

PAT 464/564 (Winter 2026)

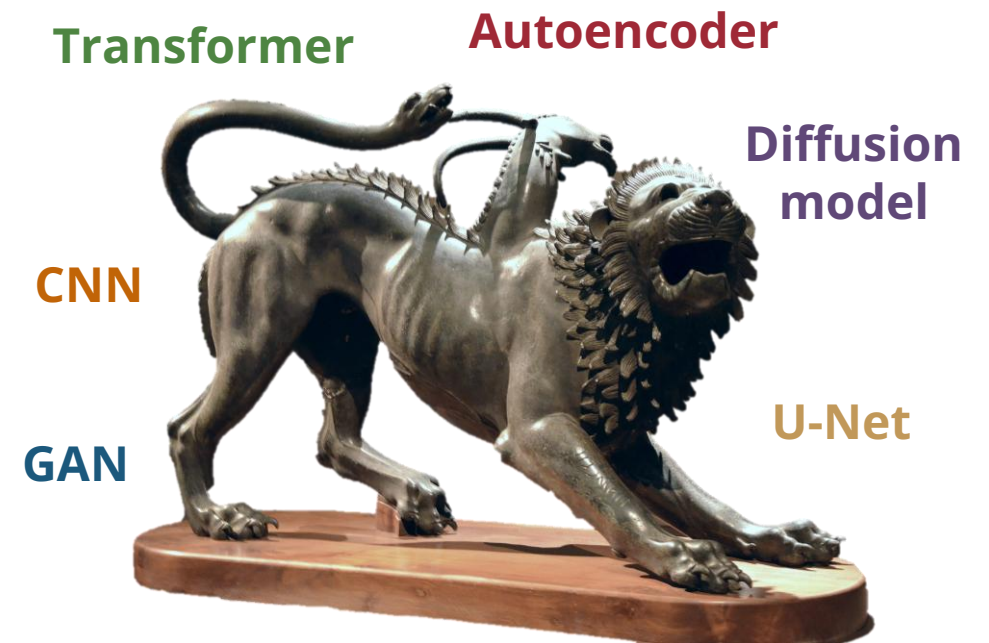
# Generative AI for Music & Audio Creation

## **Lecture 15: Latent Diffusion Models**

Instructor: Hao-Wen Dong

# Latent Diffusion Model is a Chimera

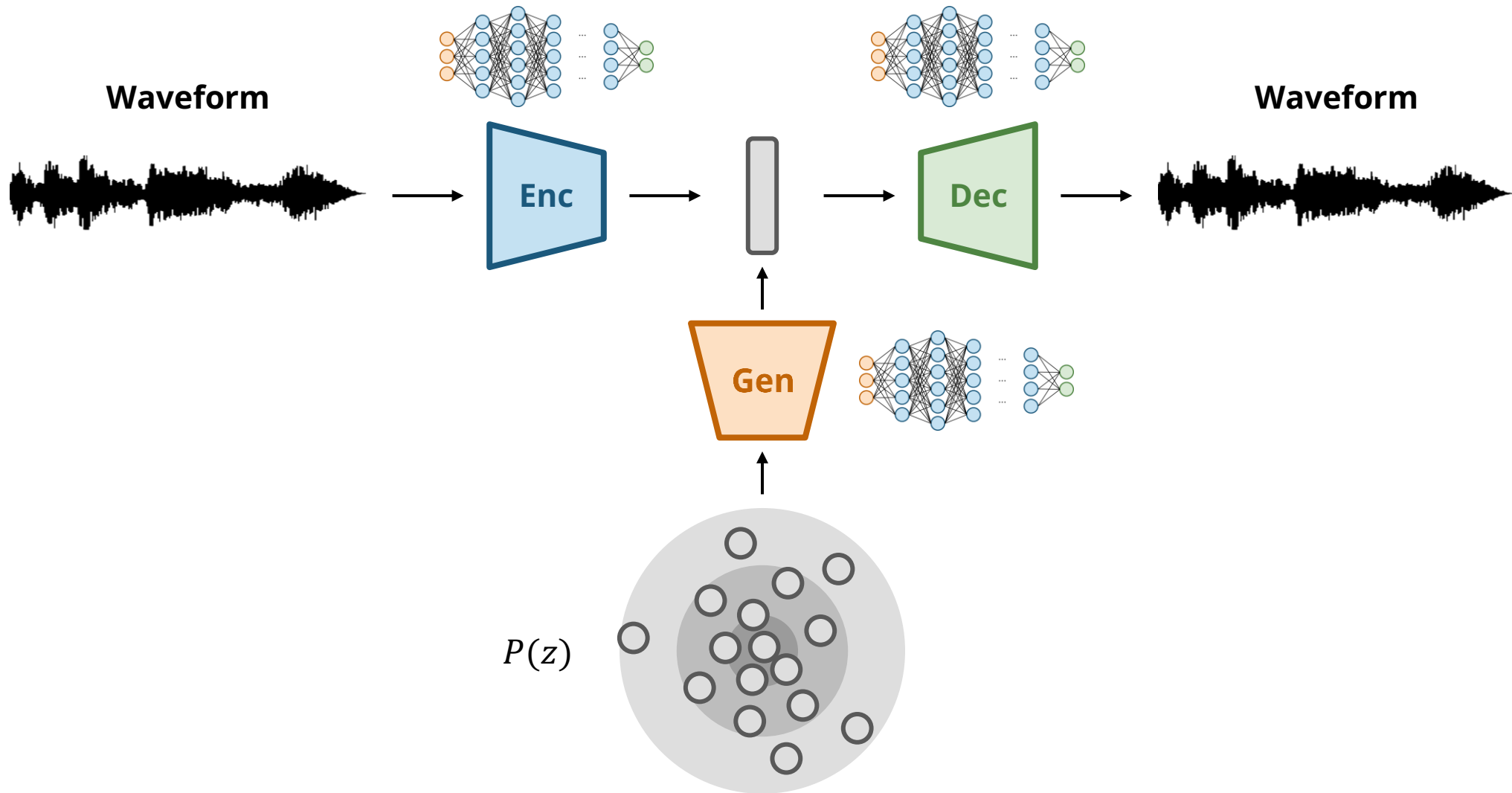
- **A neural codec**
  - An CNN-based autoencoder
  - Trained with a GAN-like adversarial loss
- **Diffusion model in the latent space**
  - A denoising U-Net
- **A conditioning module**
  - Transformer-like cross-attention mechanism



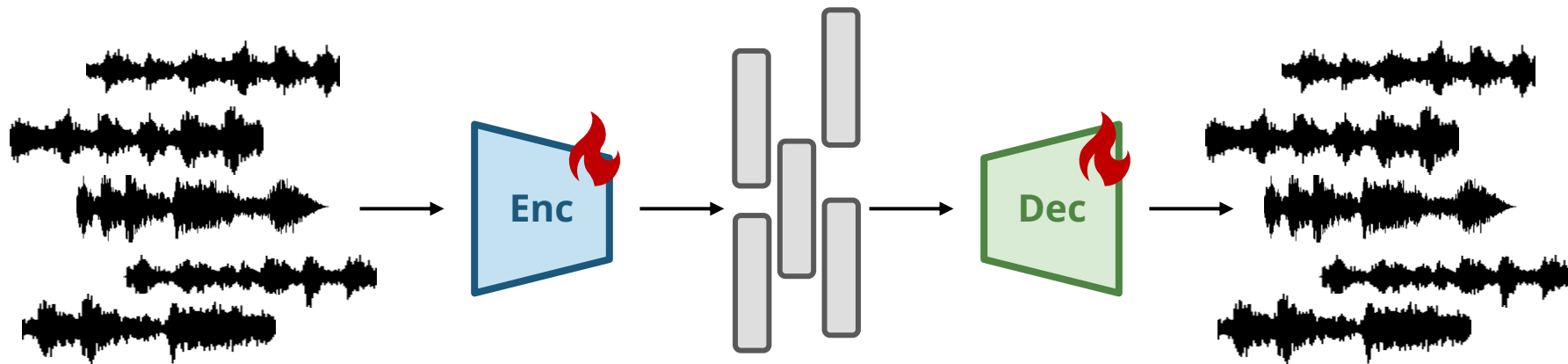
(Source: Raddato via worldhistory.org)

# Latent-based Music Synthesis

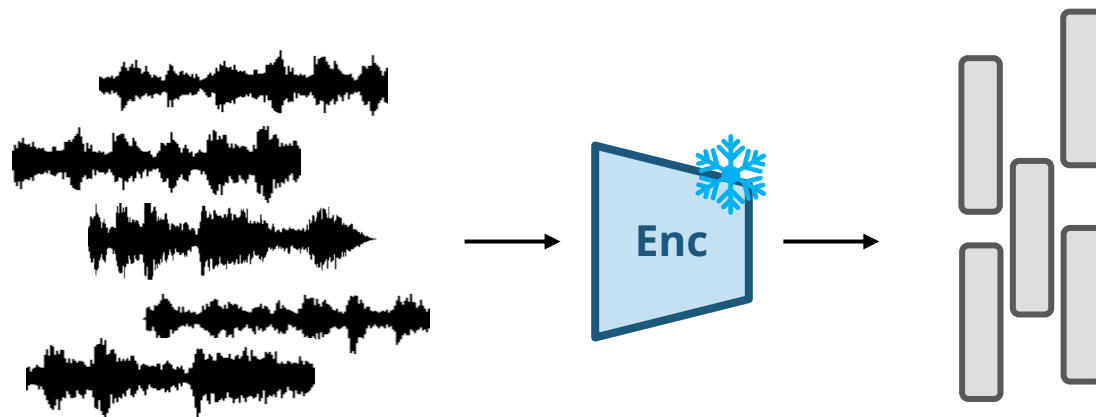
# Latent-based Audio Synthesis



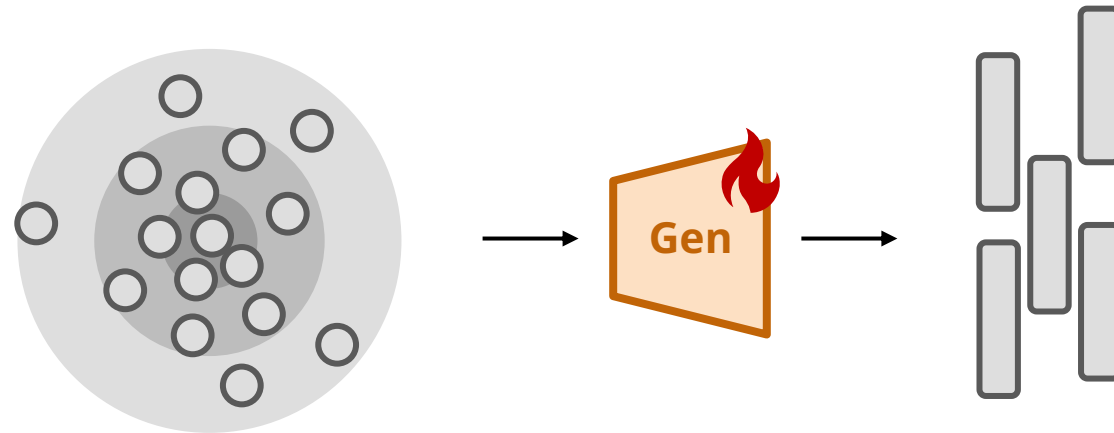
# Step 1: Train an Autoencoder



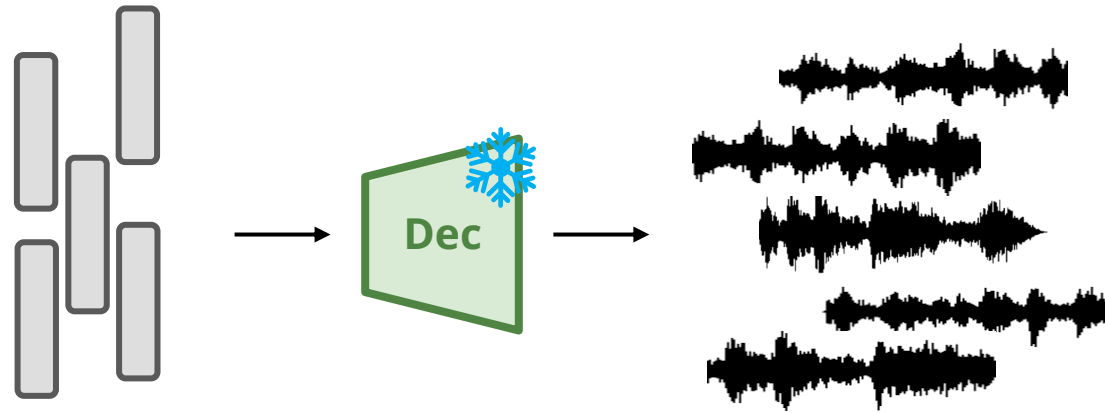
## Step 2: Compute the Latent Vectors



# Step 3: Train a Latent Generative Model

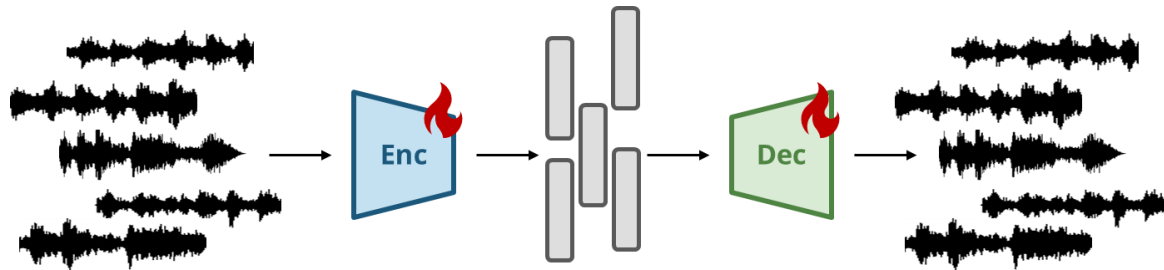


## Step 4: Decode the Latent Vectors

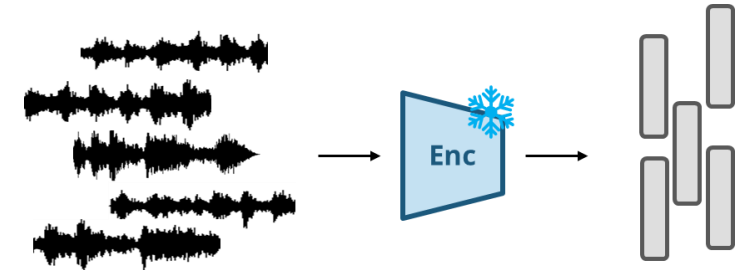


# Latent-based Audio Synthesis: Pipeline

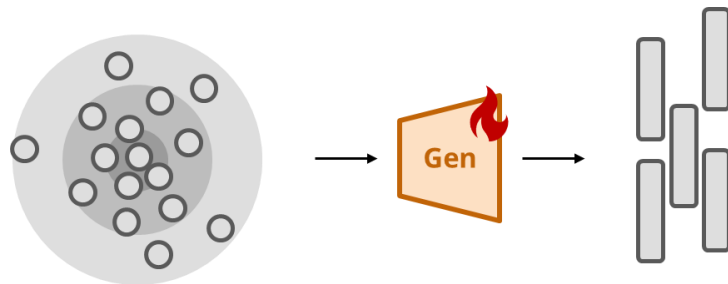
## Step 1: Train an Autoencoder



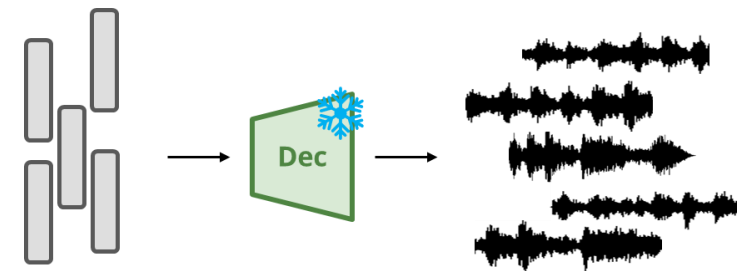
## Step 2: Compute the Latent Vectors



## Step 3: Train a Latent Generative Model

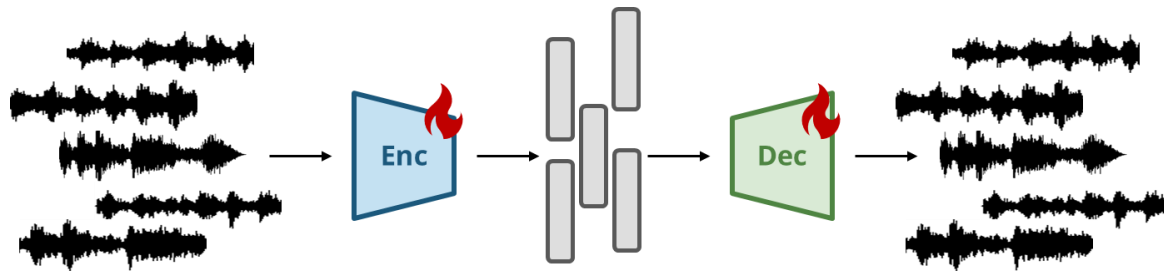


## Step 4: Decode the Latent Vectors

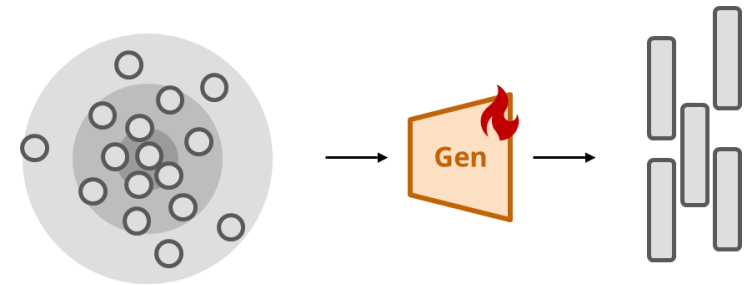


# Latent-based Audio Synthesis: Training

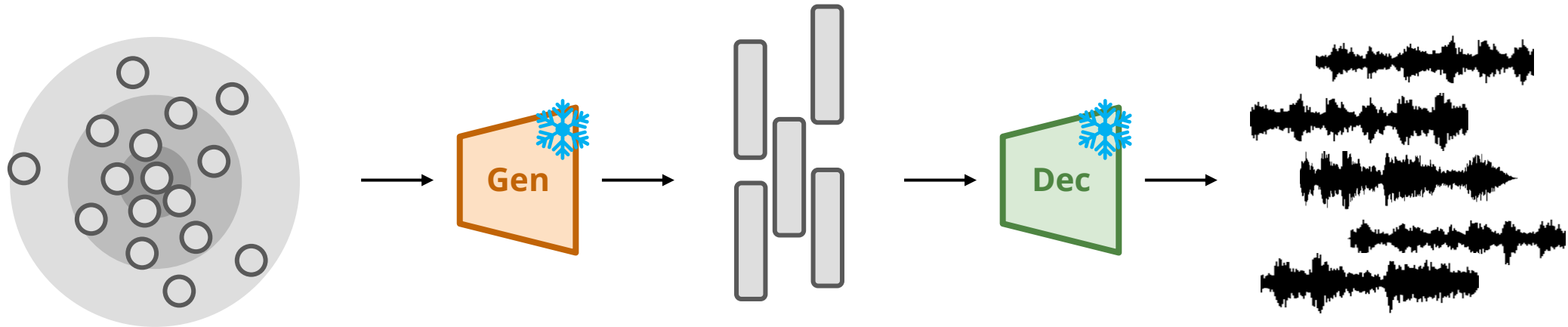
Autoencoder



Latent Generative Model

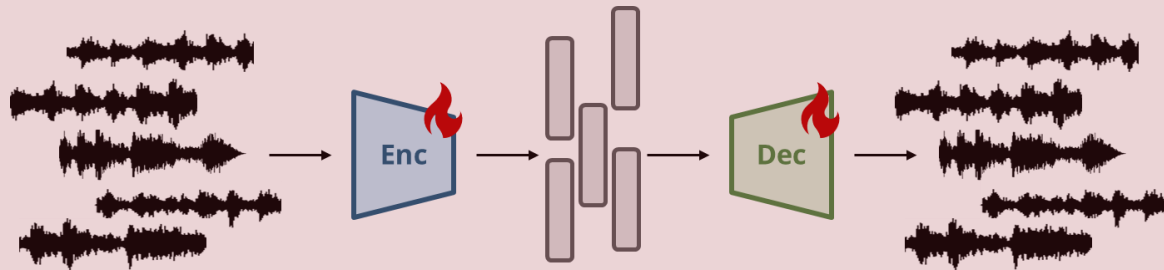


# Latent-based Audio Synthesis: Inference

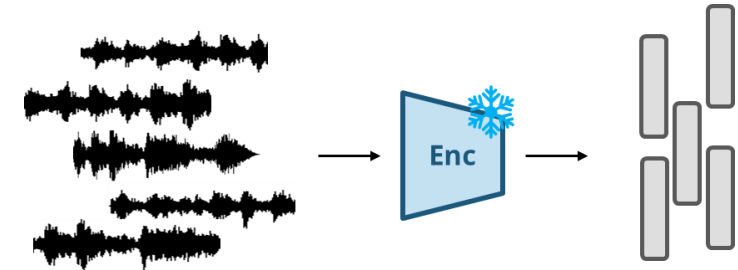


# Latent-based Audio Synthesis: Pipeline

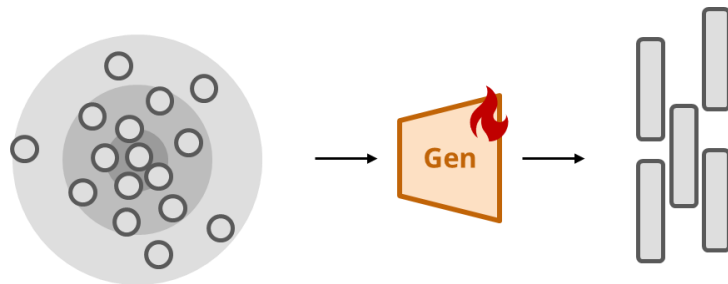
Step 1: Train an Autoencoder



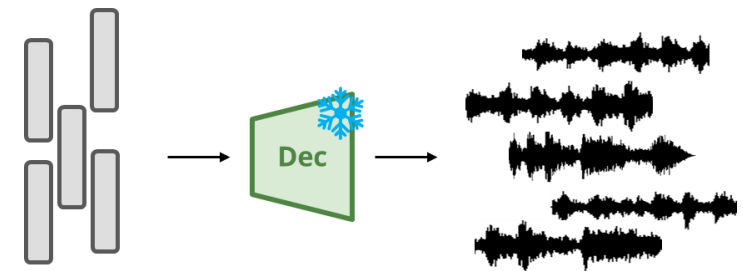
Step 2: Compute the Latent Vectors



Step 3: Train a Latent Generative Model



Step 4: Decode the Latent Vectors



# Neural Codecs

# | What is a Codec?



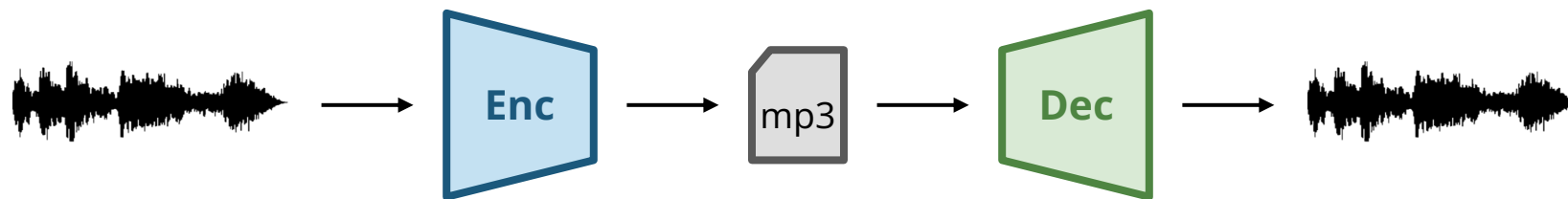
**SBC, AAC, aptX,  
aptX HD, LDAC**

# | What is a Codec?

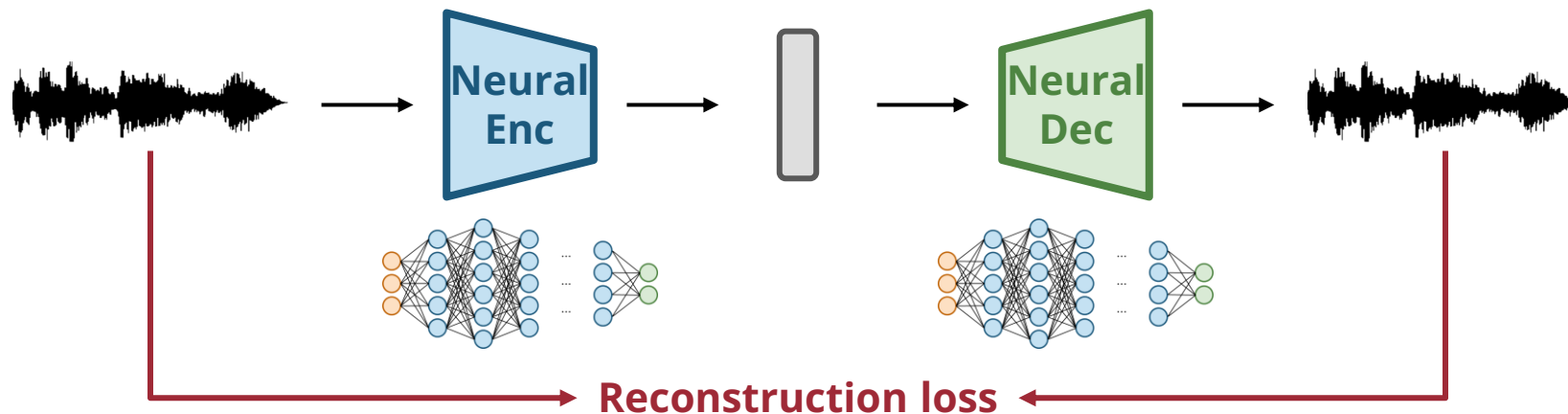


# Neural Codec

## Traditional Codec

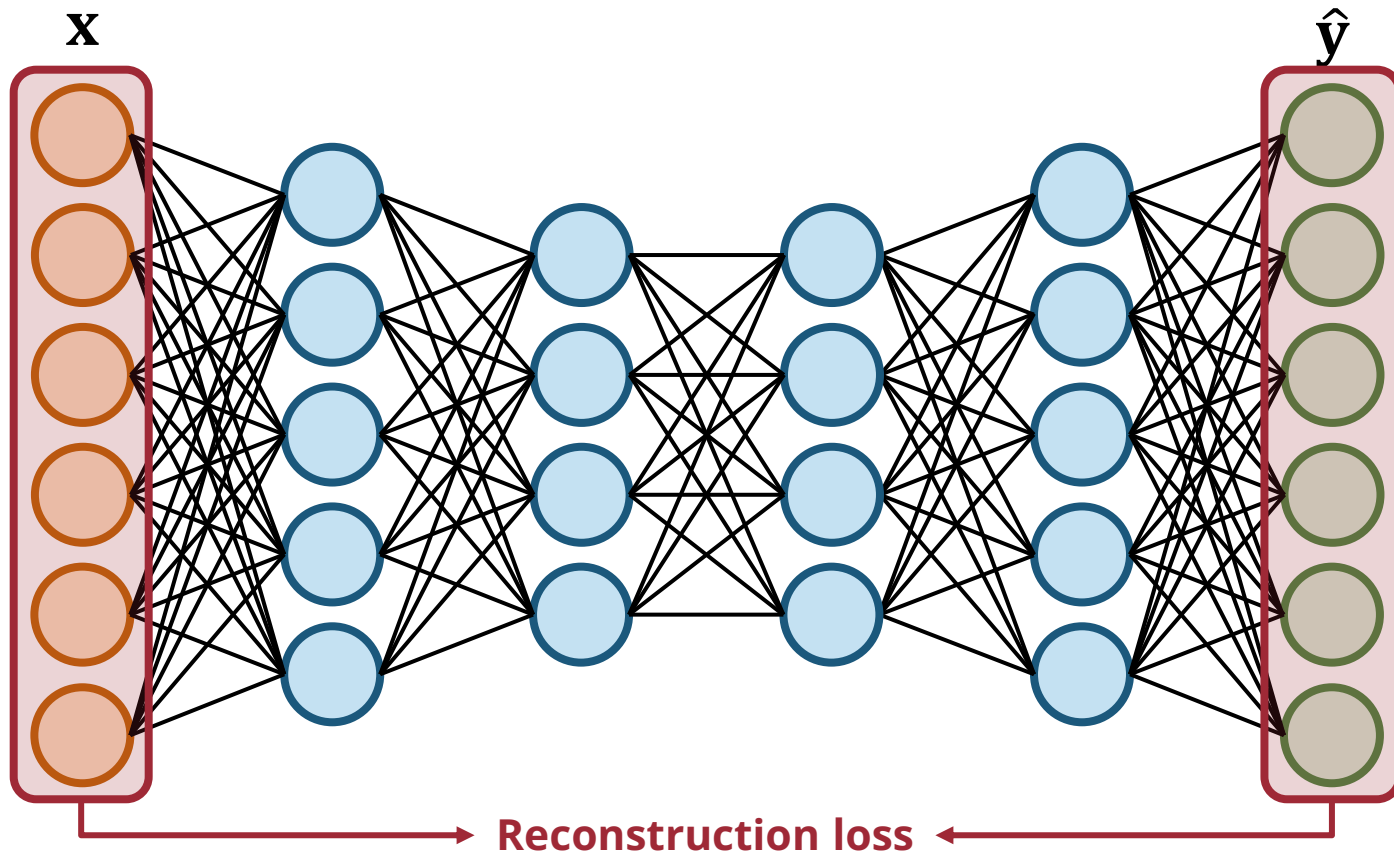


## Neural Codec

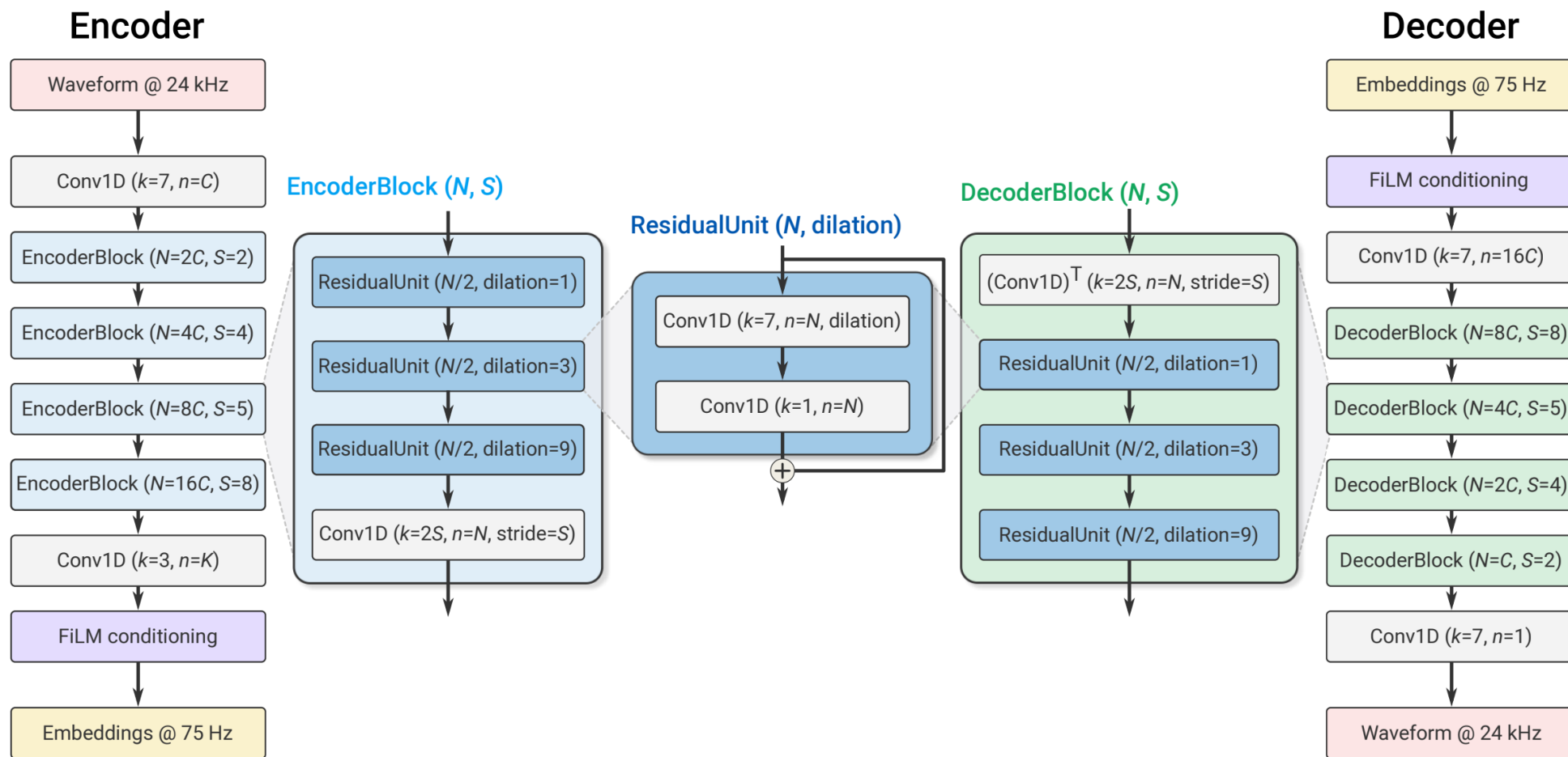


# Autoencoders

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

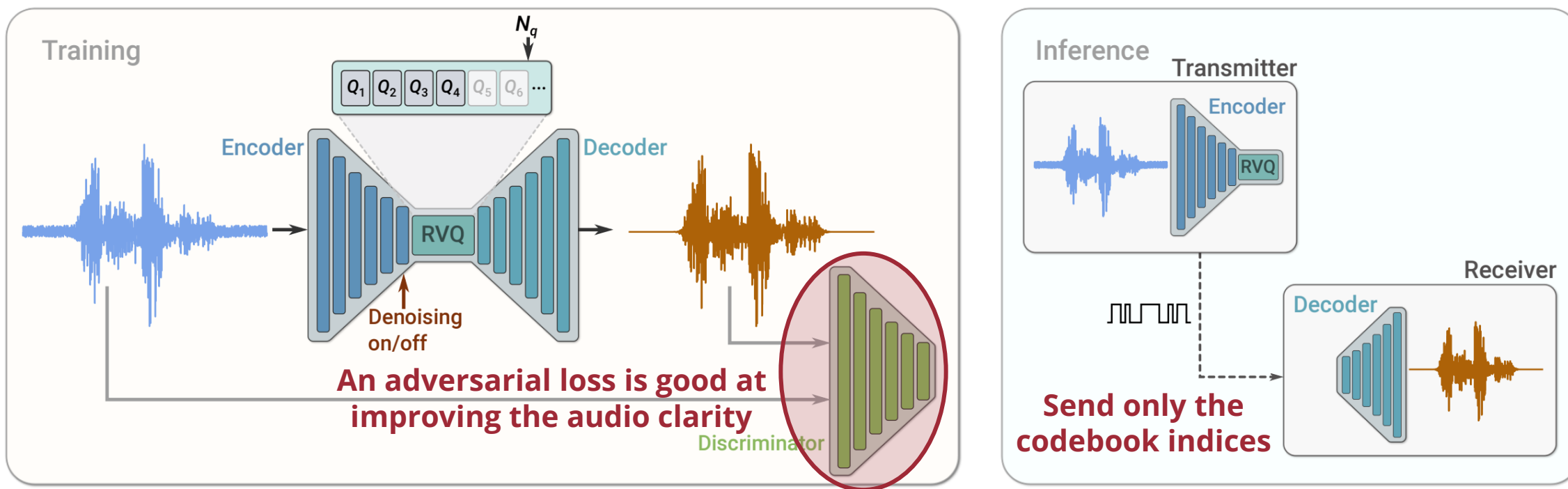


# SoundStream (Zeghidour et al., 2021)



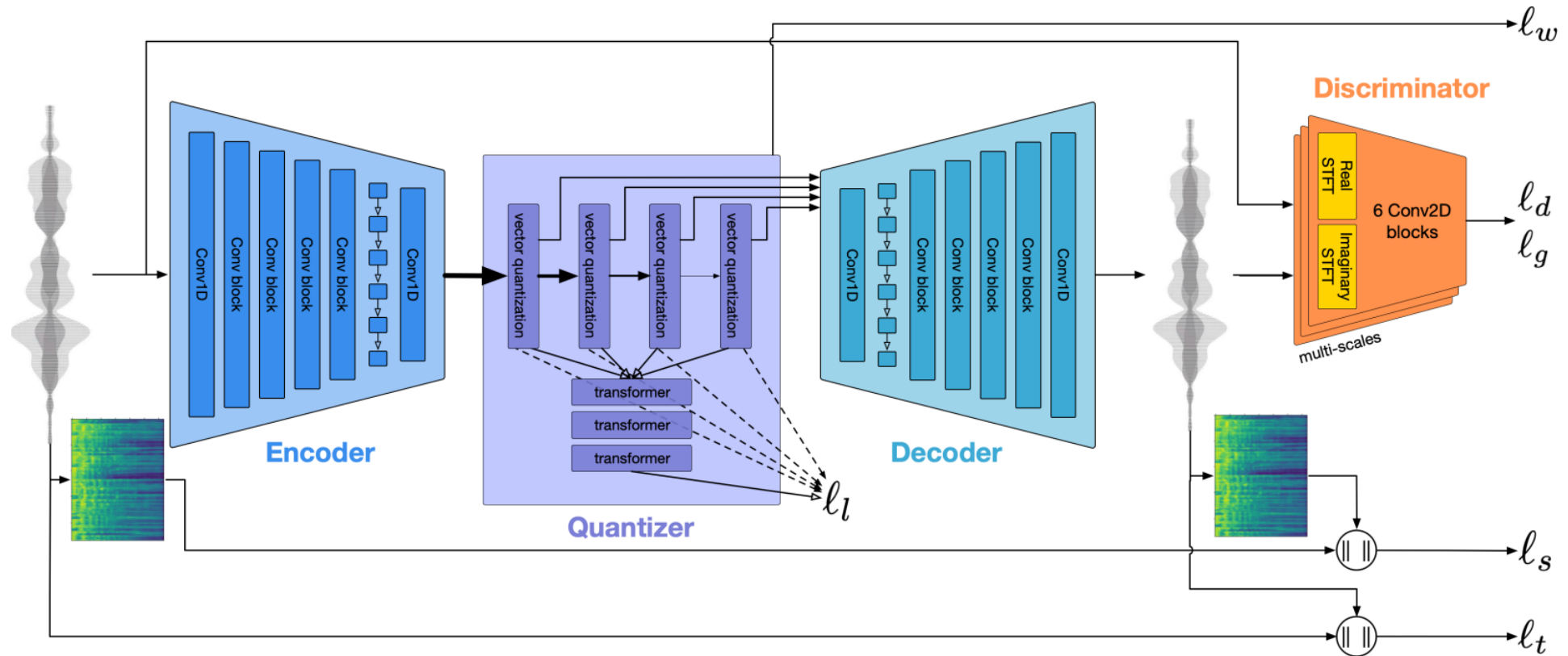
(Source: Zeghidour et al., 2021)

# SoundStream (Zeghidour et al., 2021)



(Source: Zeghidour et al., 2021)

# EnCodec (Défossez et al., 2022)

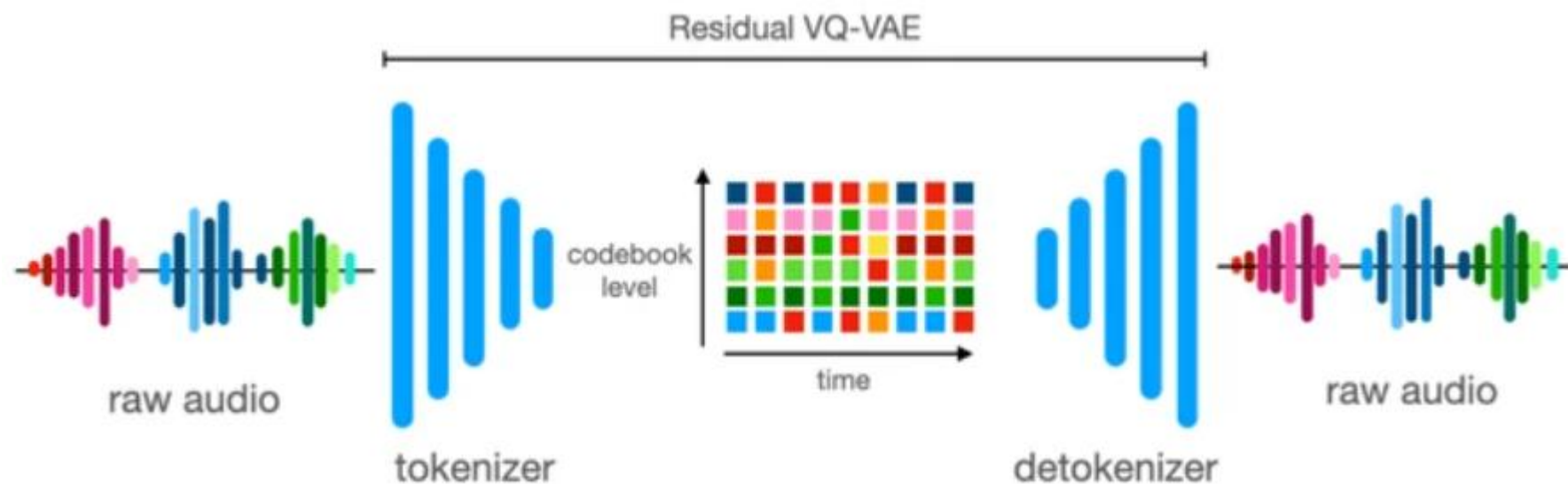


(Source: Défossez et al., 2022)

[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)



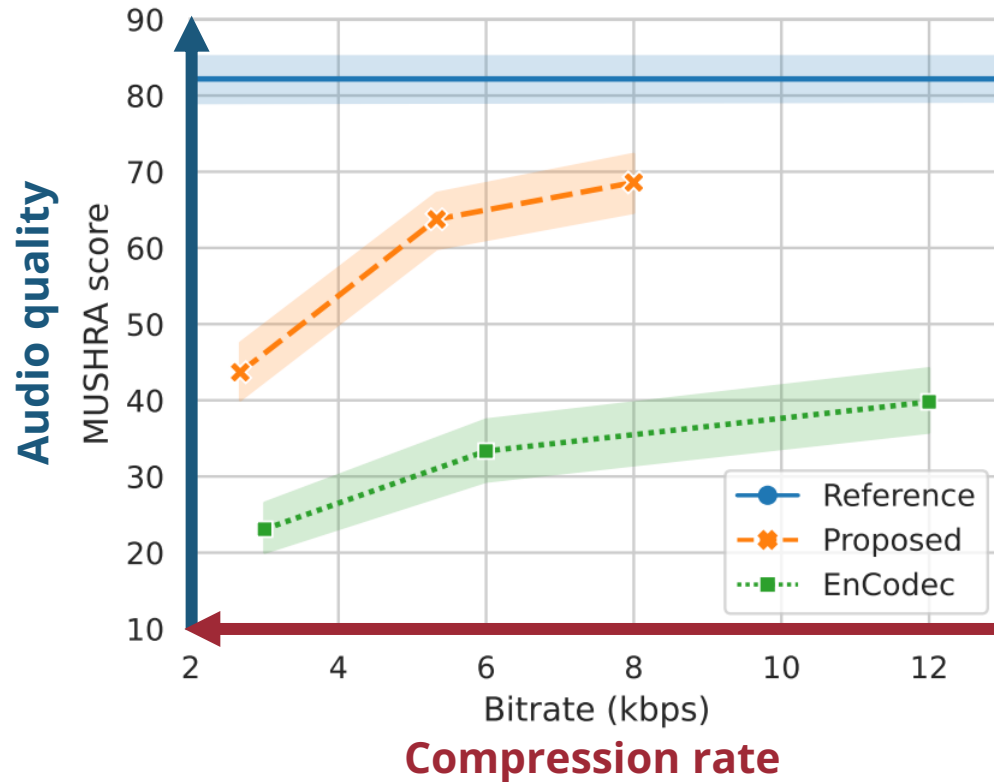
(Source: Kumar et al., 2023)

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

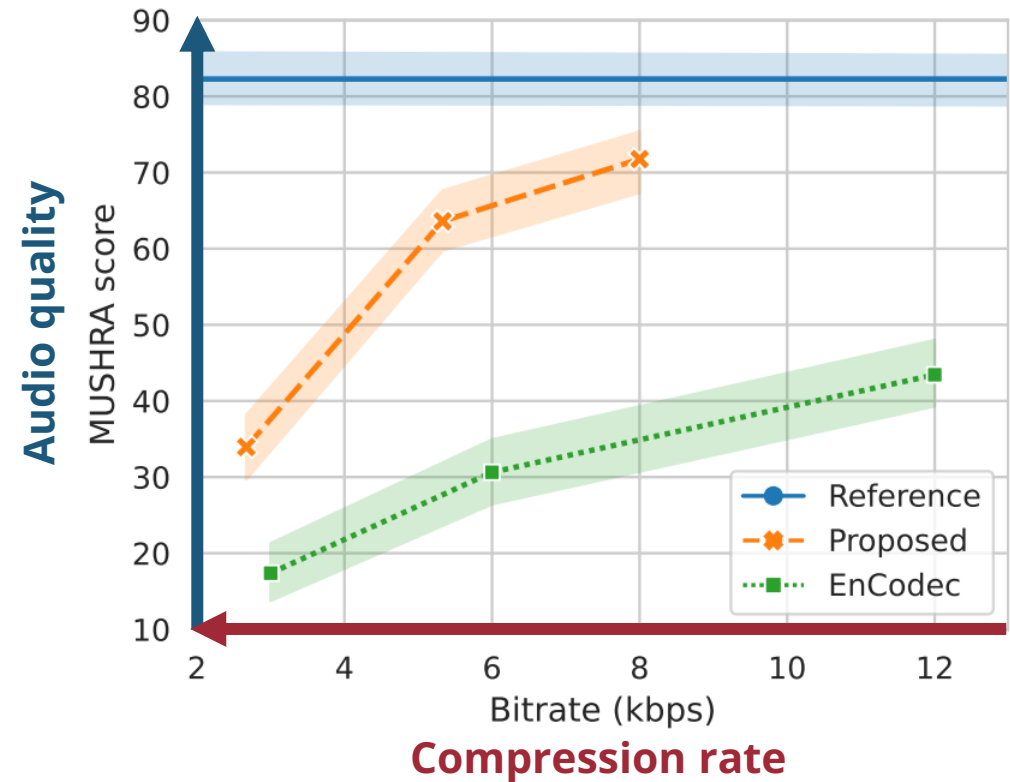
[github.com/descriptinc/descript-audio-codec](https://github.com/descriptinc/descript-audio-codec)

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

Listening Test Results @ 44.1 kHz



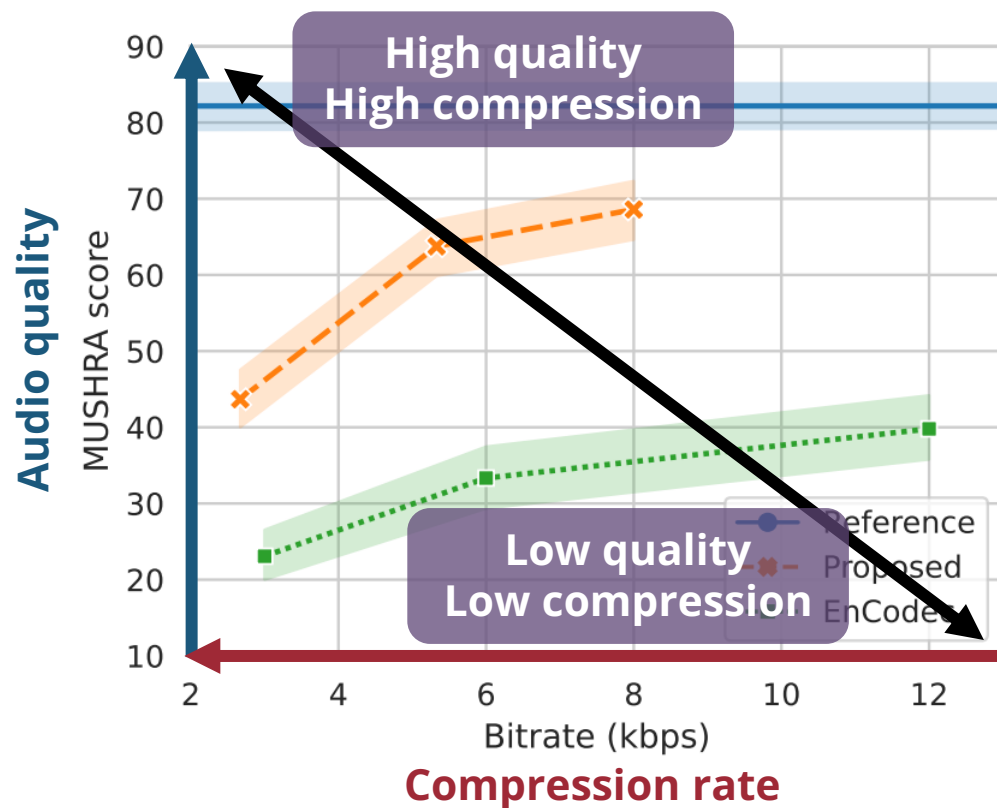
Listening Test Results @ 24 kHz



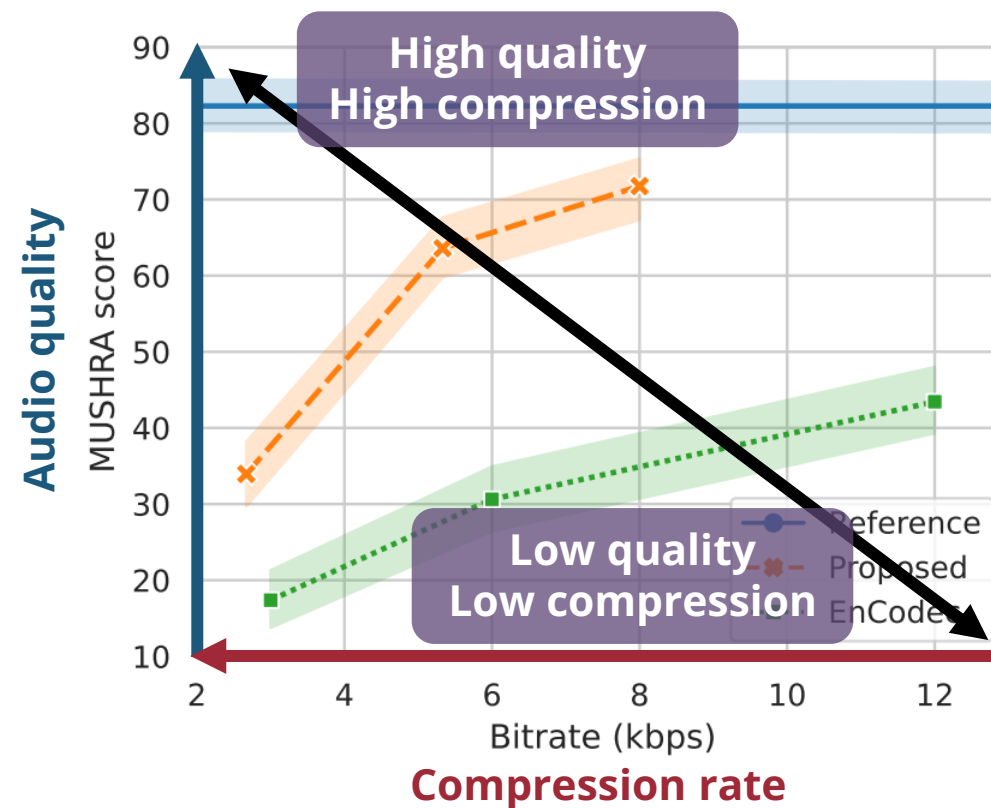
(Source: Kumar et al., 2023)

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

Listening Test Results @ 44.1 kHz



Listening Test Results @ 24 kHz

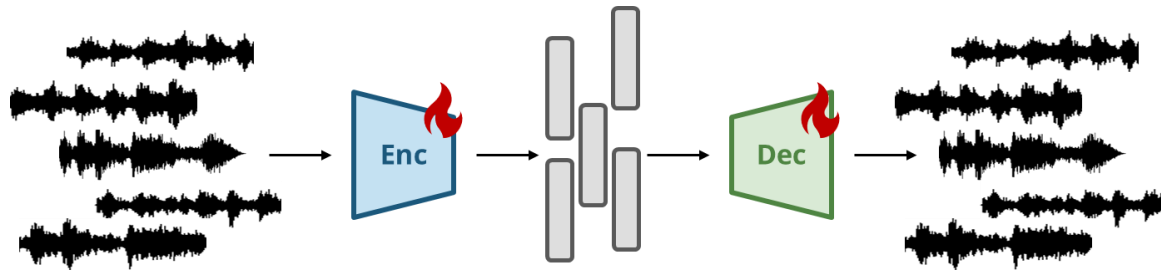


(Source: Kumar et al., 2023)

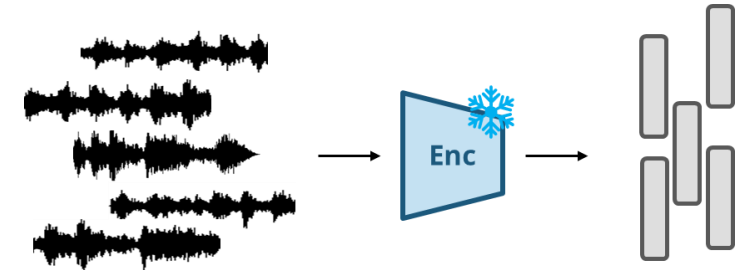
# Latent-based Text-to-Music Synthesis

# Latent-based Audio Synthesis: Pipeline

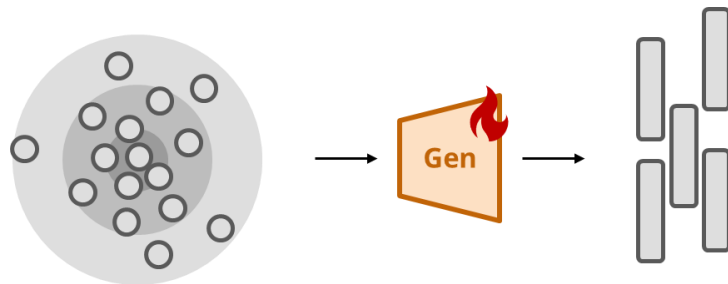
Step 1: Train an Autoencoder



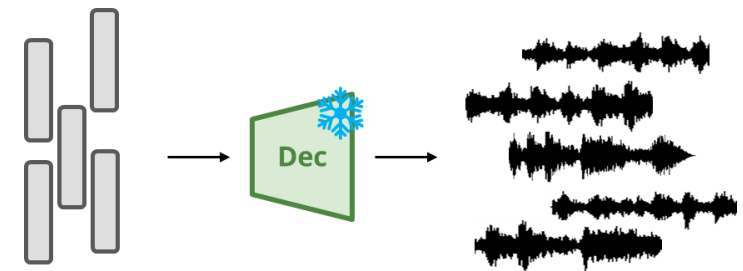
Step 2: Compute the Latent Vectors



Step 3: Train a Latent Generative Model

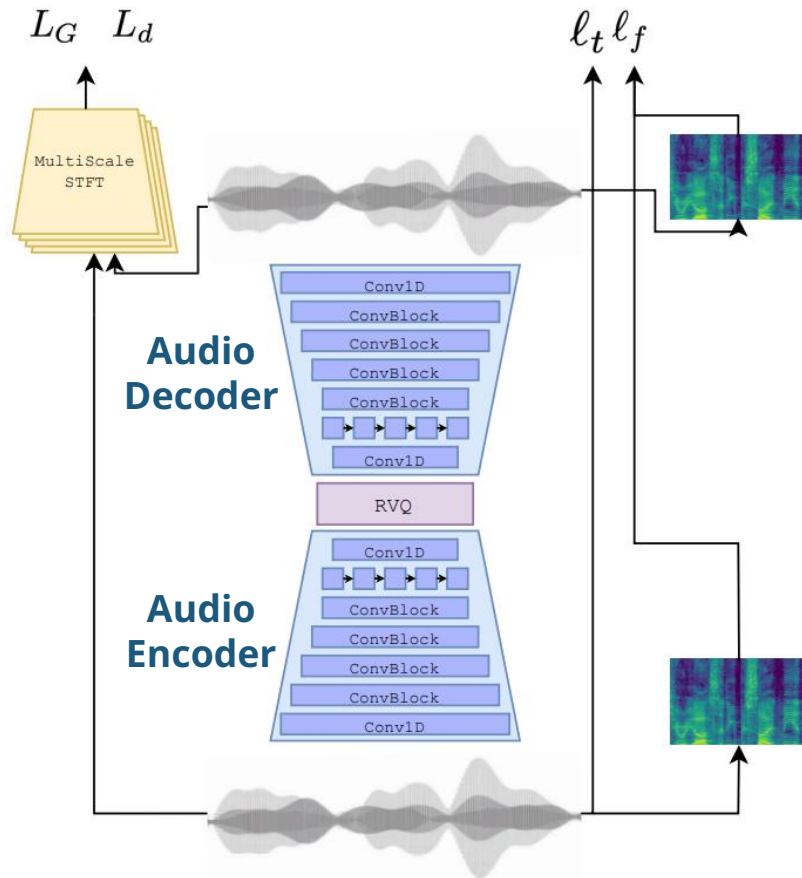


Step 4: Decode the Latent Vectors

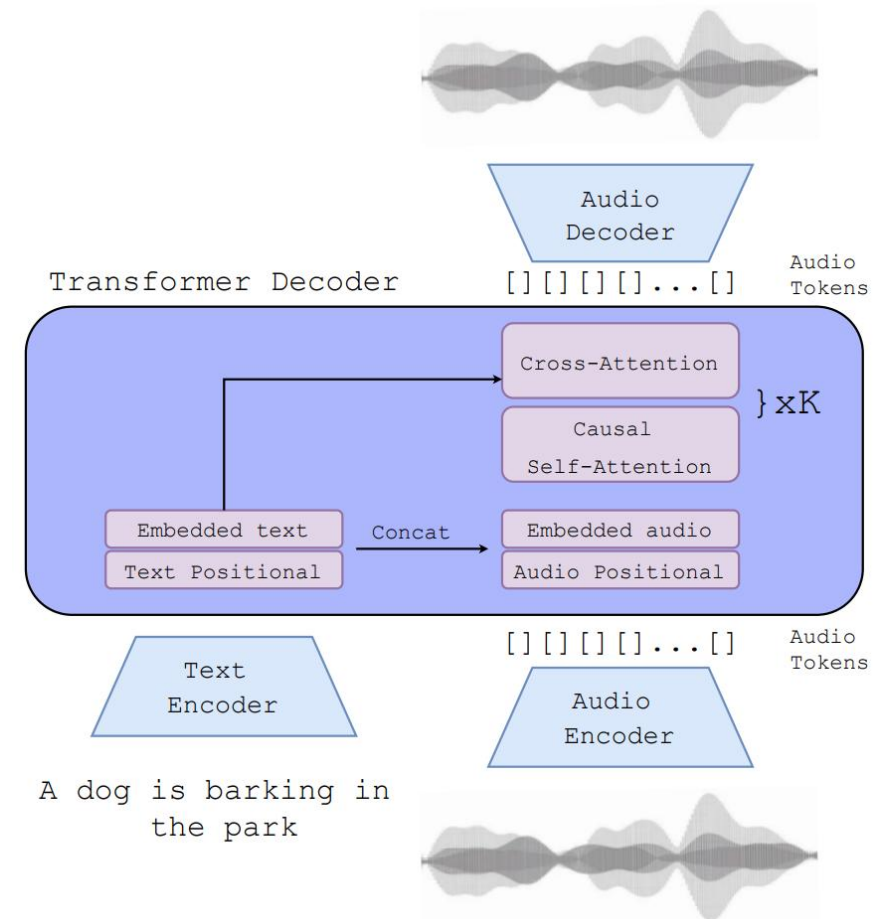


# AudioGen (Kreuk et al., 2023)

## Audio Autoencoder



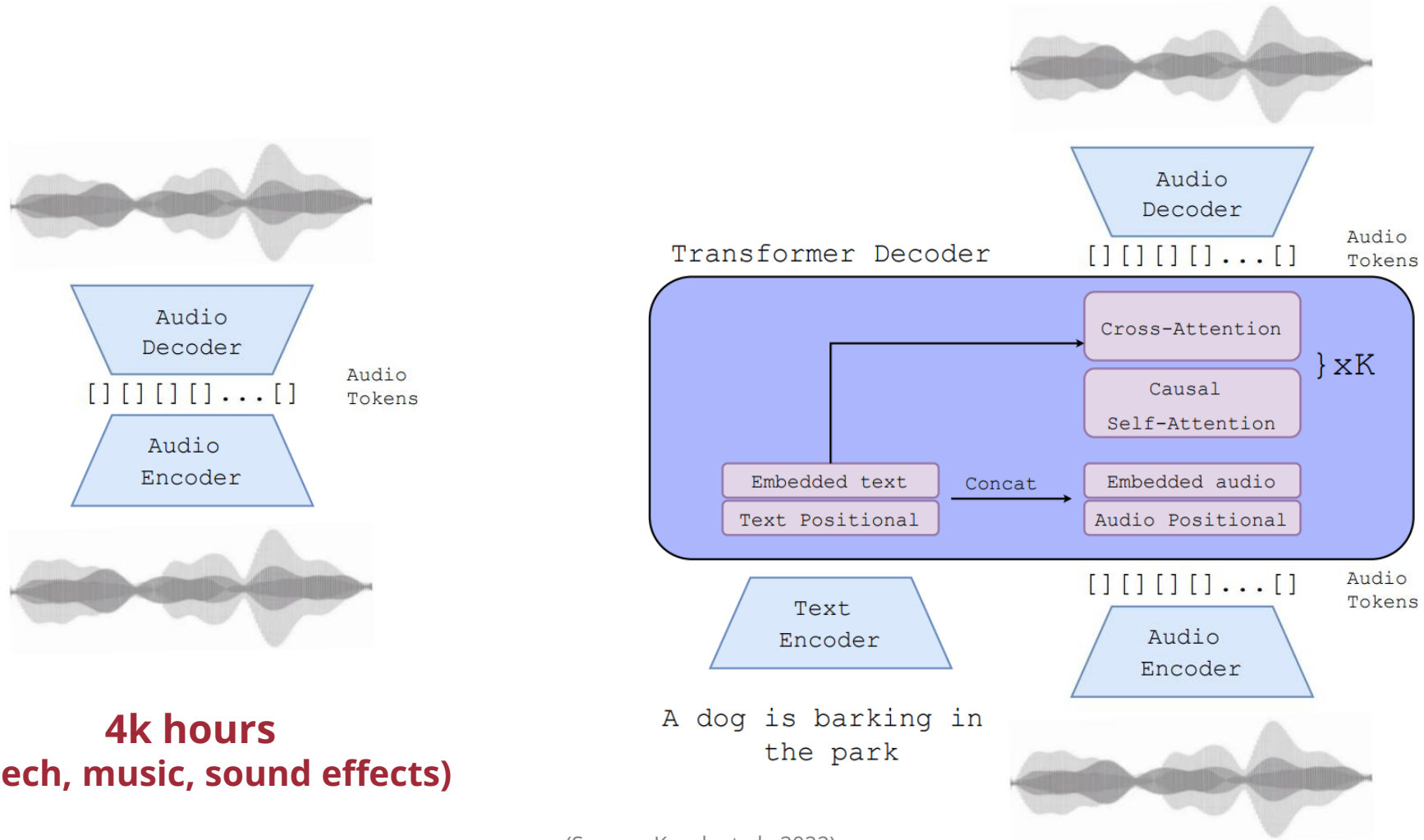
## Audio Language Model



(Source: Kreuk et al., 2022)

Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi, "AudioGen: Textually Guided Audio Generation," ICLR, 2023.

# AudioGen (Kreuk et al., 2023)



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

A dog is barking in  
the park

(Source: Kreuk et al., 2022)

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



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

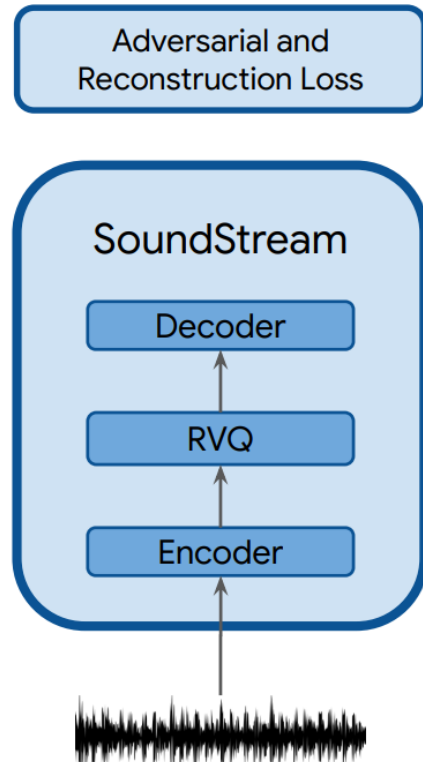
# MusicGen (Copet et al., 2023)

- 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/)

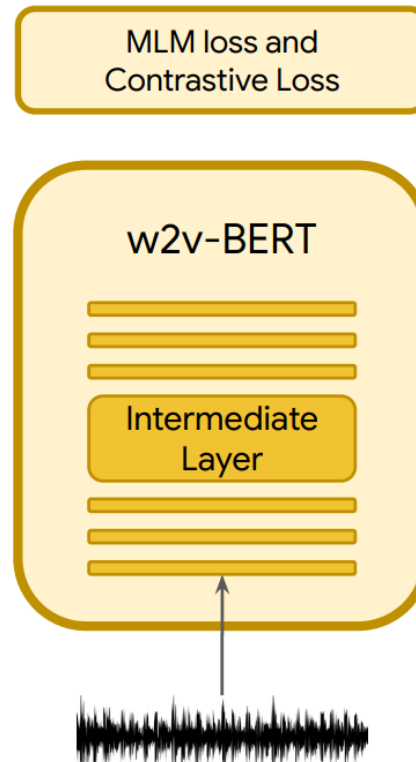
# MusicLM (Agostinelli et al., 2023)

## Audio autoencoder

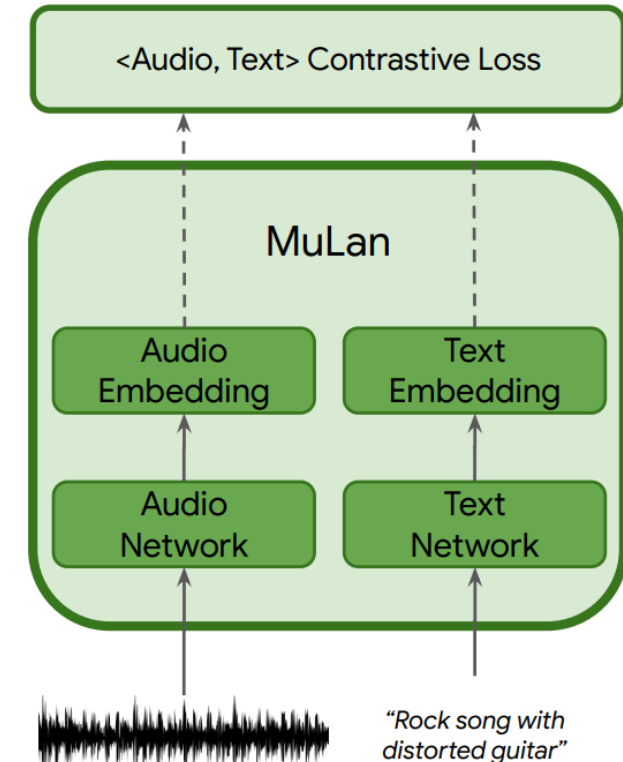


**106k songs, 8.2k hours**

## Semantic representation



## Text-music correspondence

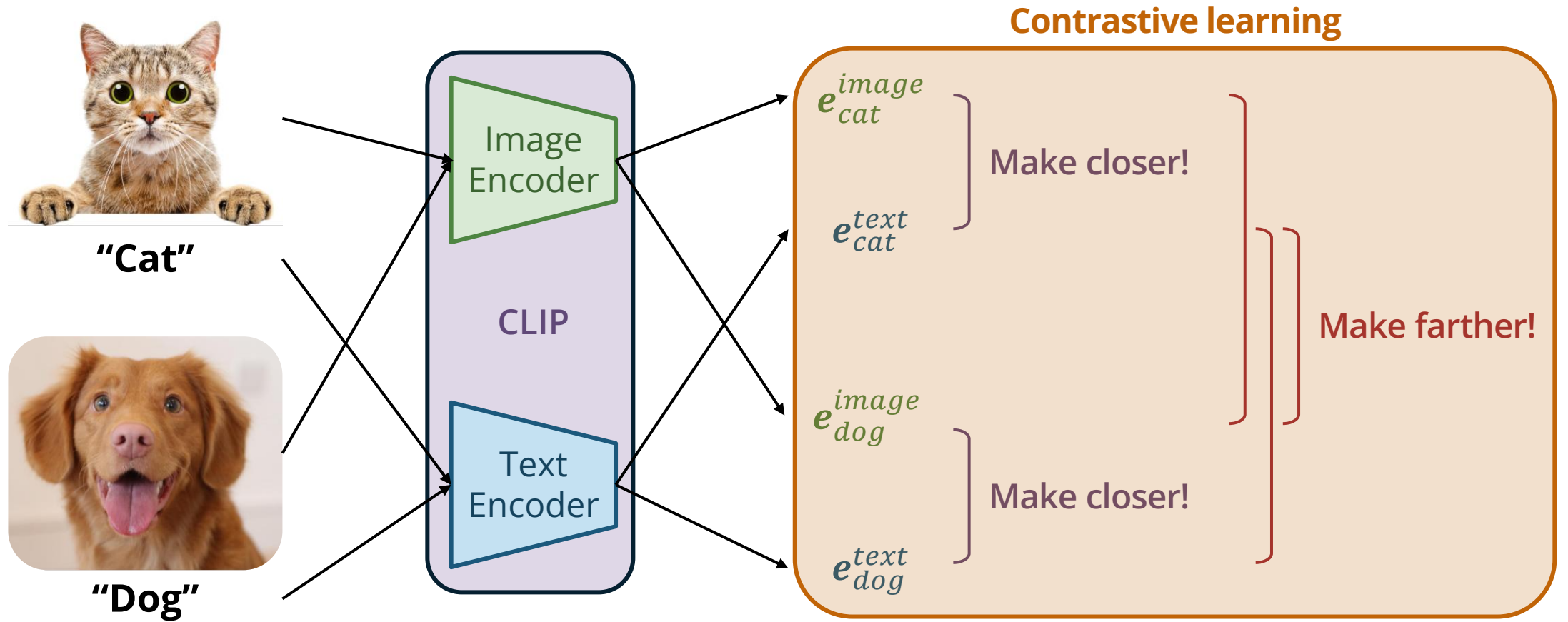


**44M 30-sec clips, 370k hours**

(Source: Agostinelli et al., 2022)

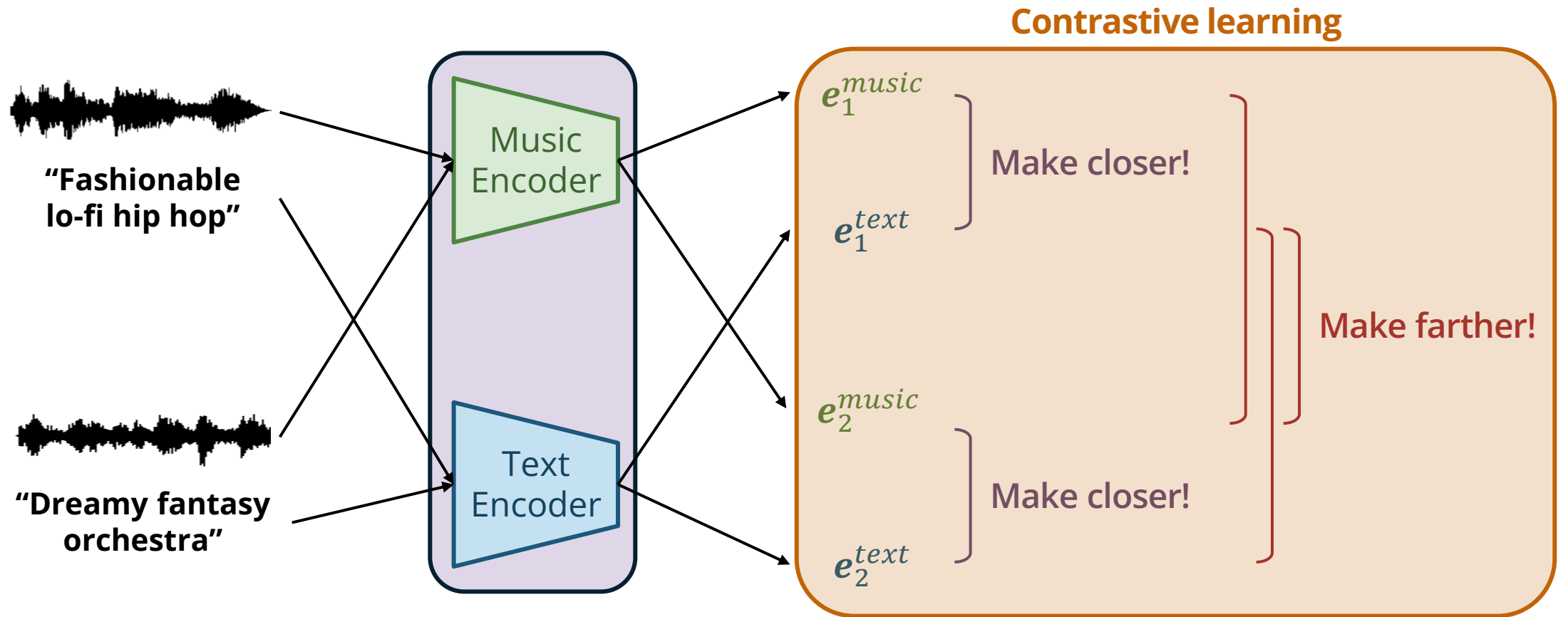
Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Christian Frank, "MusicLM: Generating Music From Text," *arXiv preprint arXiv:2301.11325*, 2023.

# Contrastive Language-Image Pretraining (CLIP)



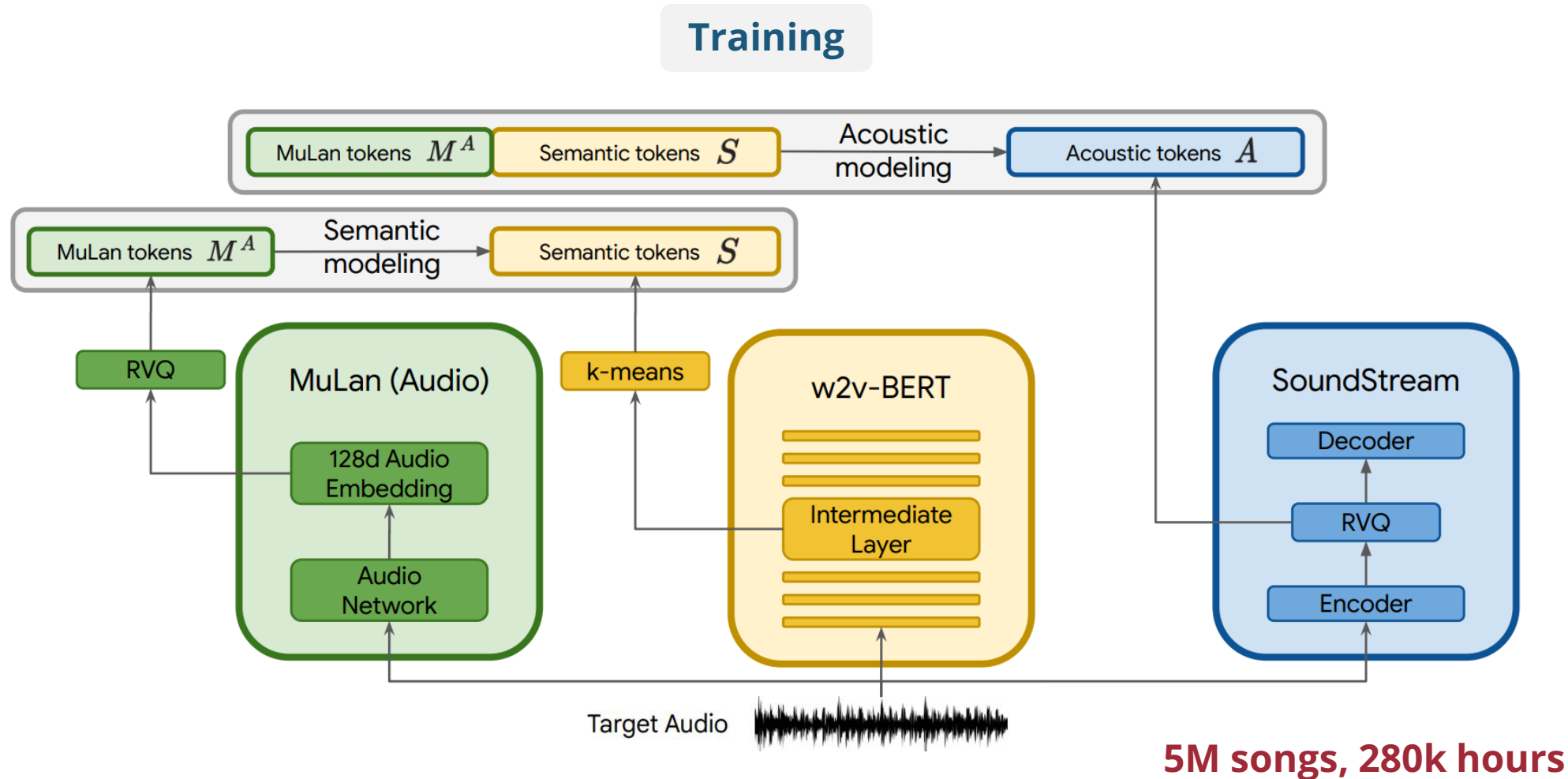
Learn a **shared embedding space** for images and language

# Contrastive Language-Music Pretraining



Learn a **shared embedding space** for music and language

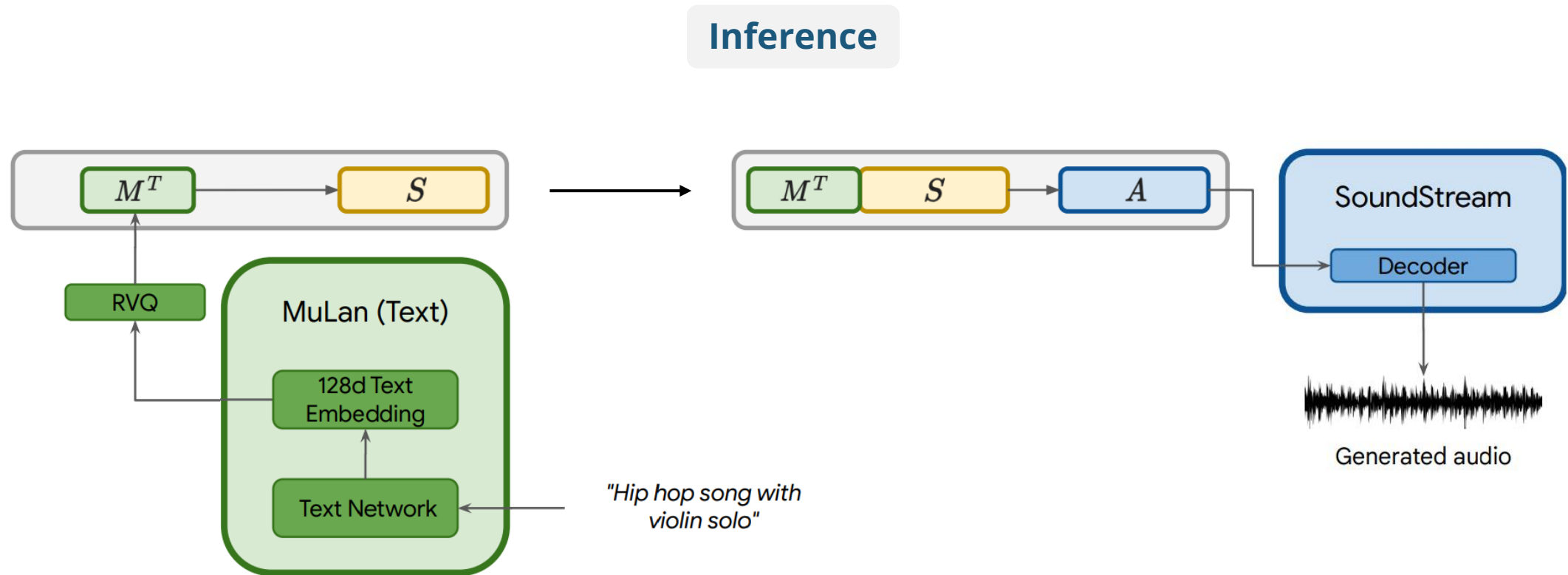
# MusicLM (Agostinelli et al., 2023)



(Source: Agostinelli et al., 2022)

Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Christian Frank, "MusicLM: Generating Music From Text," *arXiv preprint arXiv:2301.11325*, 2023.

# MusicLM (Agostinelli et al., 2023)

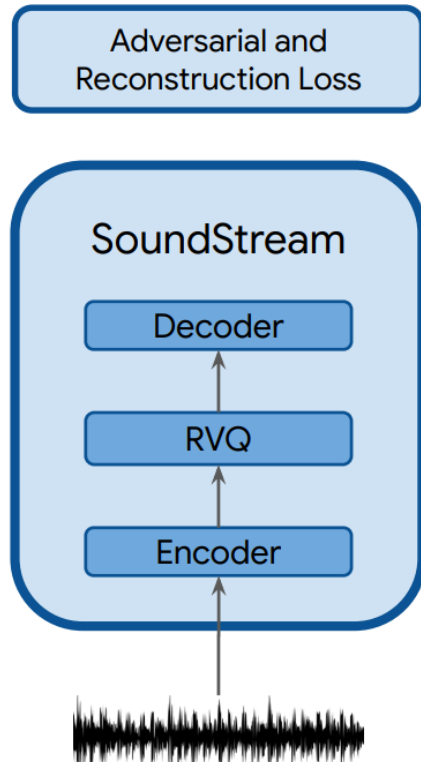


(Source: Agostinelli et al., 2022)

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

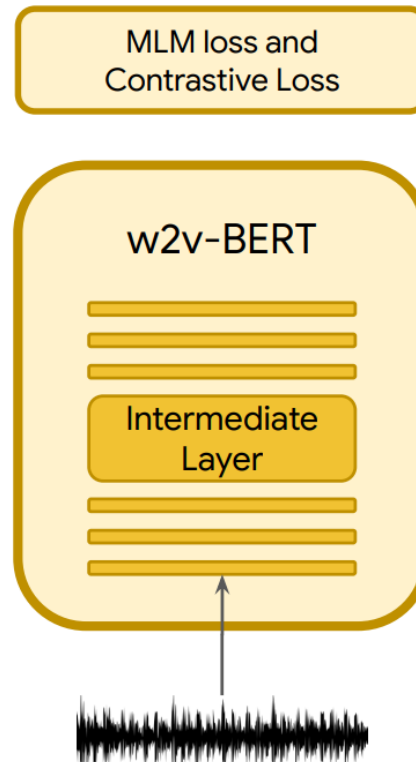
# MusicLM (Agostinelli et al., 2023)

## Audio autoencoder

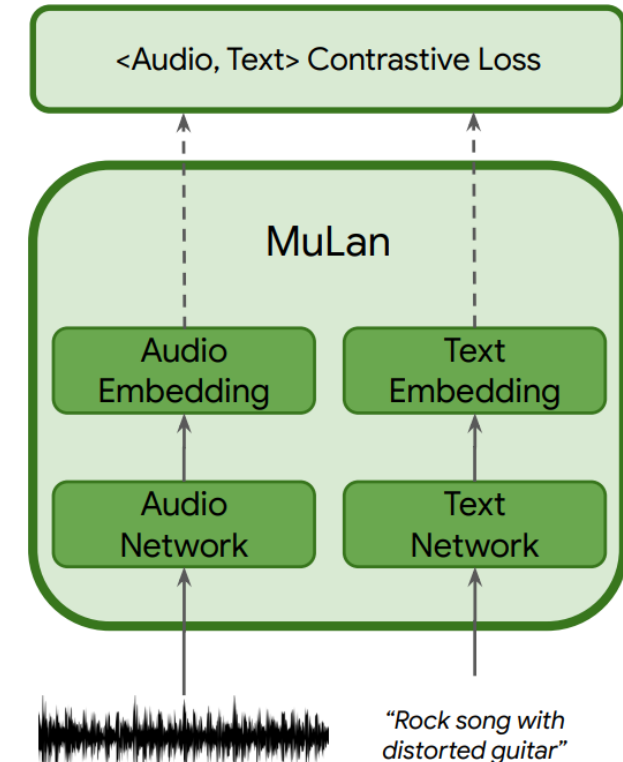


**106k songs, 8.2k hours**

## Semantic representation



## Text-music correspondence

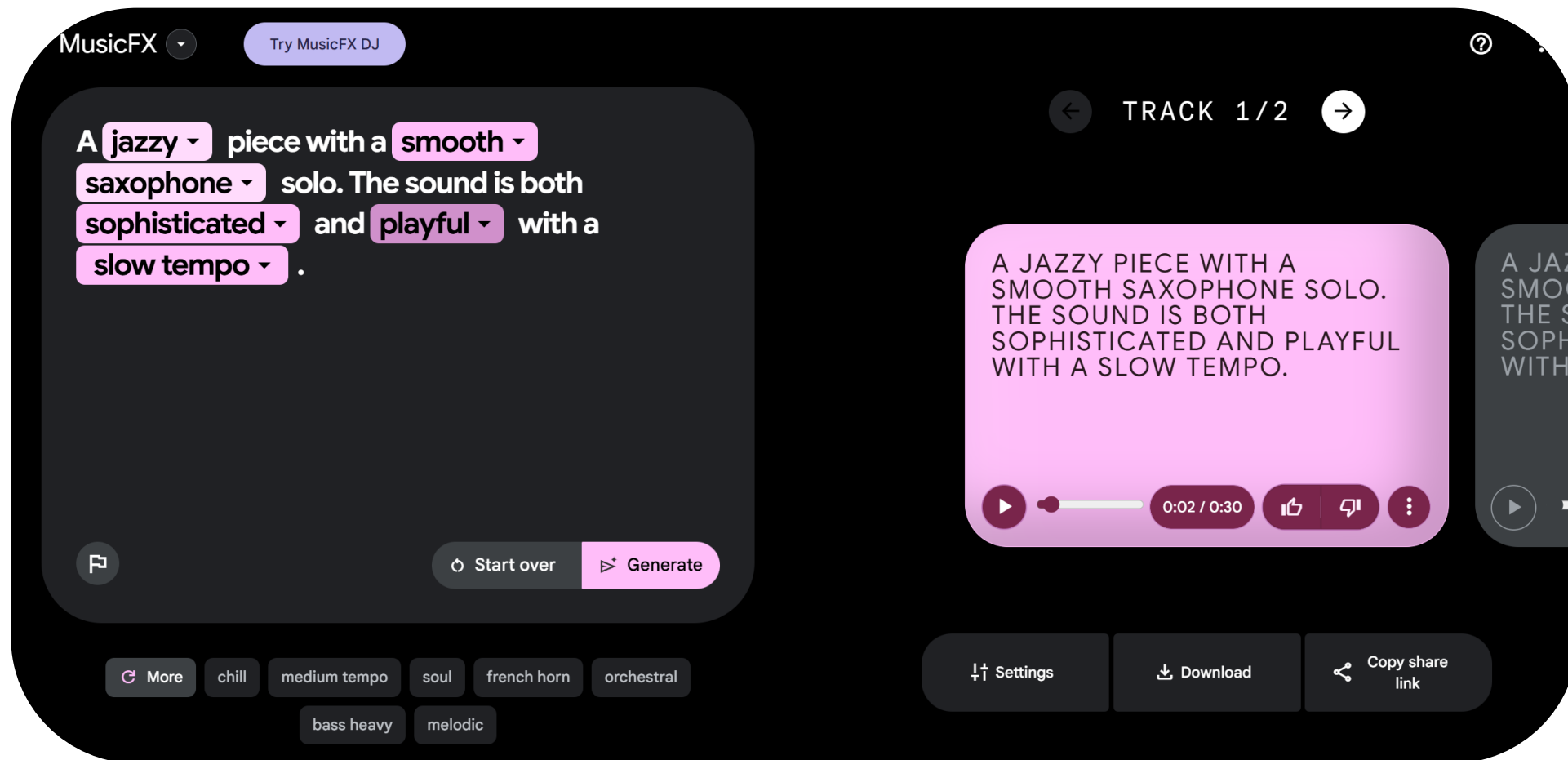


**44M 30-sec clips, 370k hours**

(Source: Agostinelli et al., 2022)

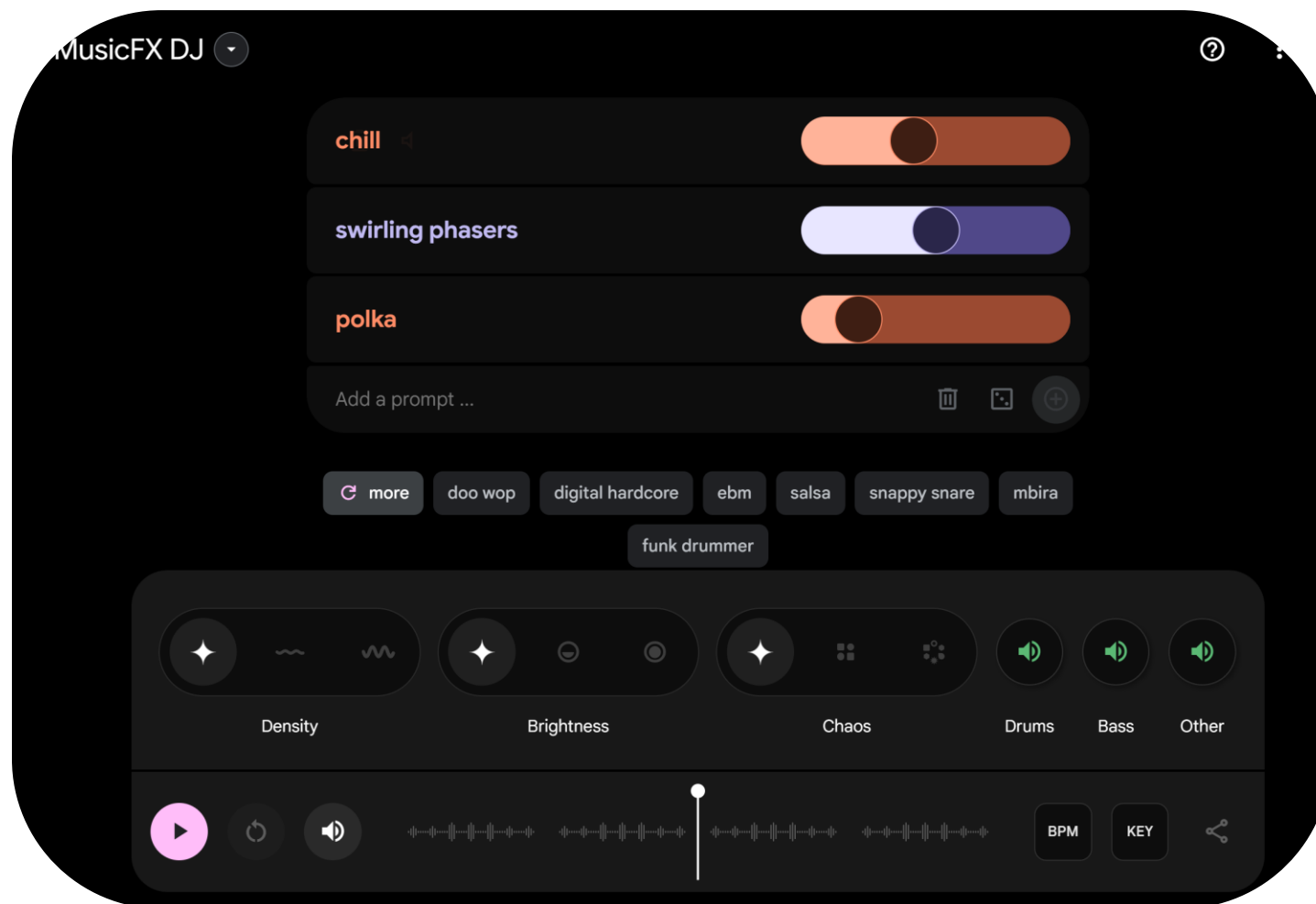
Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Christian Frank, "MusicLM: Generating Music From Text," *arXiv preprint arXiv:2301.11325*, 2023.

# Google's Music FX (2024)



[labs.google/fx/tools/music-fx](https://labs.google/fx/tools/music-fx)

# Google's Music FX DJ (2024)



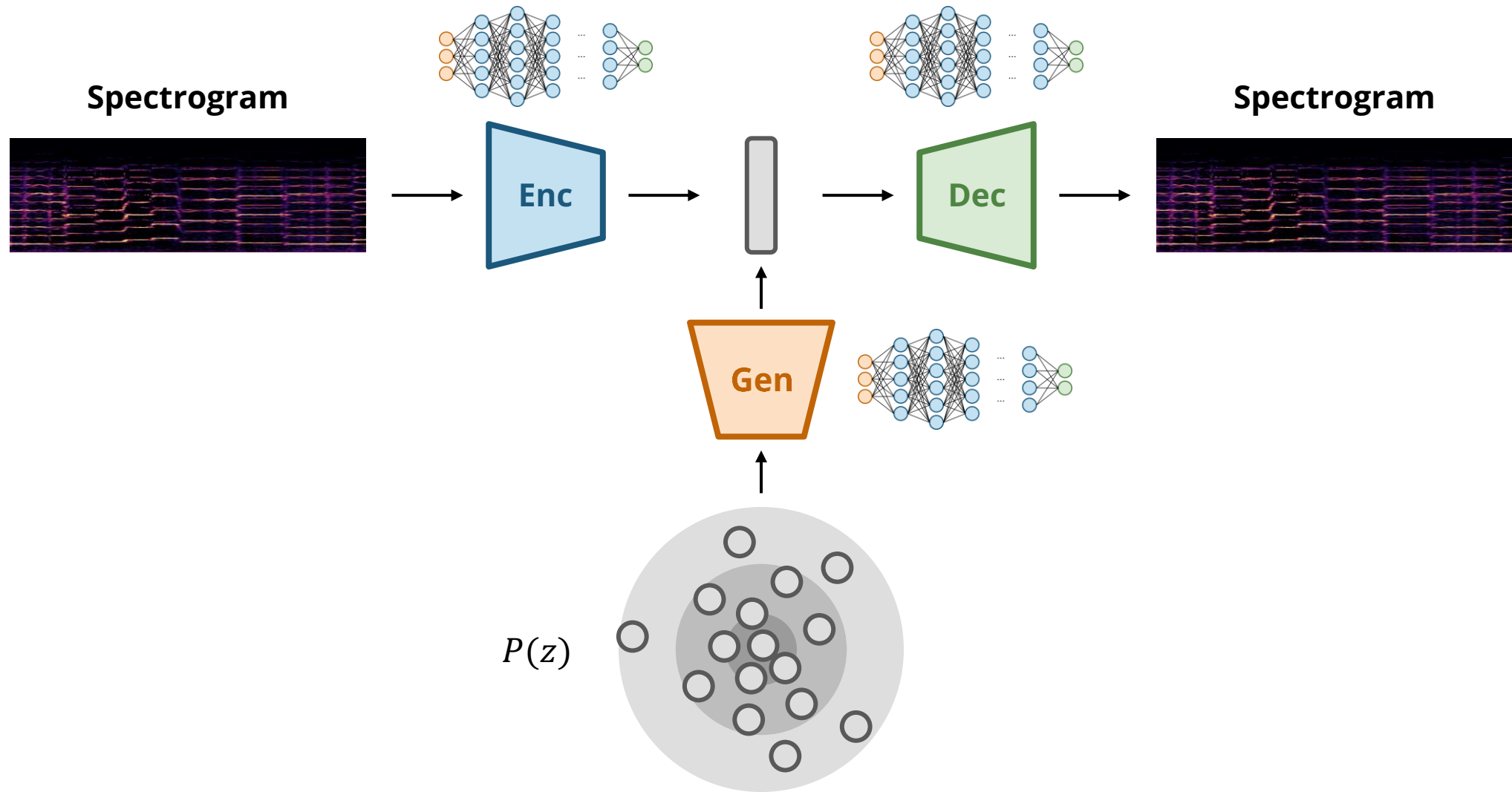
[labs.google/fx/tools/music-fx-dj](https://labs.google/fx/tools/music-fx-dj)

# Google's Music FX DJ (2024)



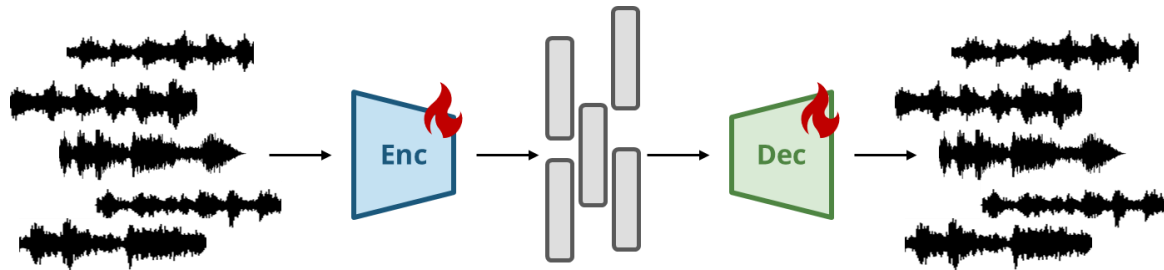
[youtube.com/live/IUQW5LgBZvQ](https://youtube.com/live/IUQW5LgBZvQ)

# Latent-based Audio Synthesis

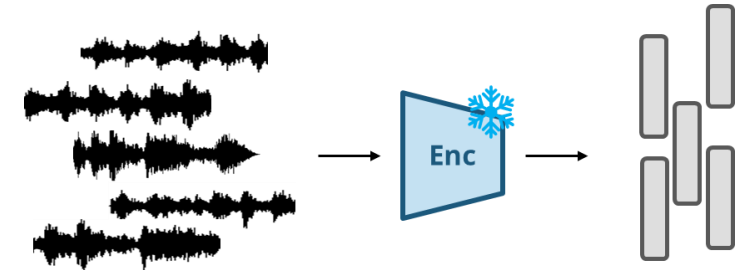


# Latent-based Audio Synthesis: Pipeline

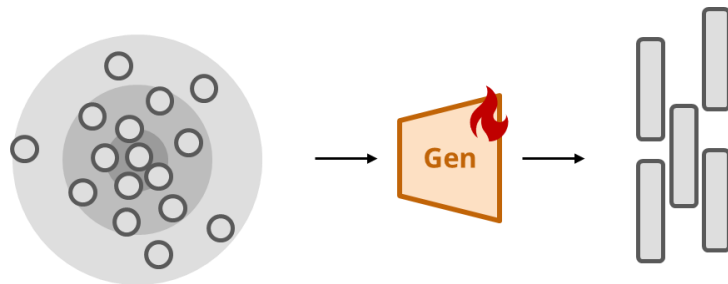
Step 1: Train an Autoencoder



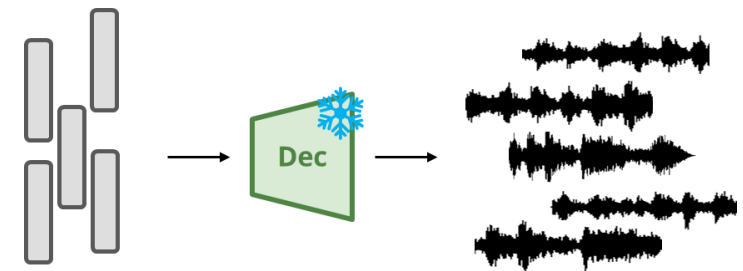
Step 2: Compute the Latent Vectors



Step 3: Train a Latent Generative Model

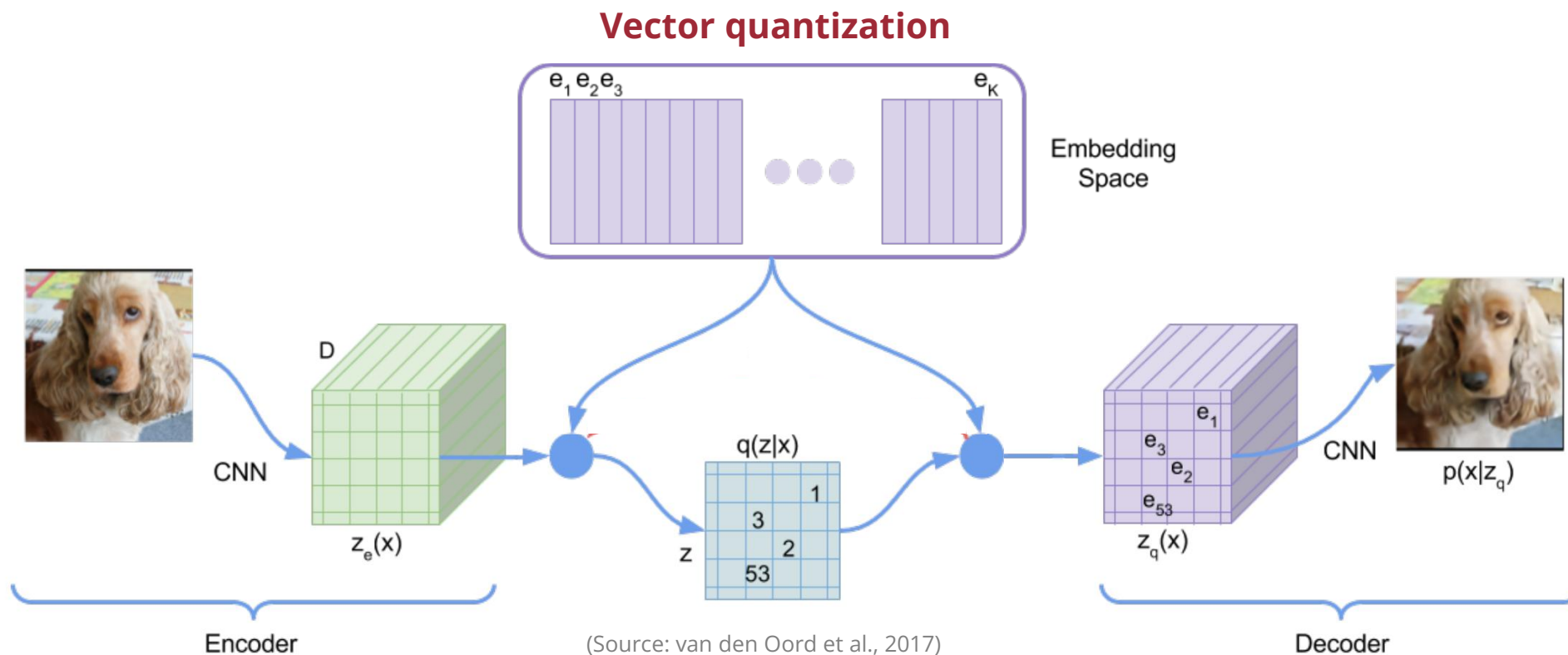


Step 4: Decode the Latent Vectors



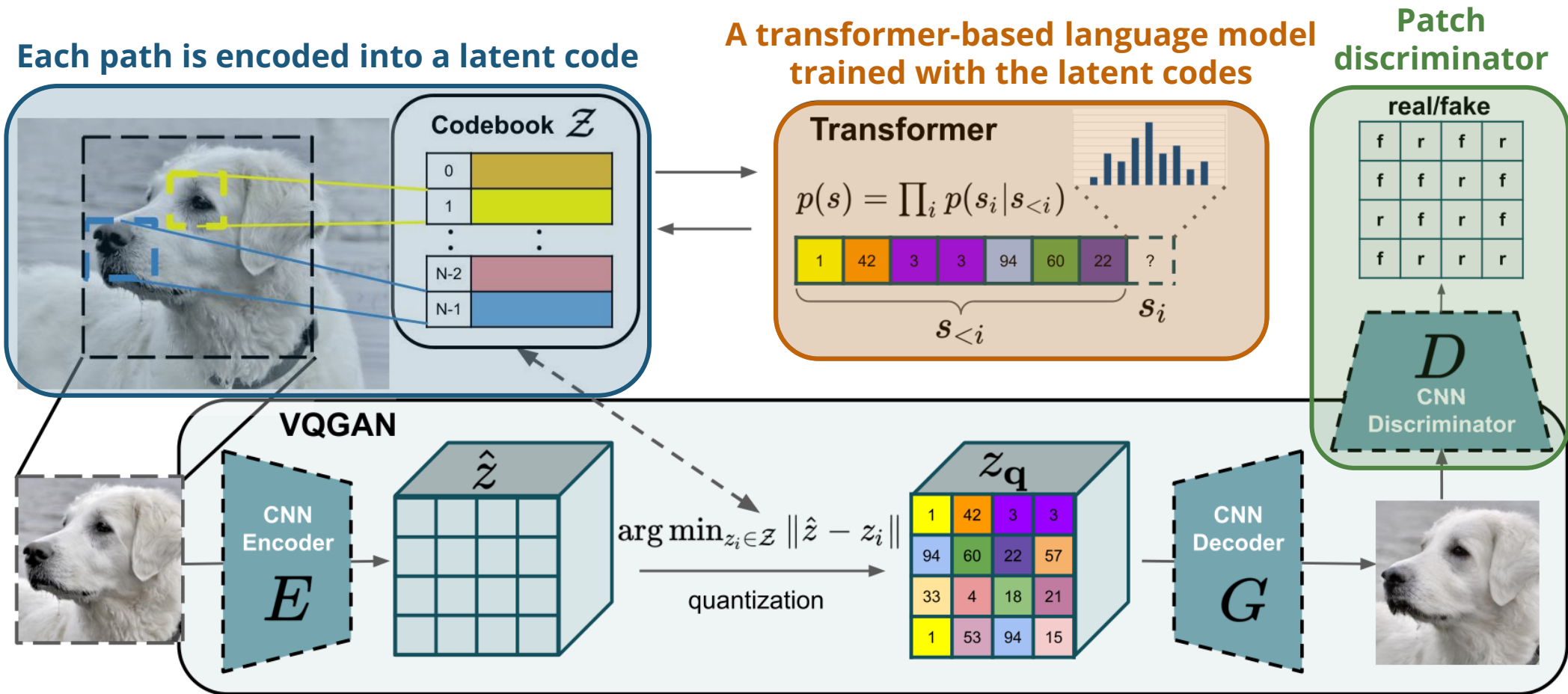
# Latent Diffusion Models

# Vector-Quantized VAE (VQVAE)



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

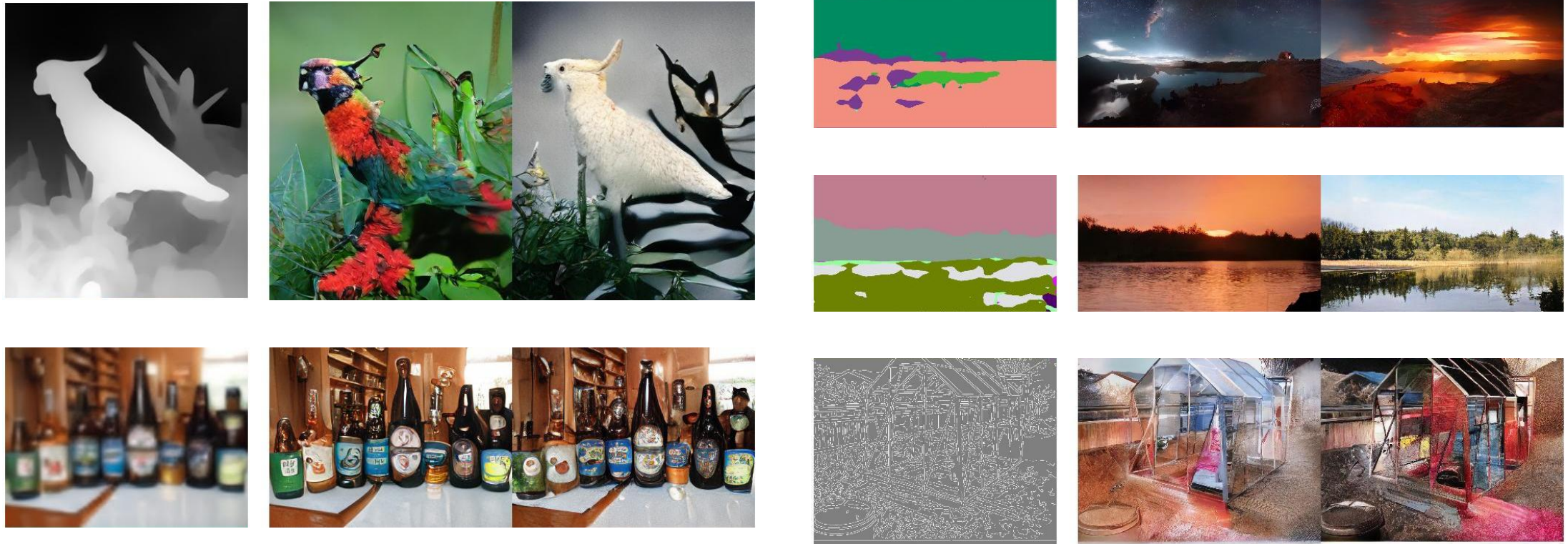
# Vector-Quantized GAN (VQGAN)



(Source: Esser et al., 2021)

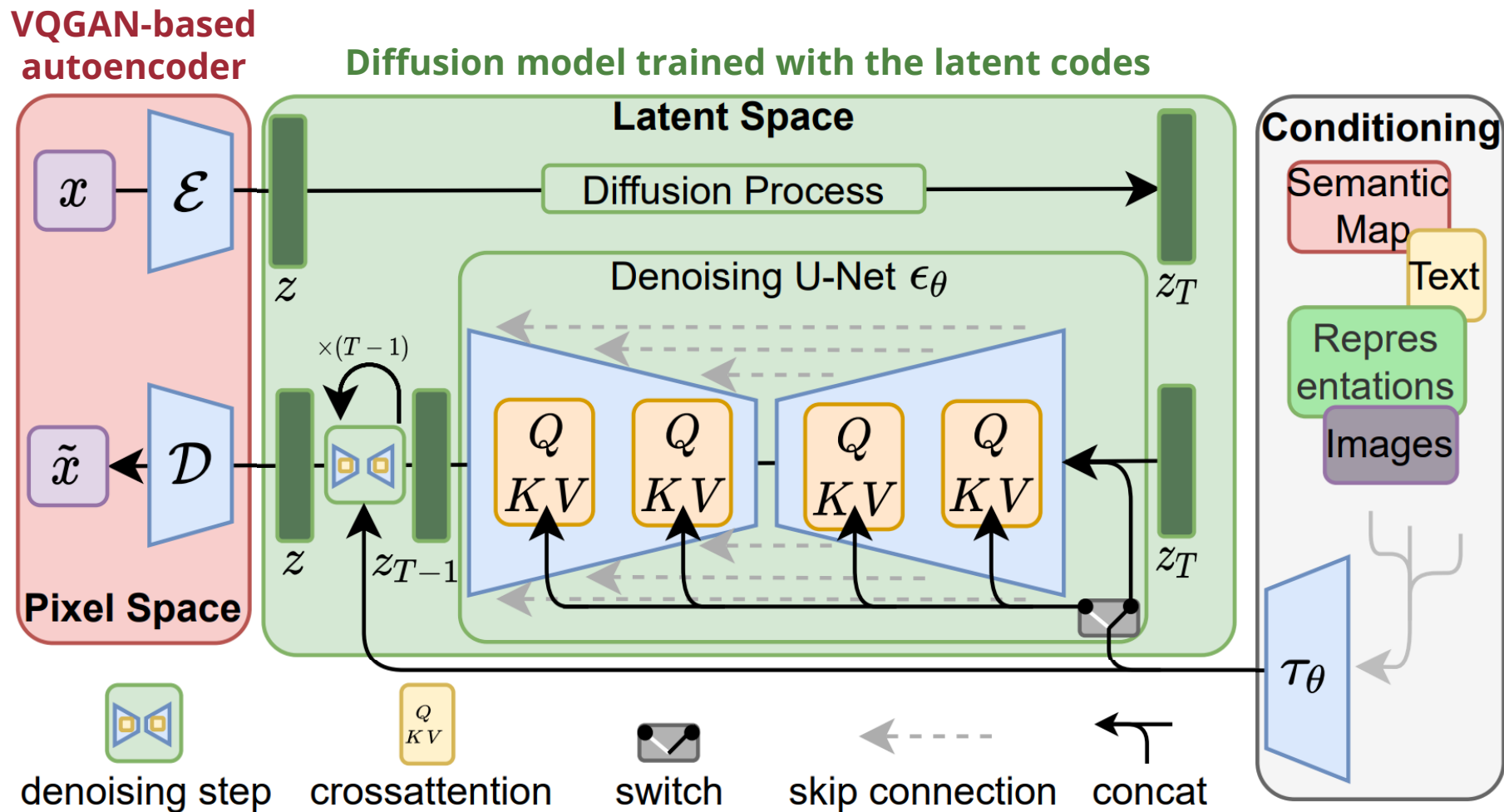
**A VQGAN is a VQVAE equipped with adversarial loss**

# VQGAN: Conditional Generation



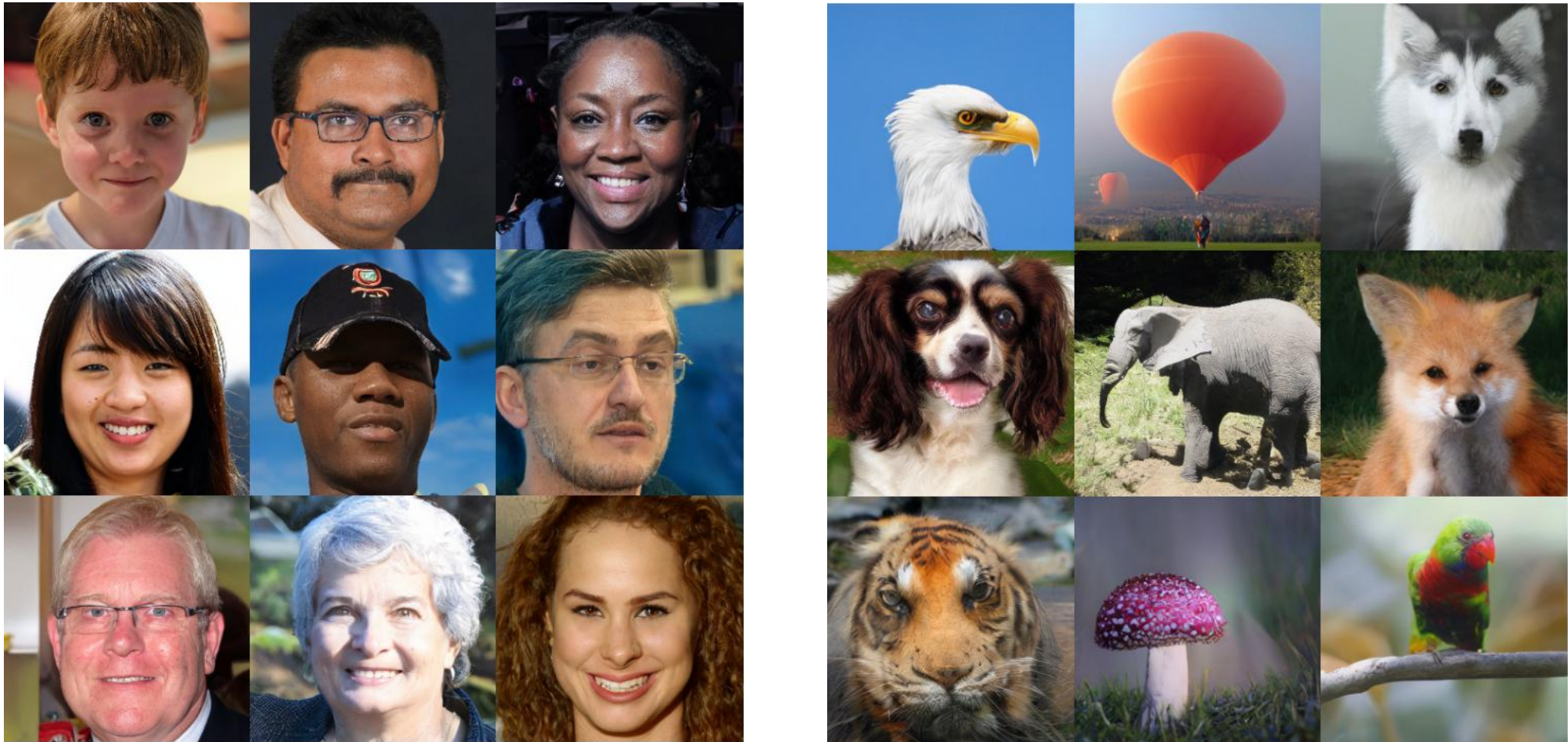
(Source: Esser et al., 2021)

# Latent Diffusion Models (LDMs)



(Source: Rombach et al., 2022)

# Latent Diffusion Models: Unconditional Generation



(Source: Rombach et al., 2022)

# Latent Diffusion Models: Unconditional Generation



(Source: Rombach et al., 2022)

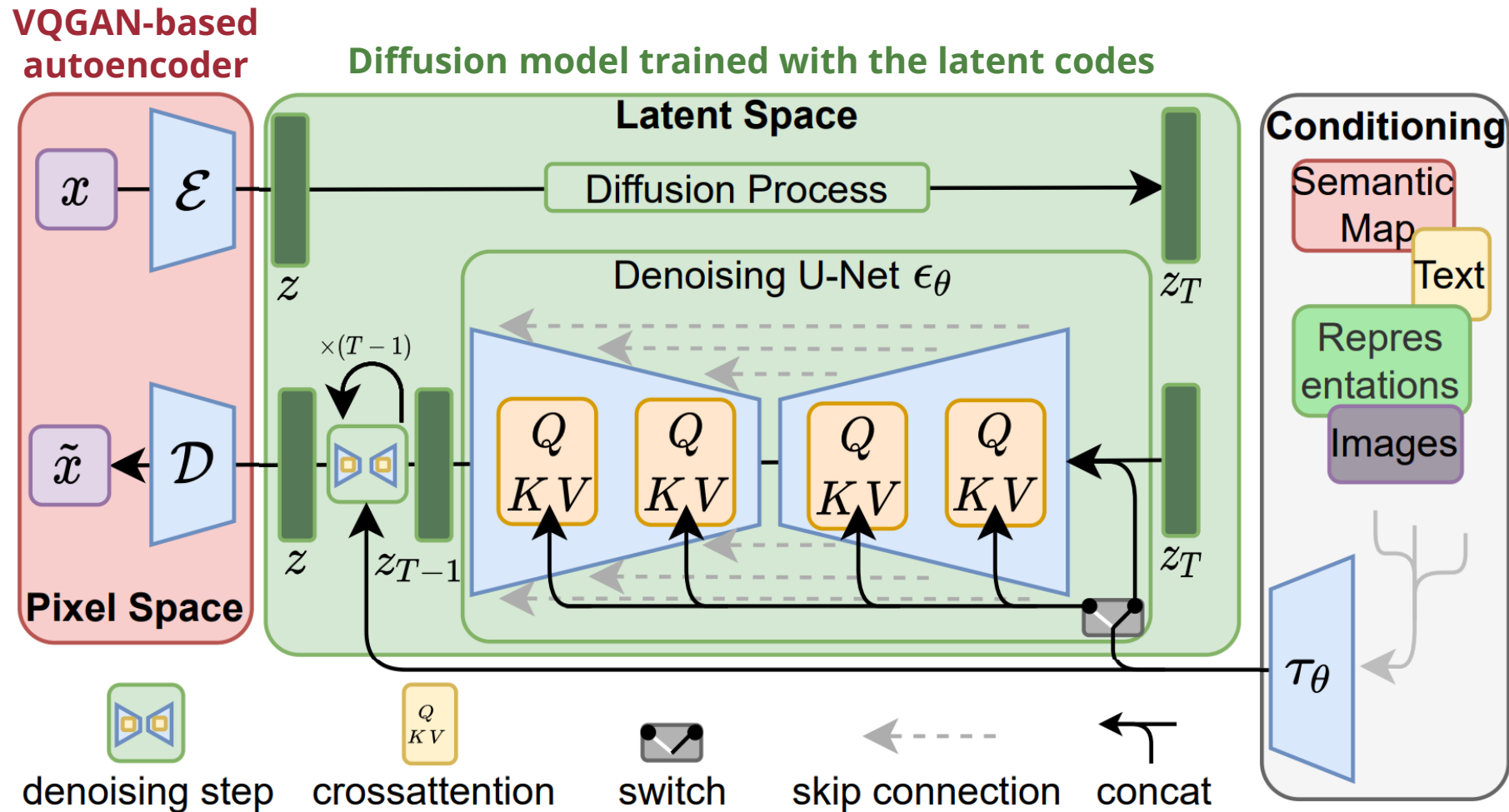
# LDMs: Semantic Synthesis



(Source: Rombach et al., 2022)



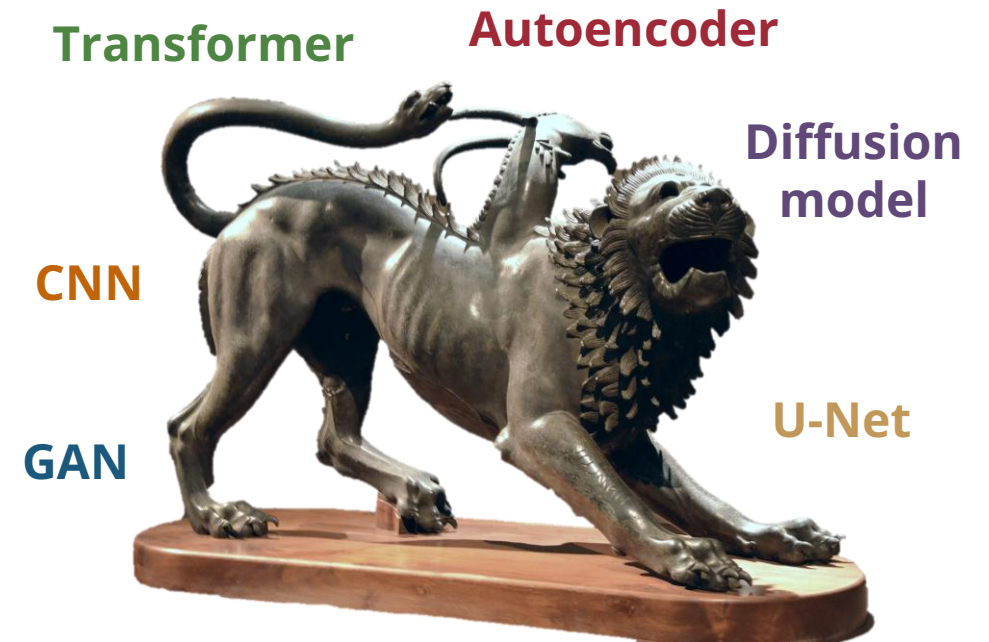
# Latent Diffusion Models (LDMs)



(Source: Rombach et al., 2022)

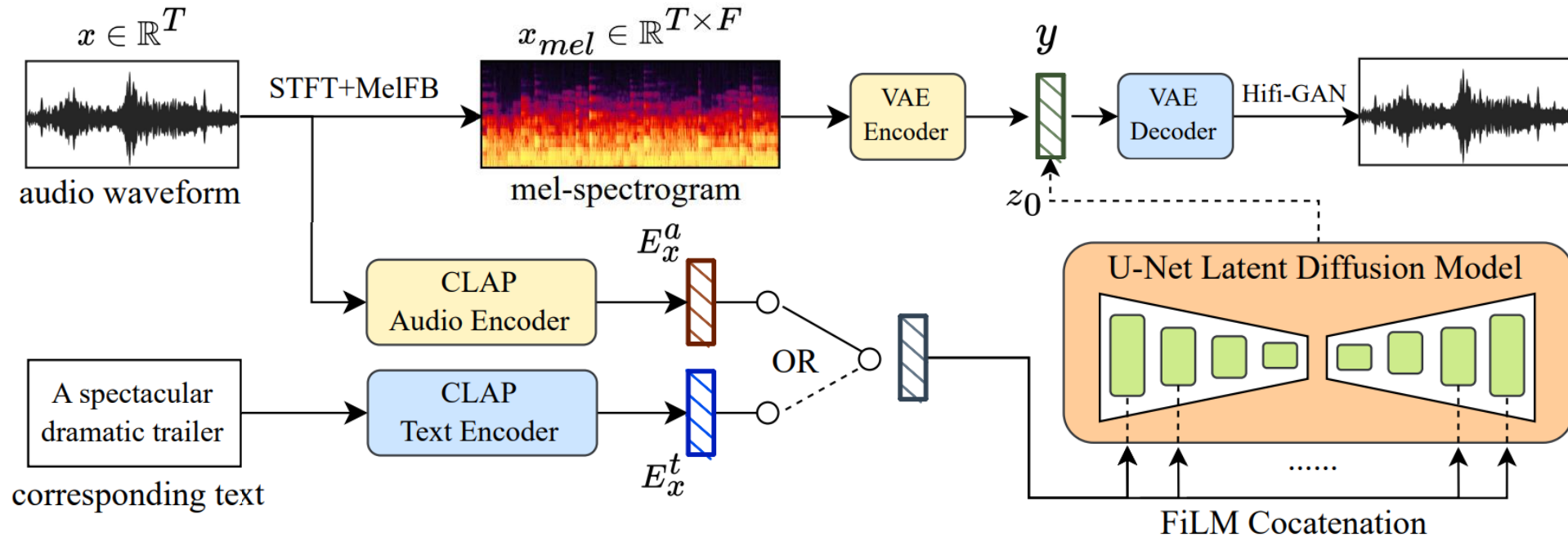
# Latent Diffusion Model is a Chimera

- **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 Diffusion for Music: MusicLDM (Chen et al., 2023)



(Source: Ke et al., 2023)

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

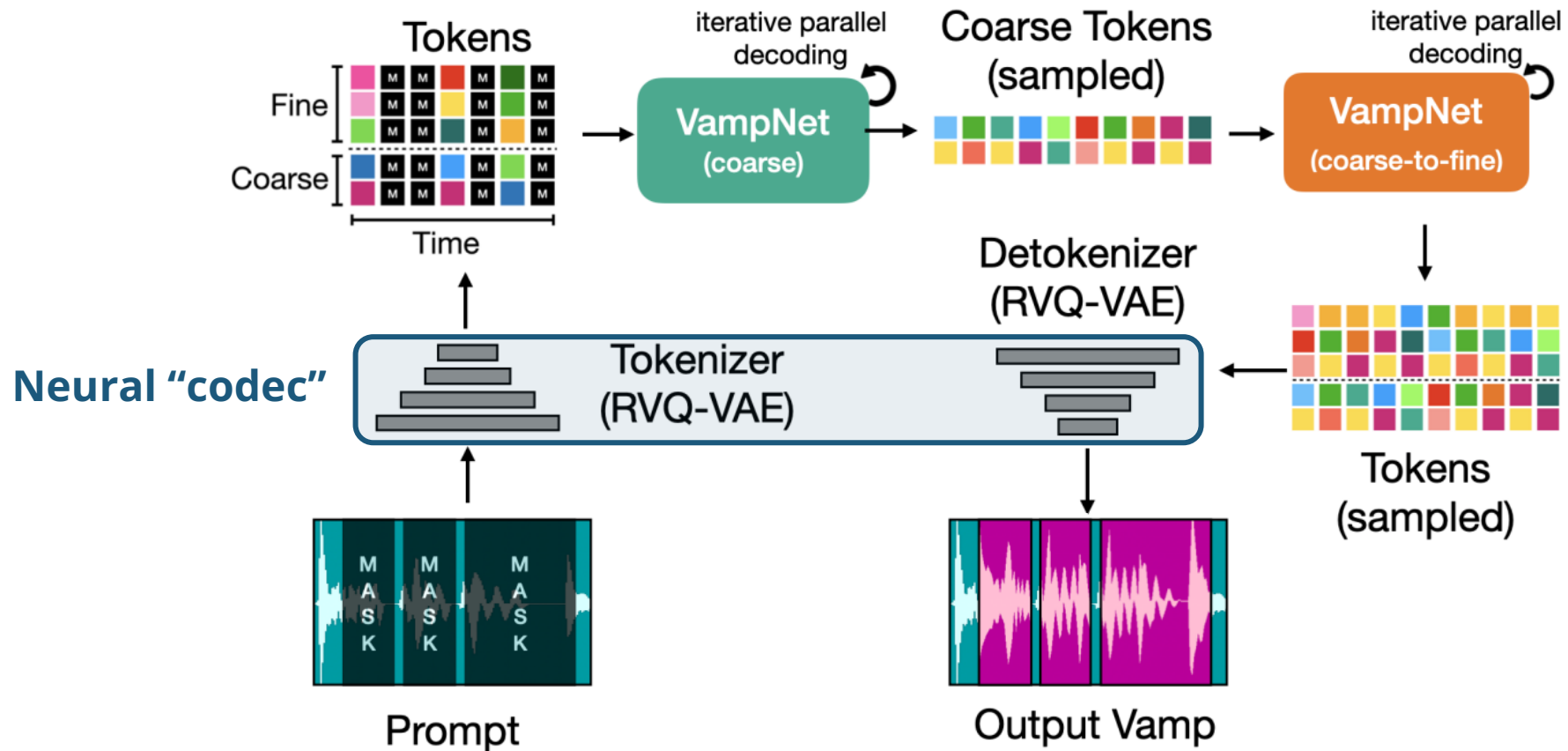
# MusicLDM: Demo (Chen et al., 2023)



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

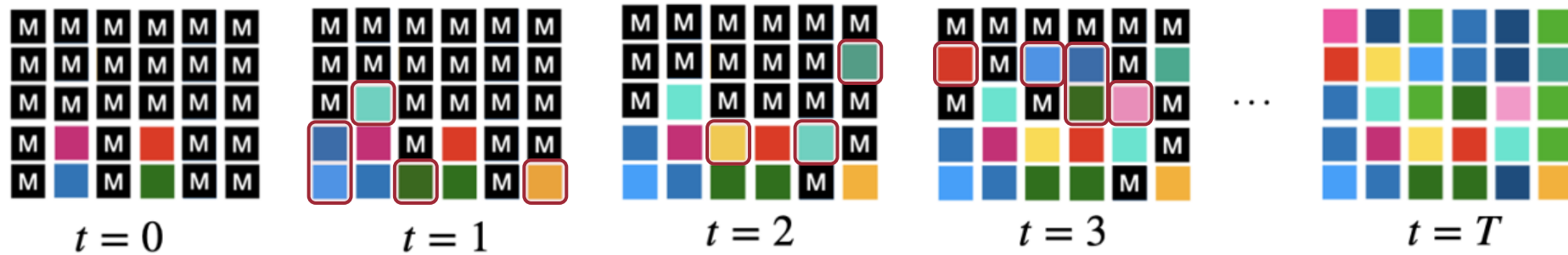
# Masked Acoustic Token Modeling

# VampNet (Garcia et al., 2023)



(Source: Garcia et al., 2023)

# VampNet (Garcia et al., 2023)

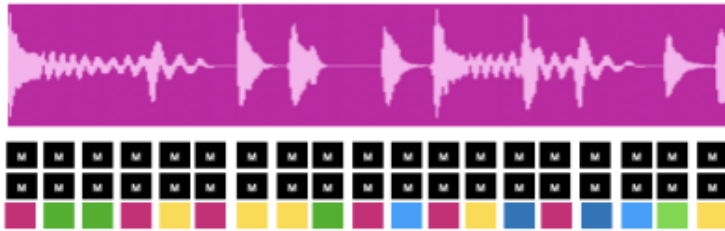


(Source: Garcia et al., 2023)

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

# VampNet (Garcia et al., 2023)

Compression



Beat Driven

= predicted beat mark



Periodic

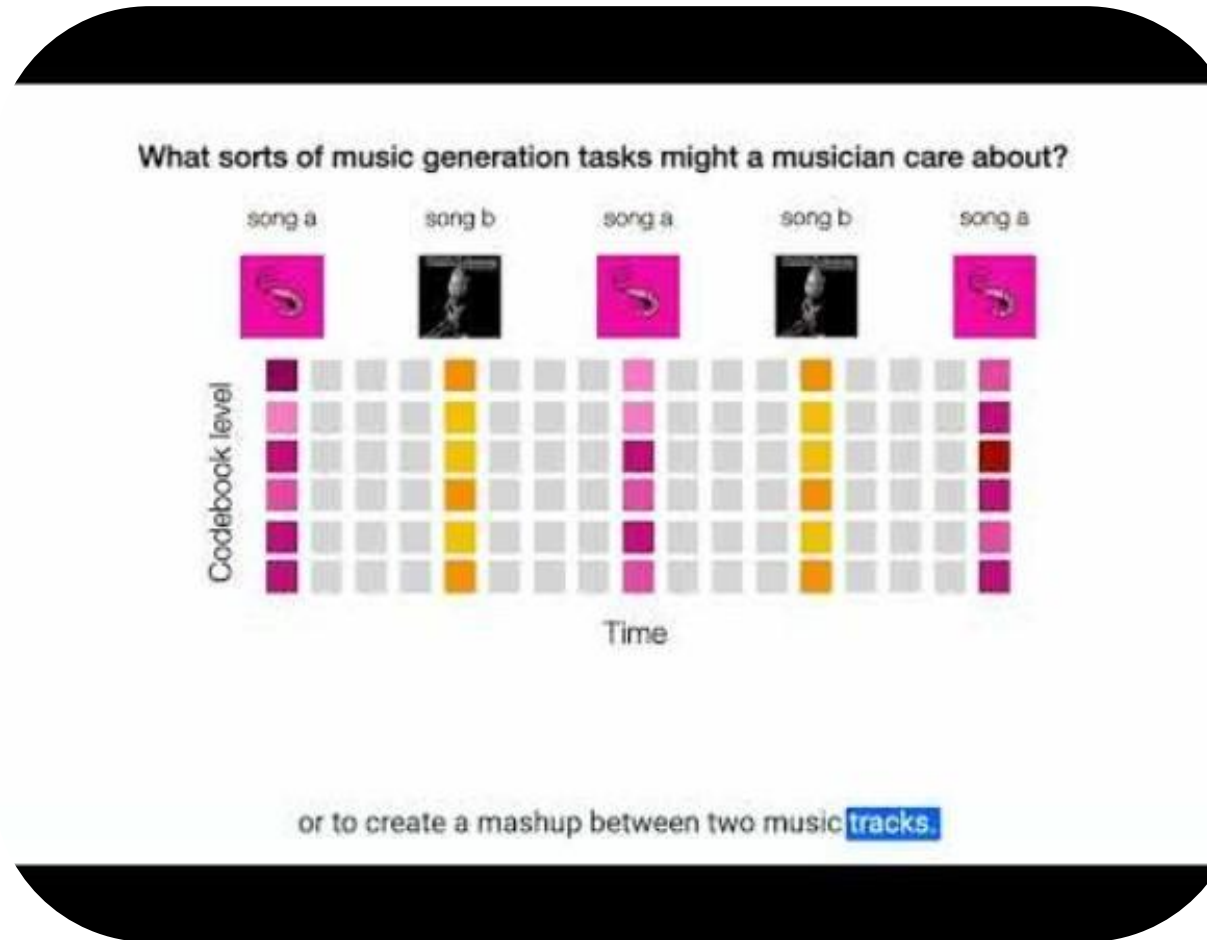


Inpainting



(Source: Garcia et al., 2023)

# VampNet (Garcia et al., 2023)



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

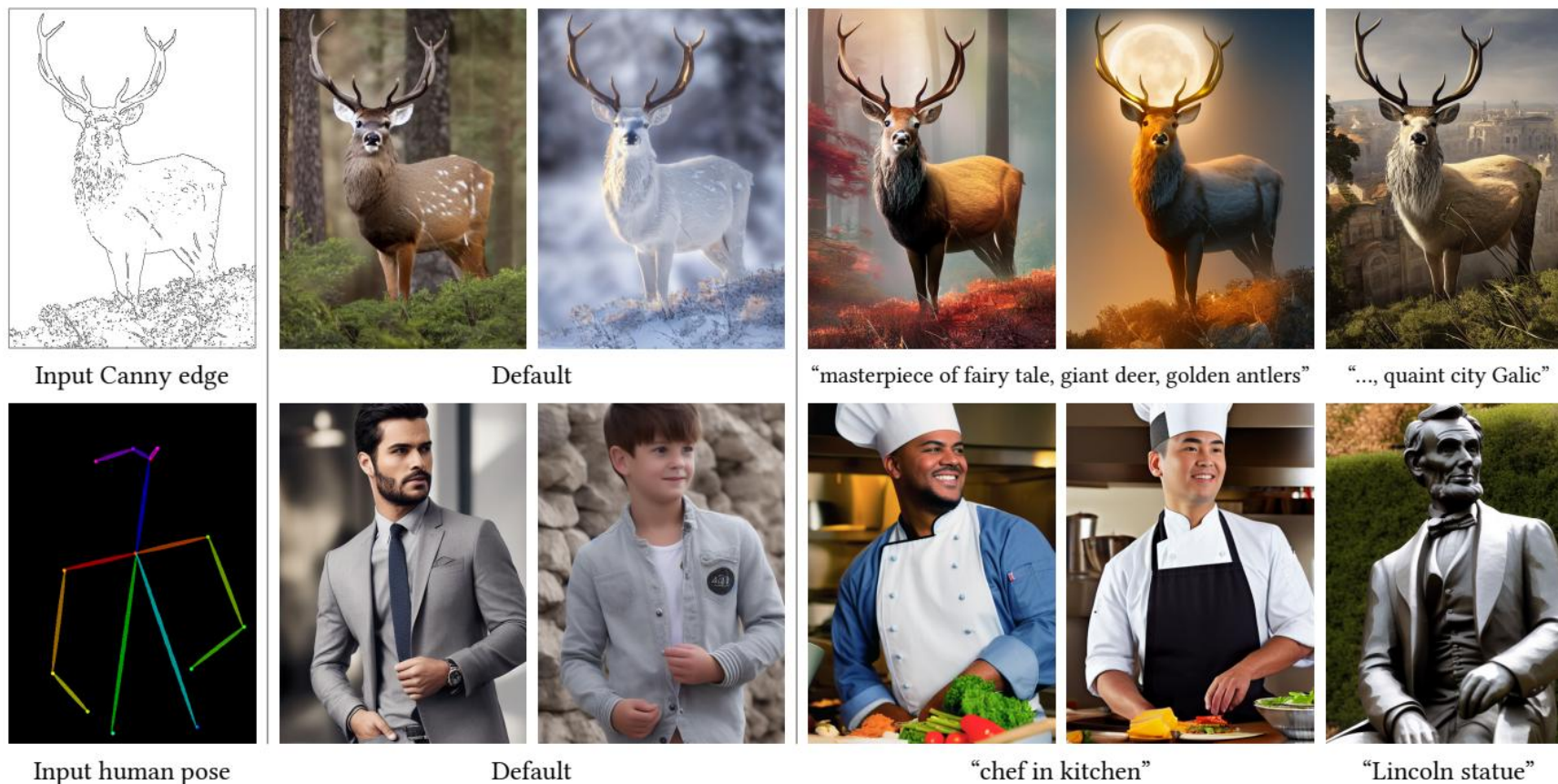
# unloop (Garcia et al., 2023)



[youtu.be/yzBl8Vcjd2s](https://youtu.be/yzBl8Vcjd2s) & [github.com/hugofloresgarcia/unloop](https://github.com/hugofloresgarcia/unloop)

# Controlling Music Generation Systems

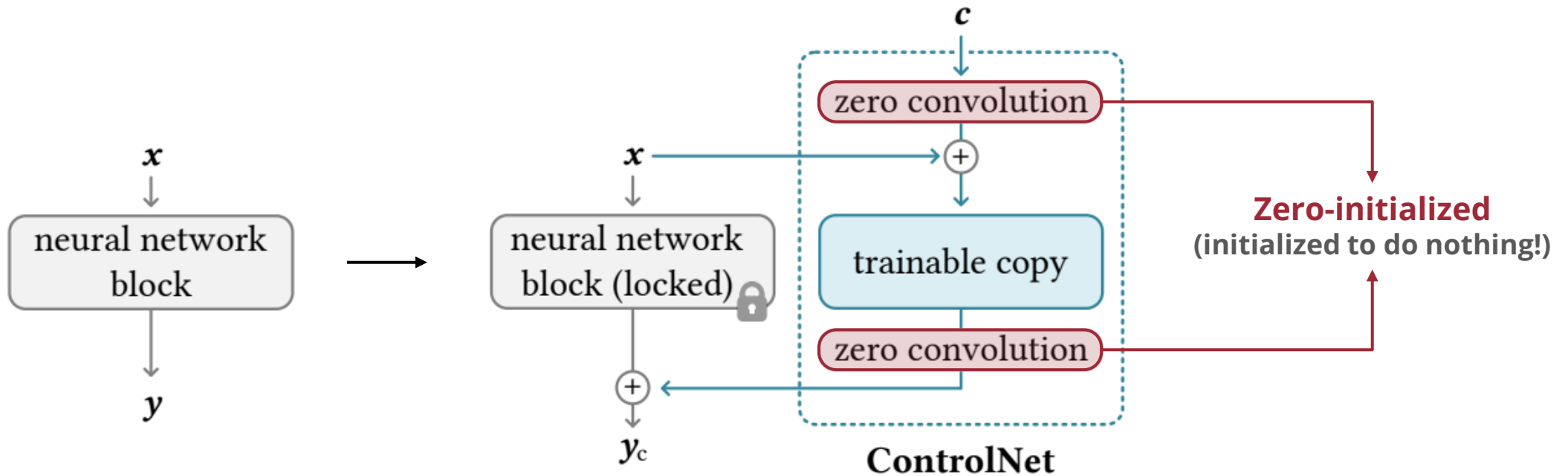
# ControlNet (Zhang et al., 2023)



(Source: Zhang et al., 2023)

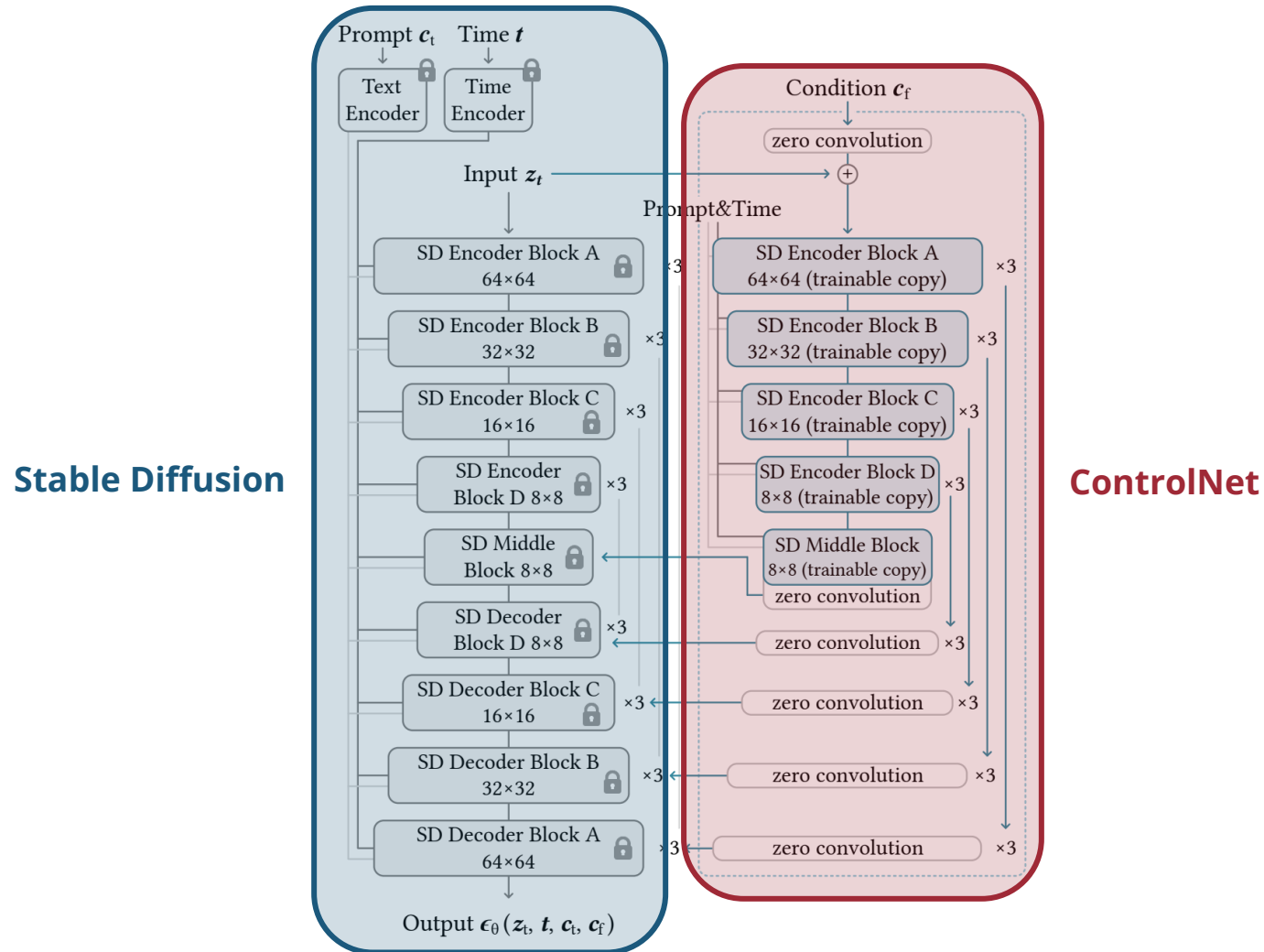
Can we **add controls** to a trained text-to-image diffusion model?

# ControlNet (Zhang et al., 2023)



(Source: Zhang et al., 2023)

# ControlNet (Zhang et al., 2023)



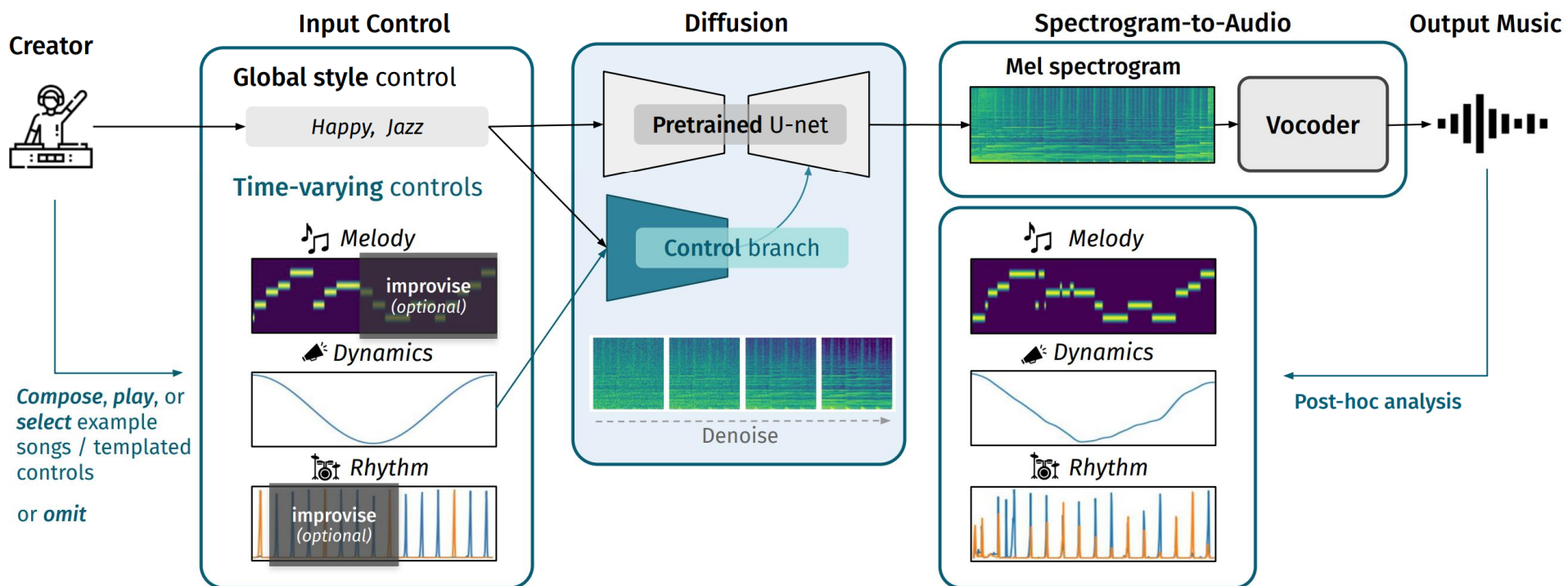
(Source: Zhang et al., 2023)

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



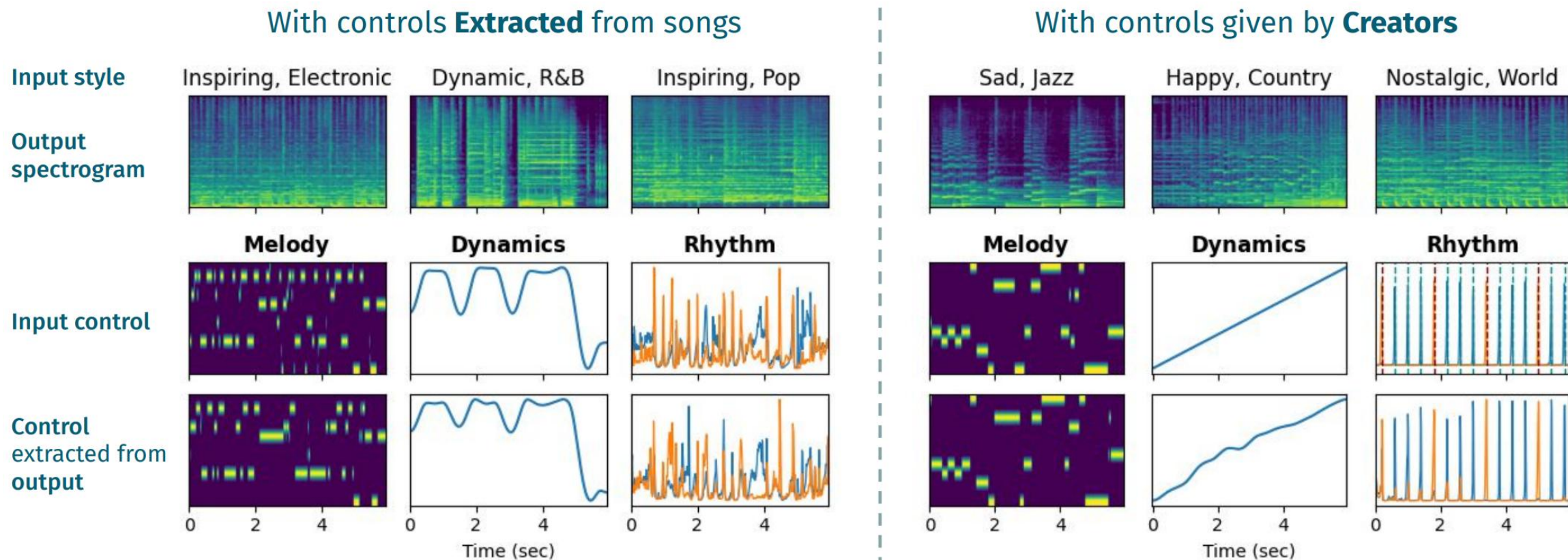
[youtu.be/O4cJ3acEGDw](https://youtu.be/O4cJ3acEGDw)

# Music ControlNet (Wu et al., 2024)



(Source: Wu et al., 2024)

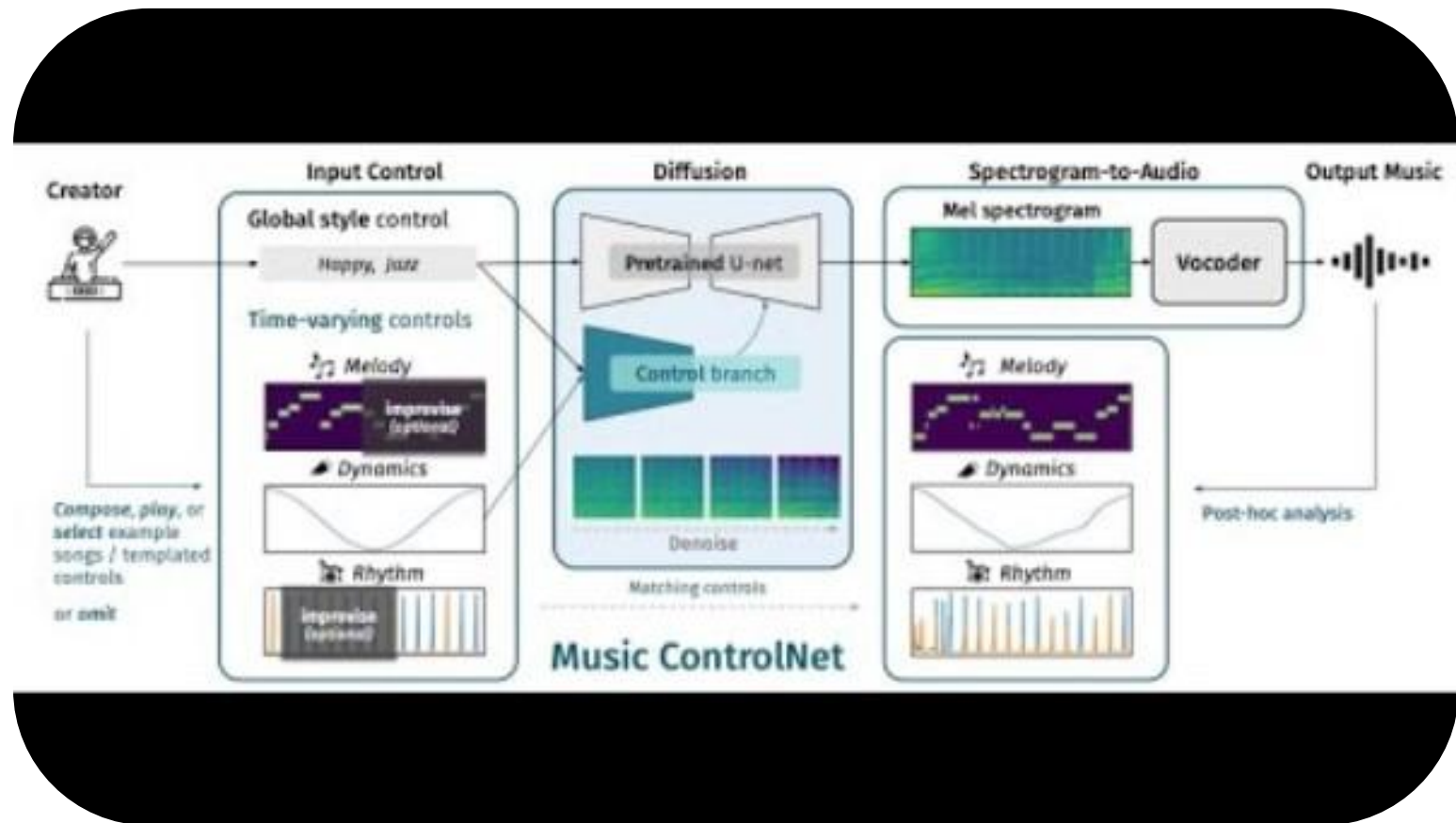
# Music ControlNet (Wu et al., 2024)



(Source: Wu et al., 2024)

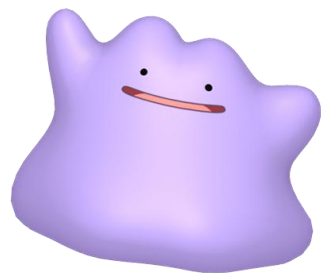
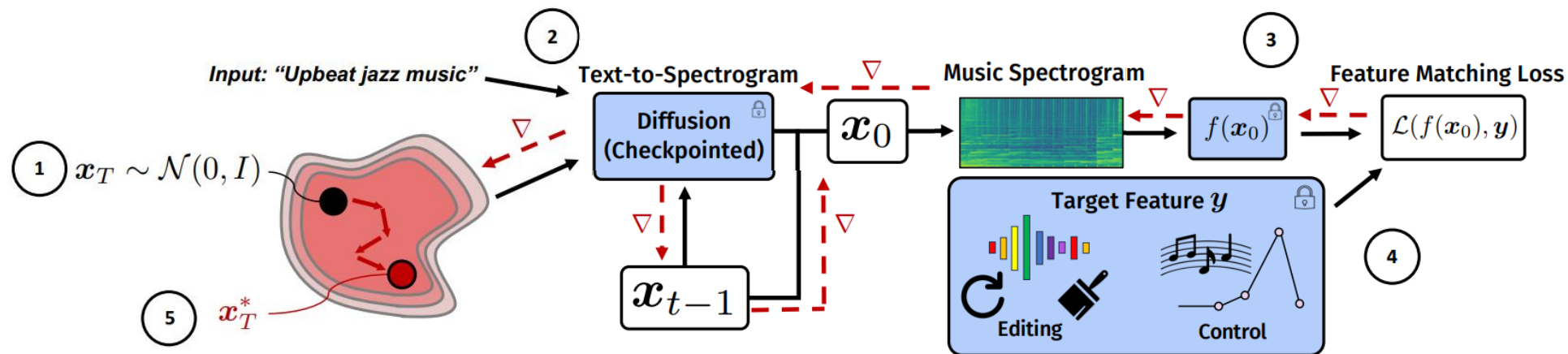
[musiccontrolnet.github.io/web](https://musiccontrolnet.github.io/web)

# Music ControlNet (Wu et al., 2024)



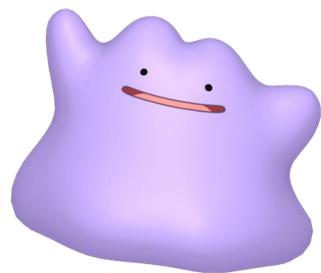
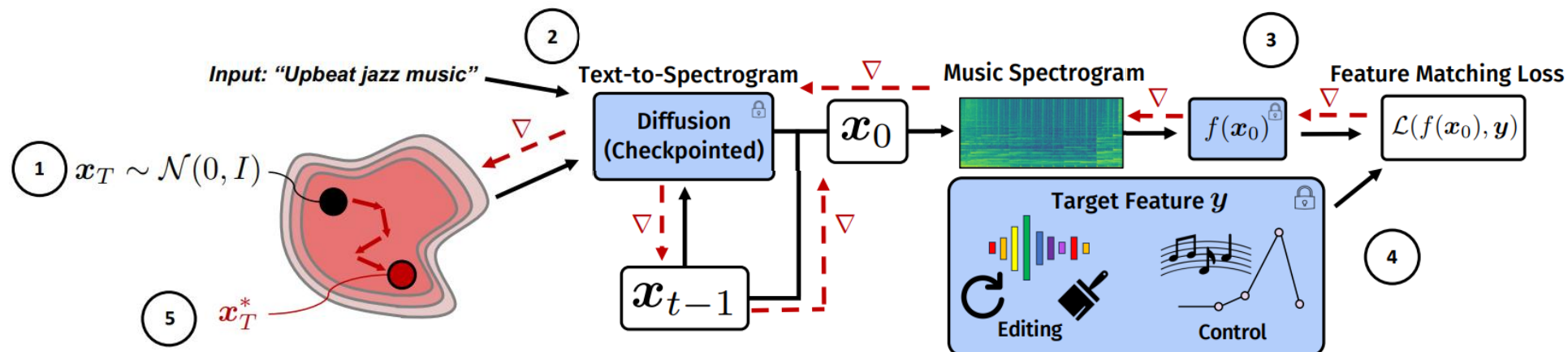
[youtu.be/QVr-S-DyccU](https://youtu.be/QVr-S-DyccU)

# Inference-time Control: DITTO (Novack et al., 2024)



(Source: Novack et al., 2024)

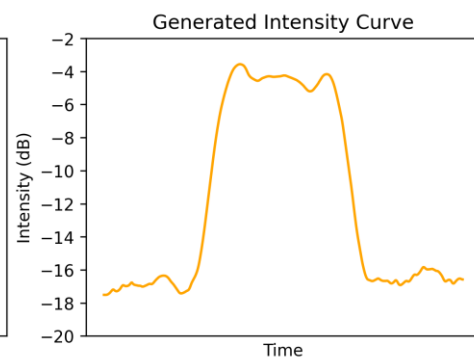
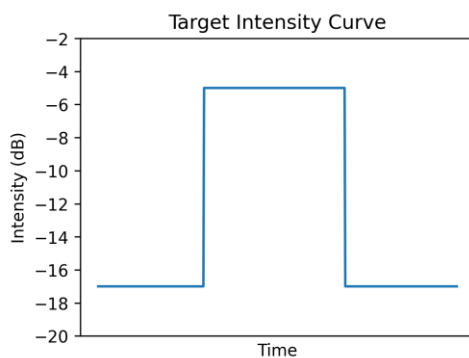
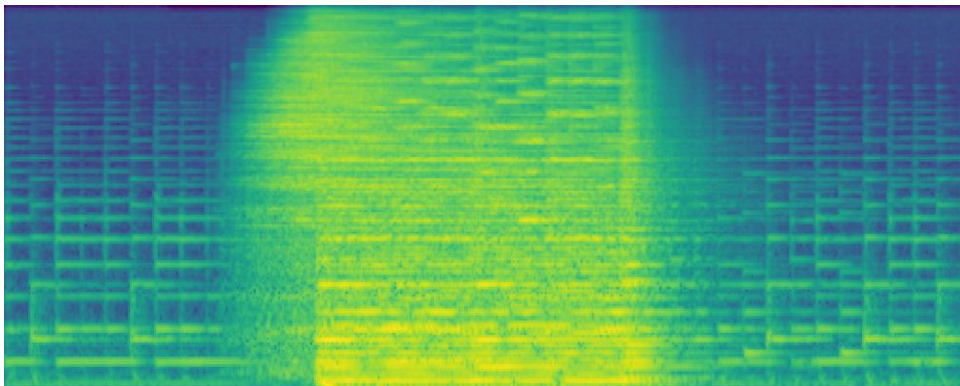
# Inference-time Control: DITTO (Novack et al., 2024)



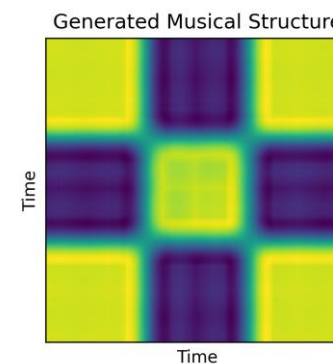
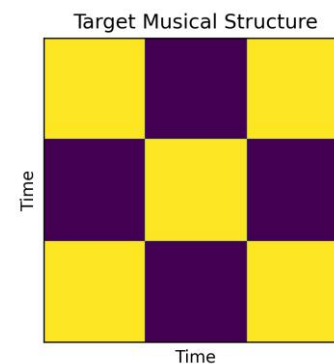
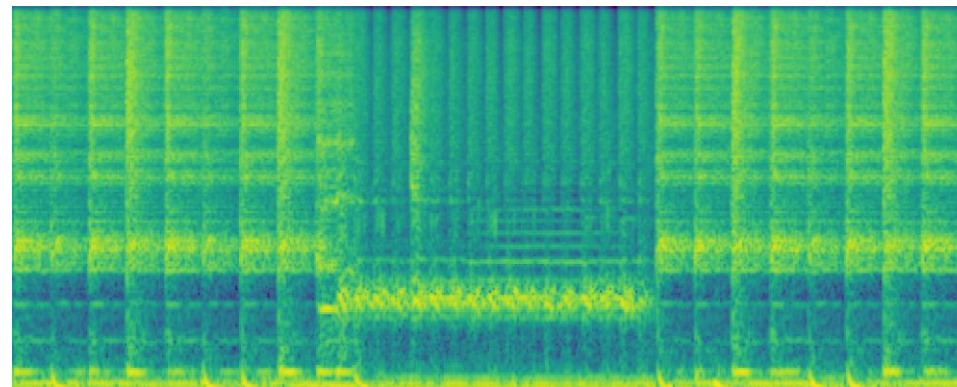
(Source: Novack et al., 2024)

# Inference-time Control: DITTO (Novack et al., 2024)

## Intensity control



## Structure control



(Source: Novack et al., 2024)

# Inference-time Control: DITTO (Novack et al., 2024)

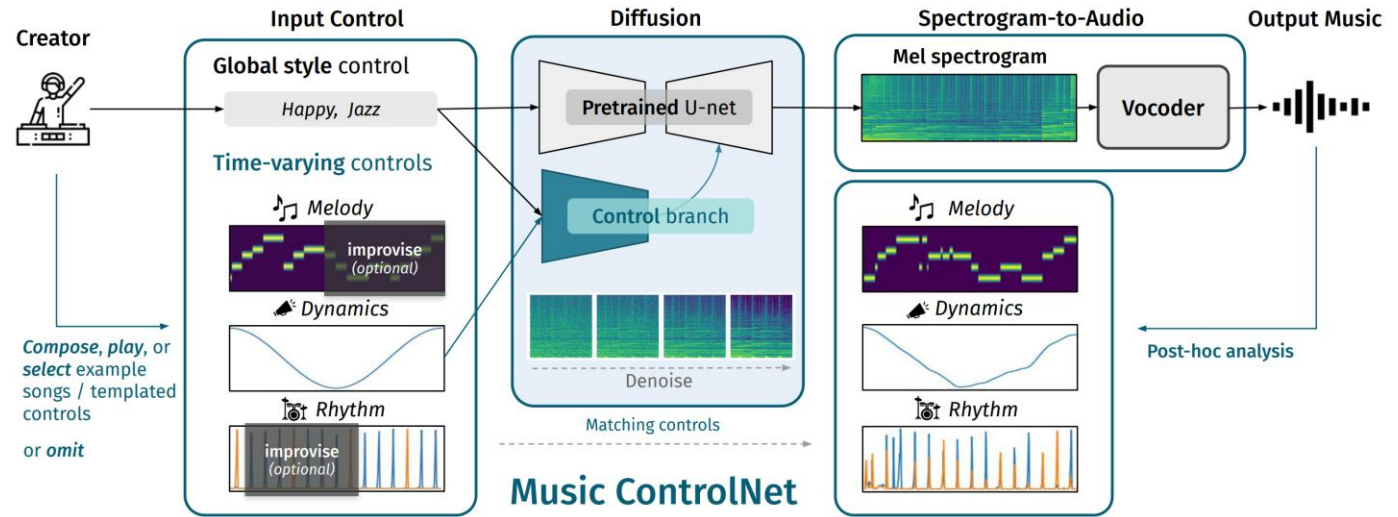


[youtu.be/KooosSNPNo8](https://youtu.be/KooosSNPNo8) & [ditto-music.github.io/web/](https://ditto-music.github.io/web/)

# Music ControlNet vs DITTO

## Music ControlNet

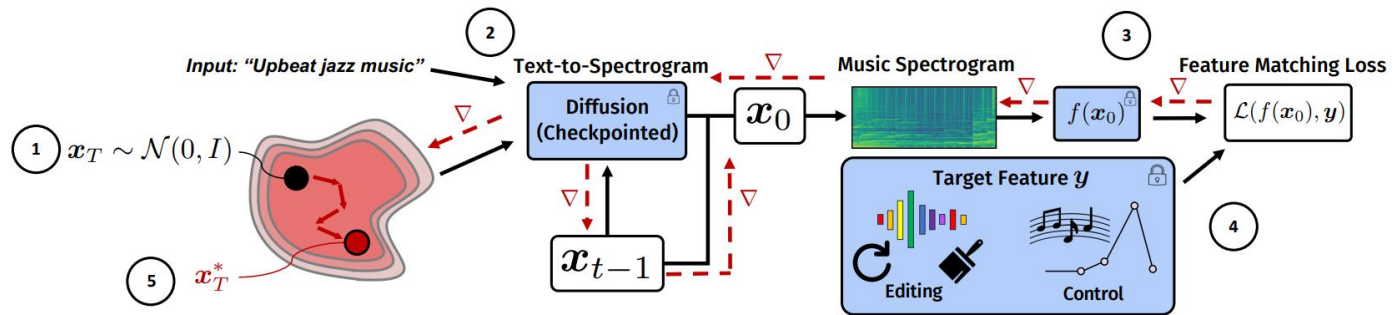
Needs some training!



(Source: Wu et al., 2024)

## DITTO

No training needed!

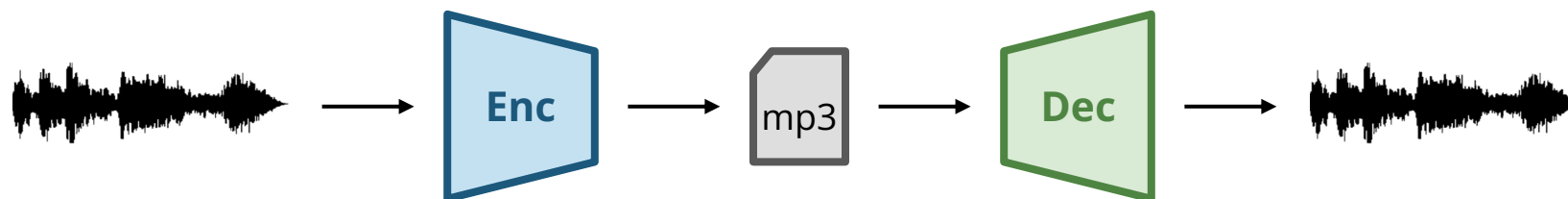


(Source: Novack et al., 2024)

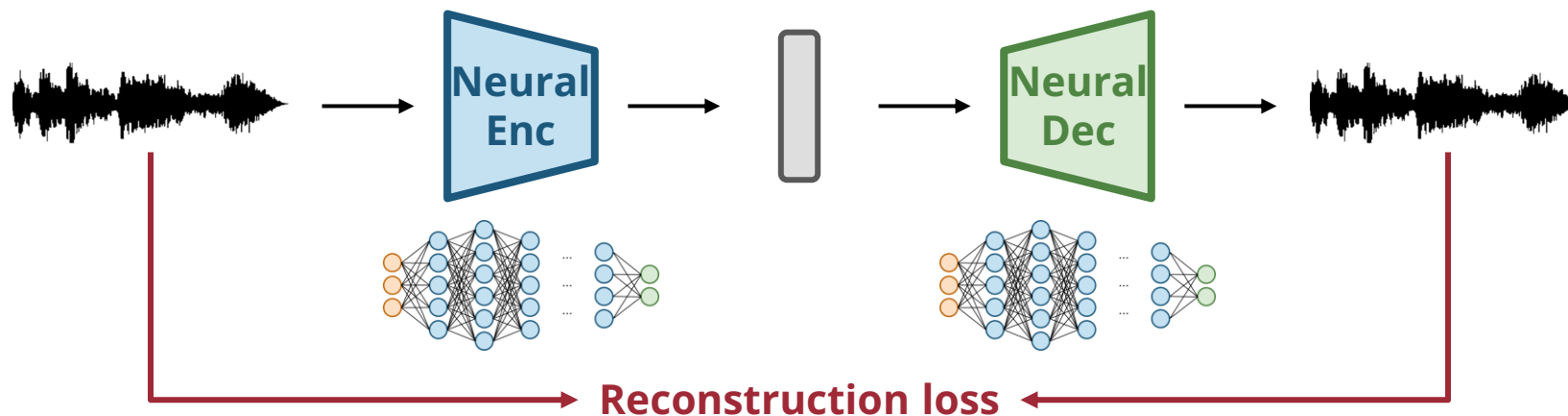
# Recap

# Neural Codec

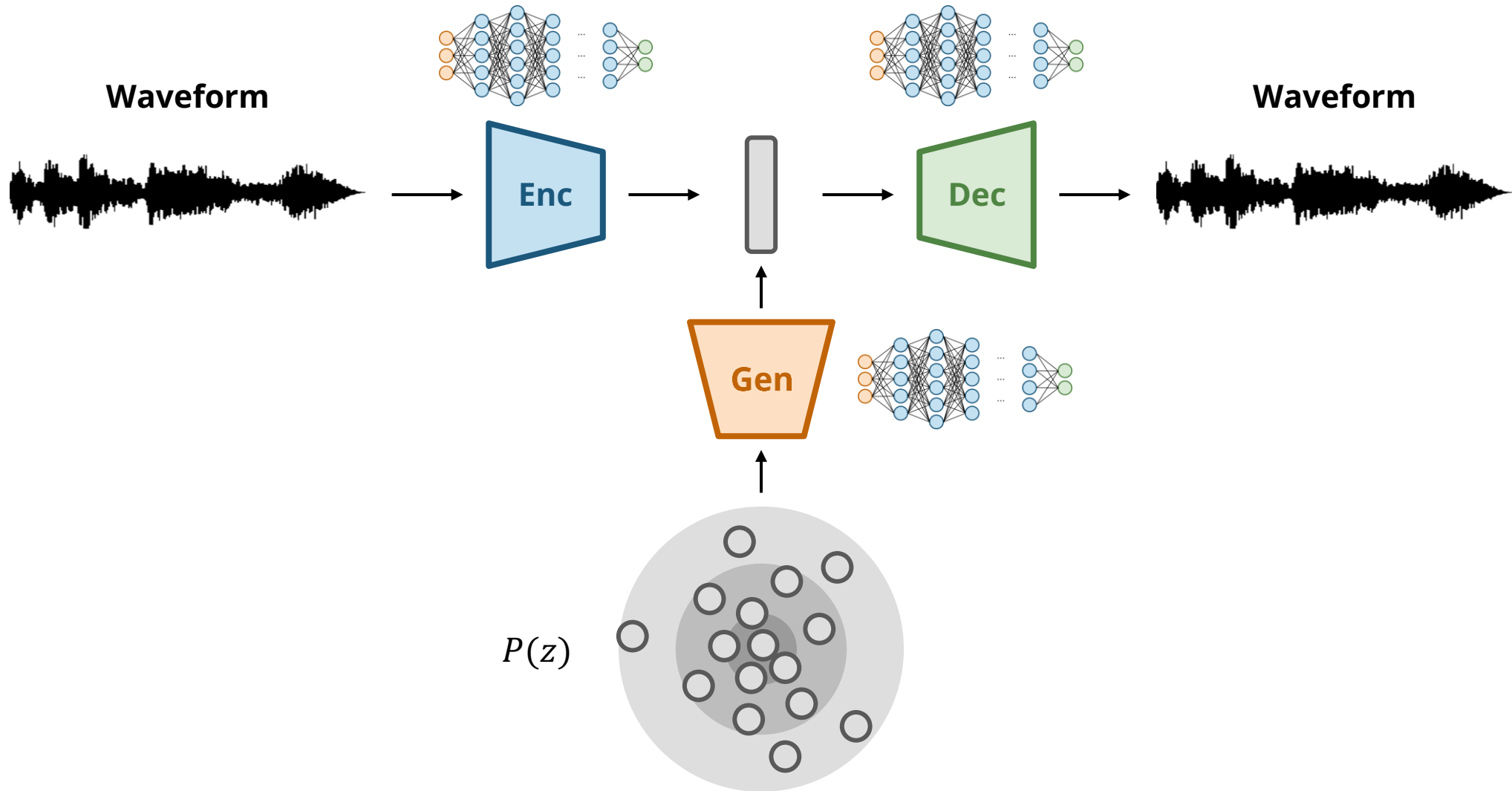
## Traditional Codec



## Neural Codec

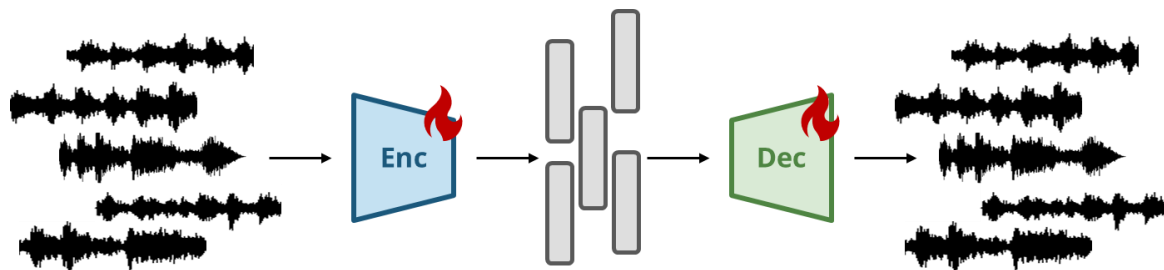


# Latent-based Audio Synthesis

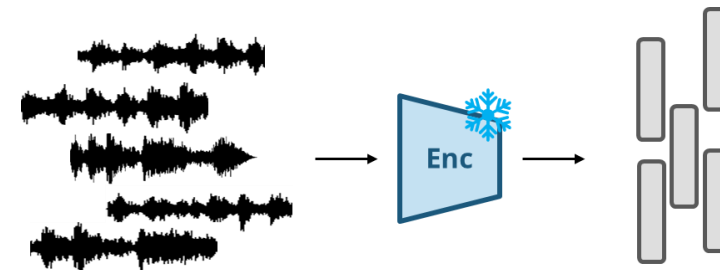


# Pipeline

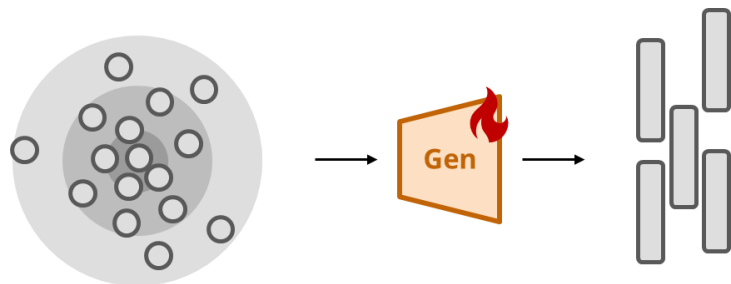
## Step 1: Train an Autoencoder



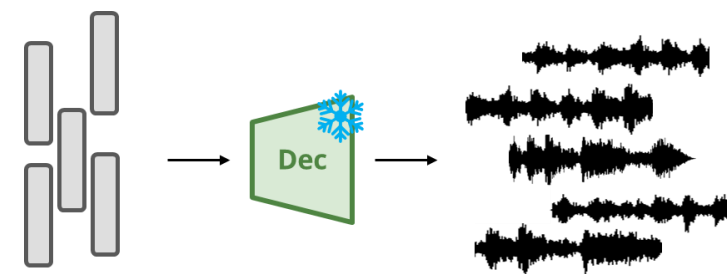
## Step 2: Compute the Latent Vectors



## Step 3: Train a Latent Generative Model

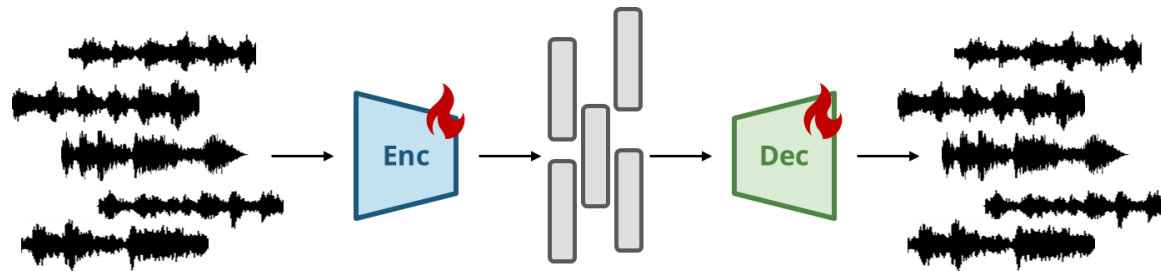


## Step 4: Decode the Latent Vectors

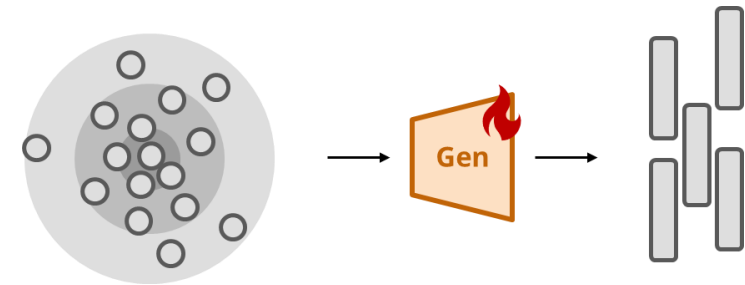


# Training

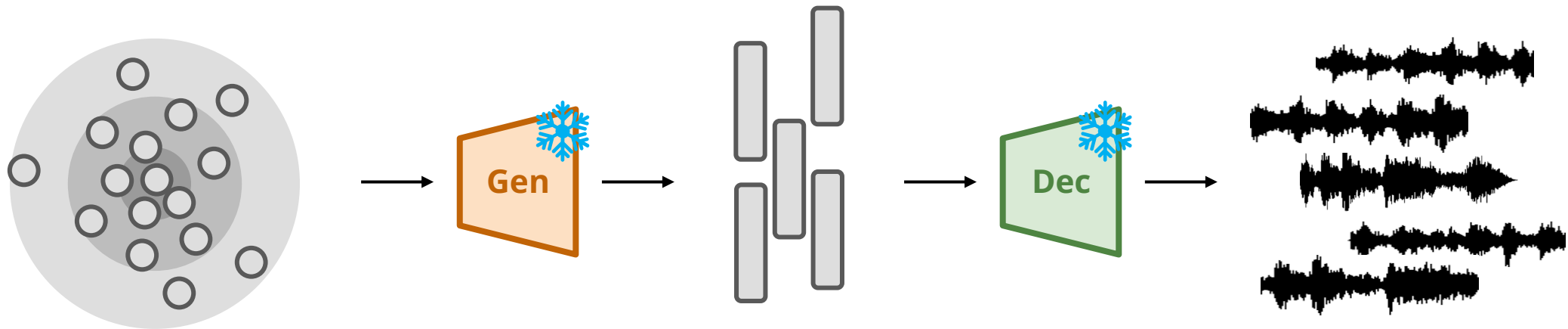
Autoencoder



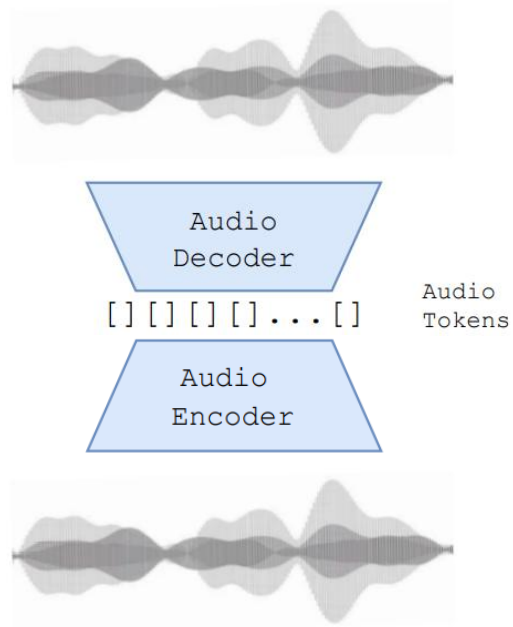
Latent Generative Model



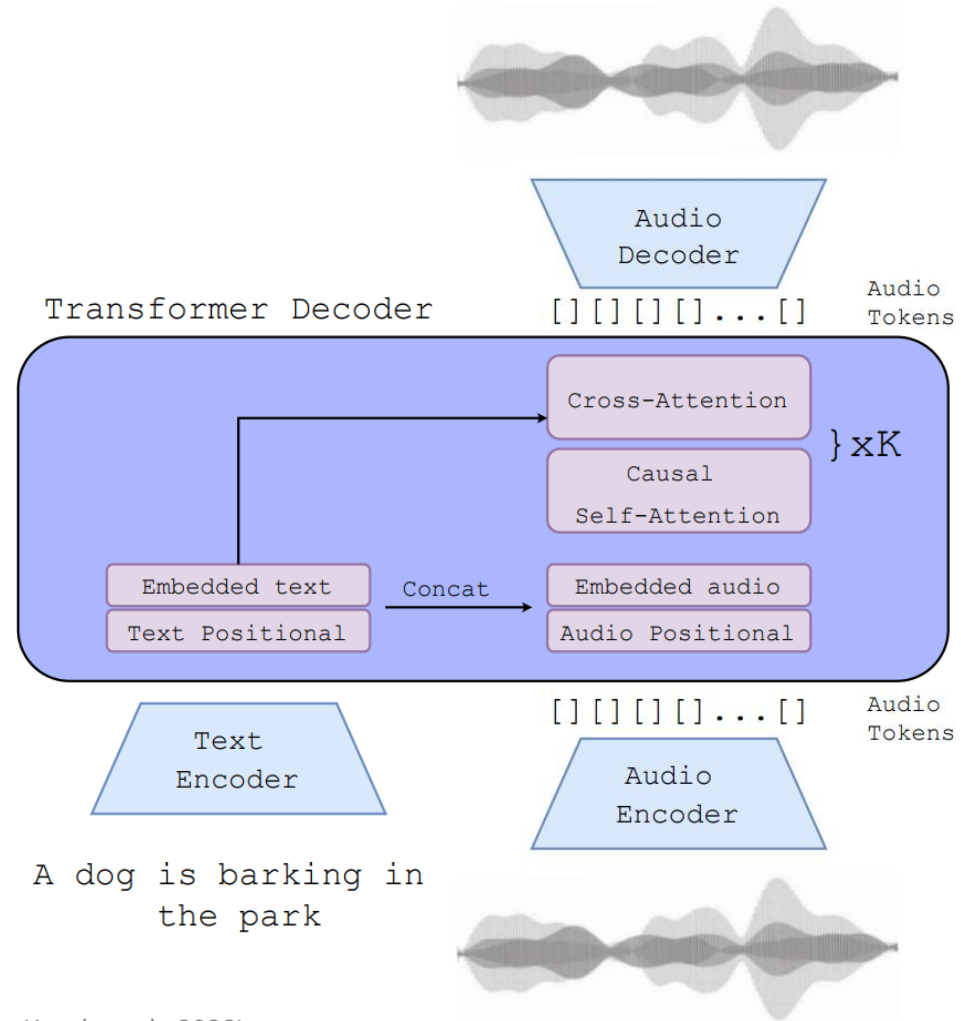
# Inference



# AudioGen (Kreuk et al., 2023)

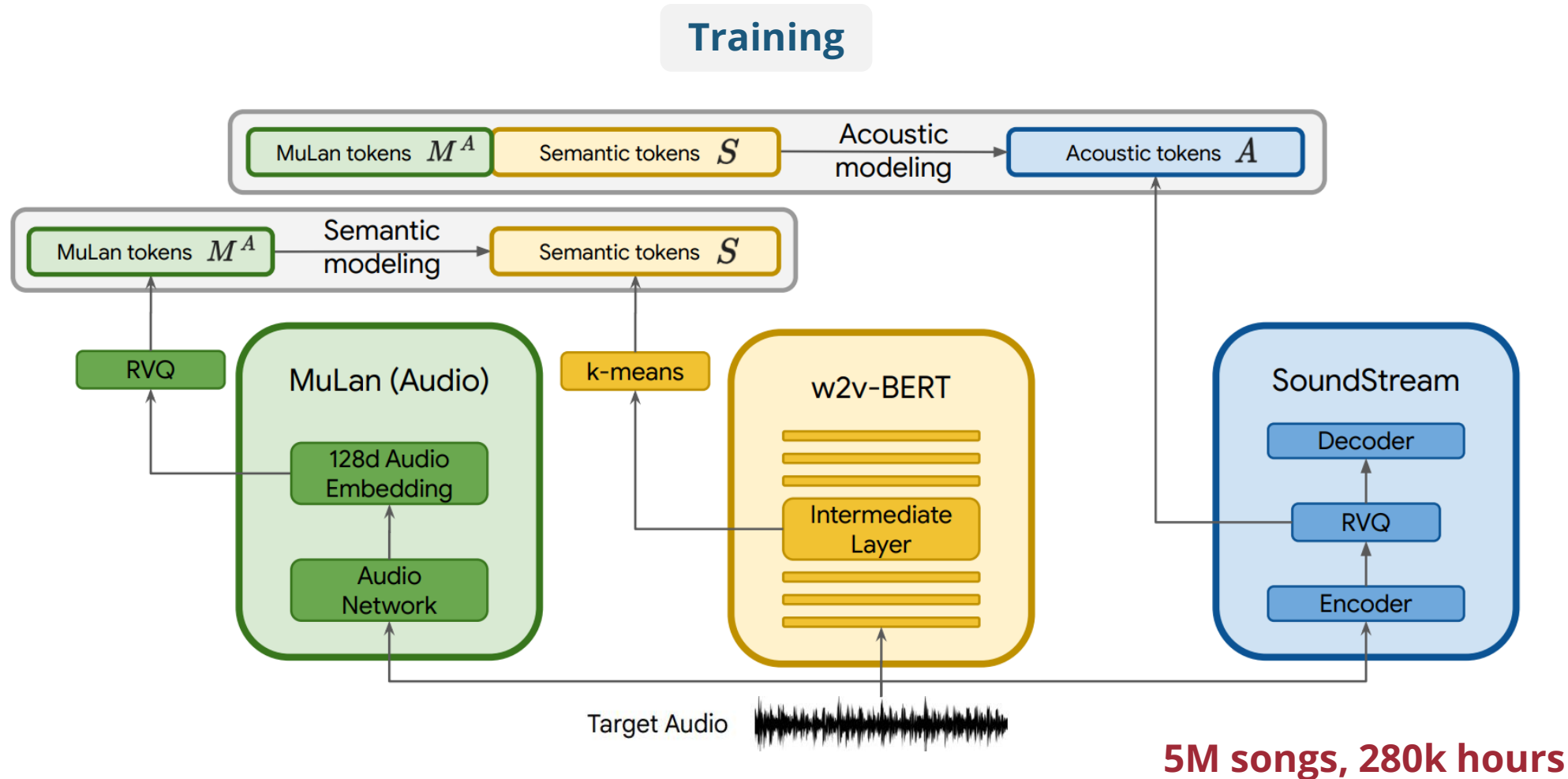


**4k hours**  
**(speech, music, sound effects)**



(Source: Kreuk et al., 2022)

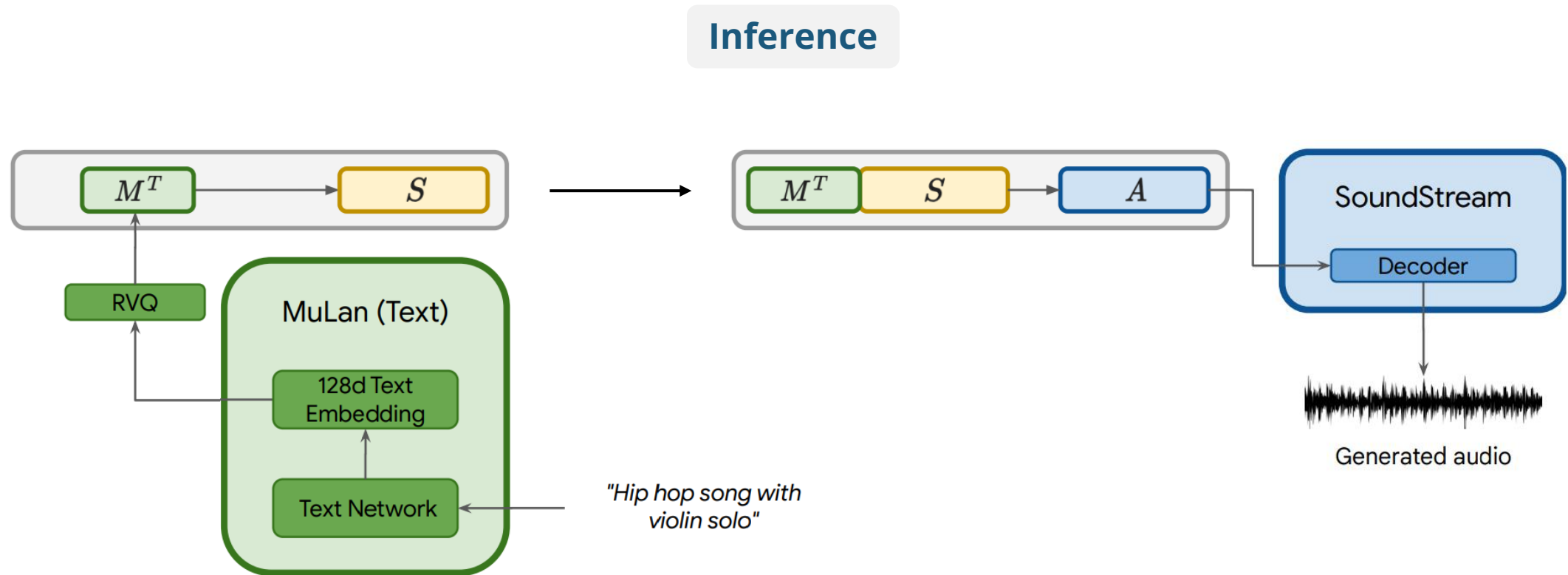
# MusicLM (Agostinelli et al., 2023)



(Source: Agostinelli et al., 2022)

Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Christian Frank, "MusicLM: Generating Music From Text," *arXiv preprint arXiv:2301.11325*, 2023.

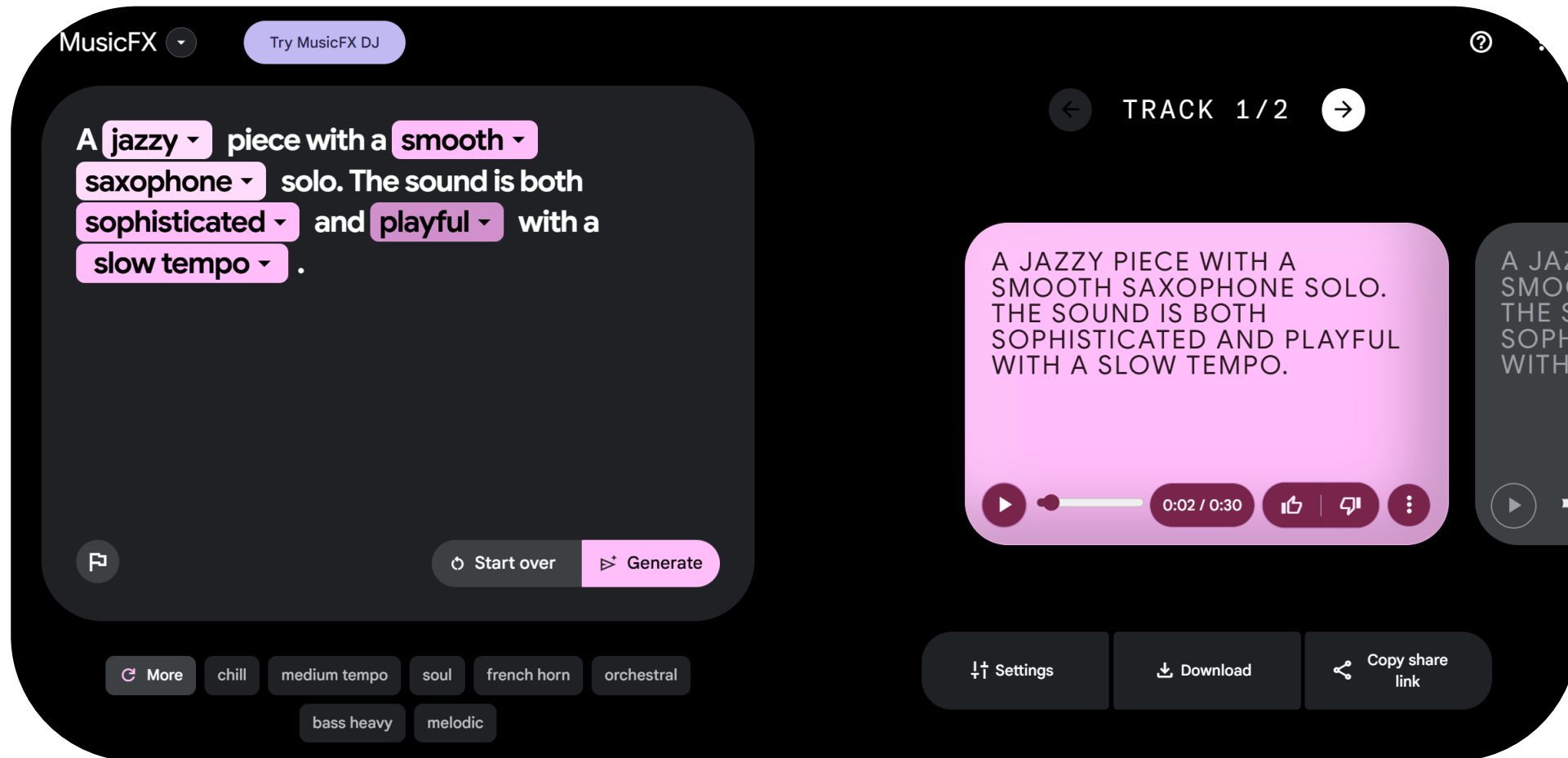
# MusicLM (Agostinelli et al., 2023)



(Source: Agostinelli et al., 2022)

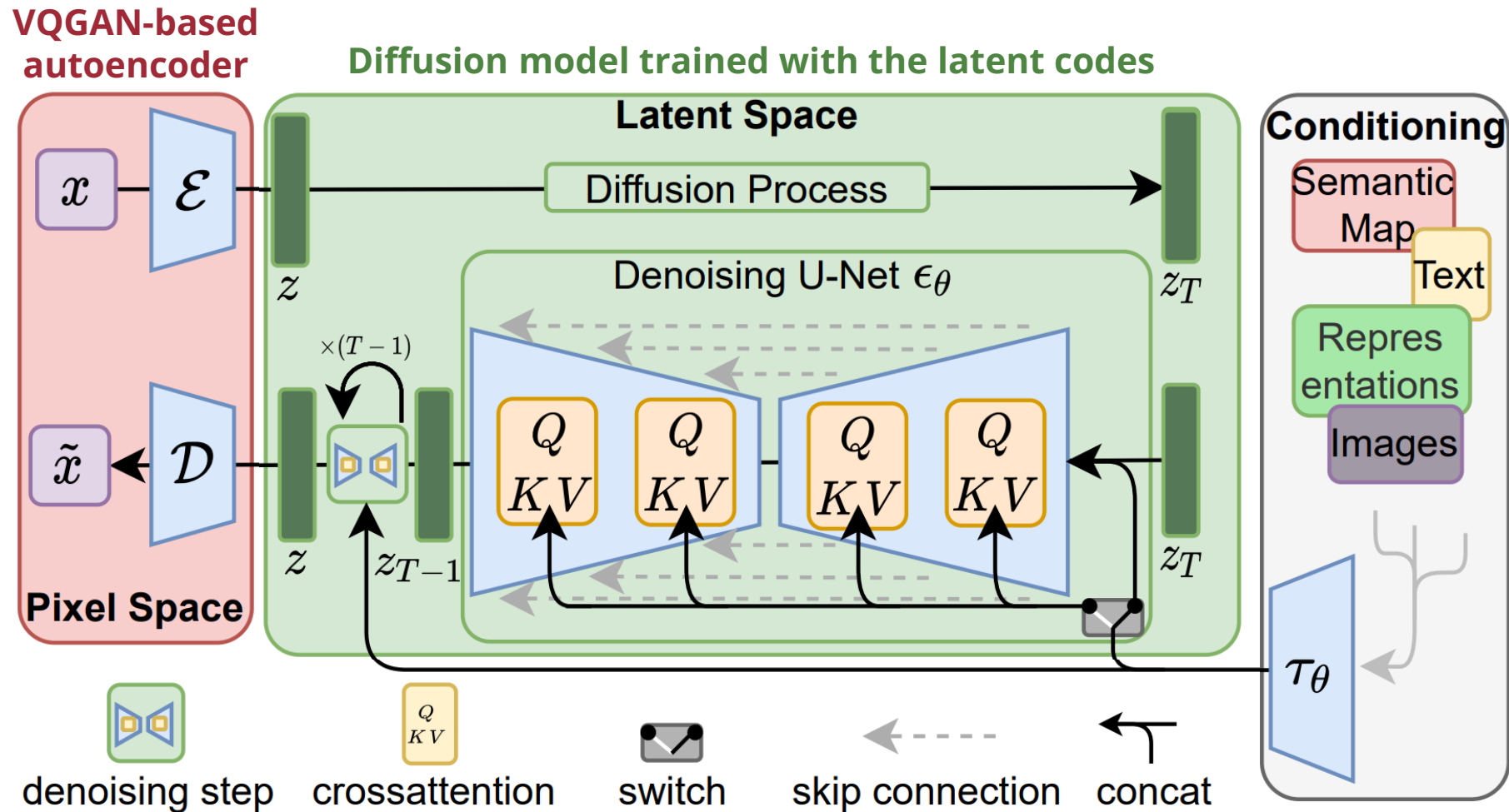
[google-research.github.io/seanet/musiclm/examples/](https://google-research.github.io/seanet/musiclm/examples/)

# Google's Music FX (2024)



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

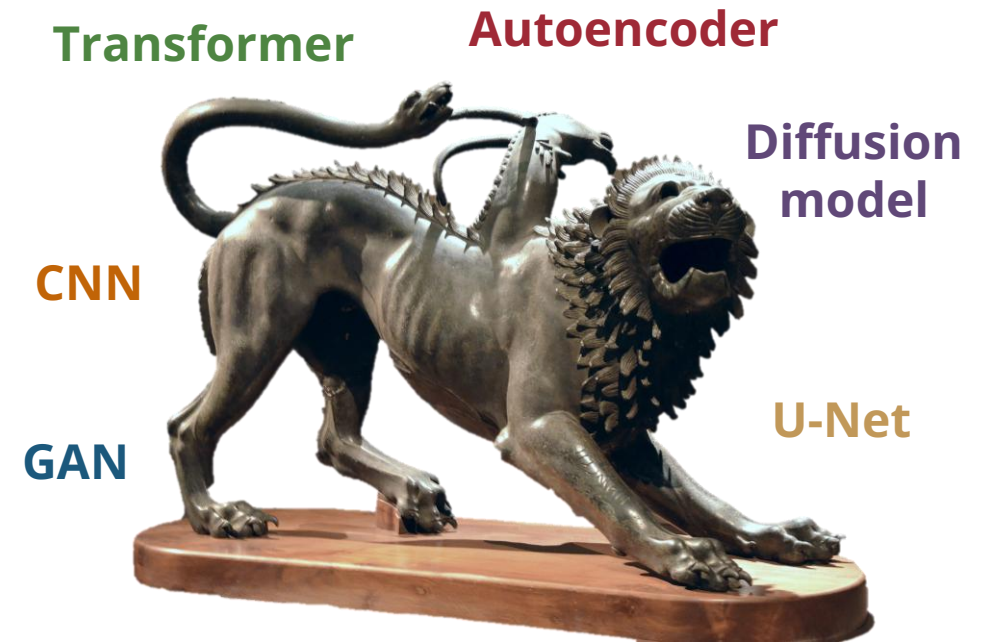
# Latent Diffusion Models (LDMs)



(Source: Rombach et al., 2022)

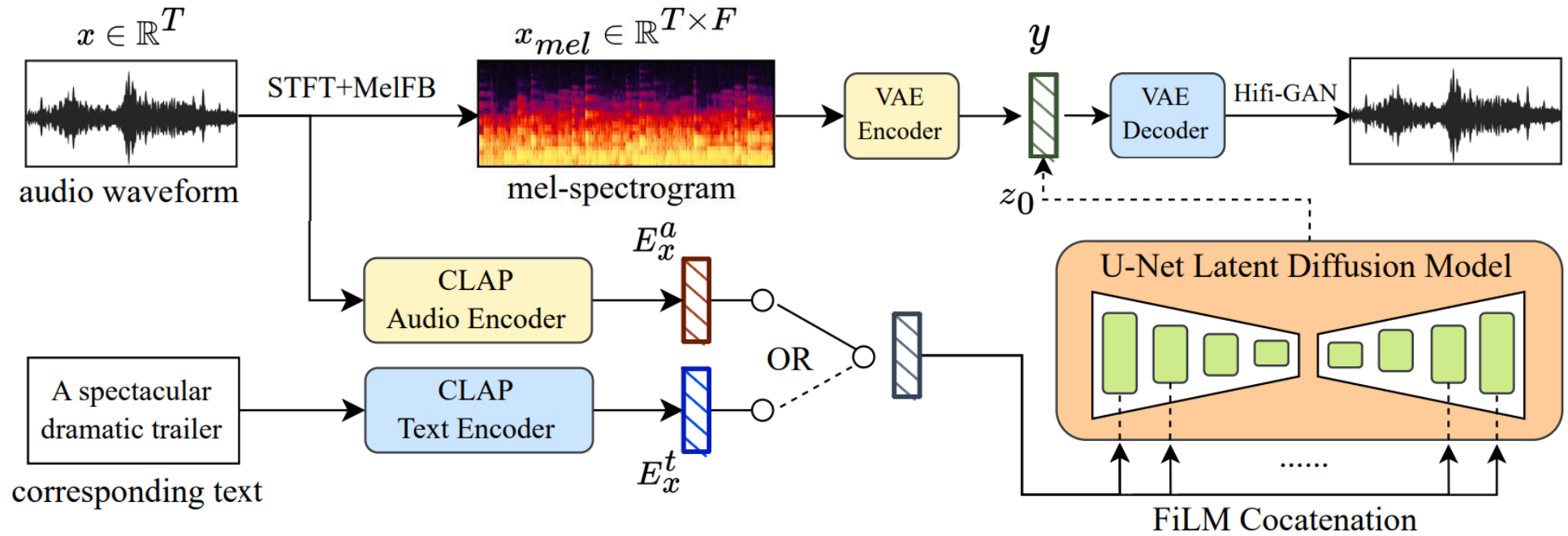
# Latent Diffusion Model is a Chimera

- **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)

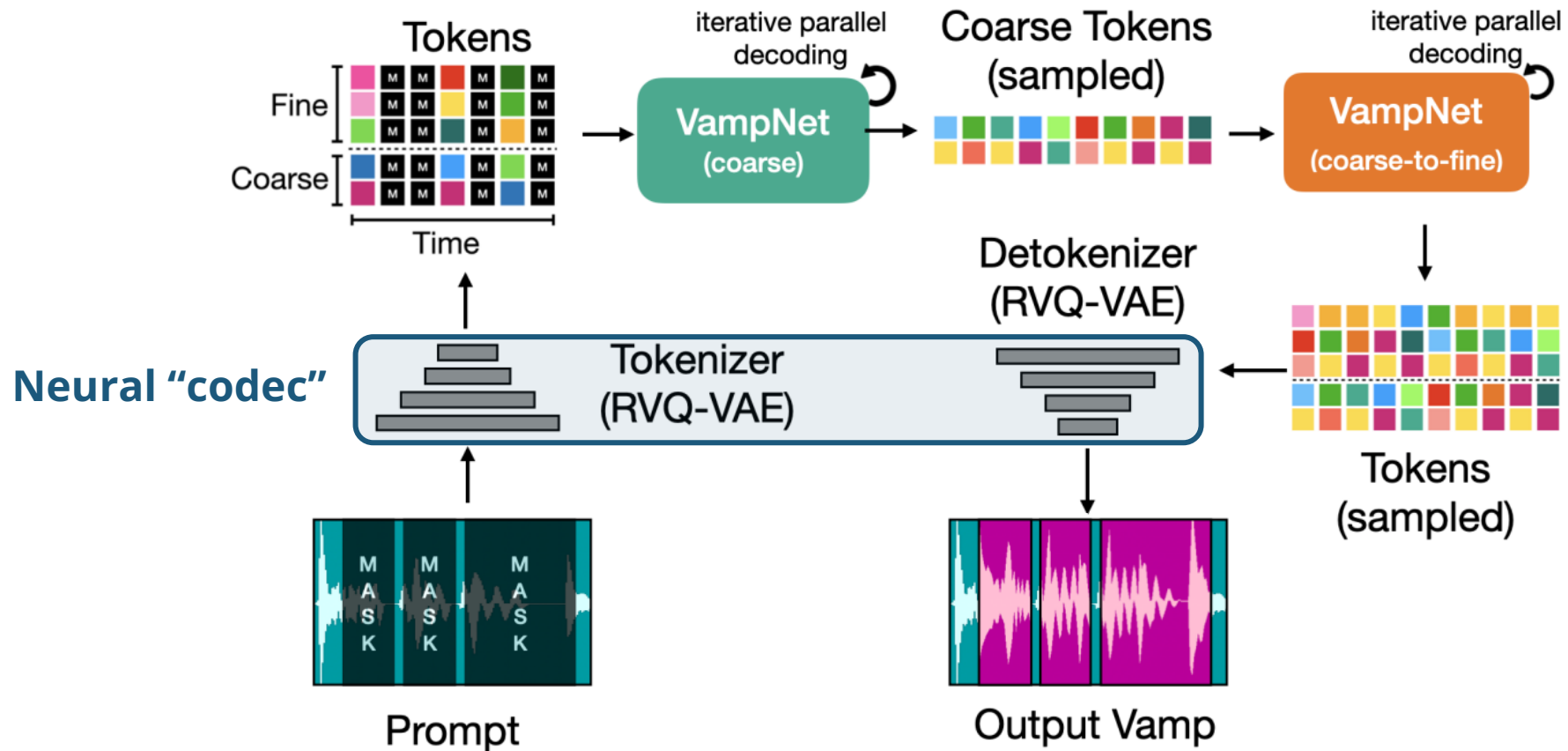
# MusicLDM (Chen et al., 2023)



(Source: Ke et al., 2023)

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

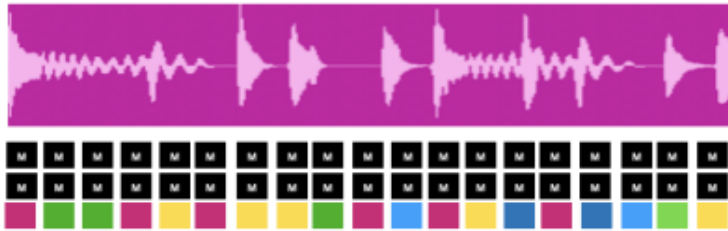
# VampNet (Garcia et al., 2023)



(Source: Garcia et al., 2023)

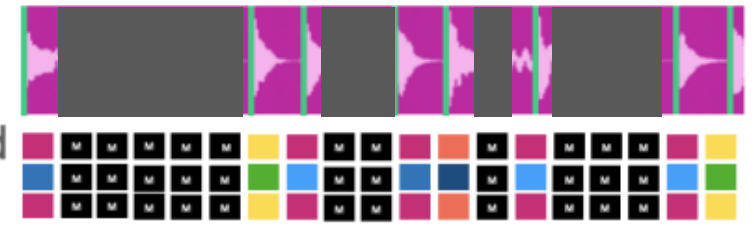
# VampNet (Garcia et al., 2023)

Compression



Beat Driven

= predicted beat mark



Periodic



Inpainting



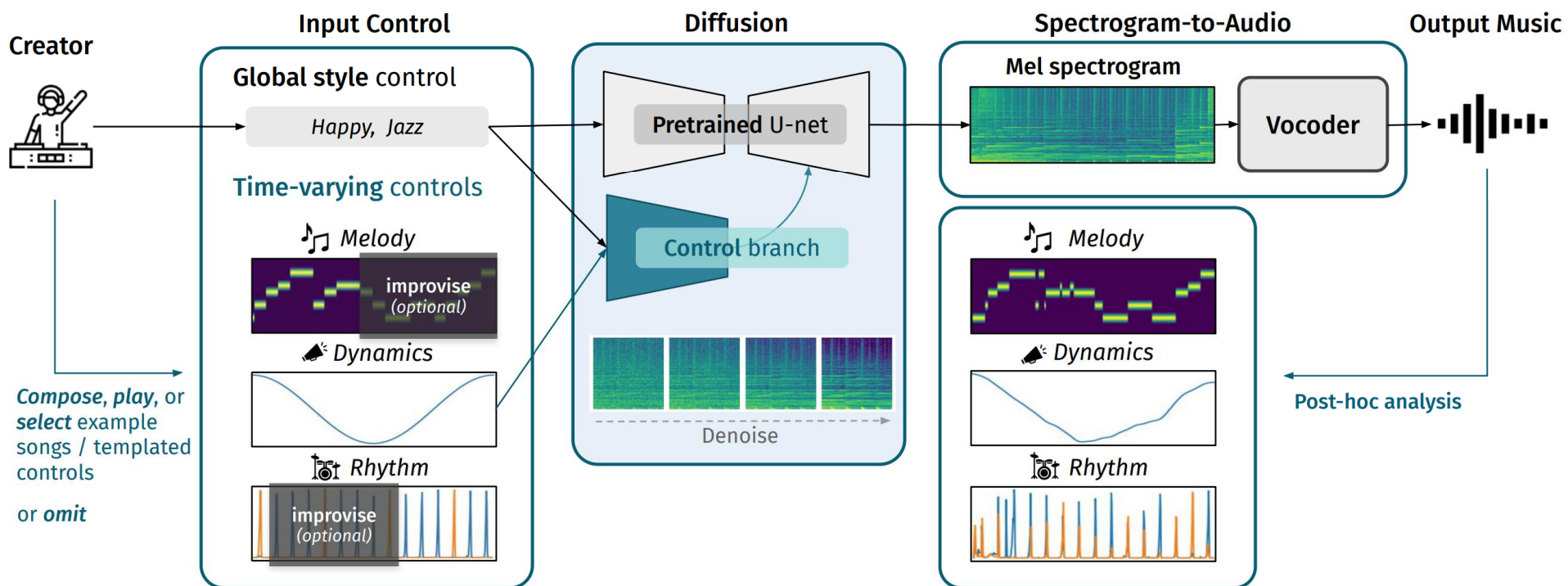
(Source: Garcia et al., 2023)

# unloop (Garcia et al., 2023)



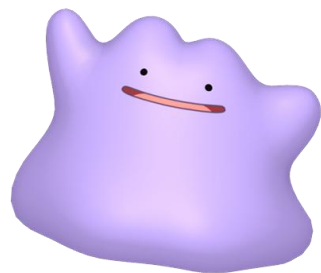
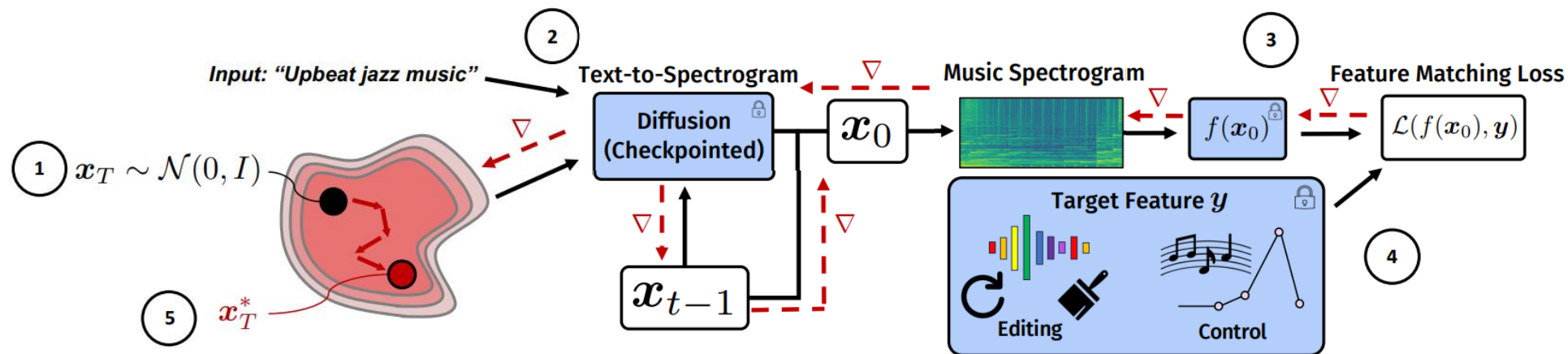
[youtu.be/yzBl8Vcjd2s](https://youtu.be/yzBl8Vcjd2s) & [github.com/hugofloresgarcia/unloop](https://github.com/hugofloresgarcia/unloop)

# Music ControlNet (Wu et al., 2024)



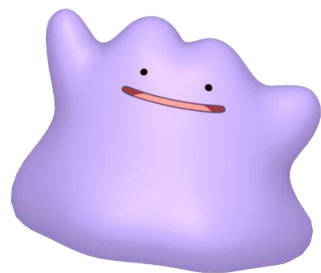
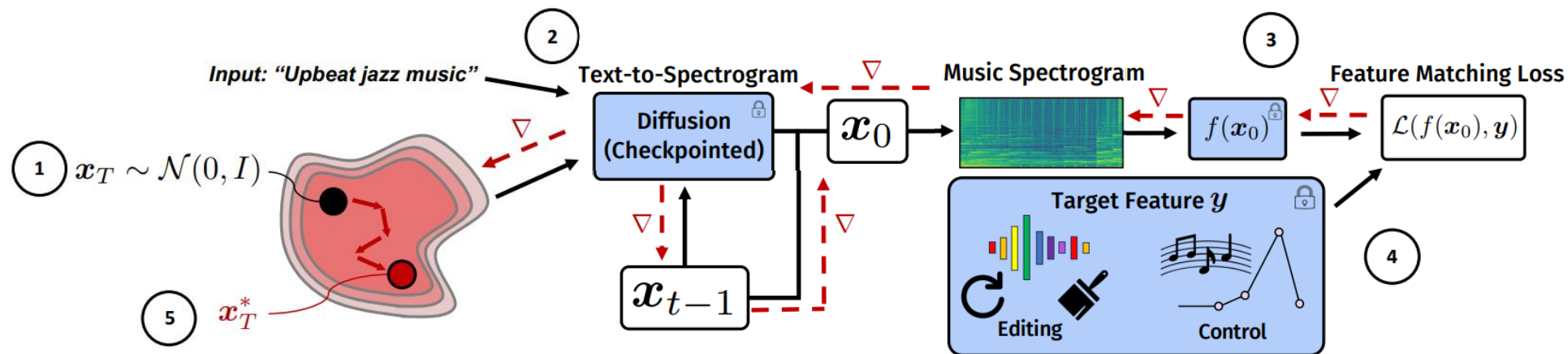
(Source: Wu et al., 2024)

# Inference-time Control: DITTO (Novack et al., 2024)



(Source: Novack et al., 2024)

# Inference-time Control: DITTO (Novack et al., 2024)

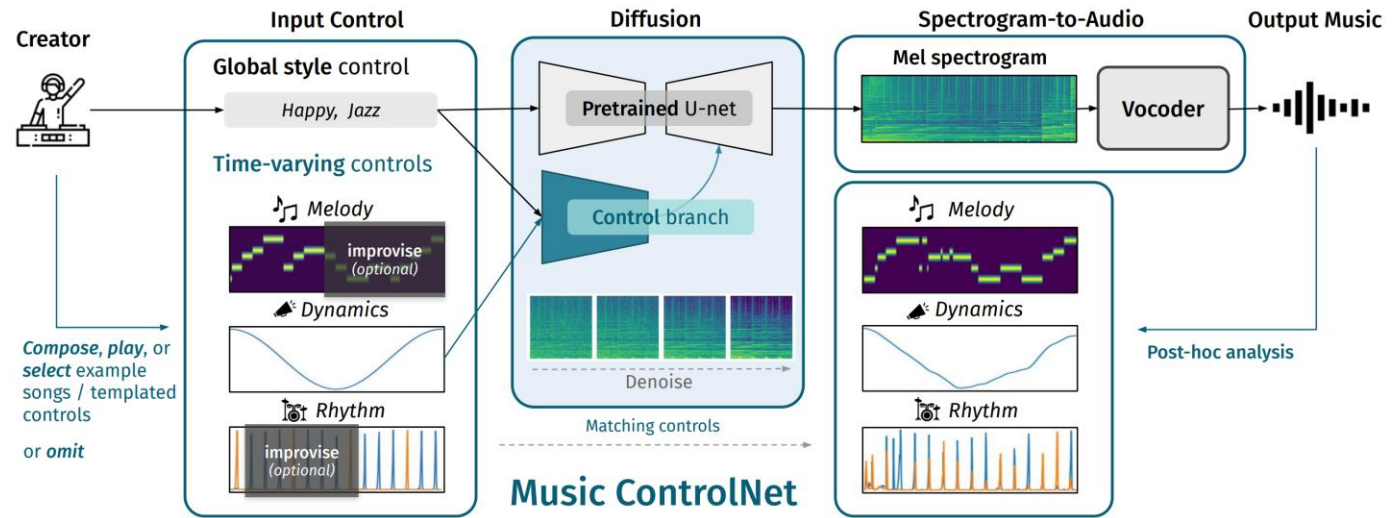


(Source: Novack et al., 2024)

# Music ControlNet vs DITTO

## Music ControlNet

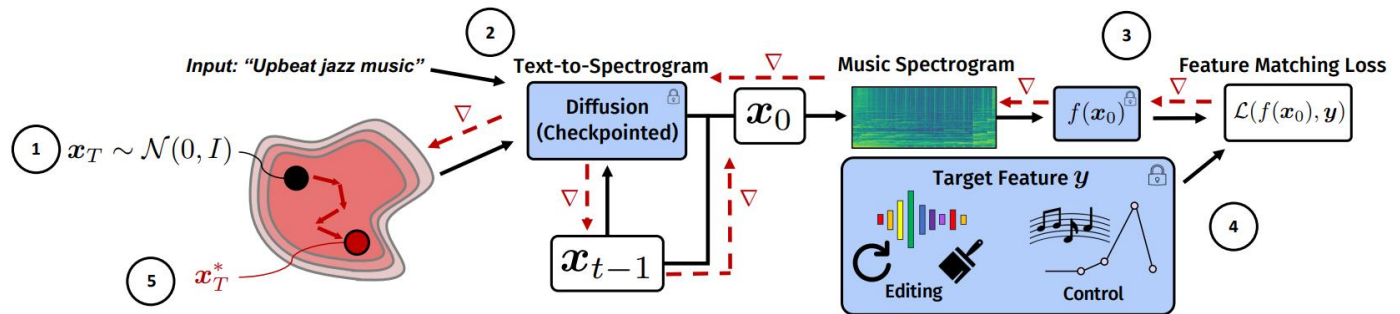
Needs some training!



(Source: Wu et al., 2024)

## DITTO

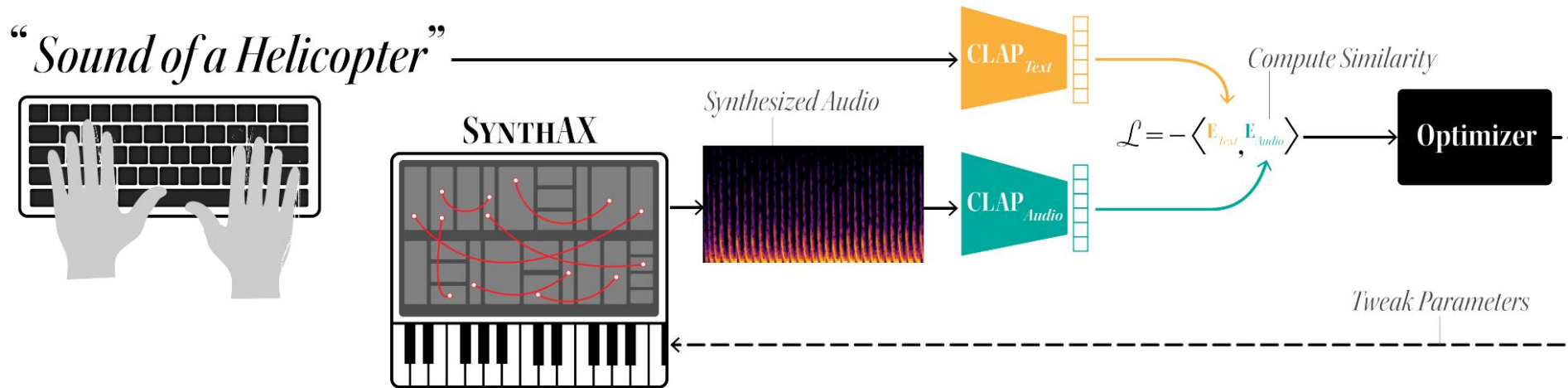
No training needed!



(Source: Novack et al., 2024)

## Next Lecture

# Multimodal Systems & Music Production



(Source: Cherep et al., 2024)

[ctag.media.mit.edu](http://ctag.media.mit.edu)