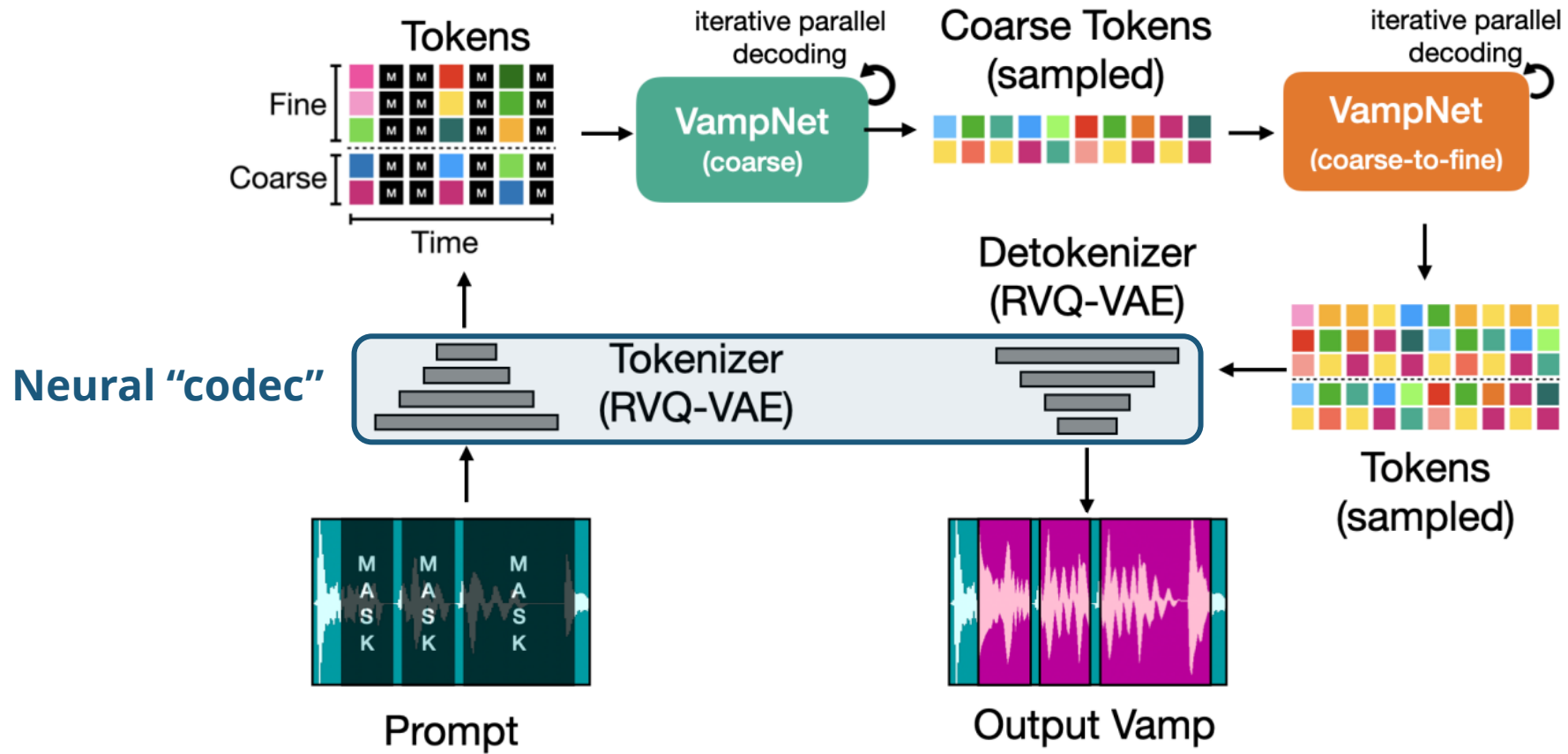


[youtu.be/yzBI8Vcjd2s](https://youtu.be/yzBI8Vcjd2s)

[github.com/hugofloresgarcia/unloop](https://github.com/hugofloresgarcia/unloop)

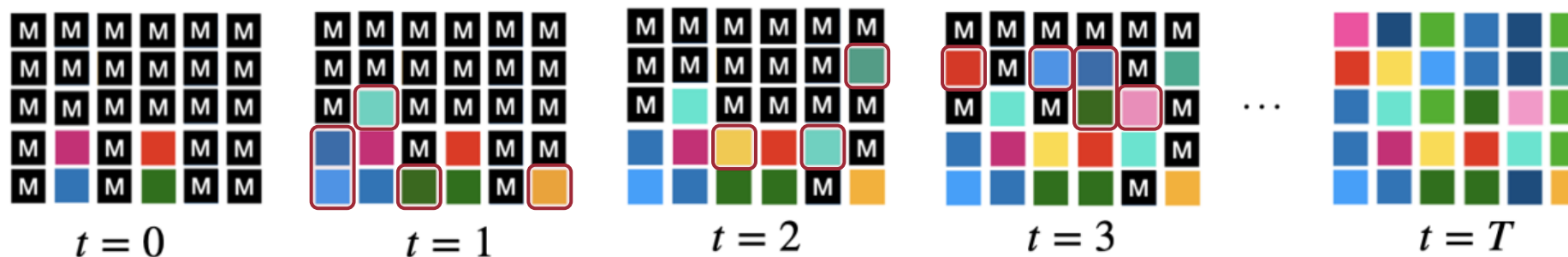


# VampNet (Garcia et al., 2023)



(Source: Garcia et al., 2023)

# VampNet (Garcia et al., 2023)

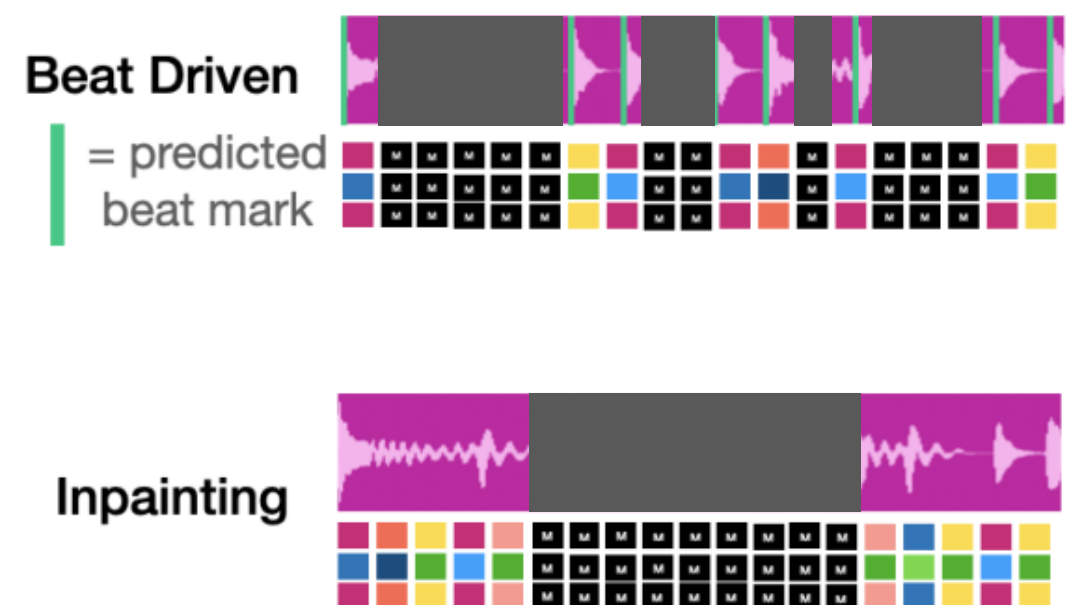
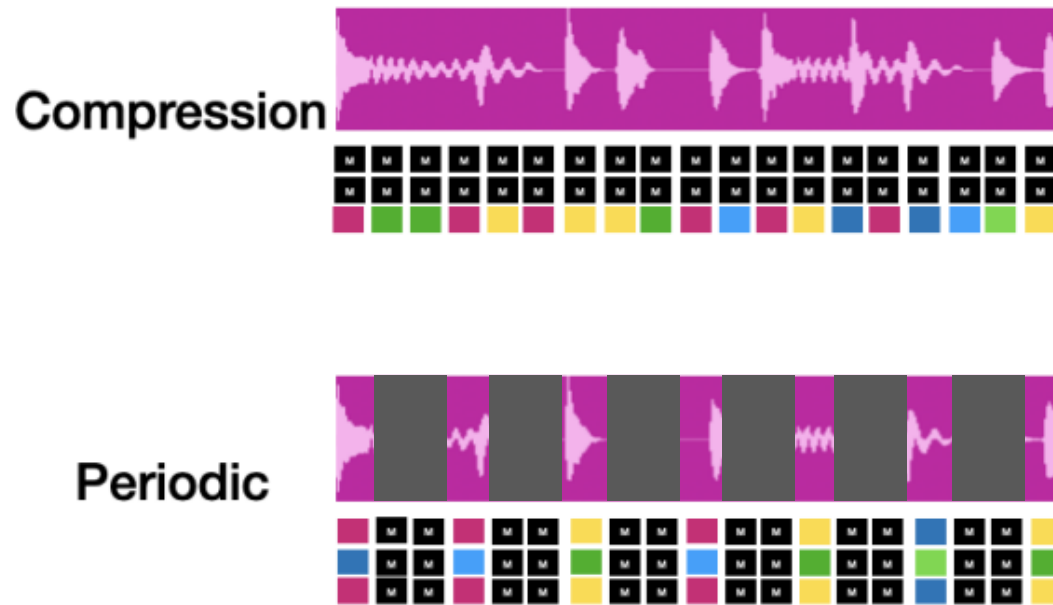


(Source: Garcia et al., 2023)

Sample a subset of the **most confident predicted tokens** in each iteration

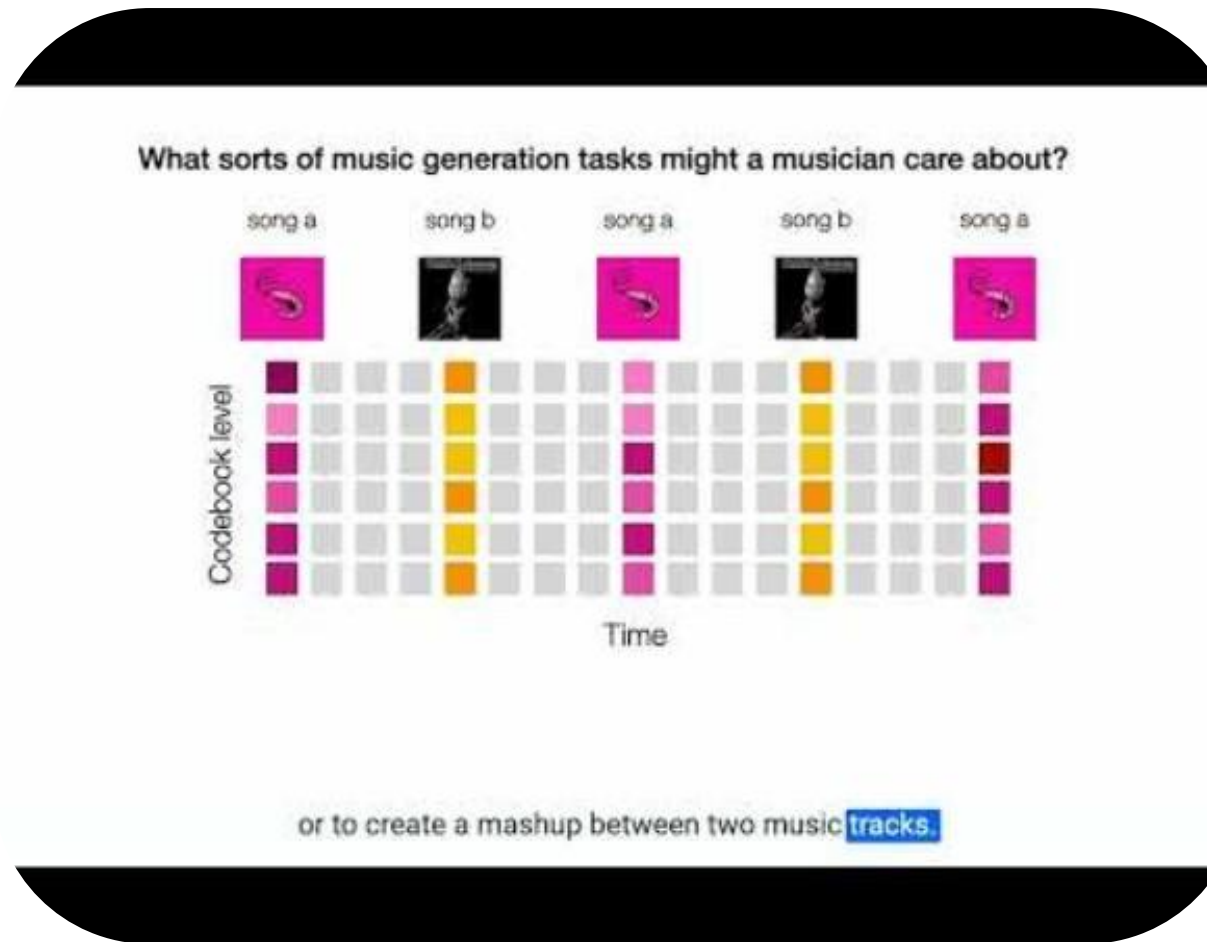


# VampNet (Garcia et al., 2023)



(Source: Garcia et al., 2023)

# VampNet (Garcia et al., 2023)



[youtu.be/3XfeWIV9Cp0](https://youtu.be/3XfeWIV9Cp0)

# Controlling Music Generation Systems

# ControlNet (Zhang et al., 2023)



Input Canny edge



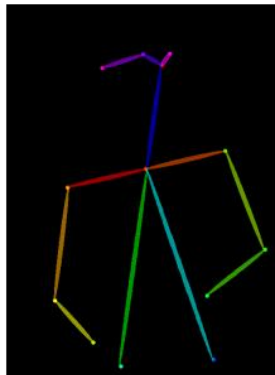
Default



“masterpiece of fairy tale, giant deer, golden antlers”



“..., quaint city Galic”



Input human pose



Default



“chef in kitchen”

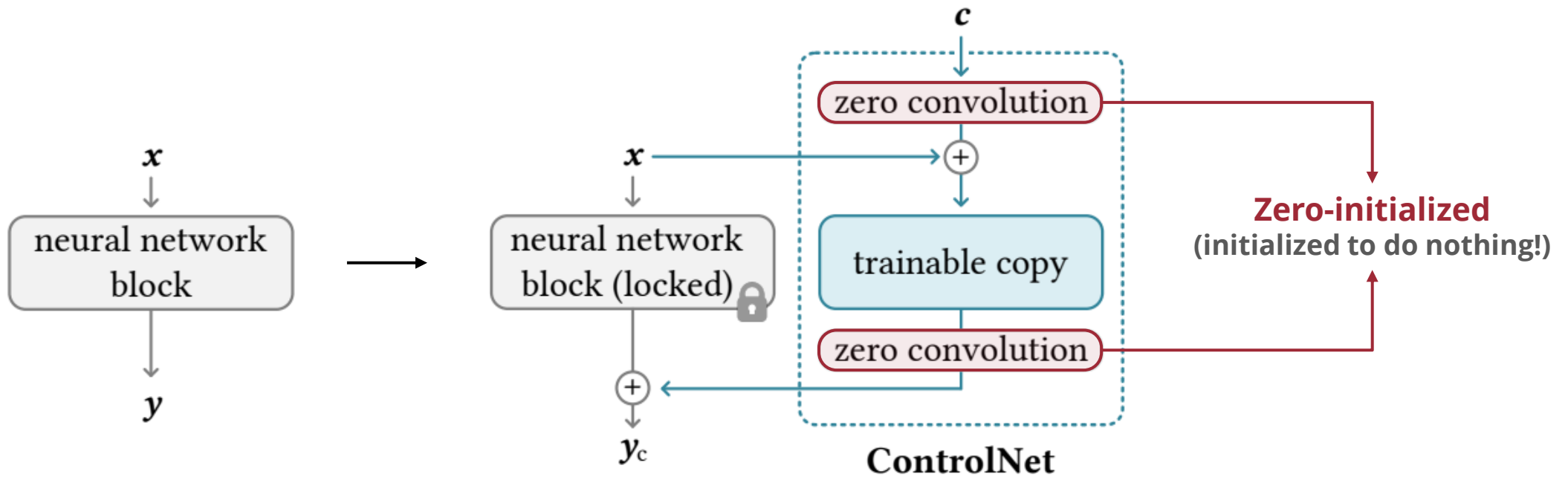


“Lincoln statue”

(Source: Zhang et al., 2023)

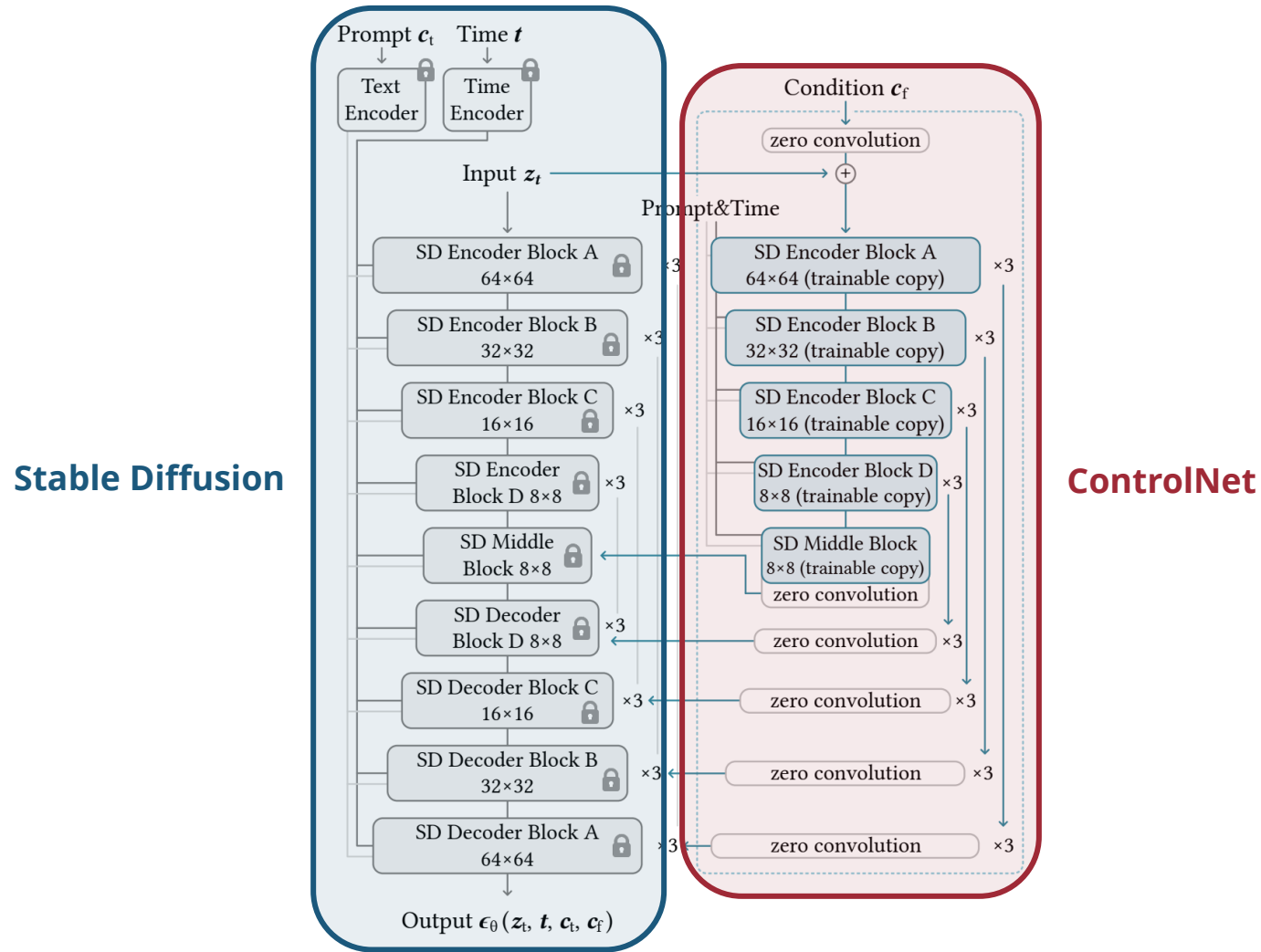
Can we **add controls** to a trained text-to-image diffusion model?

# ControlNet (Zhang et al., 2023)



(Source: Zhang et al., 2023)

# ControlNet (Zhang et al., 2023)



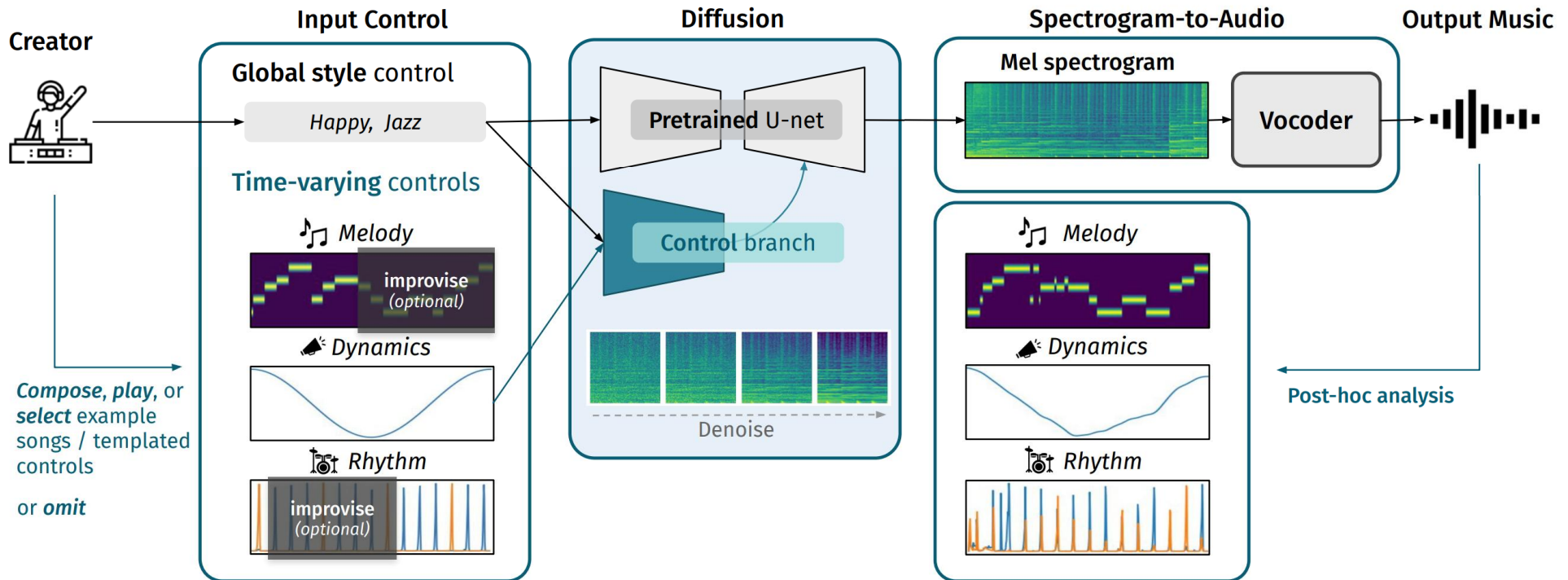
(Source: Zhang et al., 2023)

# Synthetic Beat Brigade - How would you touch me? (2023)



[youtu.be/O4cJ3acEGDw](https://youtu.be/O4cJ3acEGDw) &  
[drive.google.com/file/d/1QTQ7P3iZI6l0anlwNQ3ewf8g3JjDjesl/view](https://drive.google.com/file/d/1QTQ7P3iZI6l0anlwNQ3ewf8g3JjDjesl/view)

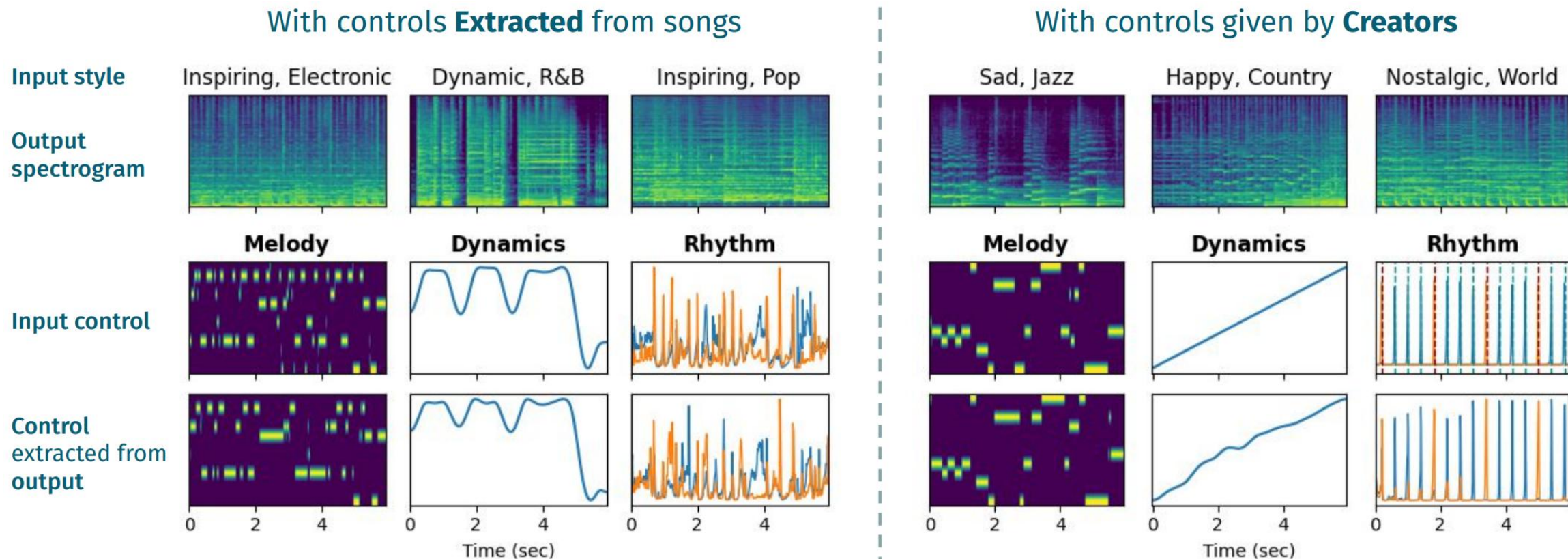
# Music ControlNet (Wu et al., 2024)



(Source: Wu et al., 2024)



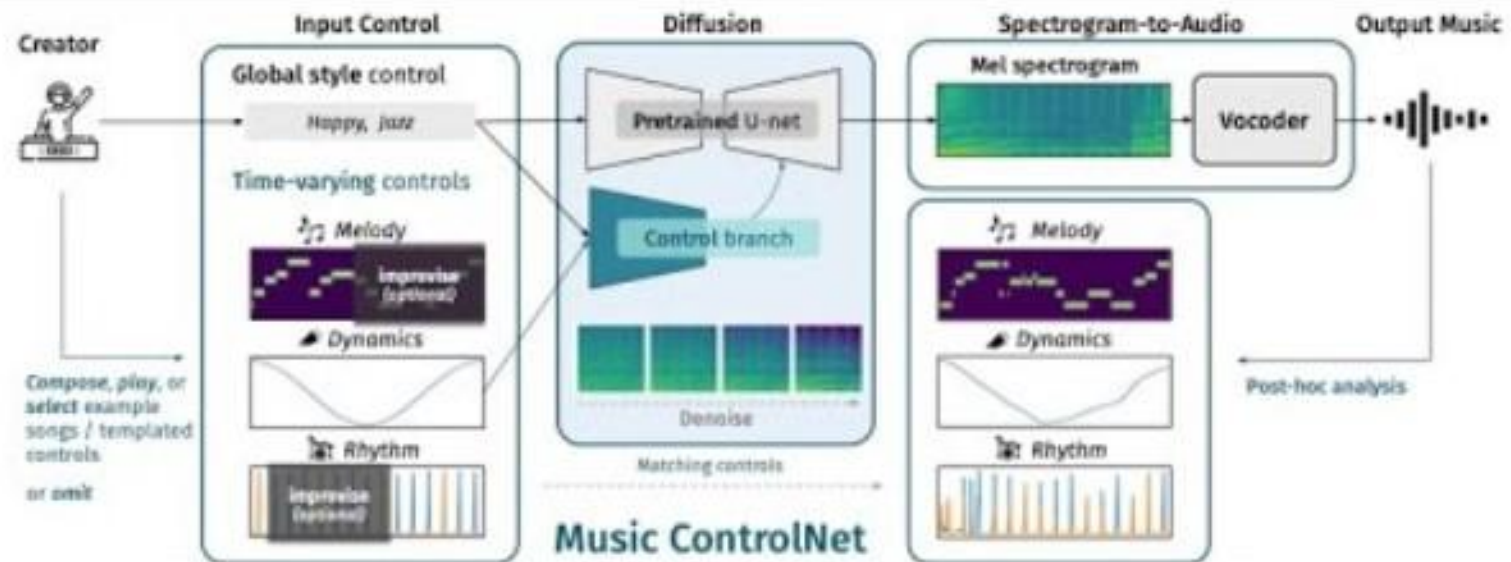
# Music ControlNet (Wu et al., 2024)



(Source: Wu et al., 2024)

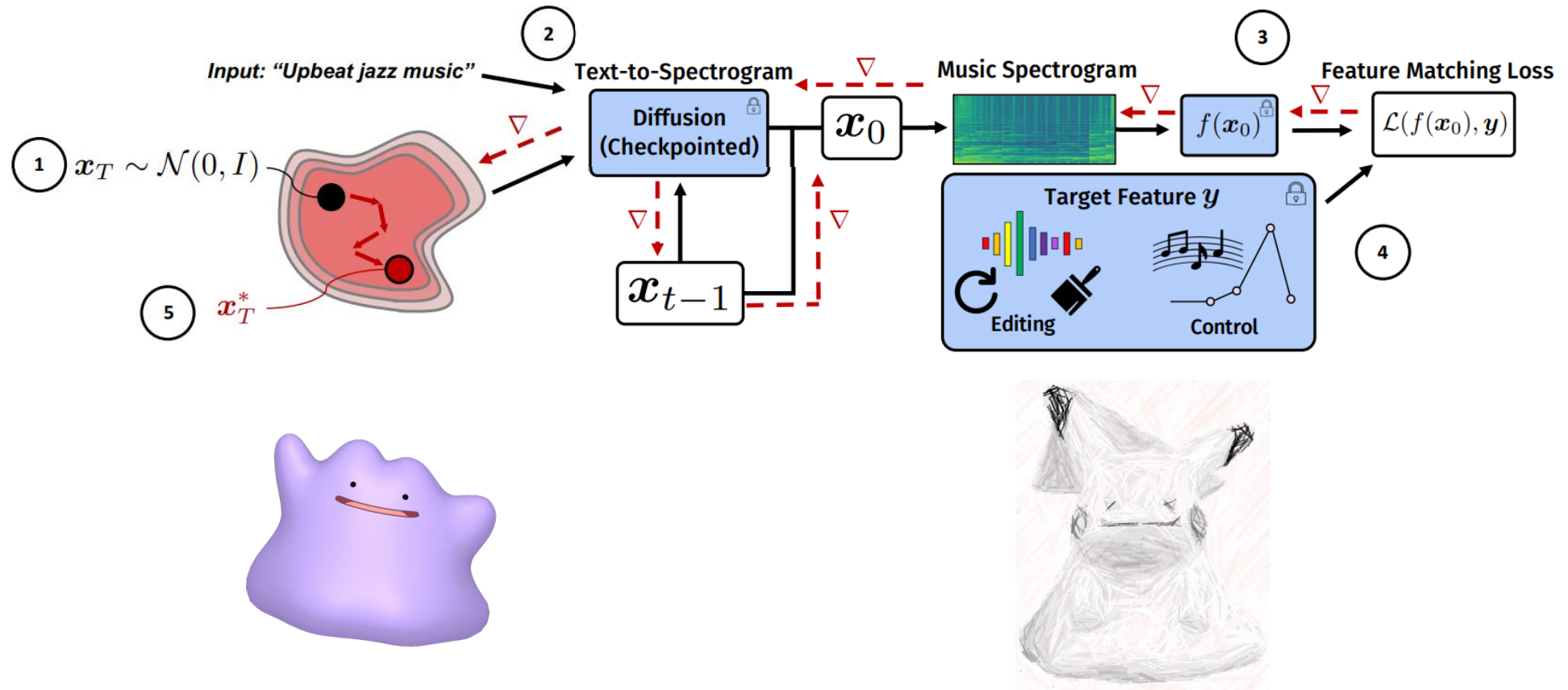
[musiccontrolnet.github.io/web](https://musiccontrolnet.github.io/web)

# Music ControlNet (Wu et al., 2024)



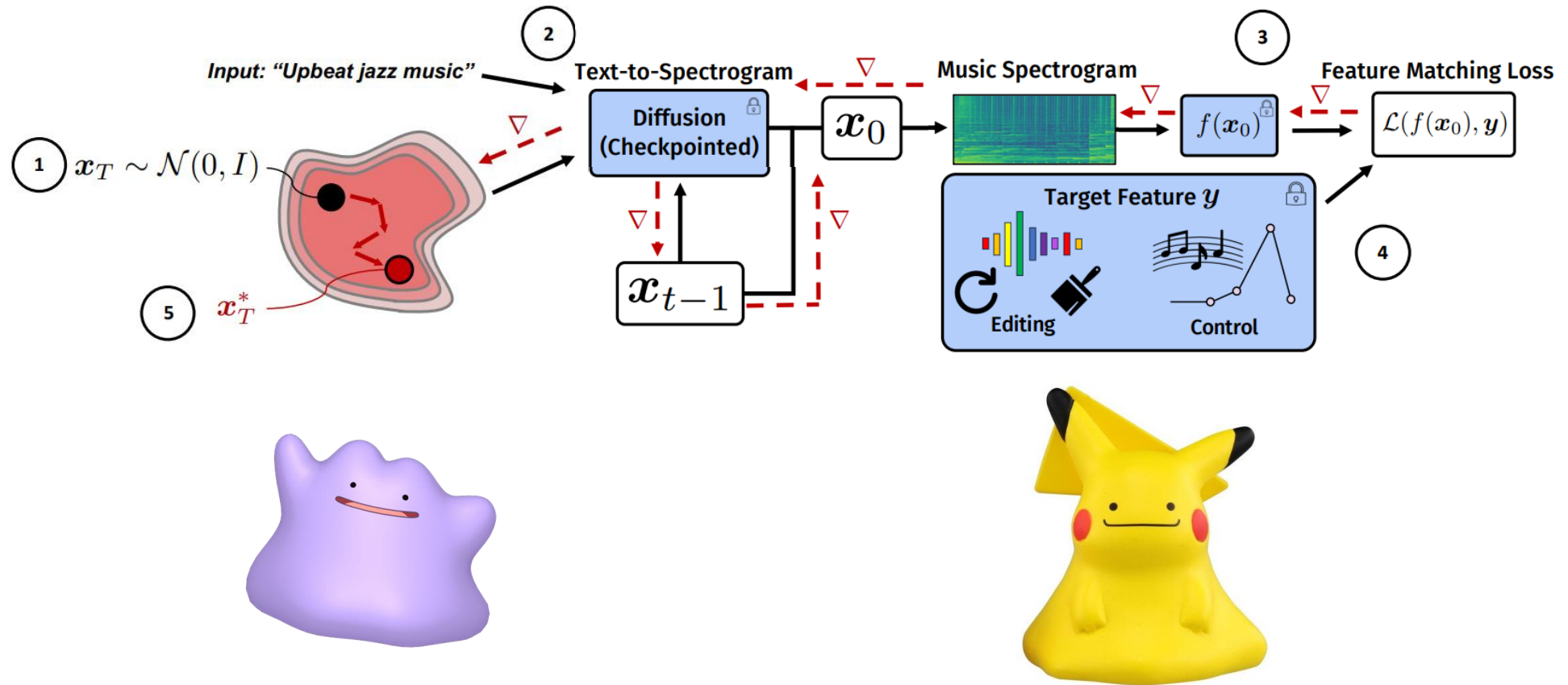
[youtu.be/QVr-S-DyccU](https://youtu.be/QVr-S-DyccU)

# DITTO (Novack et al., 2024)



(Source: Novack et al., 2024)

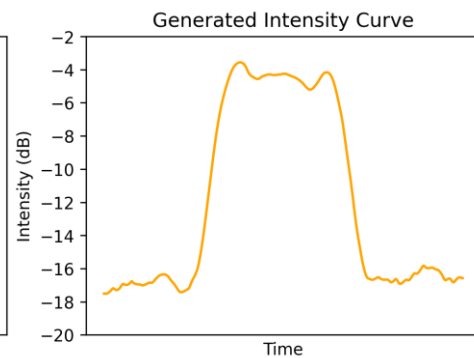
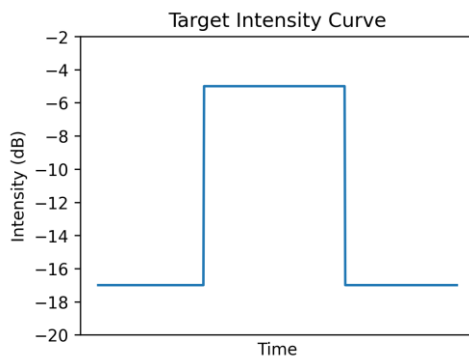
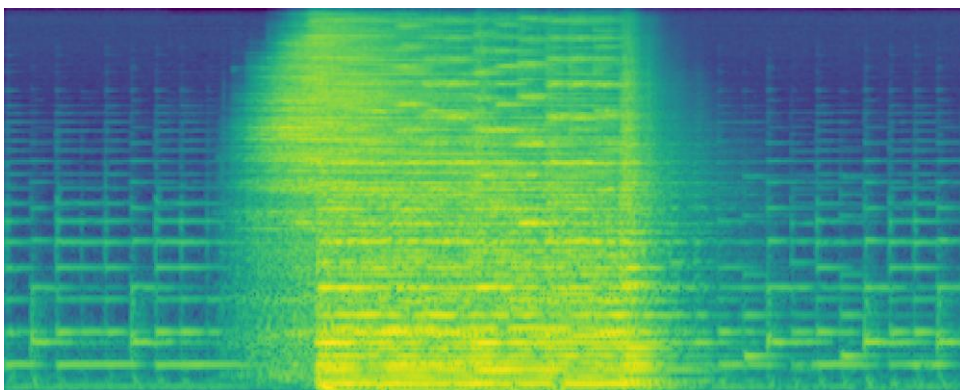
# DITTO (Novack et al., 2024)



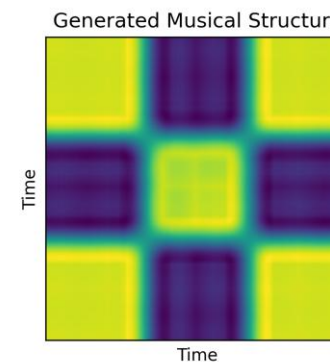
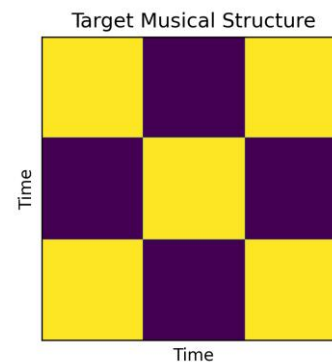
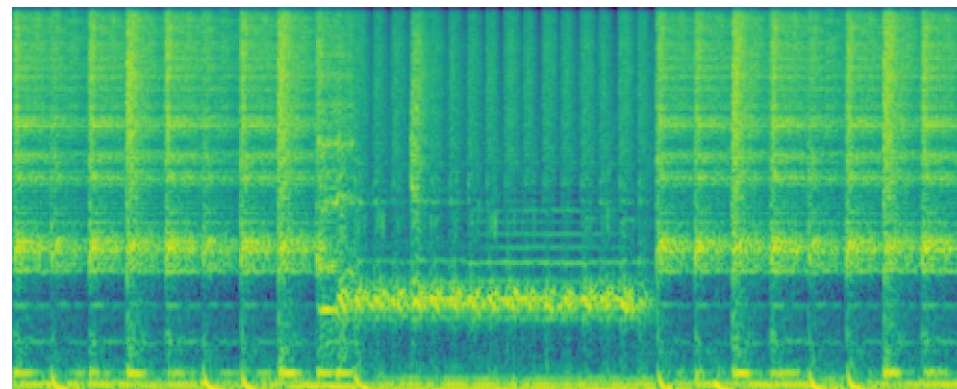
(Source: Novack et al., 2024)

# DITTO (Novack et al., 2024)

## Intensity control



## Structure control



(Source: Novack et al., 2024)

# DITTO (Novack et al., 2024)

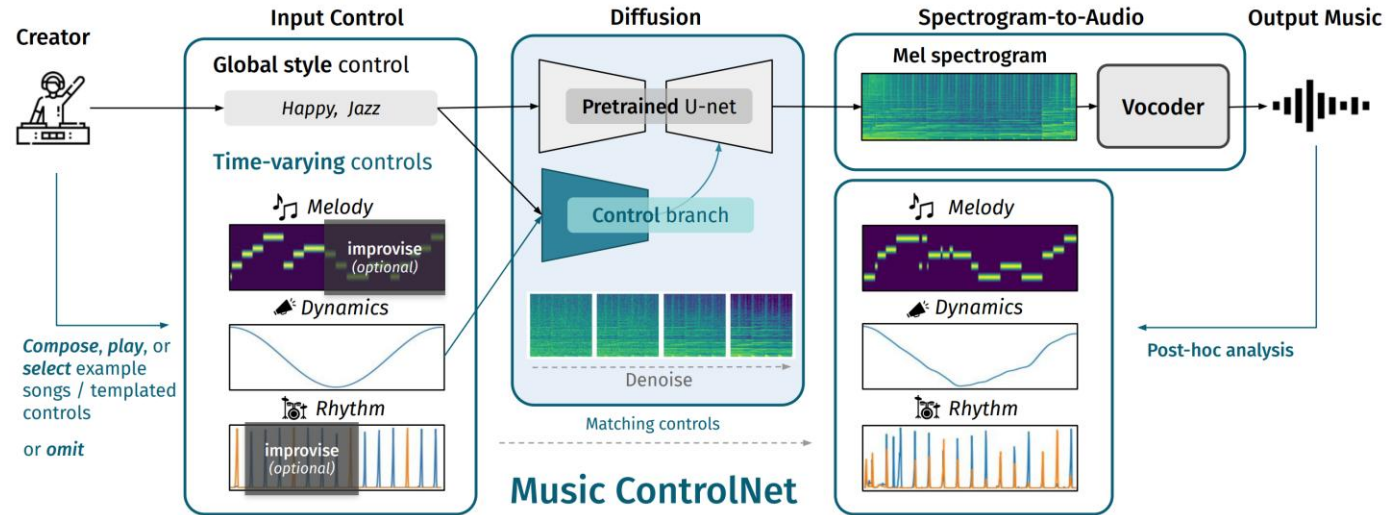


[youtu.be/KooosSNPNo8](https://youtu.be/KooosSNPNo8) & [ditto-music.github.io/web/](https://ditto-music.github.io/web/)

# Music ControlNet vs DITTO

## Music ControlNet

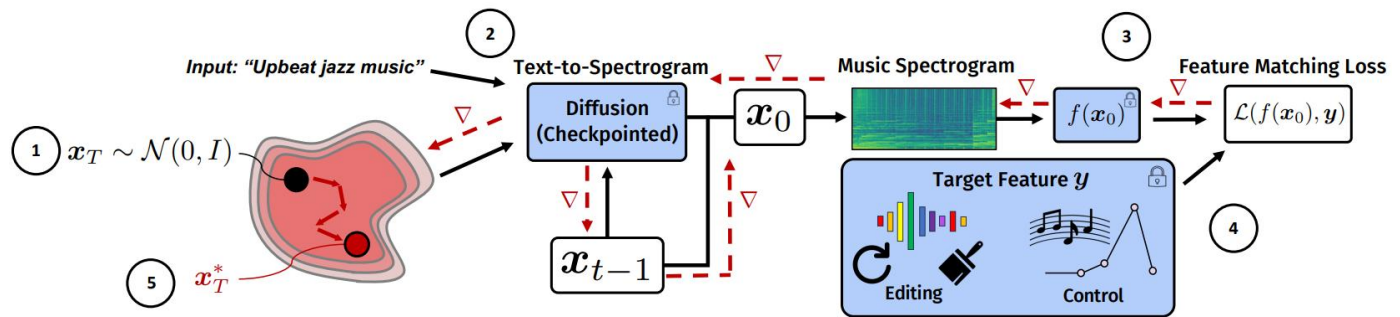
Needs some training!



(Source: Wu et al., 2024)

## DITTO

No training needed!



(Source: Novack et al., 2024)