

PAT 498/598 (Fall 2024)

# Special Topics: Generative AI for Music and Audio Creation

## Lecture 17: Latent-based Audio Synthesis

Instructor: Hao-Wen Dong



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

# Project Pitch

- Tell us briefly about
  - **What do you plan to do?**
    - Task, input/output, motivation, use cases, target users, etc.
  - **How do you plan to approach it?**
    - Papers to read, expected challenges, types of datasets/models to use, etc.
- **Send me a short paragraph or some bullet points by tomorrow** so that I can provide some feedback and suggest relevant papers to look into!

# Final Project

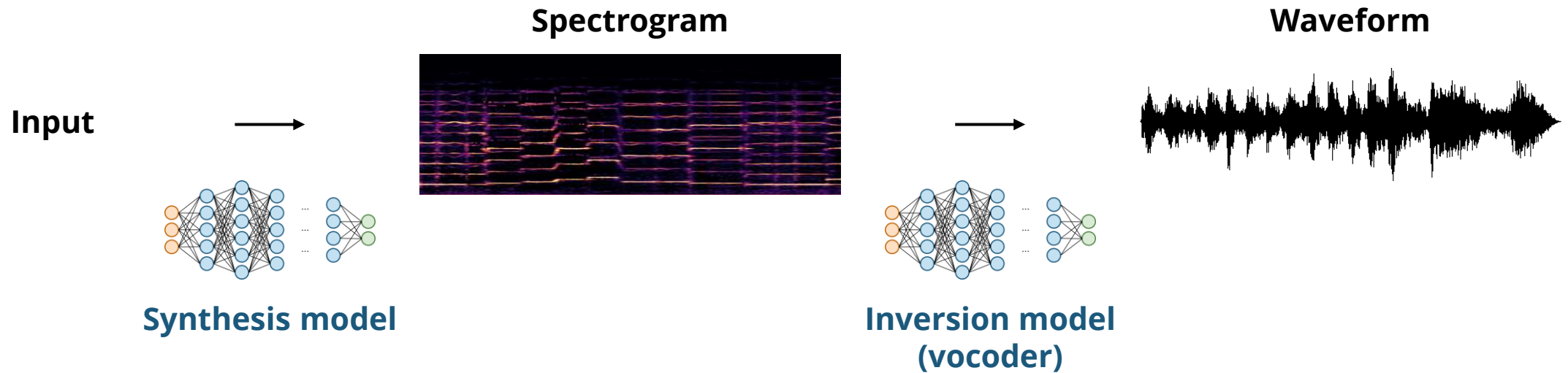
- Milestones (all due at the specified date at **11:59 PM ET**)

▪ <b>Pitch</b>	November 6	Topic & high-level plans
▪ <b>Proposal</b>	November 18	Survey & plans (1 page)
▪ <b>Presentation</b>	December 9	Showcase & report
▪ <b>Final report</b>	December 15	Full report (3-5 pages)

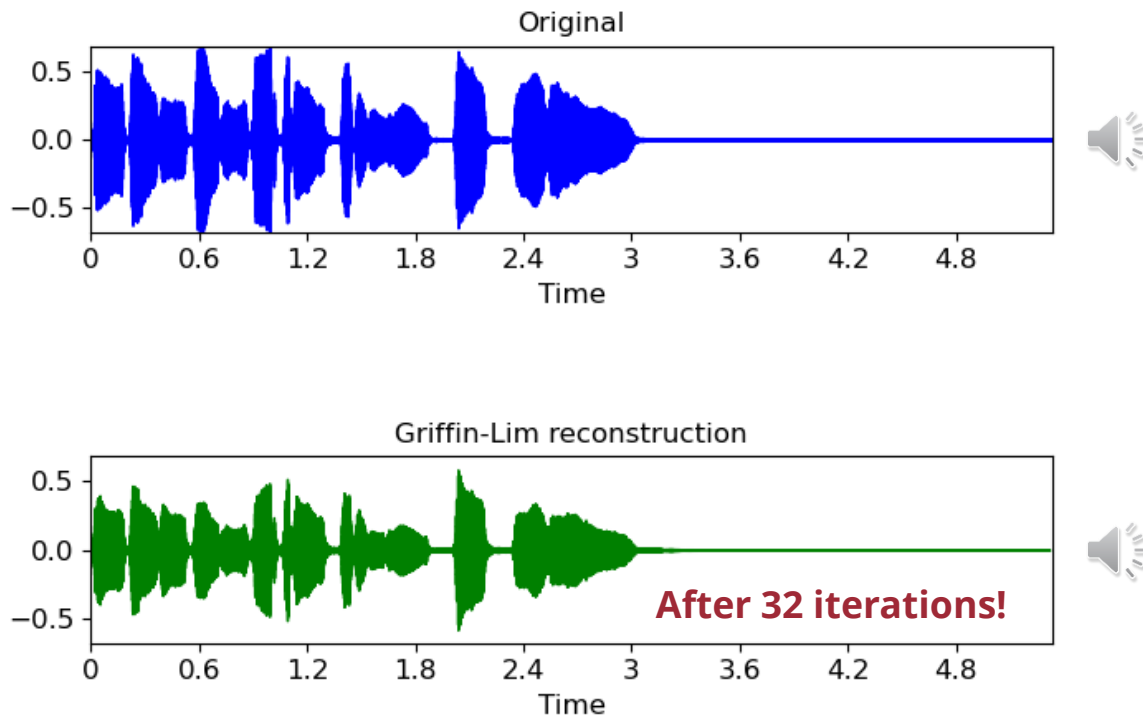
**Next milestone!**

- Instructions will be released on Gradescope
- Late submissions: **NOT accepted**

# (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  $\text{STFT}(y)$

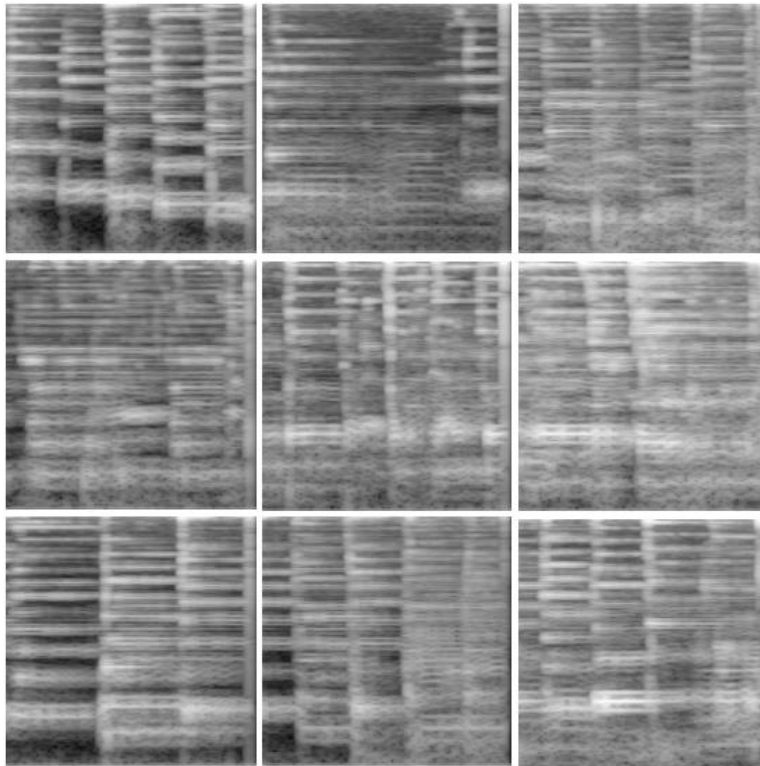


$$M' = \text{STFT}(y')$$

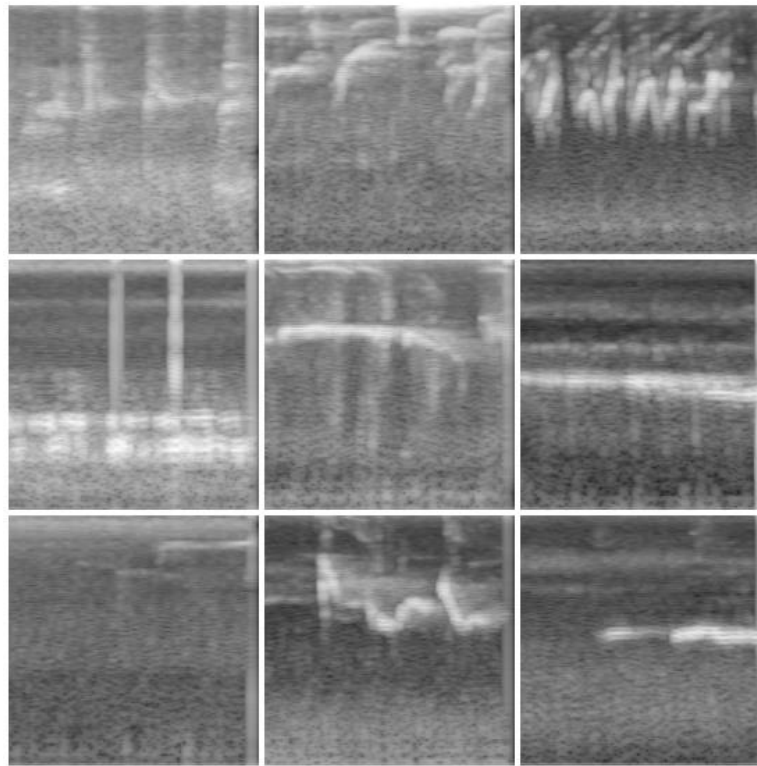
Find the STFT of the signal  $y$

# (Recap) Example: SpecGAN (Donahue et al., 2019)

Piano sounds



Bird sounds

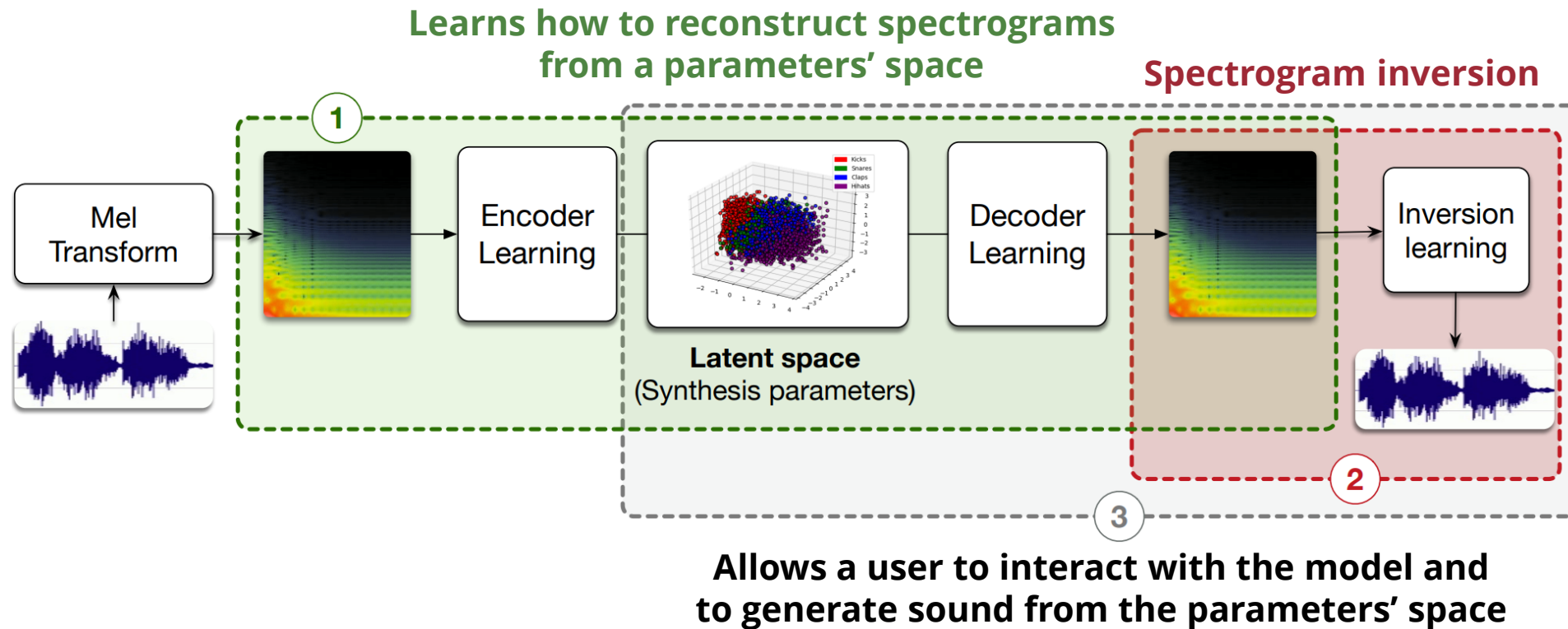


Example of generated music



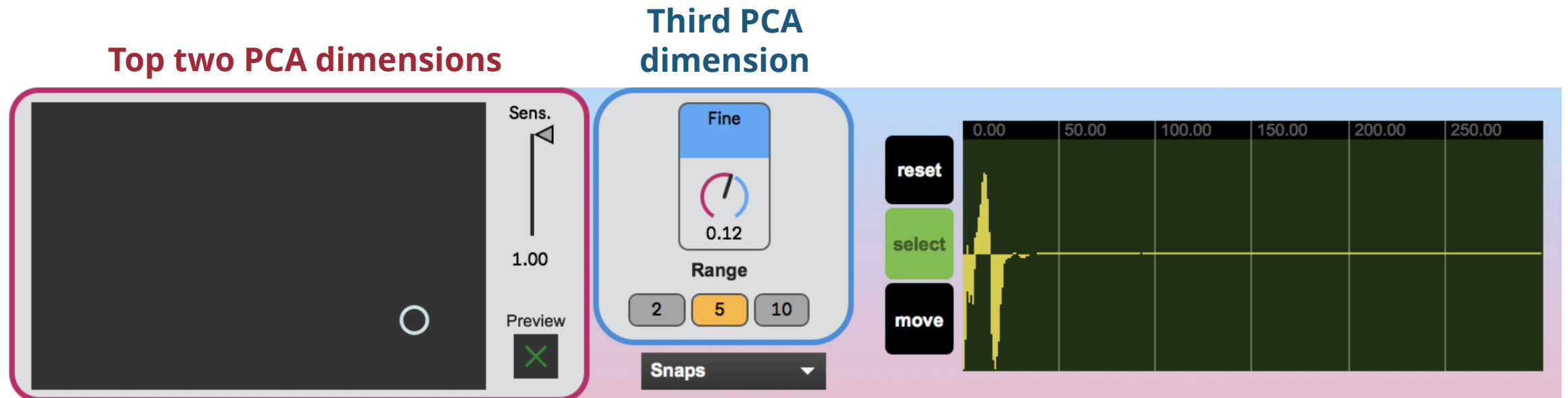
(Source: Donahue et al., 2019)

# (Recap) Example: Neural Drum Machine (Aouameur et al., 2019)



(Source: Aouameur et al., 2019)

# (Recap) Example: Neural Drum Machine (Aouameur et al., 2019)

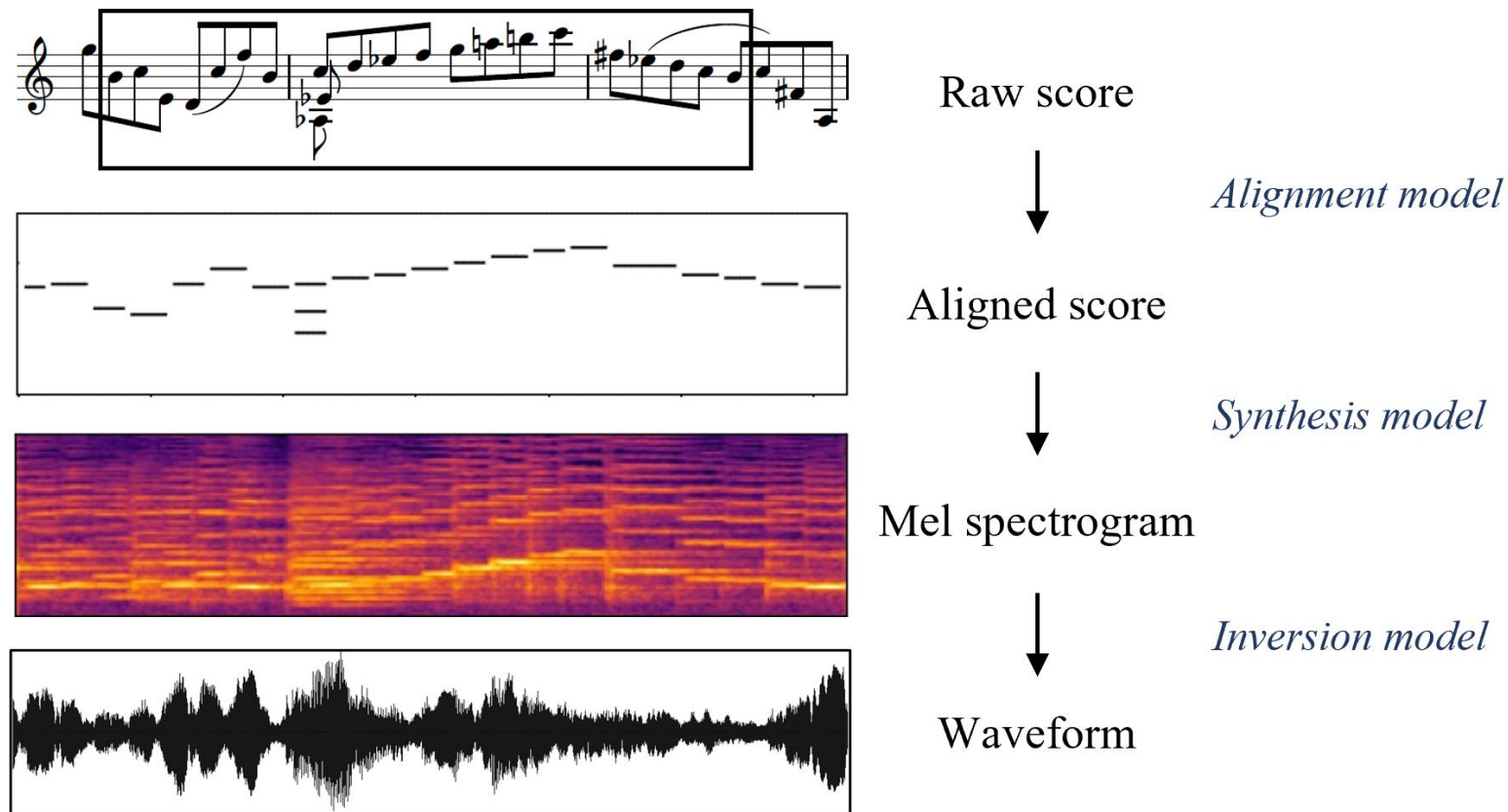


(Source: Aouameur et al., 2019)

[drive.google.com/file/d/1DDo0\\_KnwkWirCM4t0PT8cp6uotsfuufj/view](https://drive.google.com/file/d/1DDo0_KnwkWirCM4t0PT8cp6uotsfuufj/view)

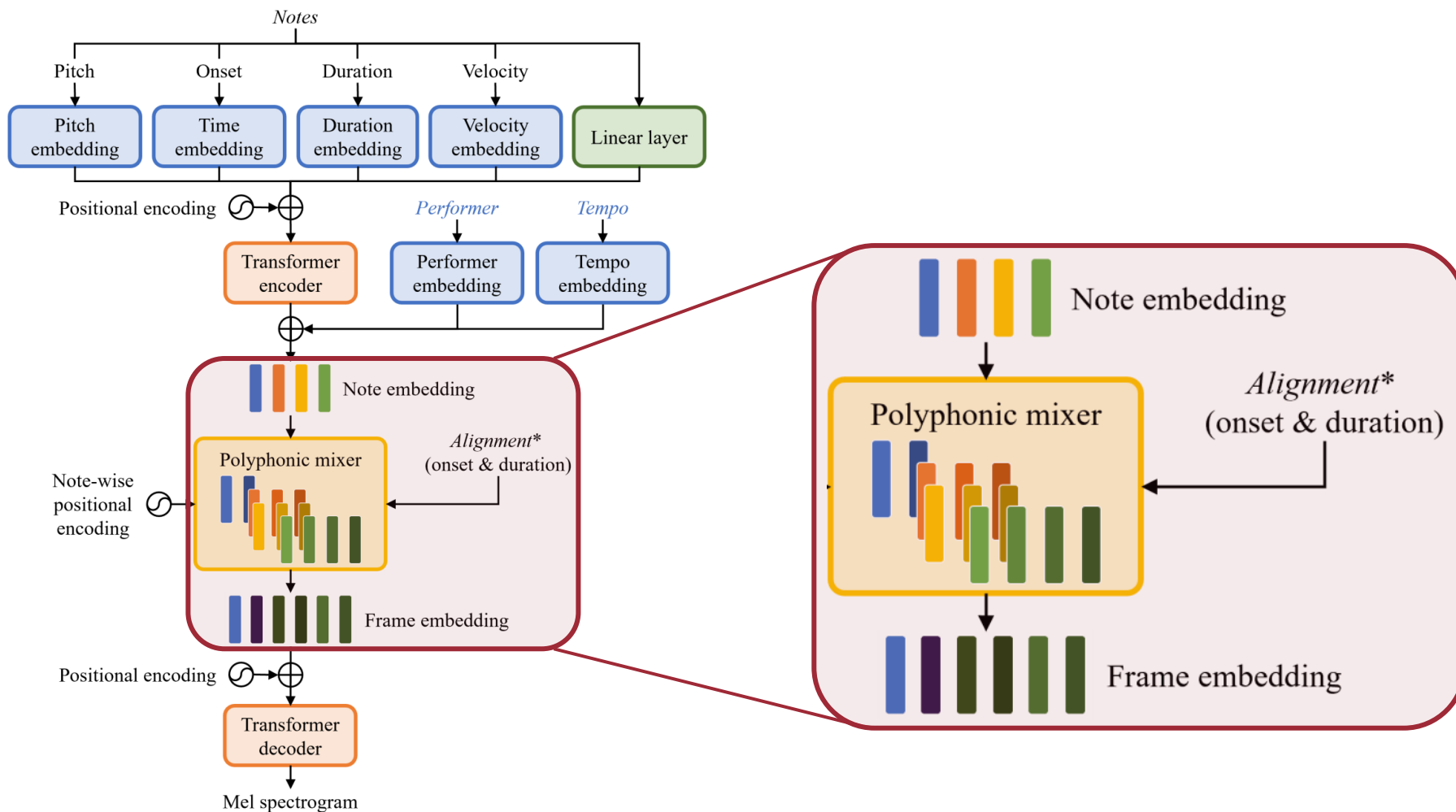


# (Recap) Example: DeepPerformer (Dong et al., 2022)



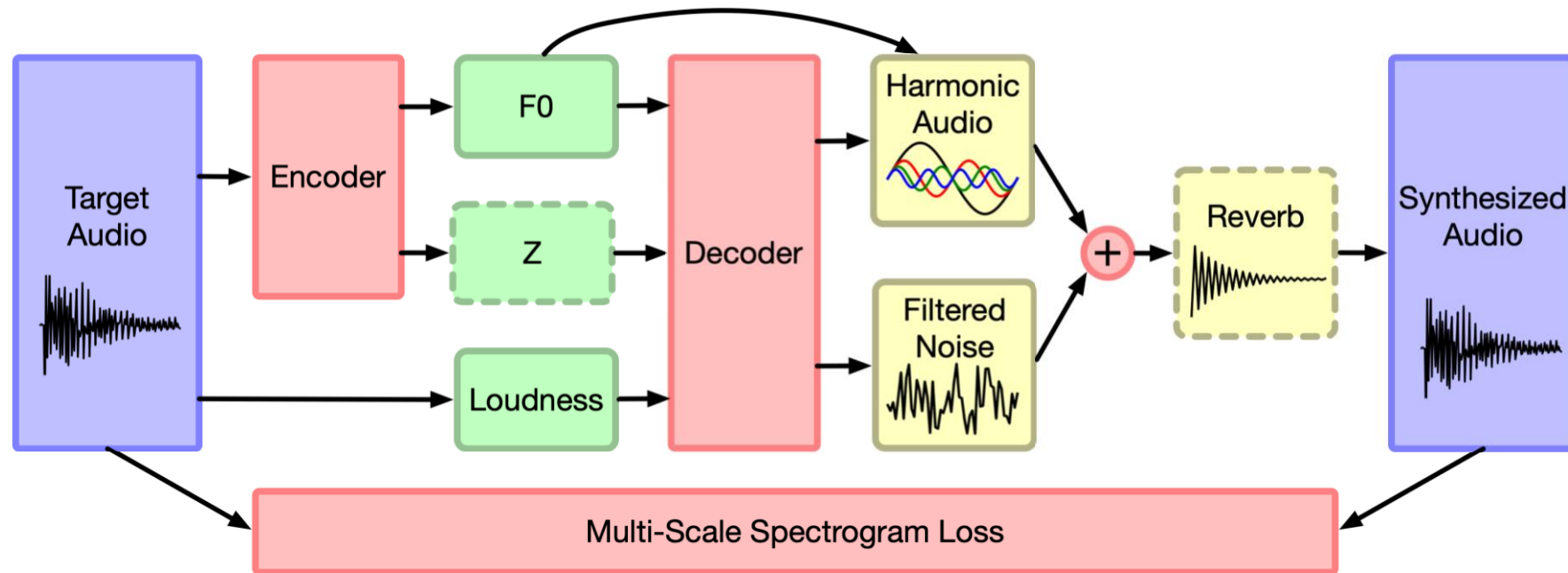
(Source: Dong et al., 2022)

# (Recap) Example: DeepPerformer (Dong et al., 2022)



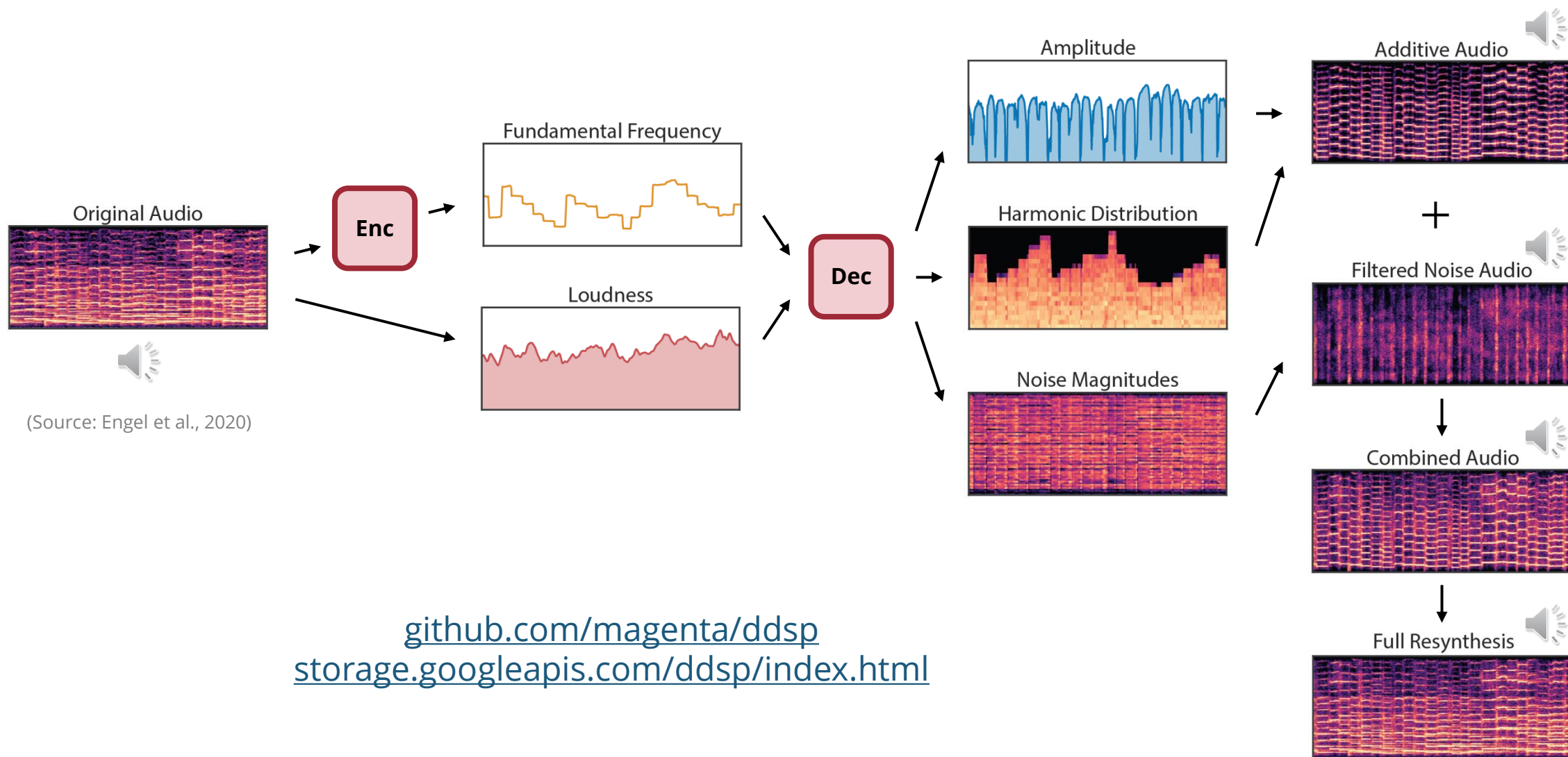
(Source: Dong et al., 2022)

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



(Source: Engel et al., 2020)

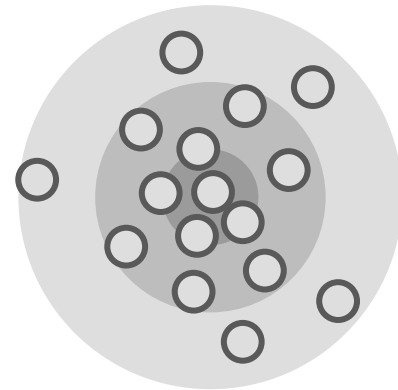
# (Recap) Example: Differentiable DSP (DDSP) (Engel et al., 2020)



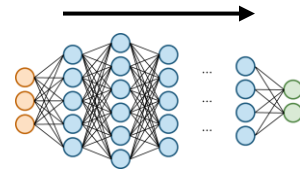
# Latent-based Audio Synthesis

# (Recap) Generating Data from a Random Distribution

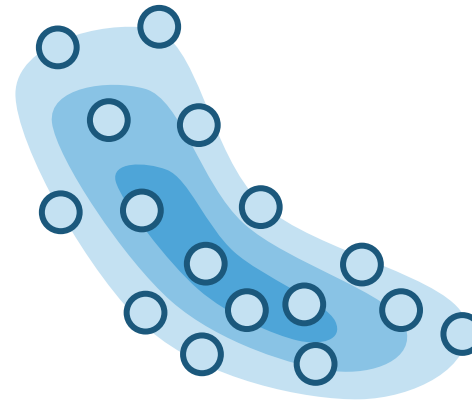
Random distribution



$P(z)$



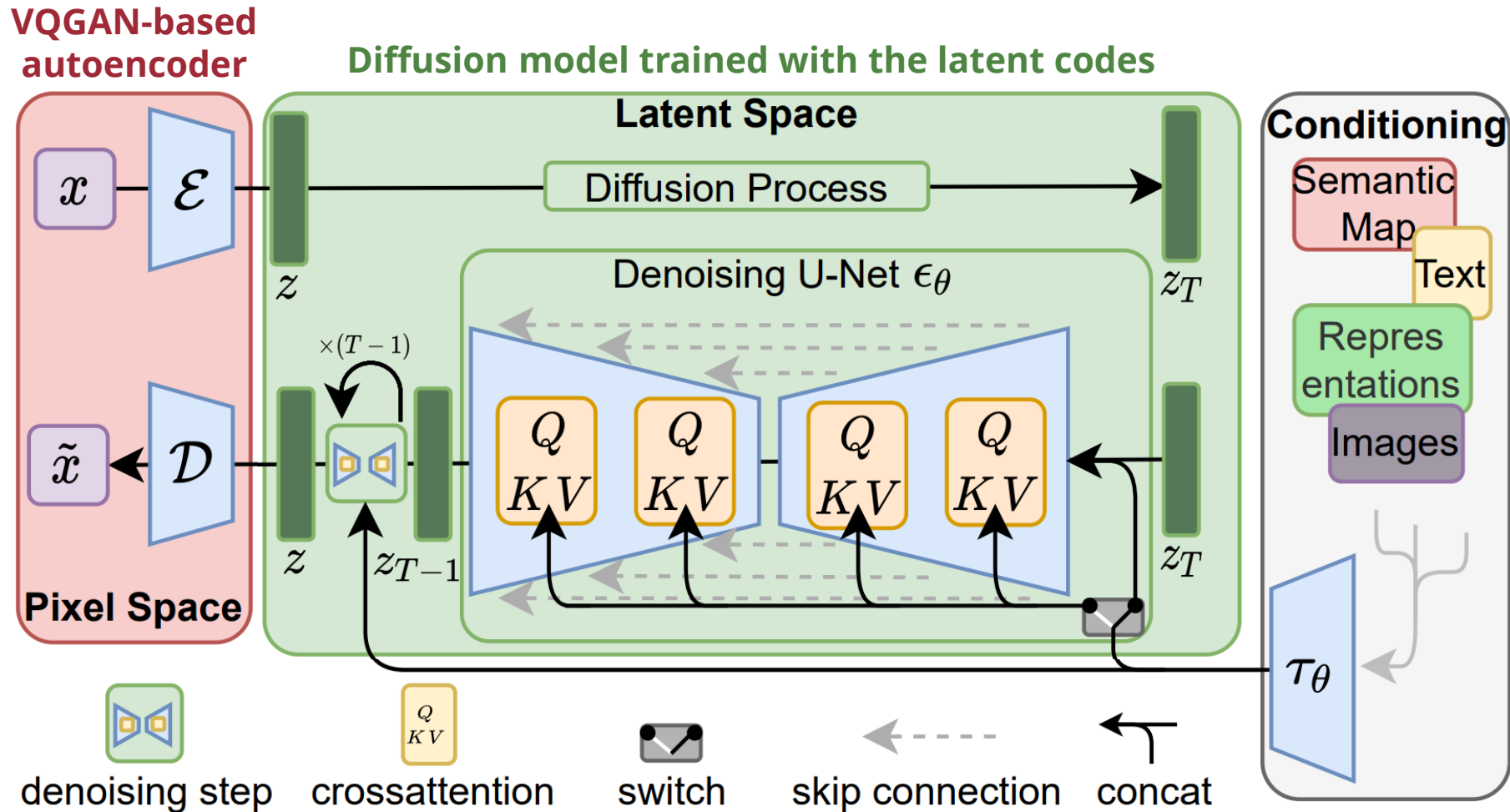
Data distribution



$P(x)$

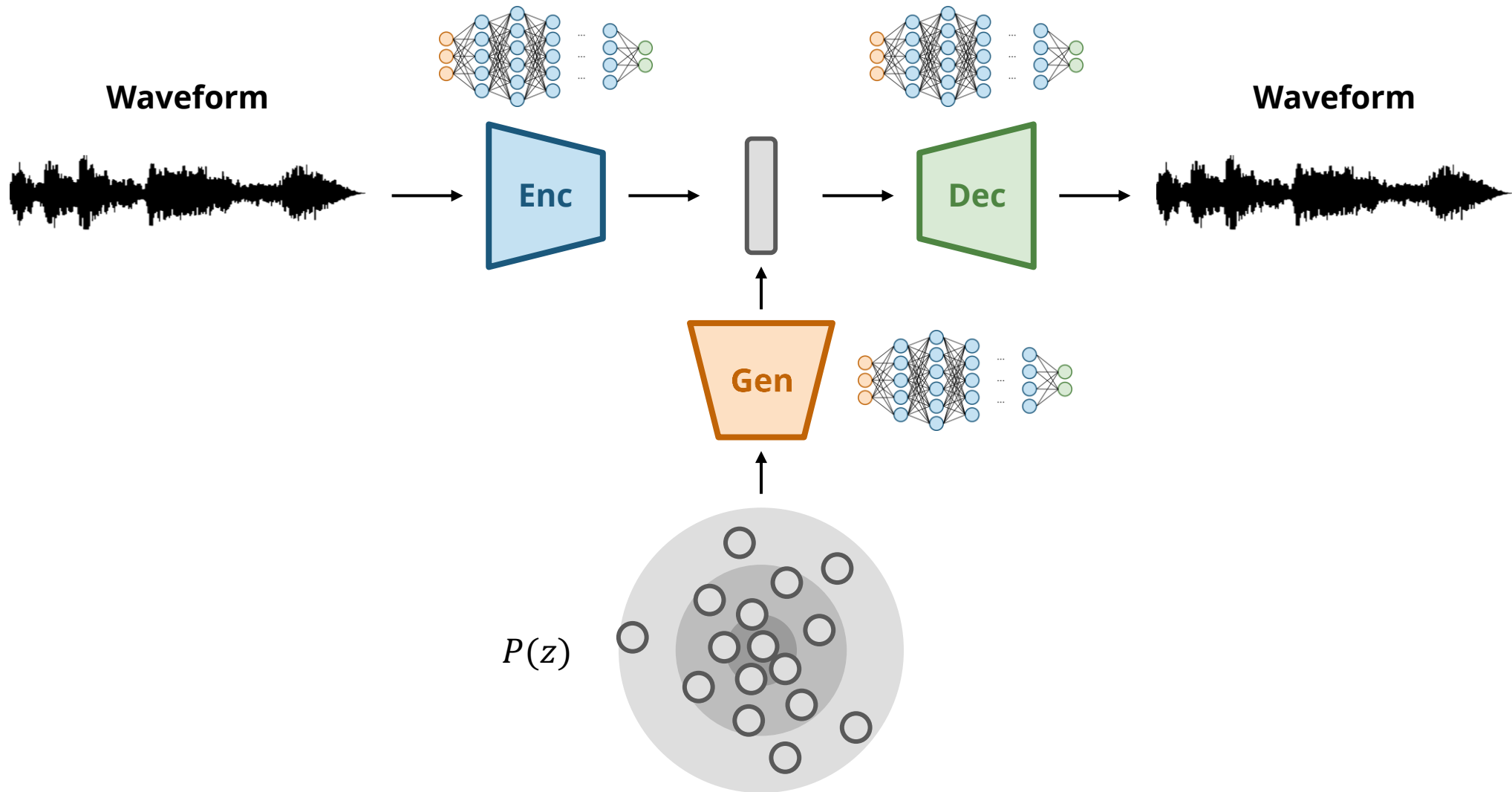
**If we can learn this mapping, we can easily generate new samples from the data distribution**

# (Recap) Latent Diffusion Models (LDMs)



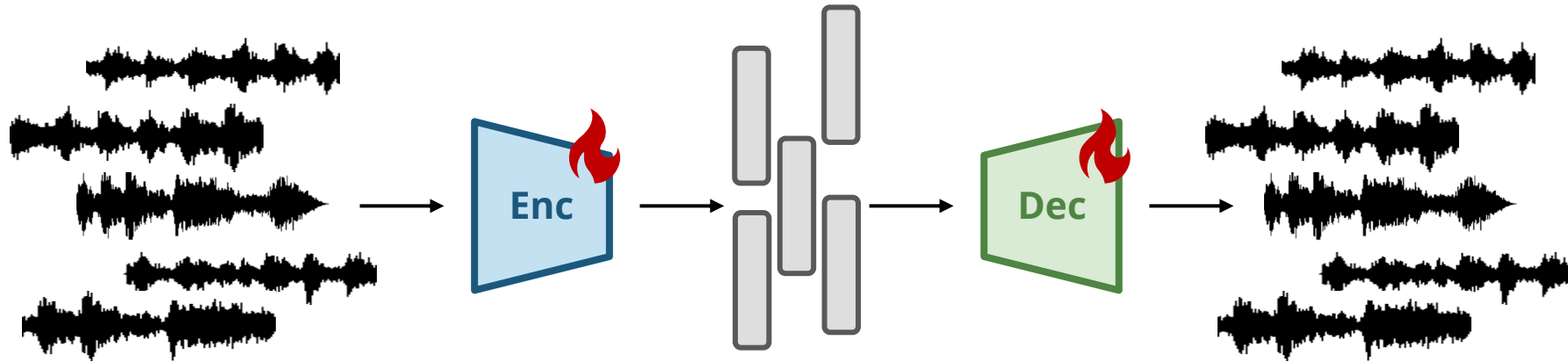
(Source: Rombach et al., 2022)

# Latent-based Audio Synthesis

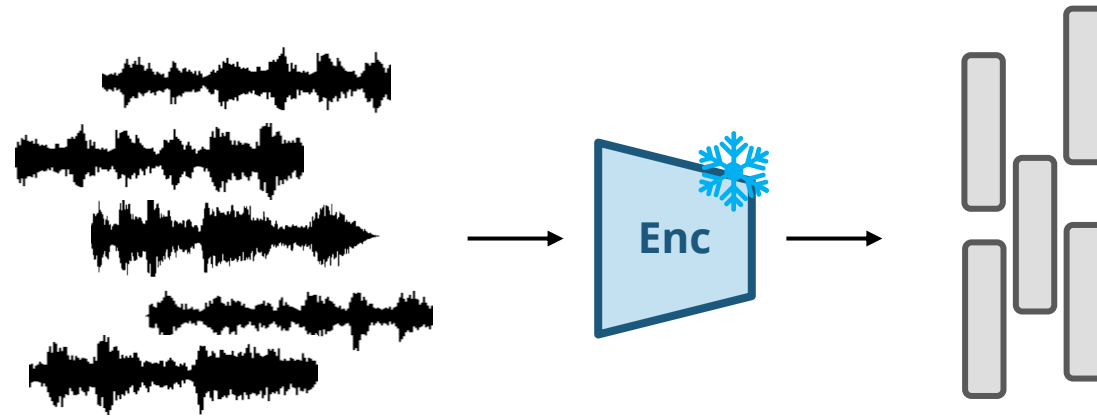




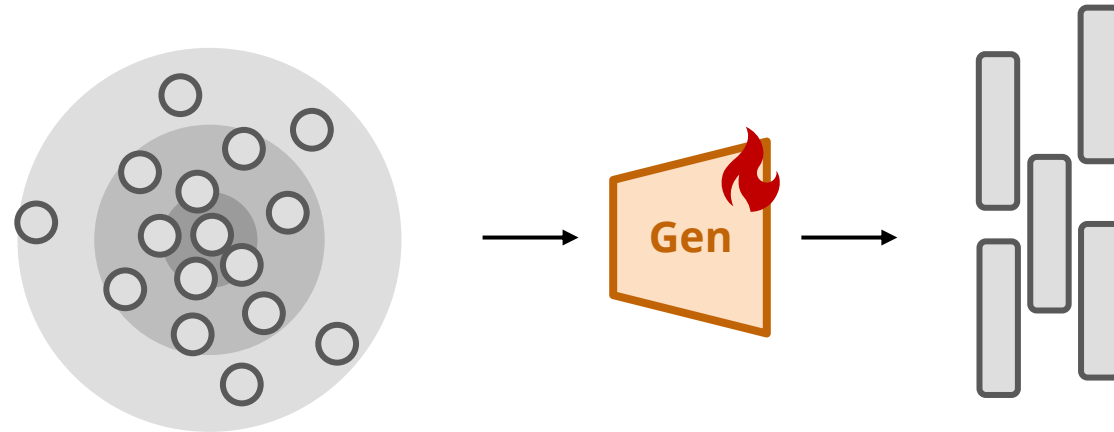
# Step 1: Train an Autoencoder



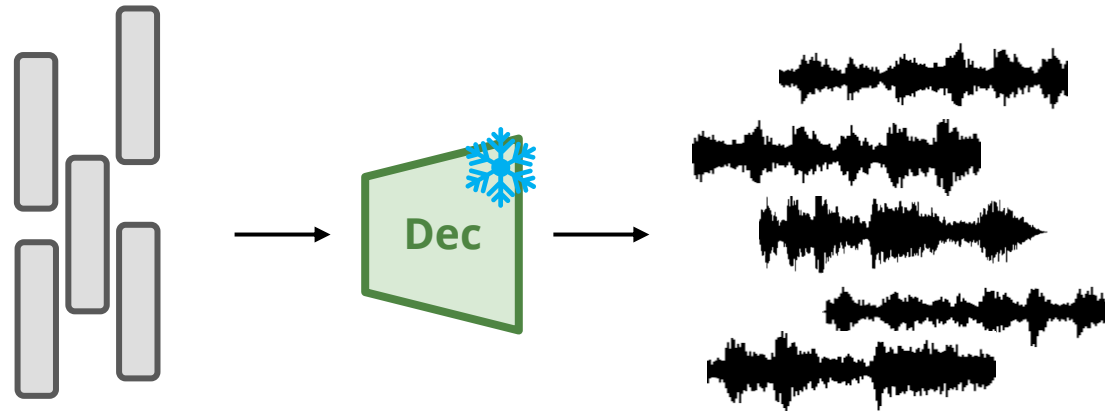
## Step 2: Compute the Latent Vectors



## Step 3: Train a Latent Generative Model

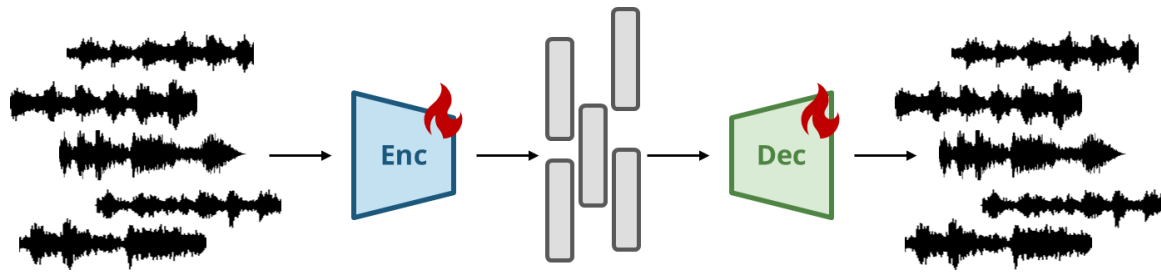


## Step 4: Decode the Latent Vectors

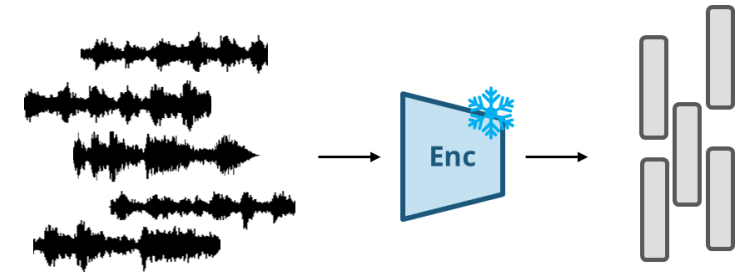


# Pipeline

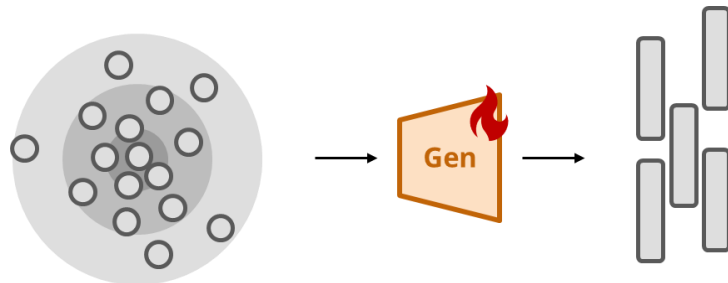
Step 1: Train an Autoencoder



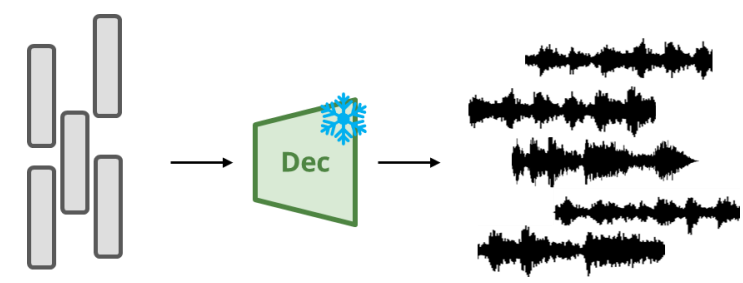
Step 2: Compute the Latent Vectors



Step 3: Train a Latent Generative Model

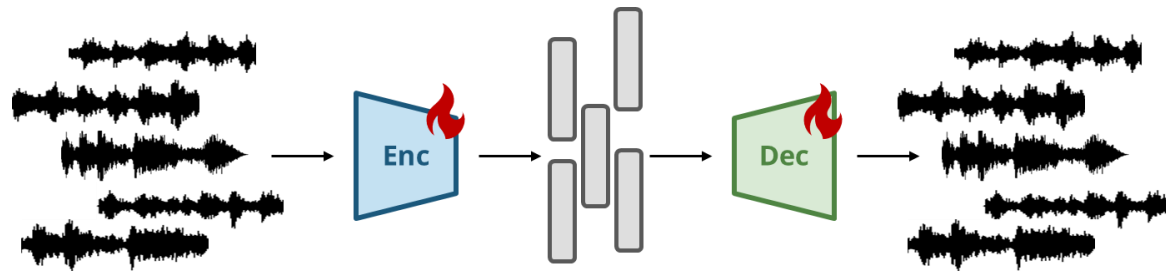


Step 4: Decode the Latent Vectors

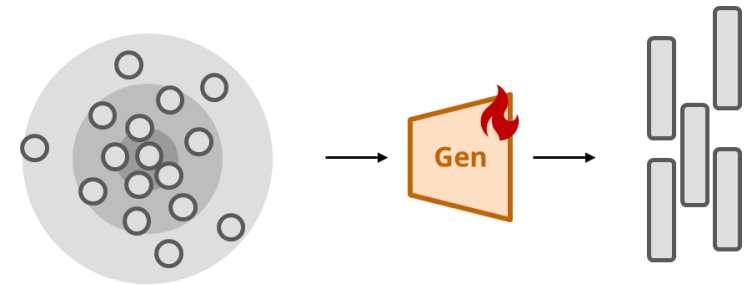


# Training

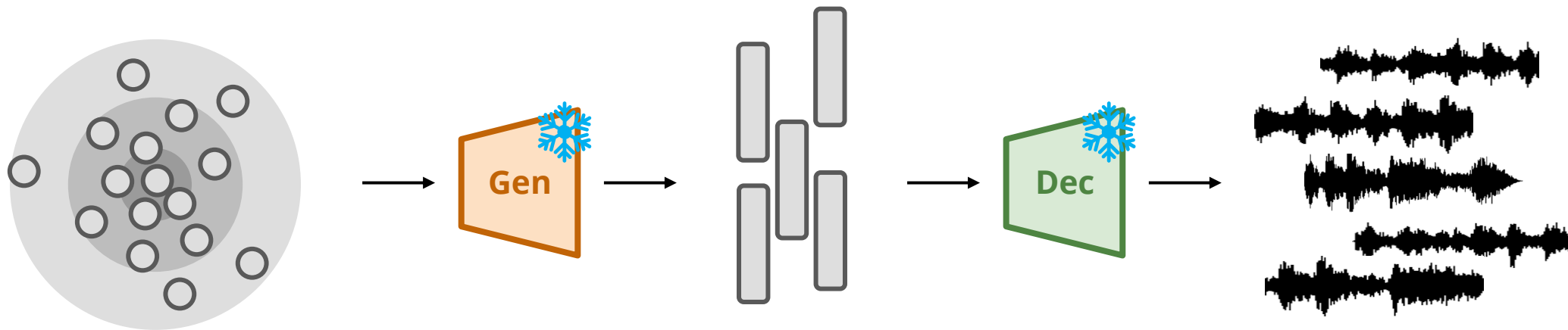
Autoencoder



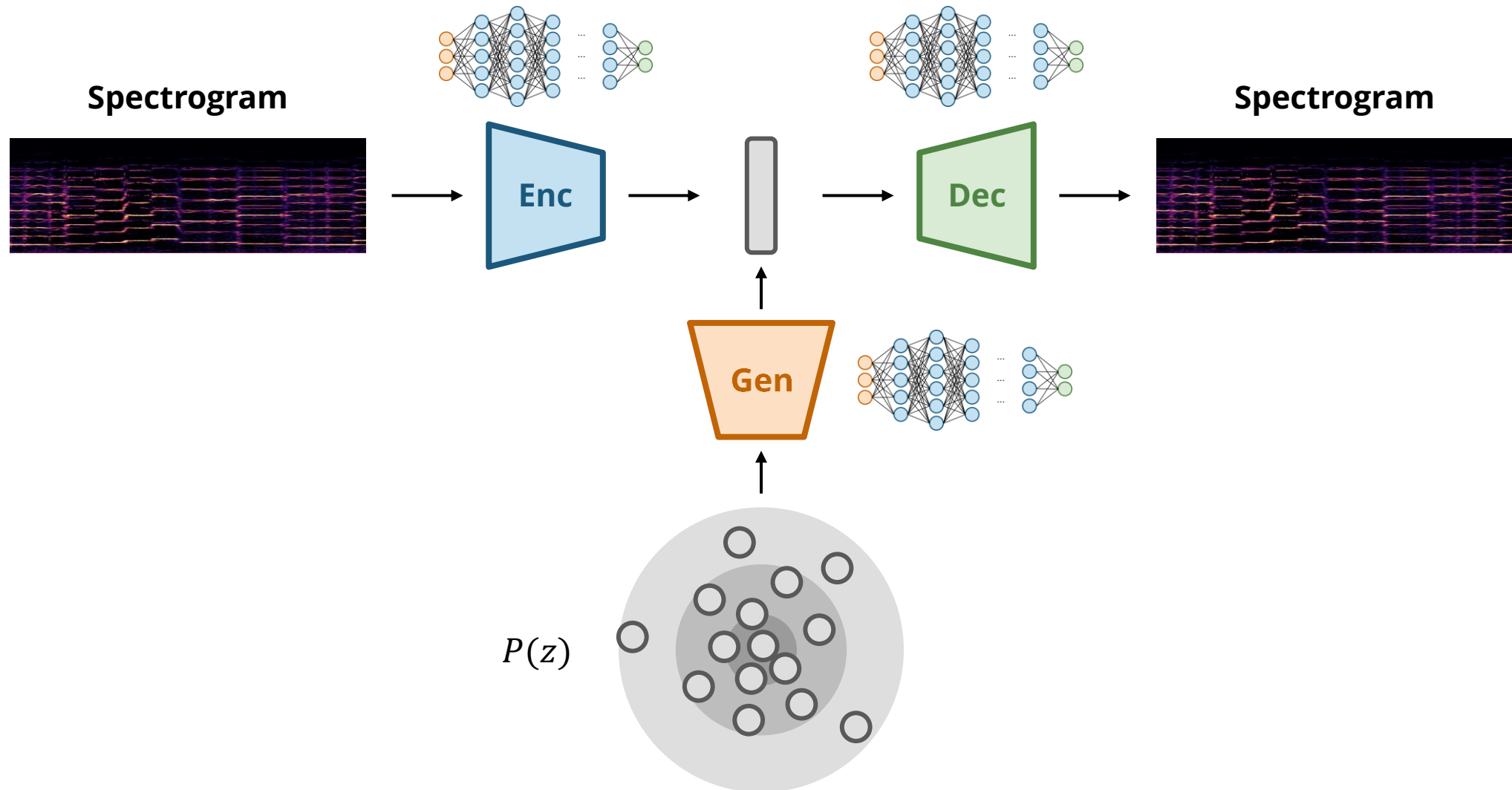
Latent Generative Model



# Inference



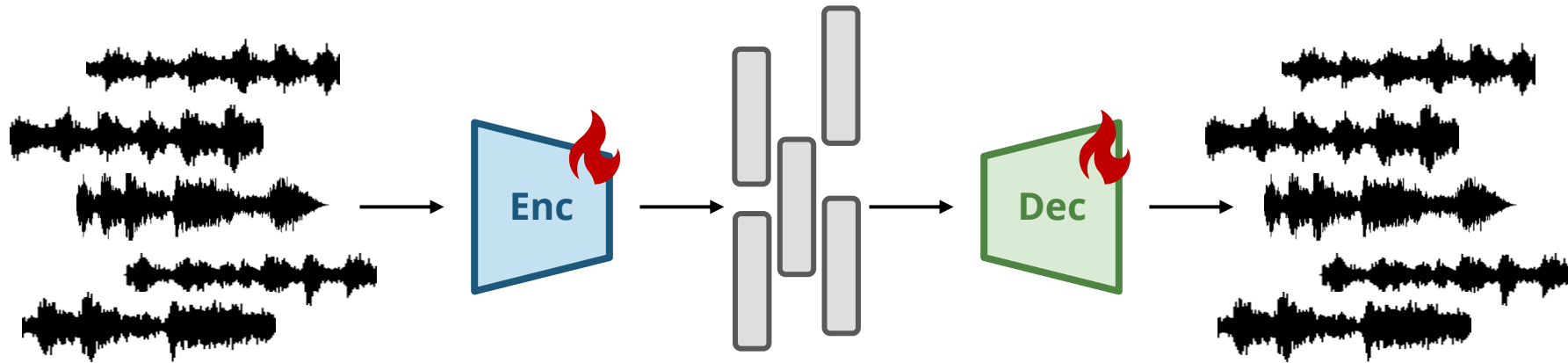
# Latent-based Audio Synthesis



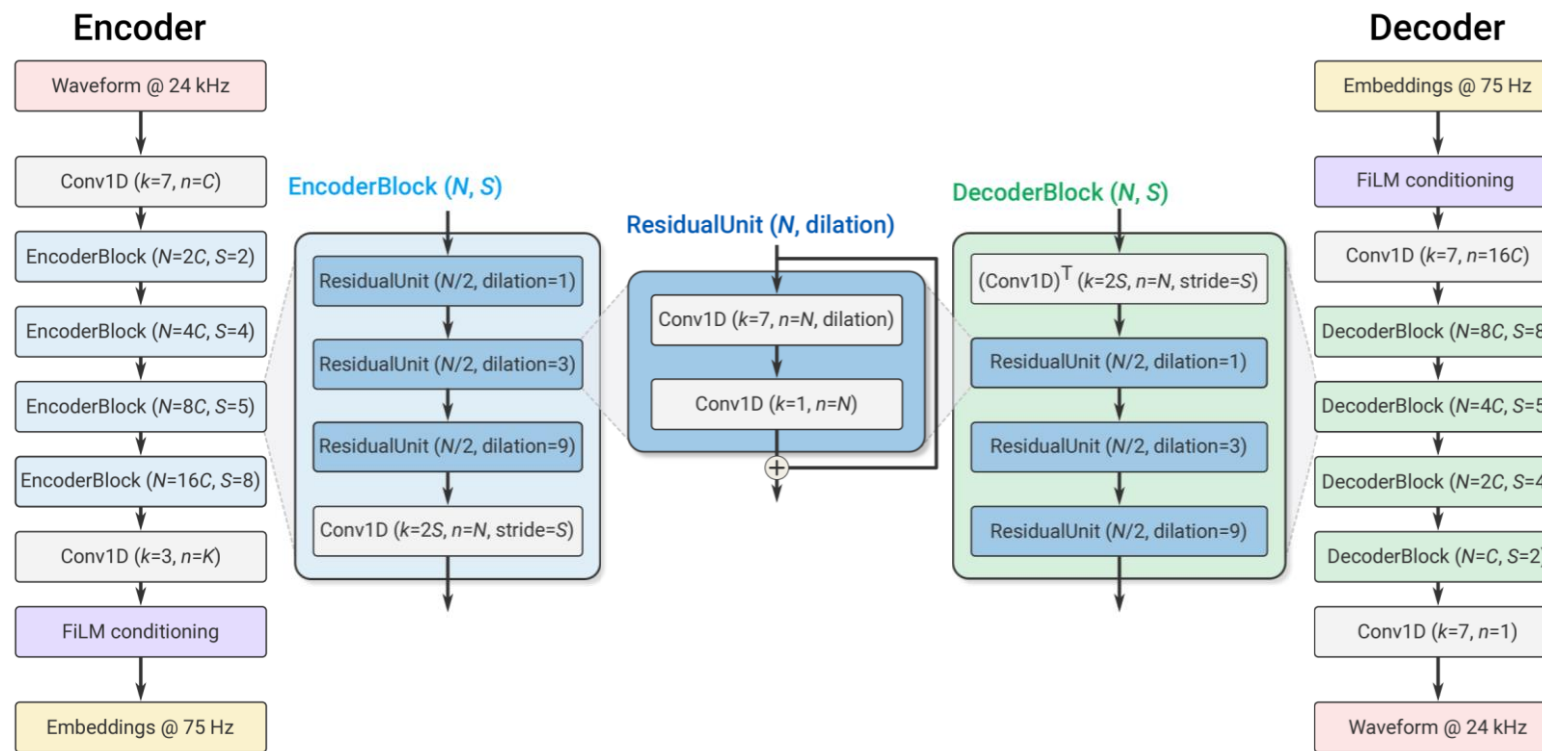


# Neural Codecs

# (Recap) Step 1: Train an Autoencoder

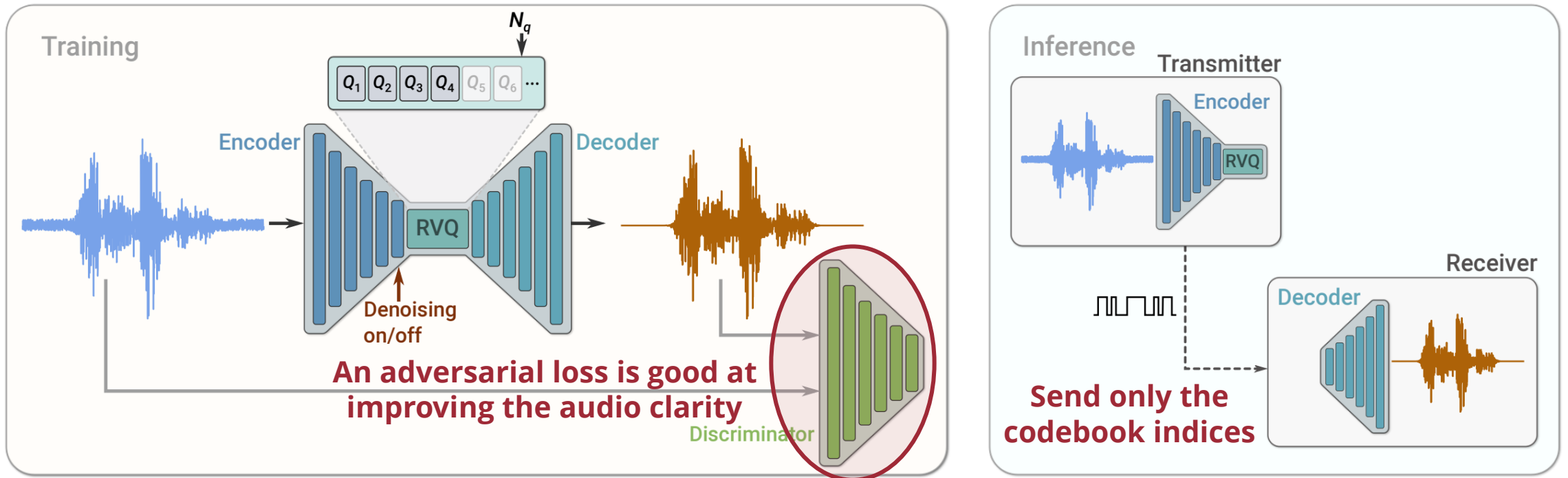


# Example: SoundStream (Zeghidour et al., 2021)



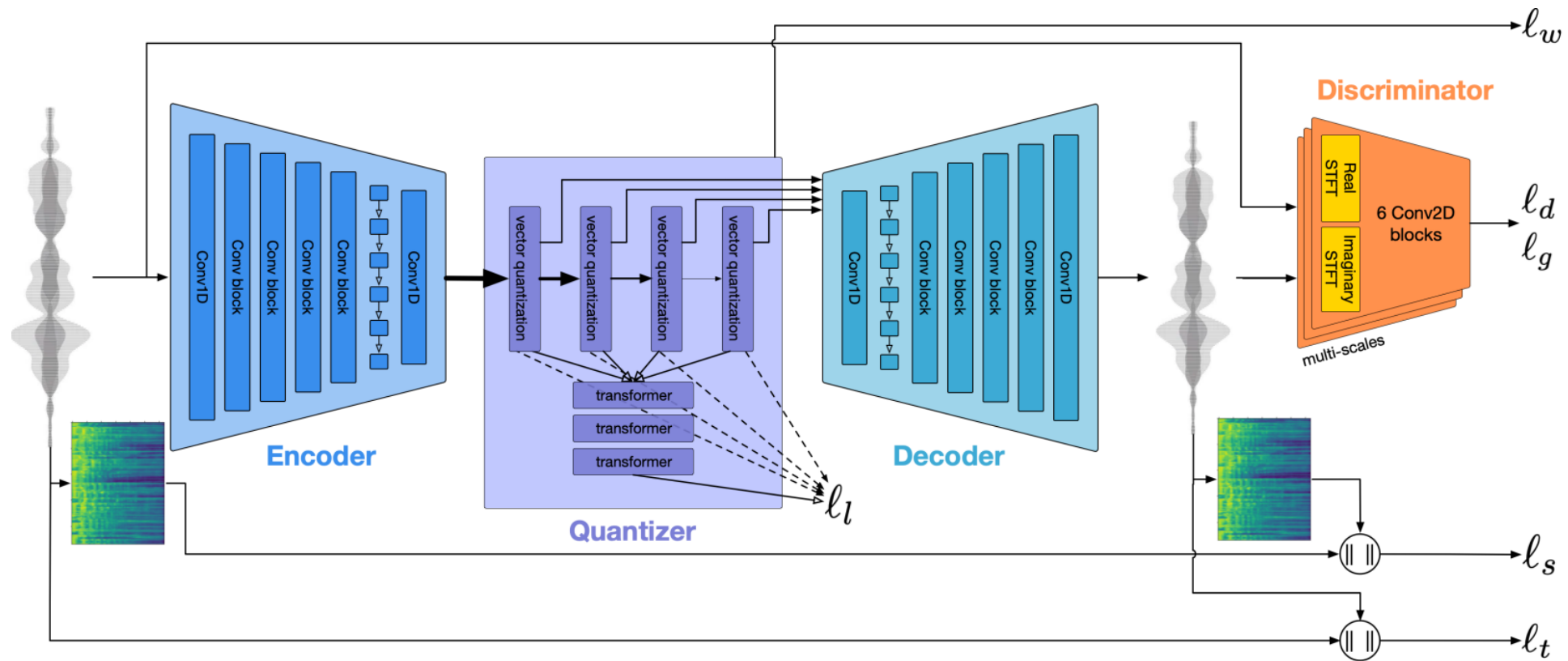
(Source: Zeghidour et al., 2021)

# Example: SoundStream (Zeghidour et al., 2021)



(Source: Zeghidour et al., 2021)

# Example: EnCodec (Défossez et al., 2022)

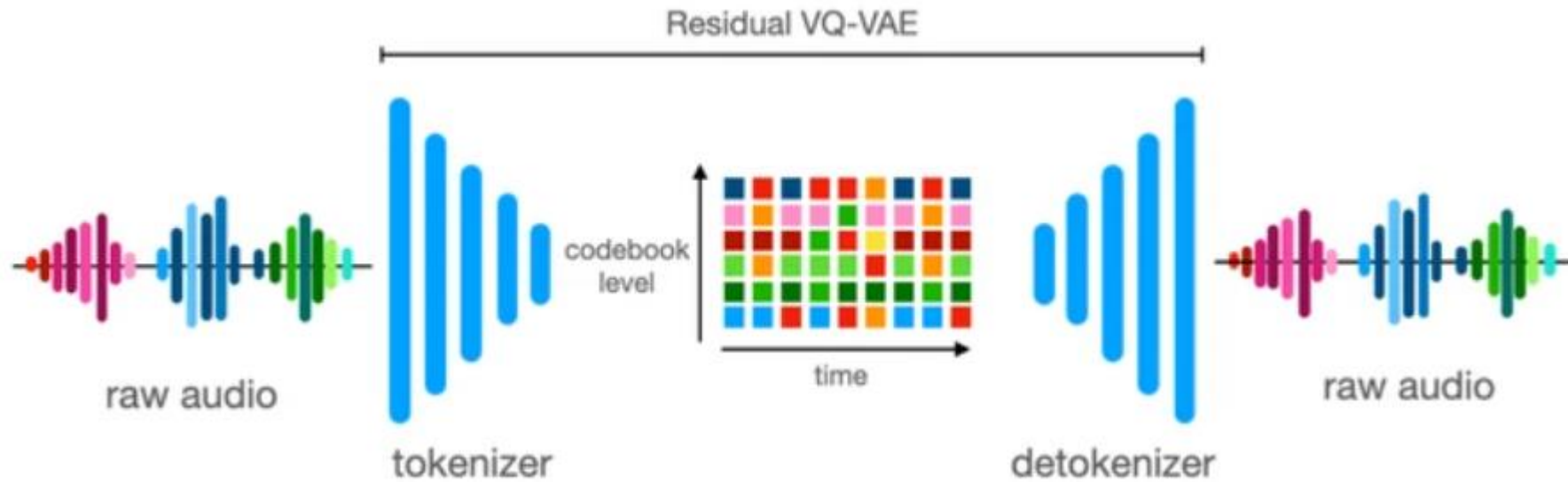


(Source: Défossez et al., 2022)

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

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

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



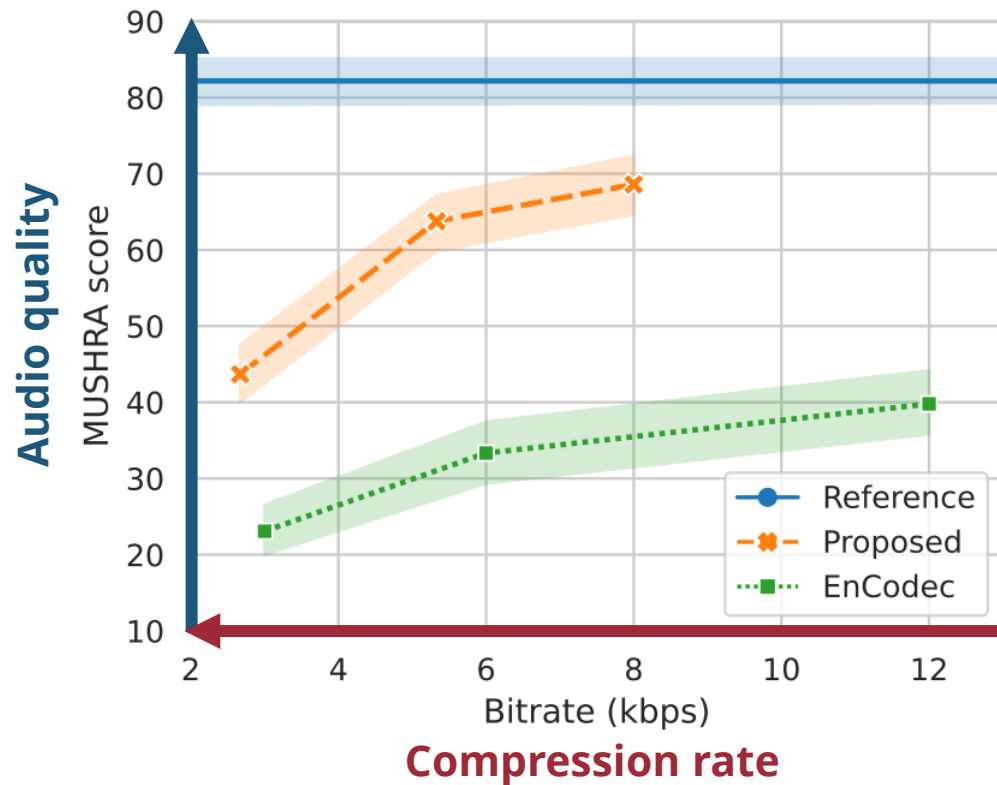
(Source: Kumar et al., 2023)

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

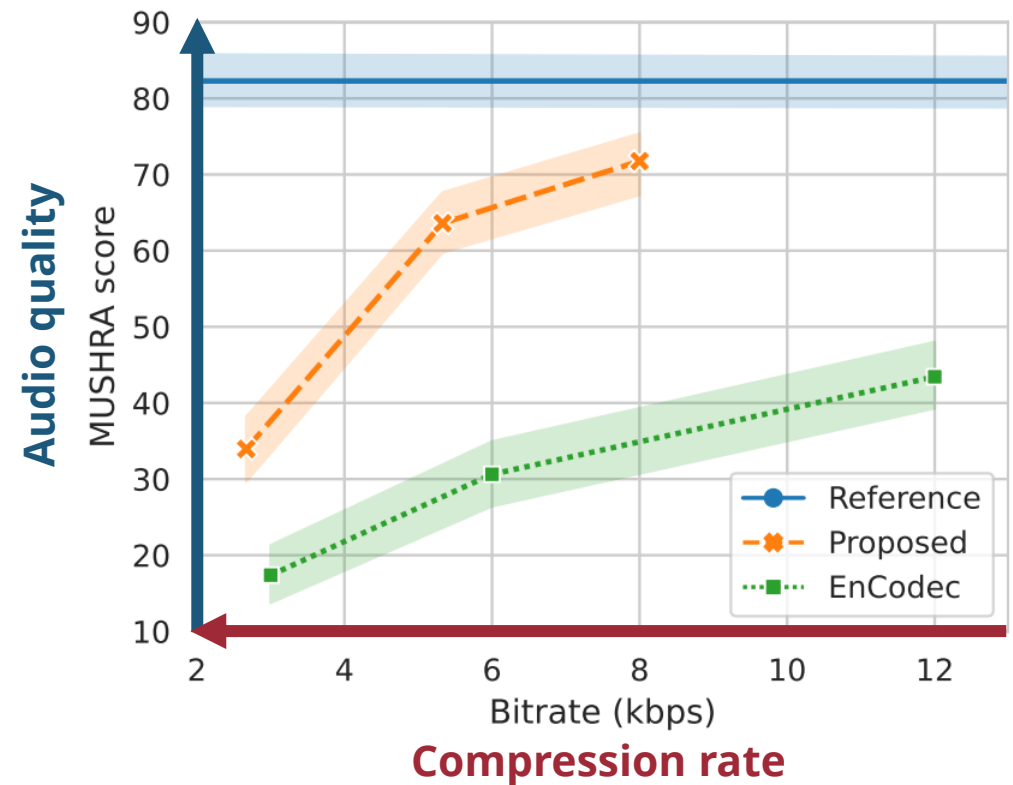
[github.com/descriptinc/descript-audio-codec](https://github.com/descriptinc/descript-audio-codec)

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

Listening Test Results @ 44.1 kHz



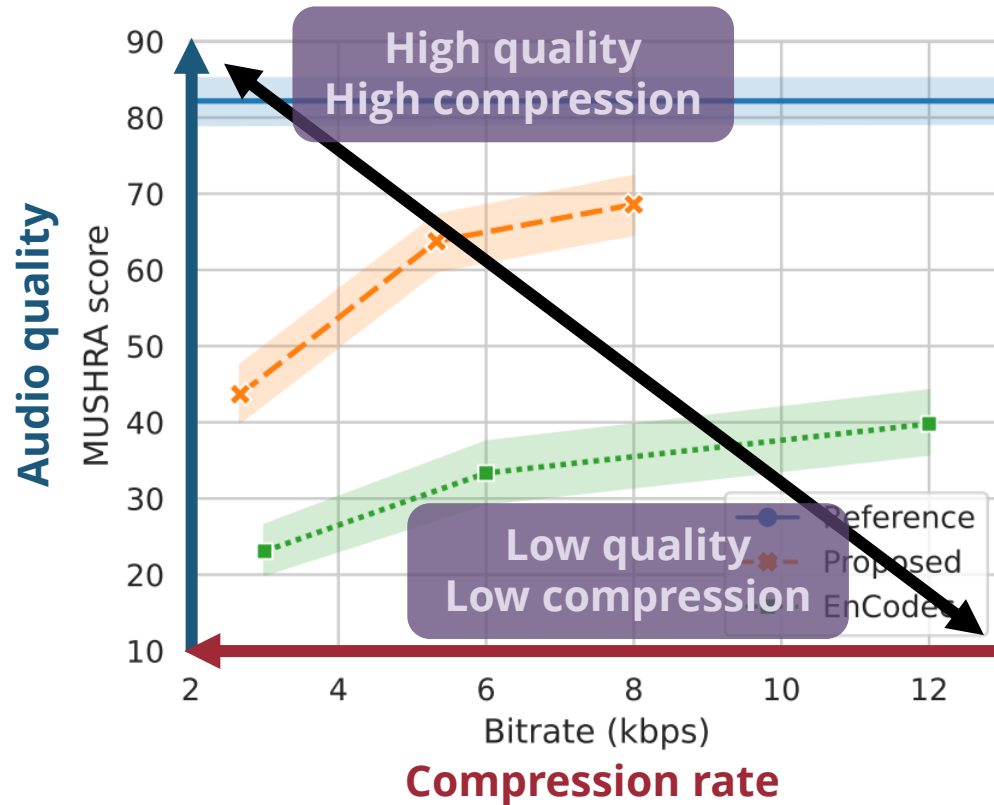
Listening Test Results @ 24 kHz



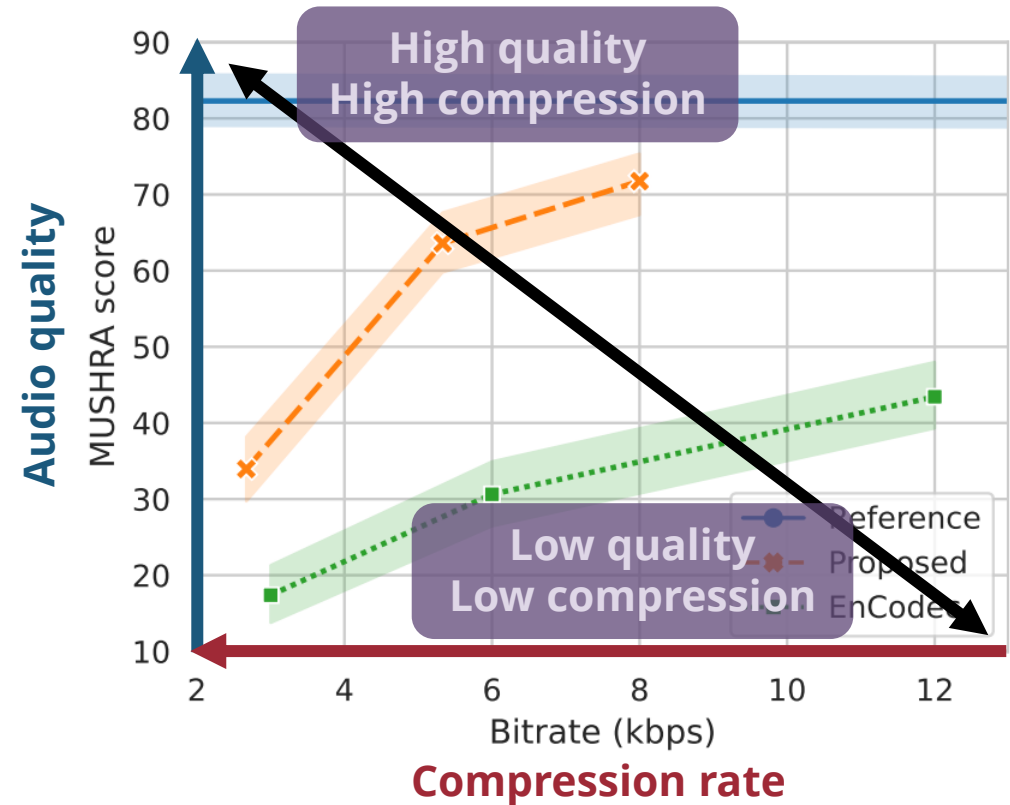
(Source: Kumar et al., 2023)

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

Listening Test Results @ 44.1 kHz



Listening Test Results @ 24 kHz

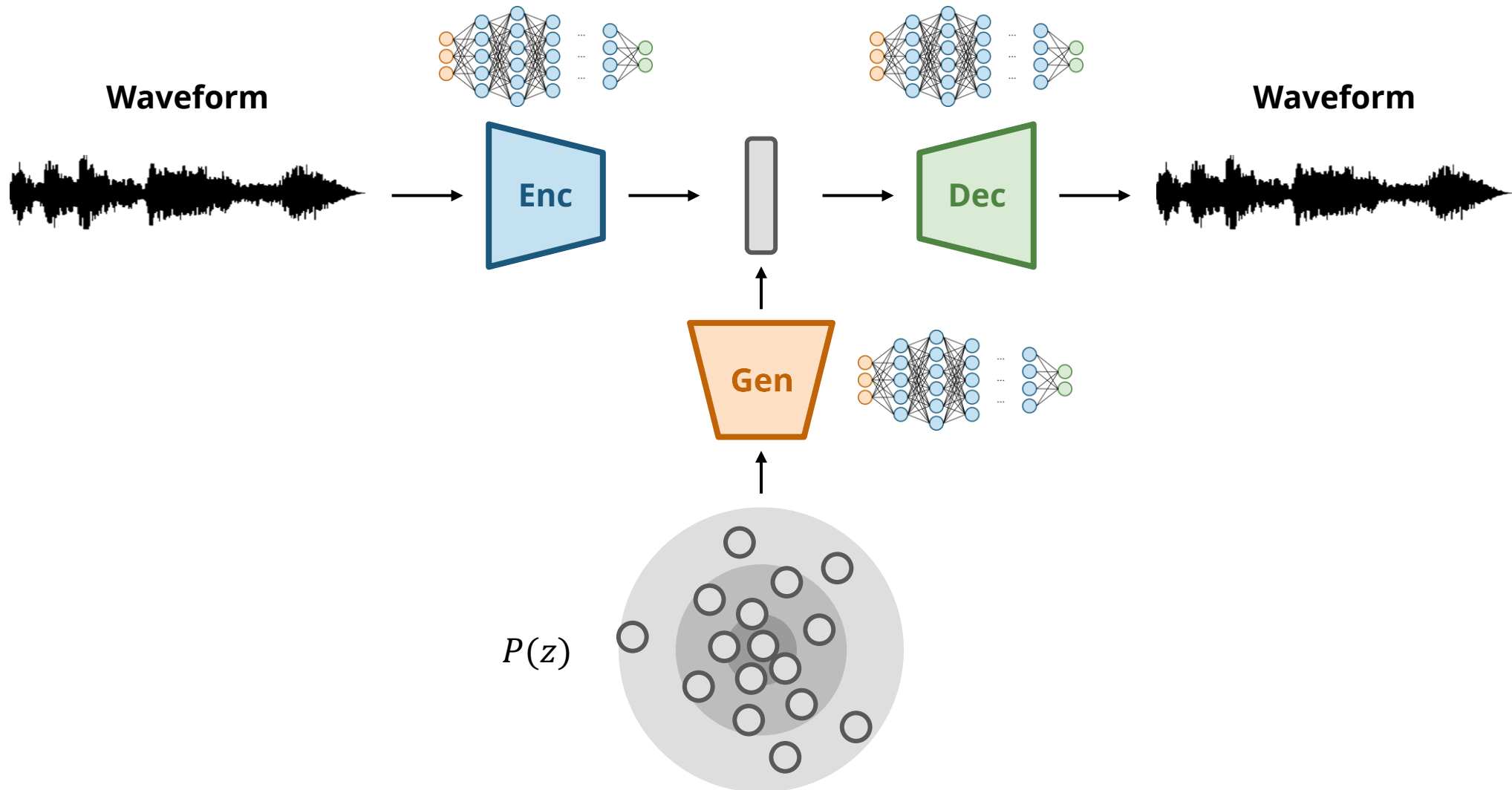


(Source: Kumar et al., 2023)



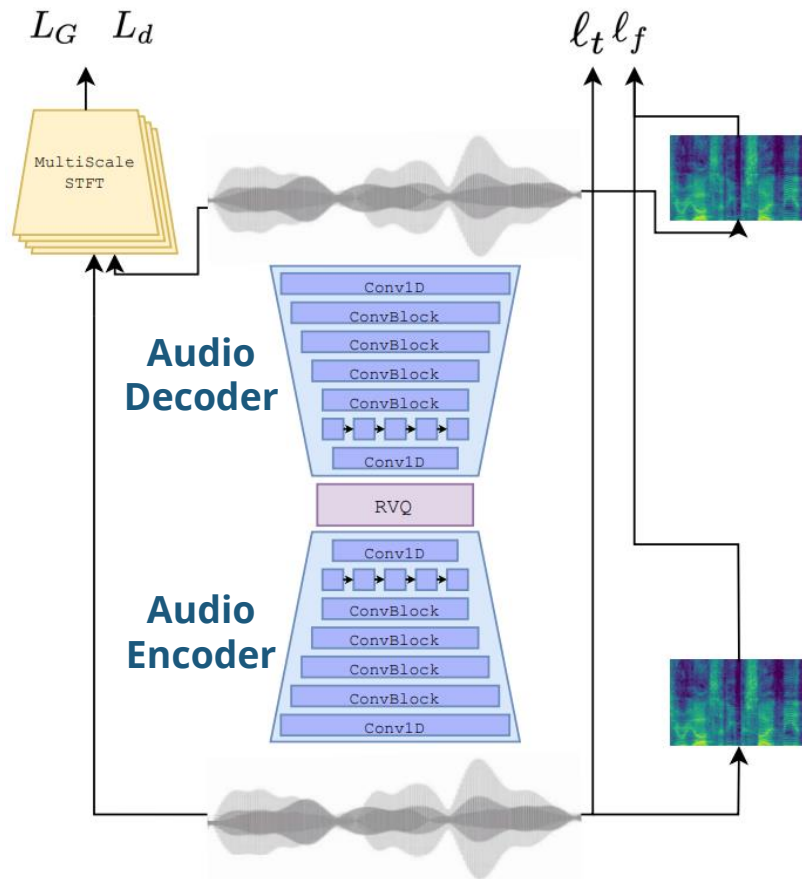
# Latent-based Audio Synthesis

# (Recap) Latent-based Audio Synthesis

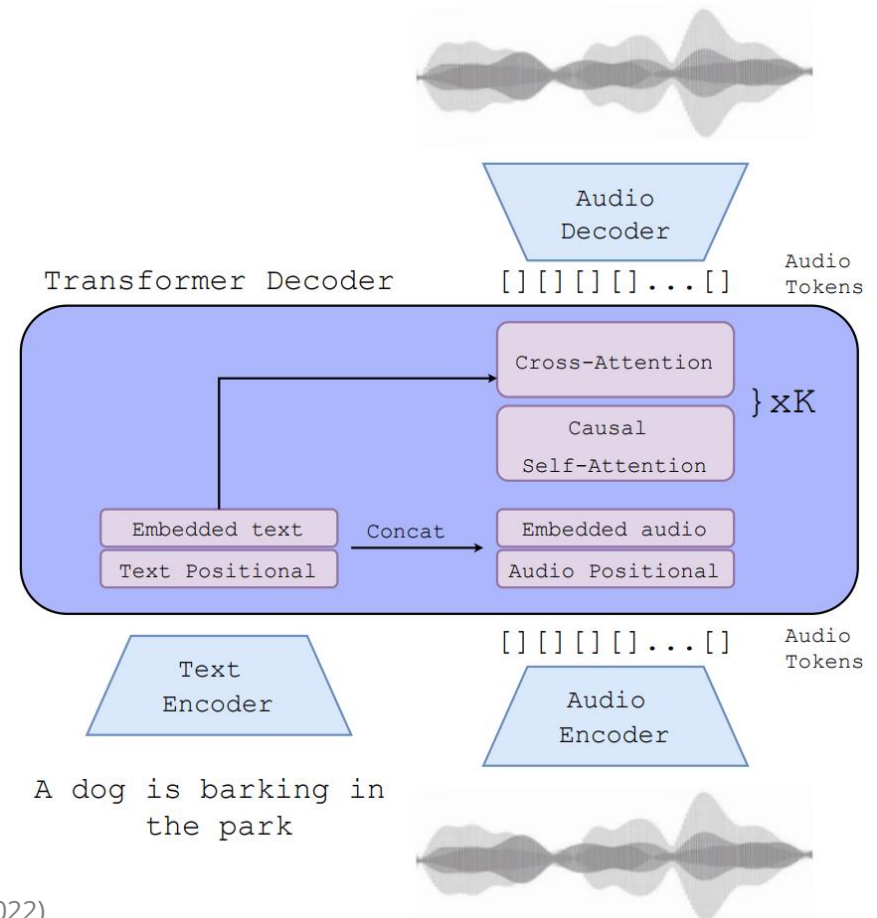


# Example: AudioGen (Kreuk et al., 2023)

## Audio Autoencoder

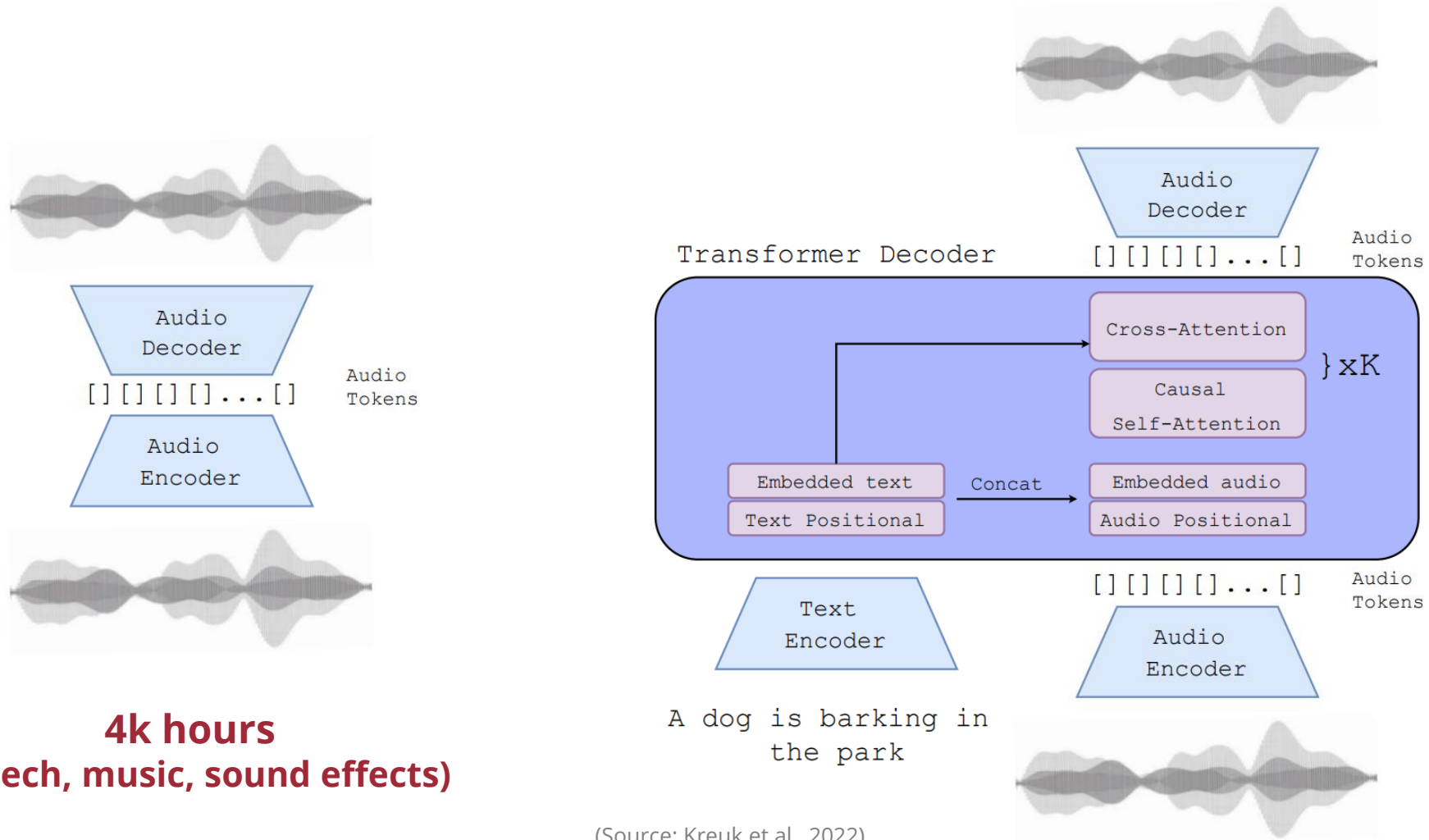


## Audio Language Model



(Source: Kreuk et al., 2022)

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



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

(Source: Kreuk et al., 2022)

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



(Source: Kreuk et al., 2022)

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

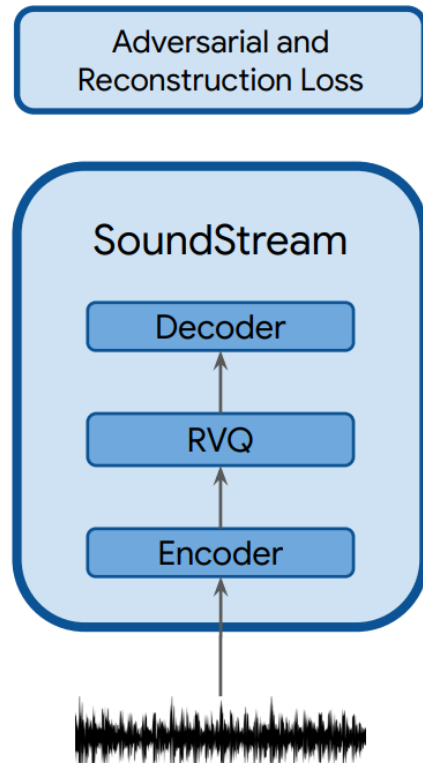
## Example: MusicGen (Copet et al., 2023)

- 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/)

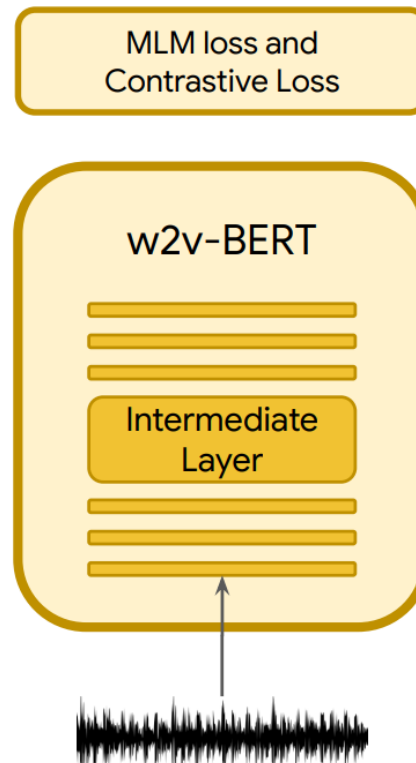
# Example: MusicLM (Agostinelli et al., 2023)

## Audio autoencoder

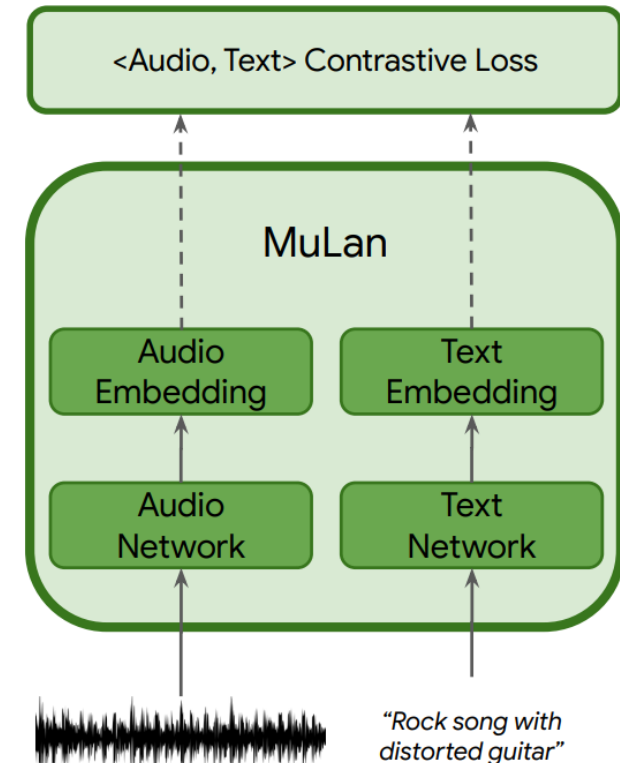


**106k songs, 8.2k hours**

## Semantic representation



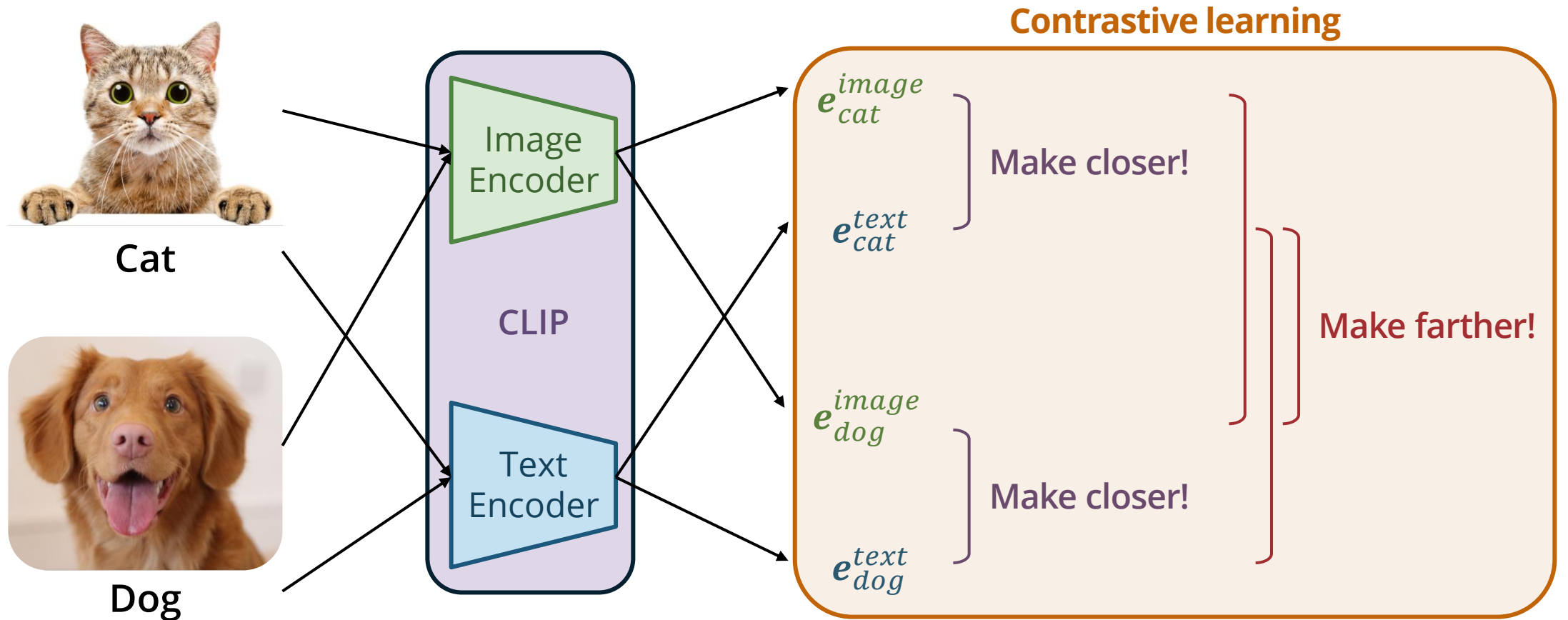
## Text-music correspondence



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

(Source: Agostinelli et al., 2022)

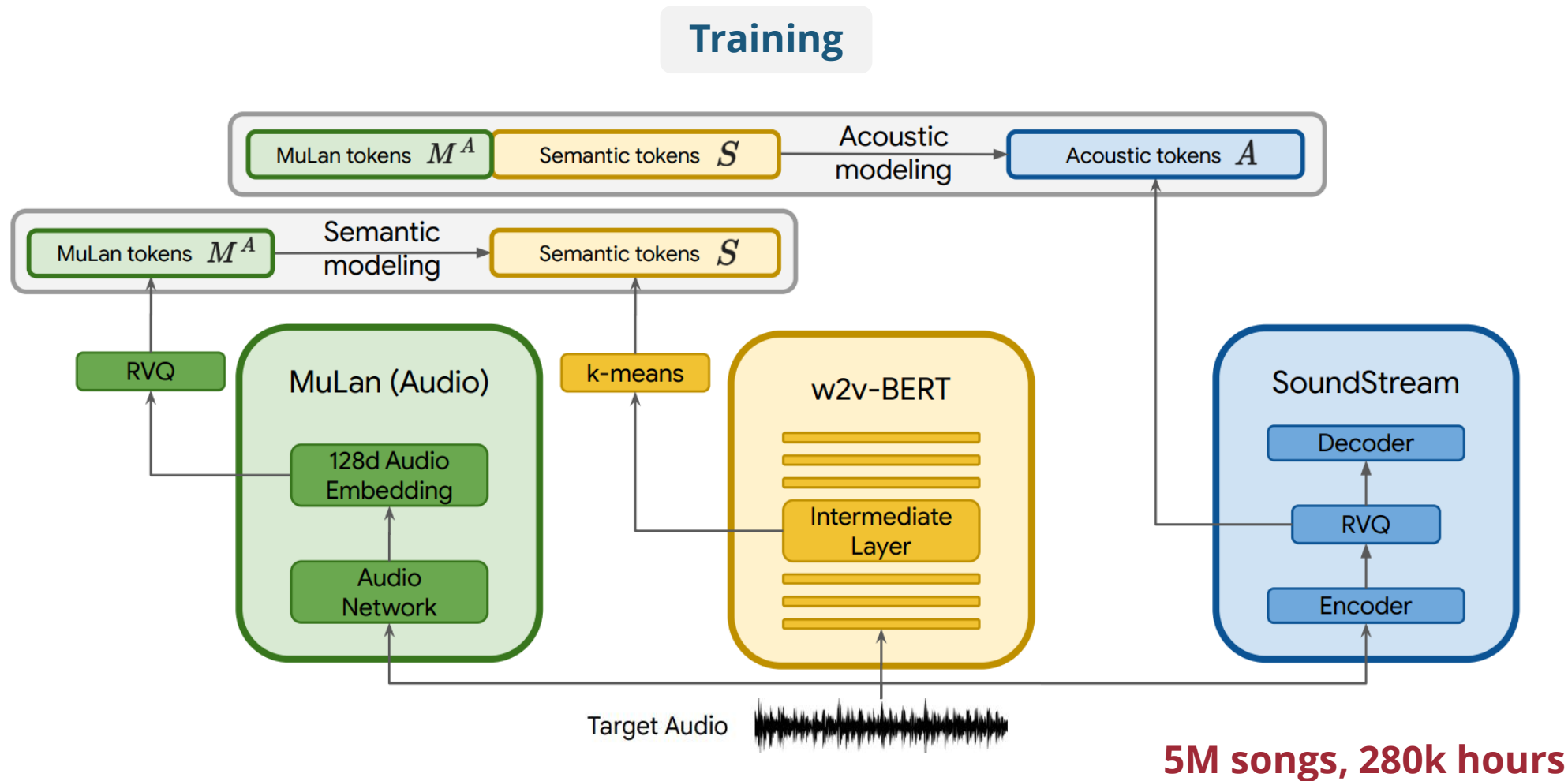
# Contrastive Language-Image Pretraining (CLIP)



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

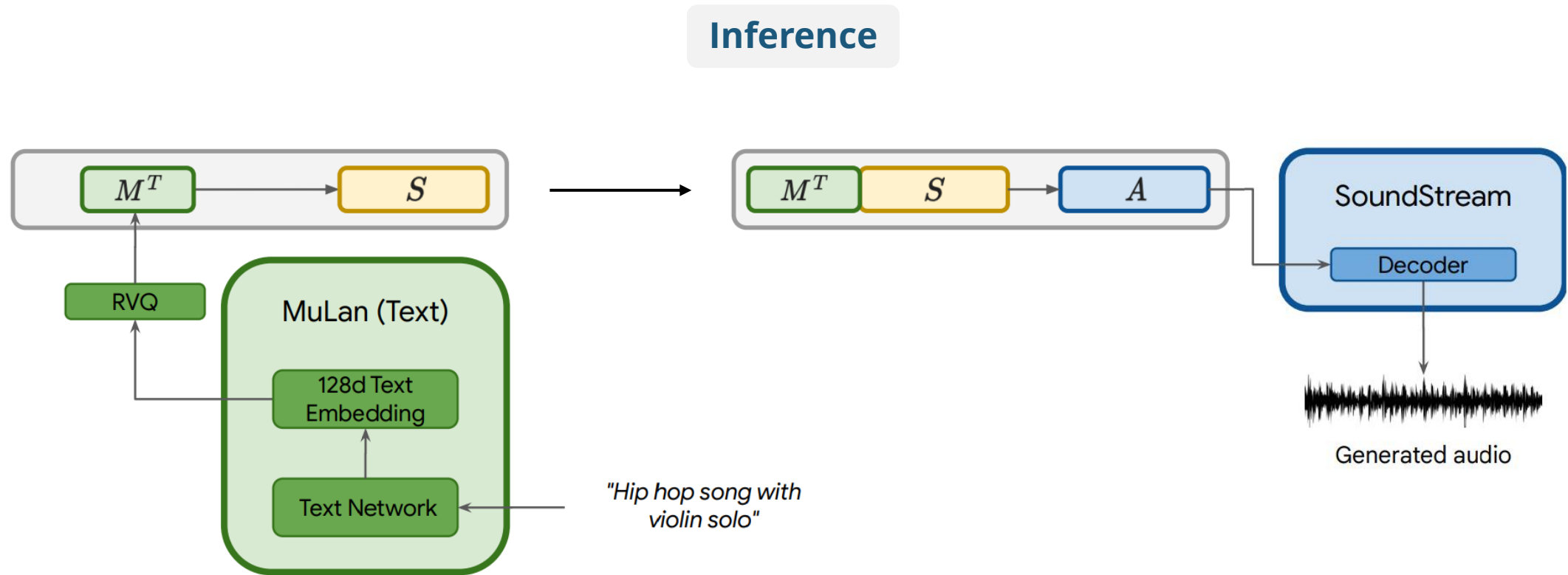


# Example: MusicLM (Agostinelli et al., 2023)



(Source: Agostinelli et al., 2022)

# Example: MusicLM (Agostinelli et al., 2023)



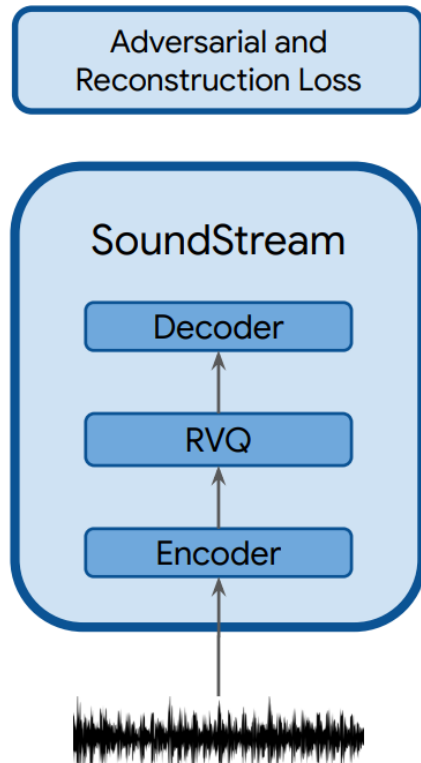
(Source: Agostinelli et al., 2022)

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

# Example: MusicLM (Agostinelli et al., 2023)

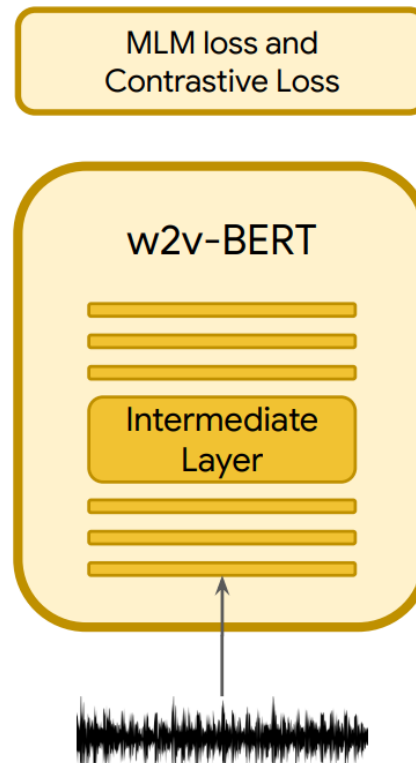
Which is the most challenging component?

## Audio autoencoder

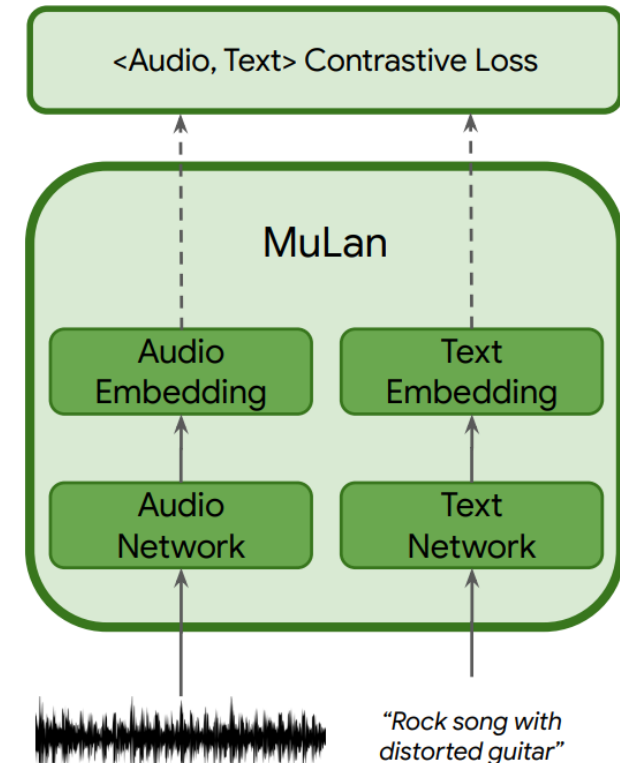


106k songs, 8.2k hours

## Semantic representation



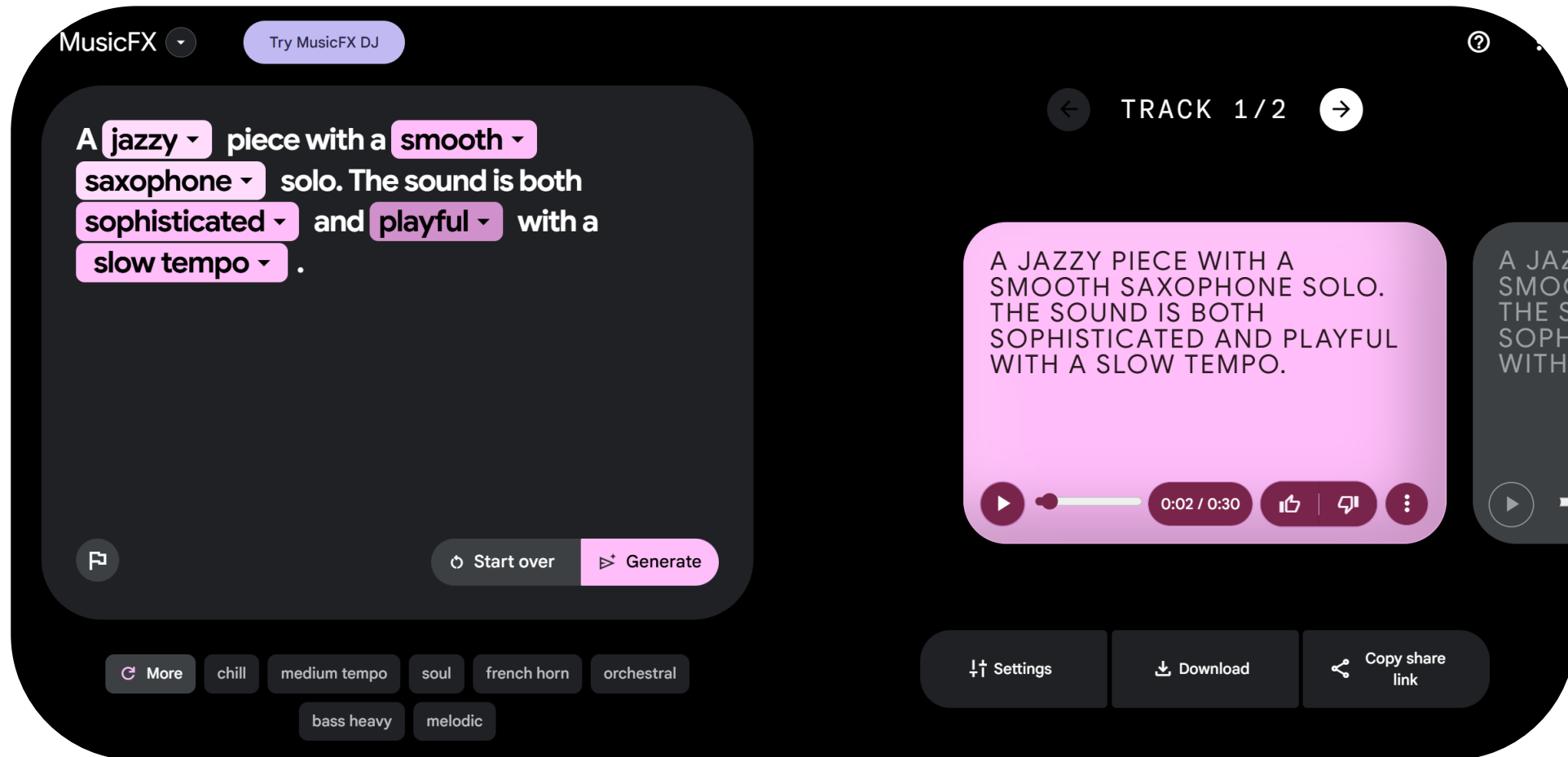
## Text-music correspondence



44M 30-sec clips, 370k hours

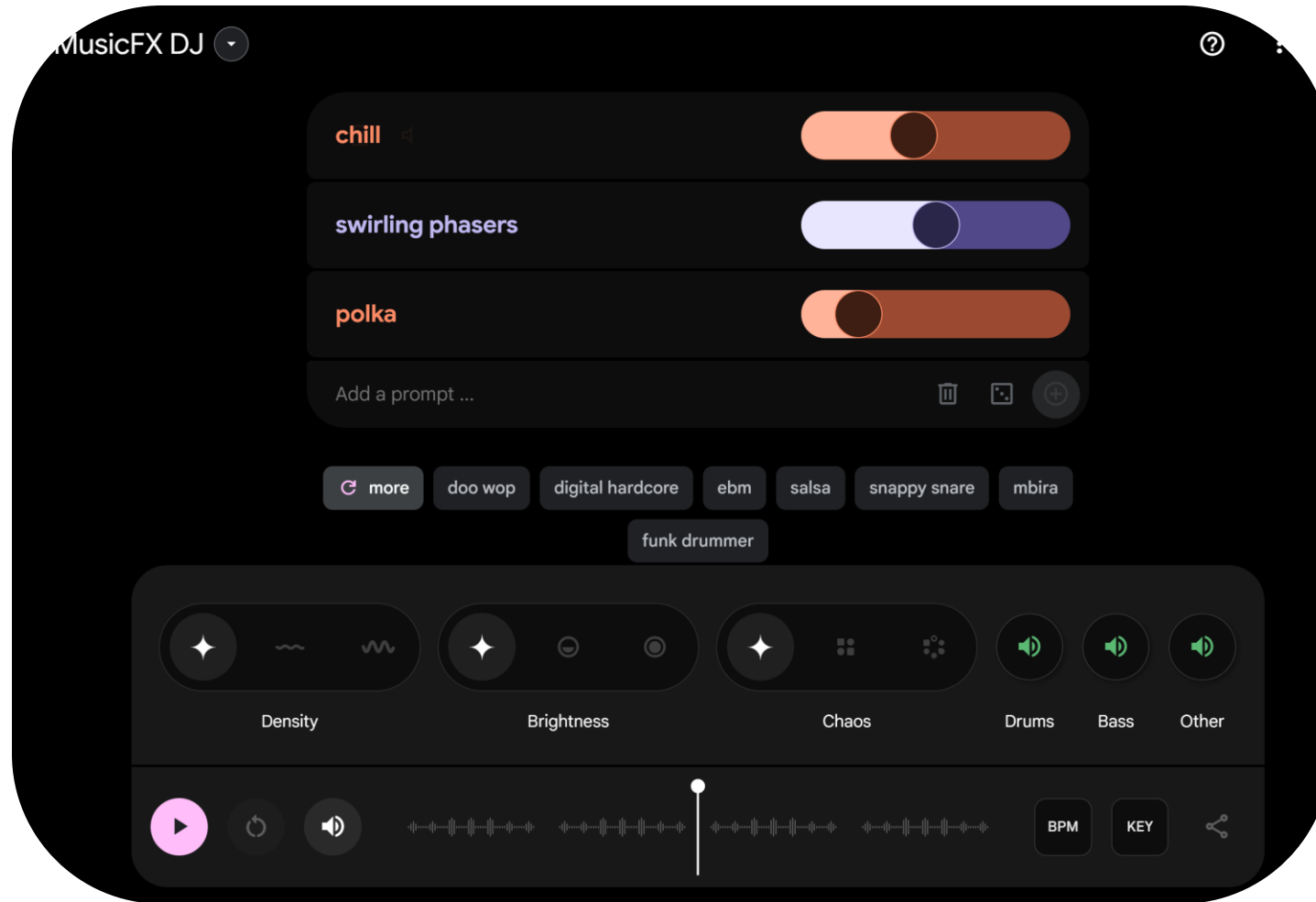
(Source: Agostinelli et al., 2022)

# Music FX (2024)



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

# Music FX DJ (2024)



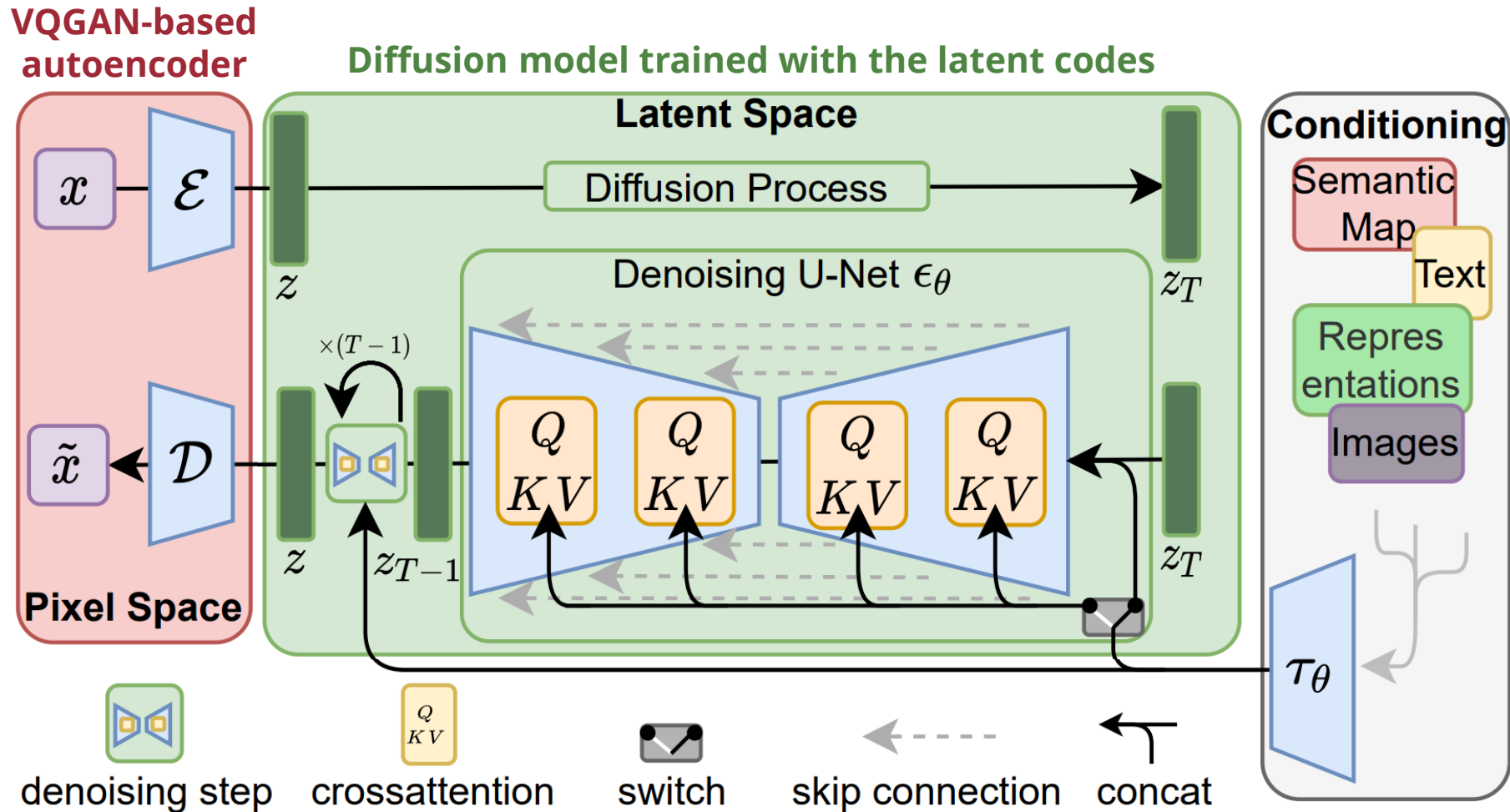
[aitestkitchen.withgoogle.com/tools/music-fx-dj](https://aitestkitchen.withgoogle.com/tools/music-fx-dj)

# Music FX DJ (2024)



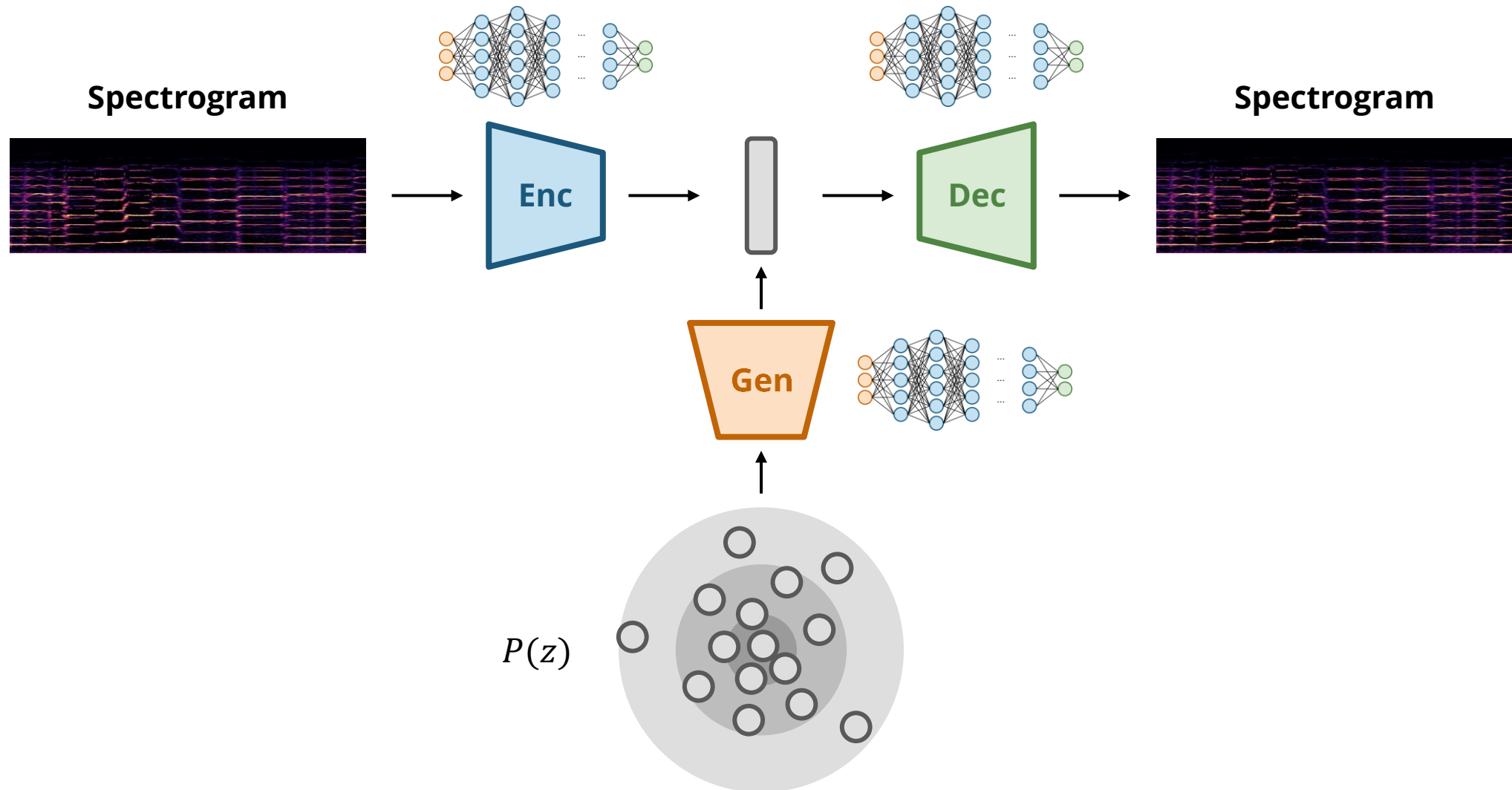
[youtube.com/live/IUQW5LgBZvQ](https://youtube.com/live/IUQW5LgBZvQ)

# (Recap) Latent Diffusion Models (LDMs)



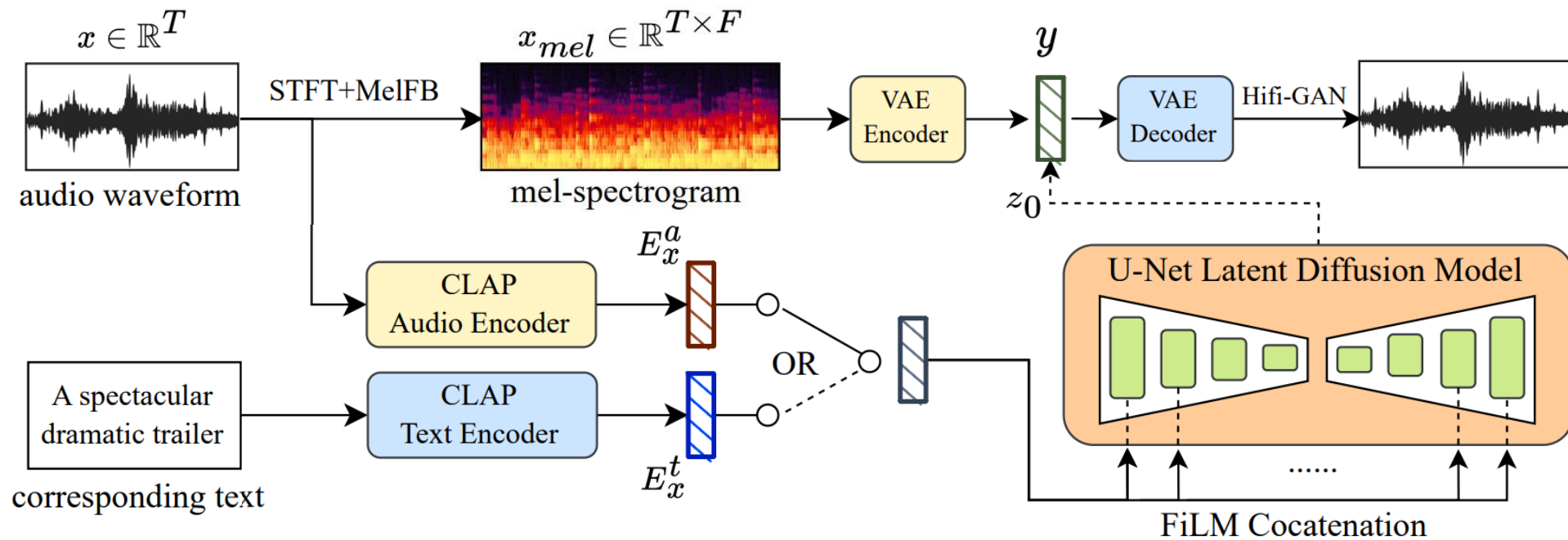
(Source: Rombach et al., 2022)

# (Recap) Latent-based Audio Synthesis





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



(Source: Ke et al., 2023)

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

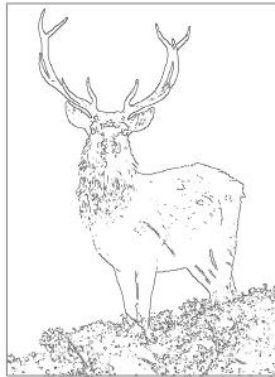
# Example: MusicLDM (Chen et al., 2023)



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

# Controlling Music Generation Systems

# ControlNet (Zhang et al., 2023)



Input Canny edge



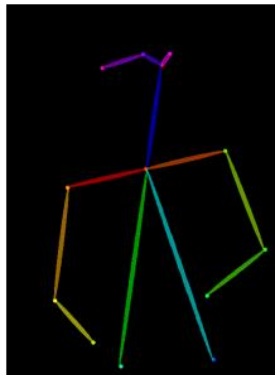
Default



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



“..., quaint city Galic”



Input human pose



Default



“chef in kitchen”

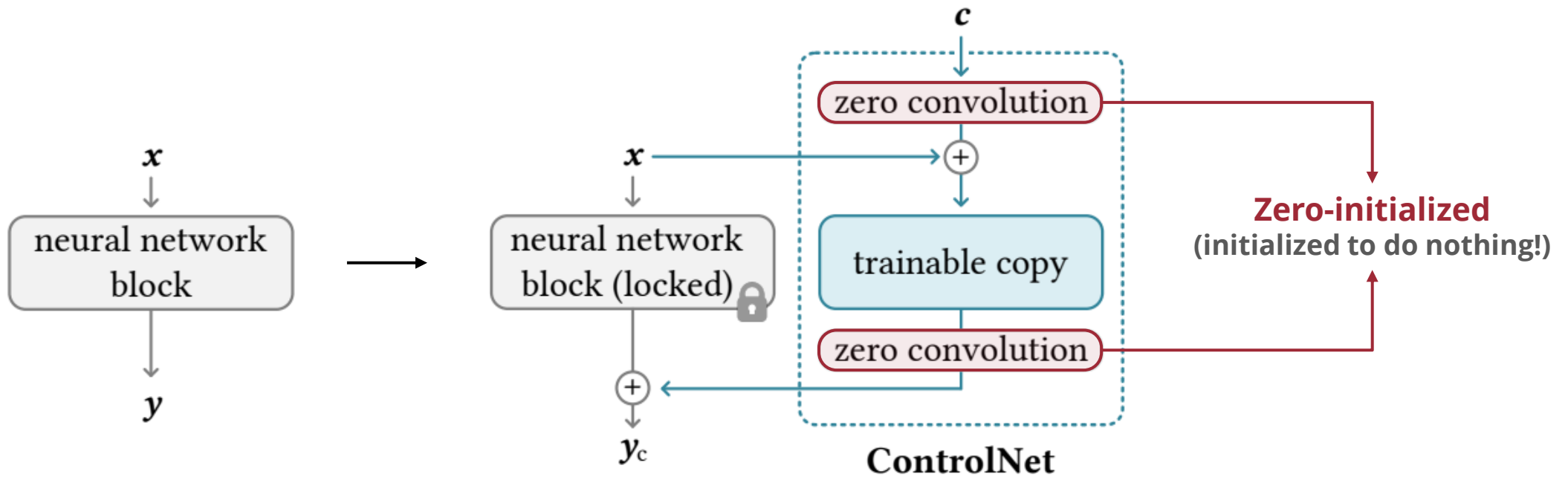


“Lincoln statue”

(Source: Zhang et al., 2023)

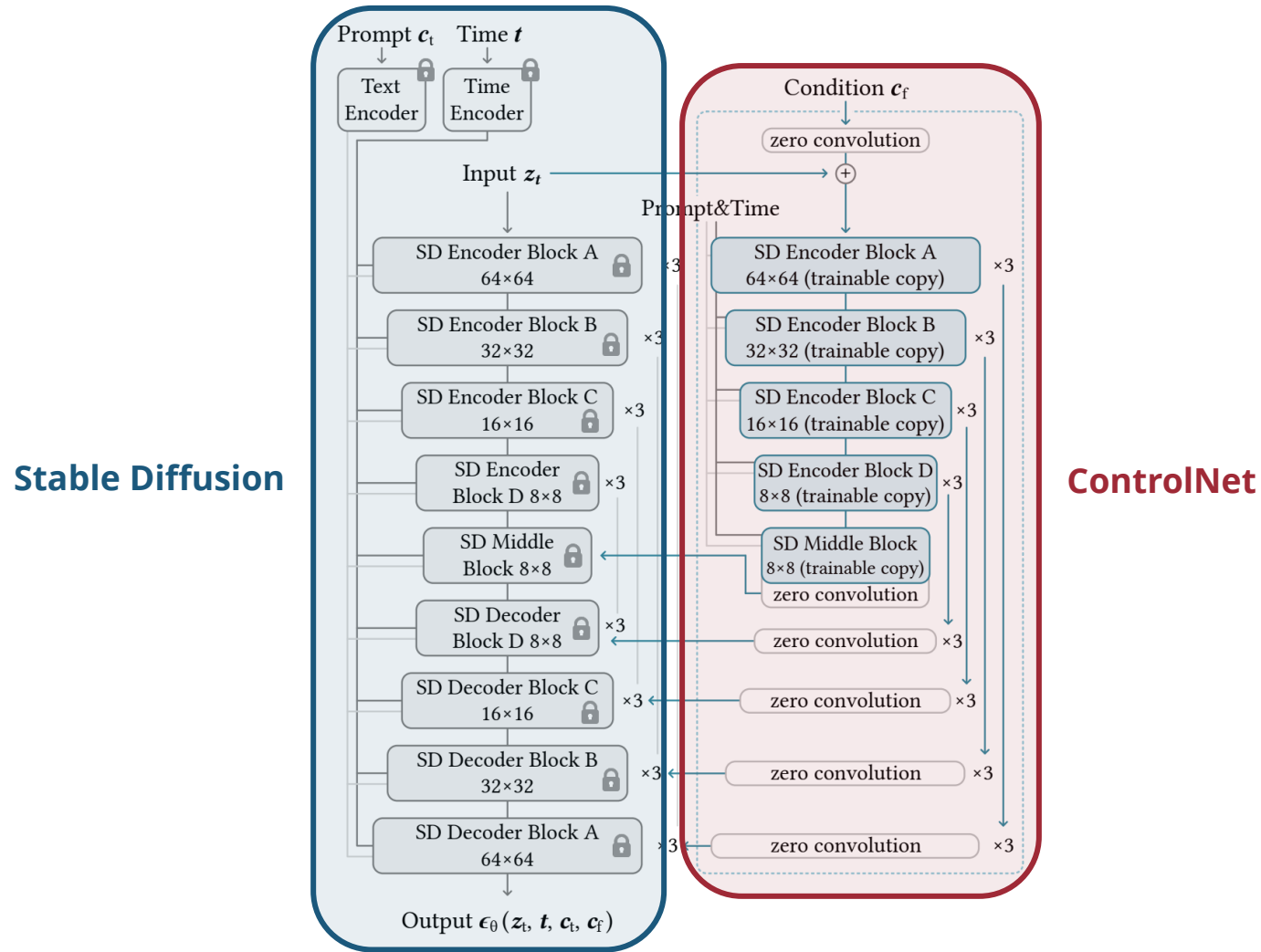
Can we **add controls** to a trained text-to-image diffusion model?

# ControlNet (Zhang et al., 2023)



(Source: Zhang et al., 2023)

# ControlNet (Zhang et al., 2023)



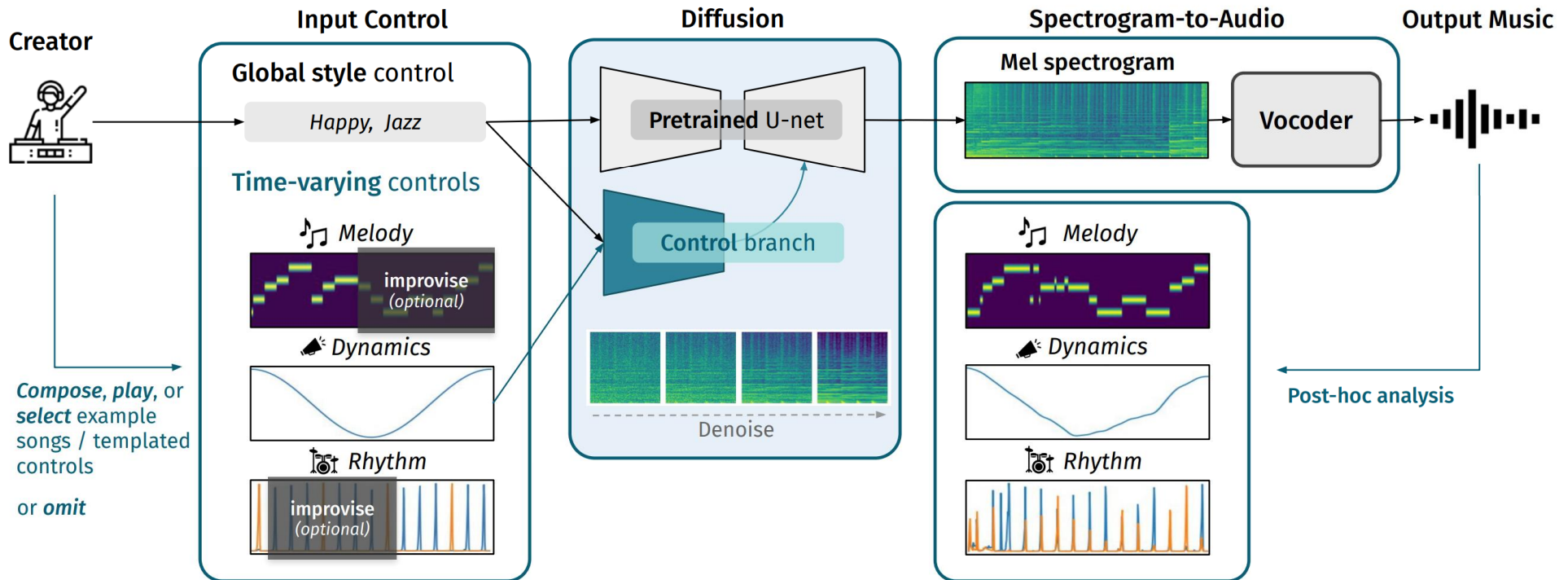
(Source: Zhang et al., 2023)

# 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)

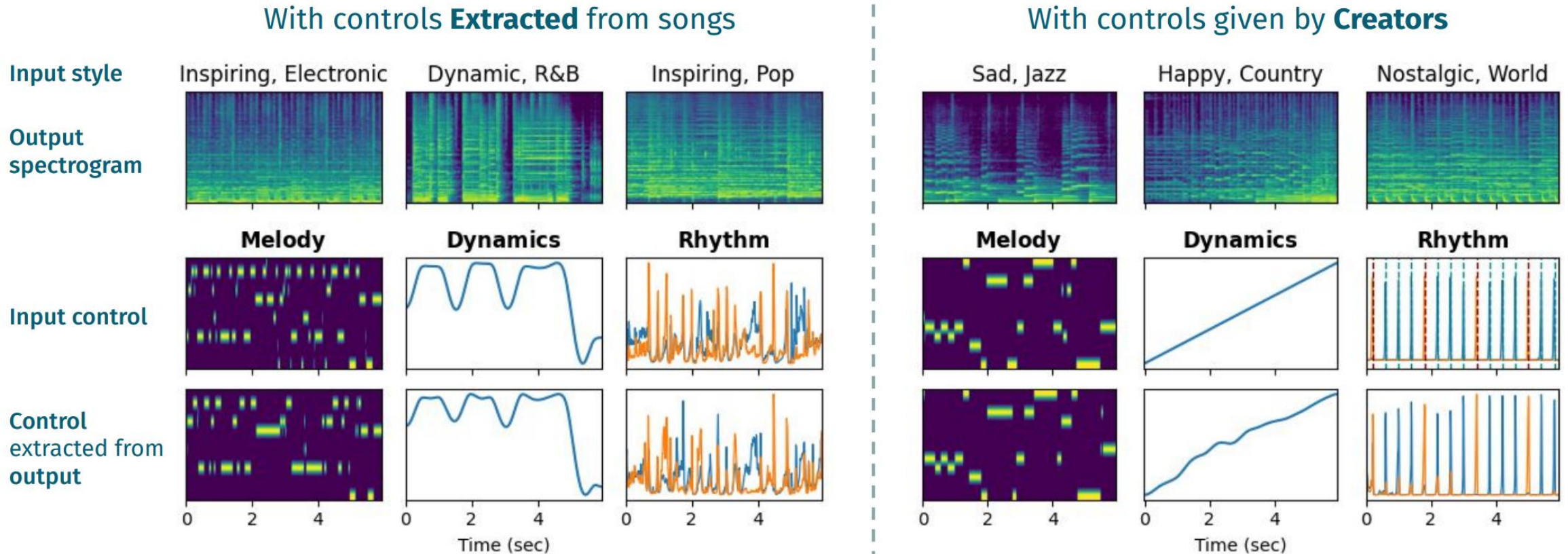
# Example: Music ControlNet (Wu et al., 2024)



(Source: Wu et al., 2024)



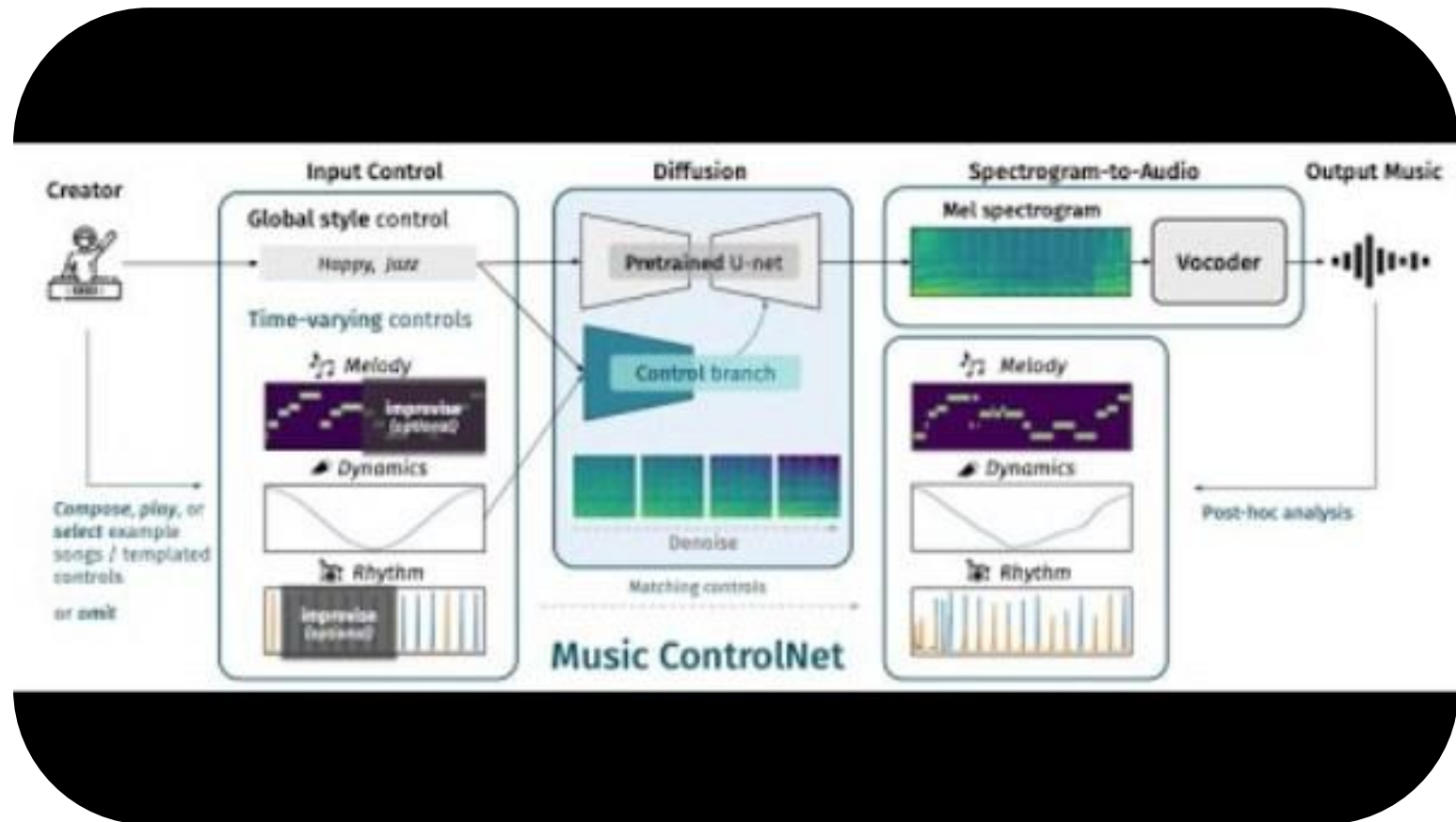
# Example: Music ControlNet (Wu et al., 2024)



(Source: Wu et al., 2024)

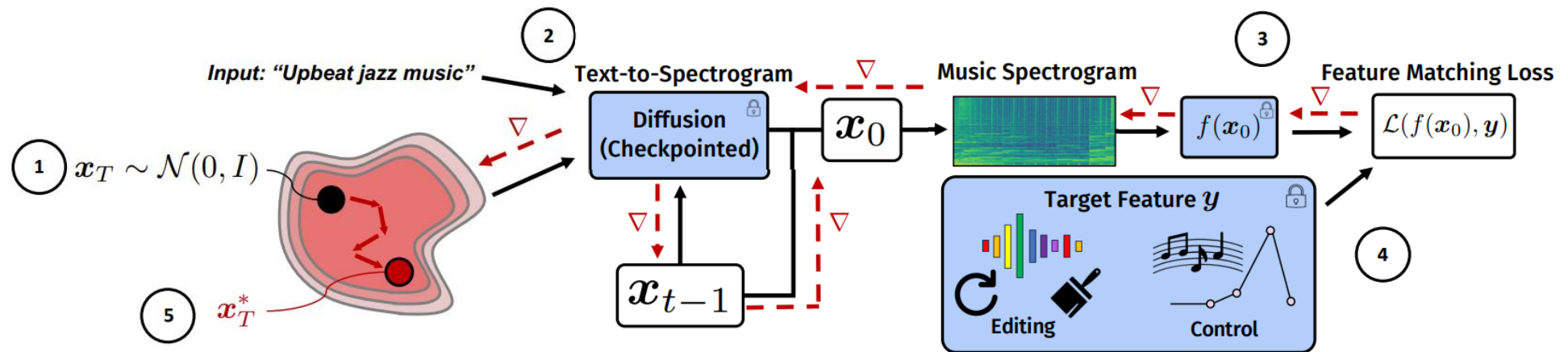
[musiccontrolnet.github.io/web](https://musiccontrolnet.github.io/web)

# Example: Music ControlNet (Wu et al., 2024)



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

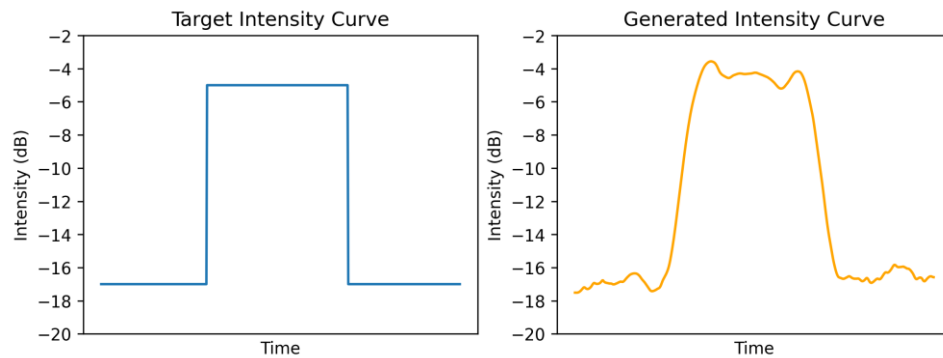
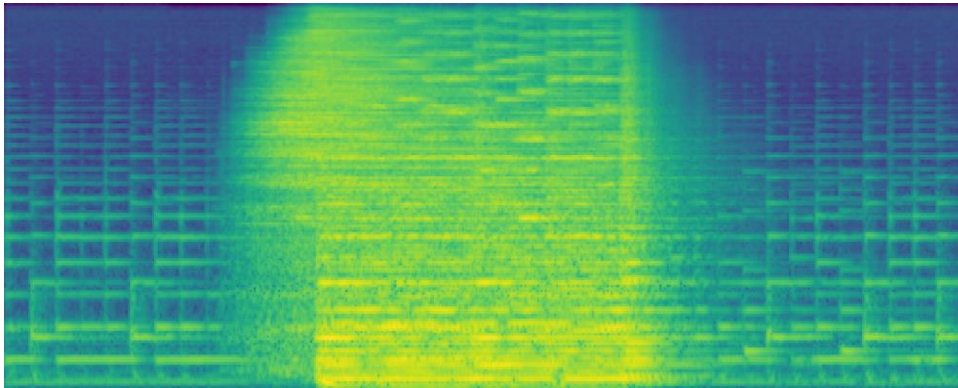
# Example: DITTO (Novack et al., 2024)



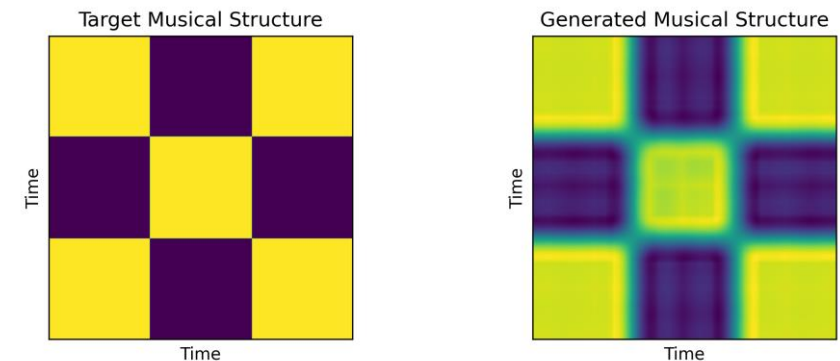
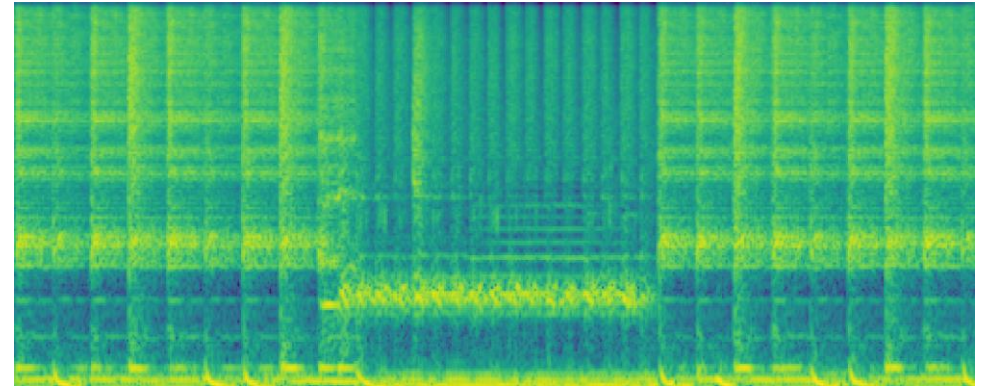
(Source: Novack et al., 2024)

# Example: DITTO (Novack et al., 2024)

## Intensity control



## Structure control



(Source: Novack et al., 2024)

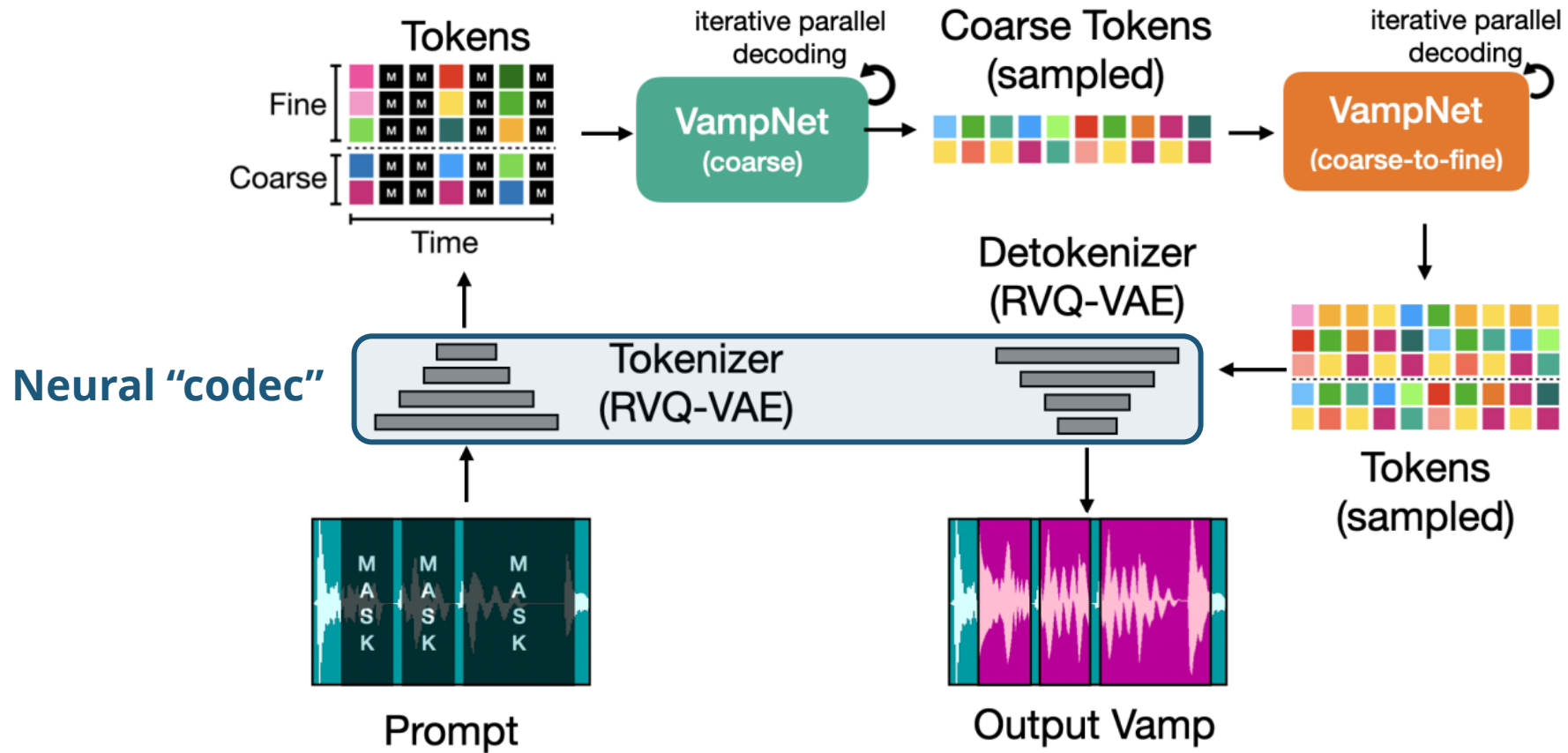
# DITTO: Diffusion for Music Generation (2024)



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

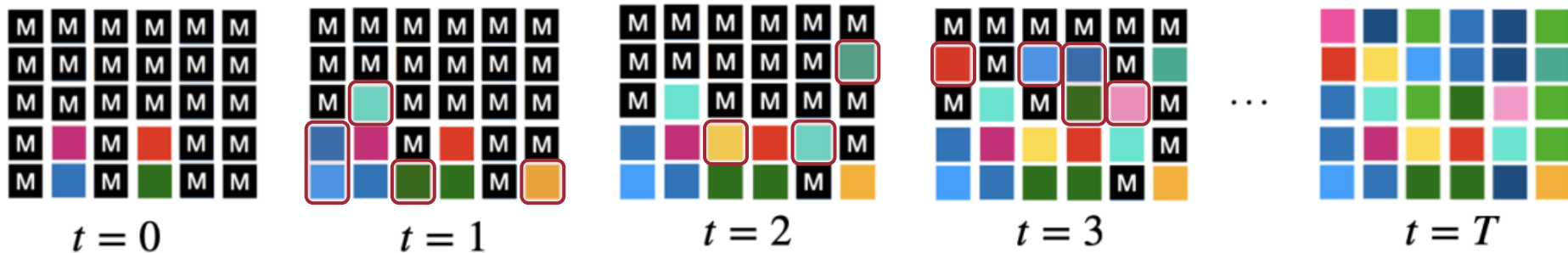
# Creative Applications of Music Generation Systems

# Example: VampNet (Garcia et al., 2023)



(Source: Garcia et al., 2023)

# Example: VampNet (Garcia et al., 2023)

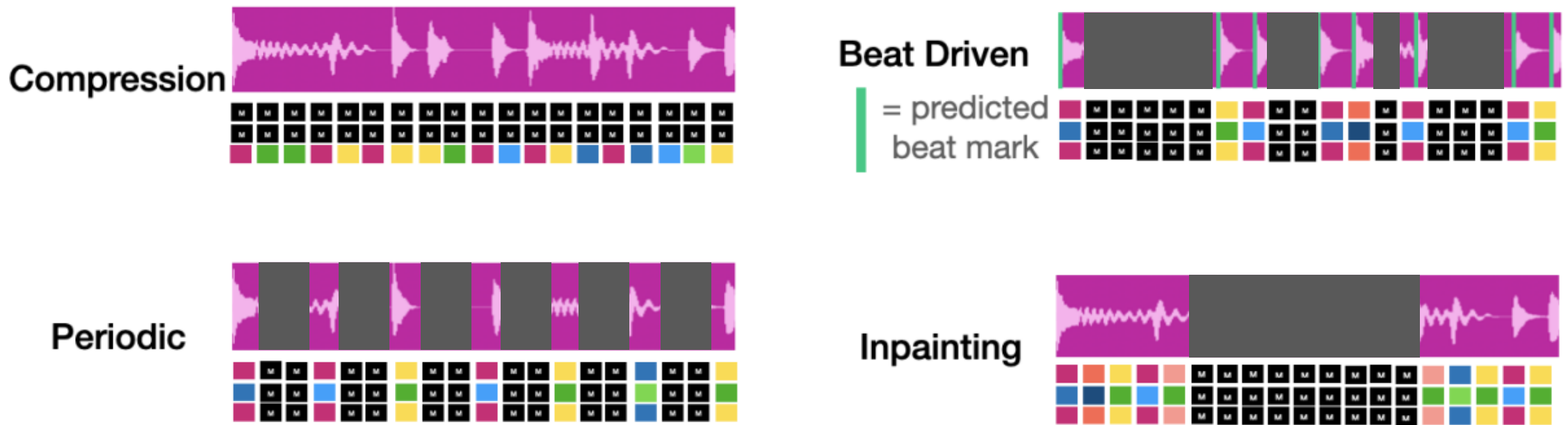


(Source: Garcia et al., 2023)

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

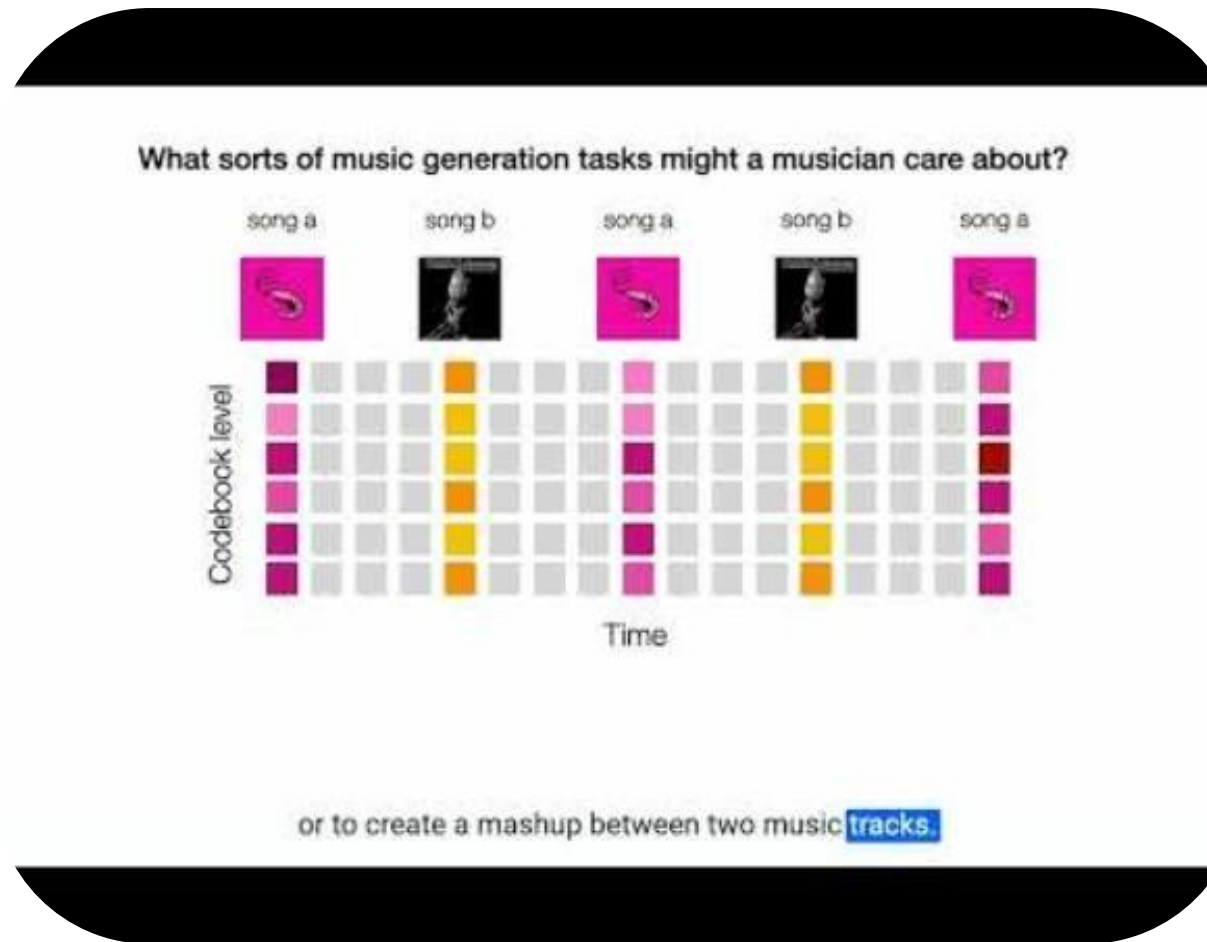


# Example: VampNet (Garcia et al., 2023)



(Source: Garcia et al., 2023)

# Example: VampNet (Garcia et al., 2023)



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

Example: **unloop** (Garcia et al., 2023)



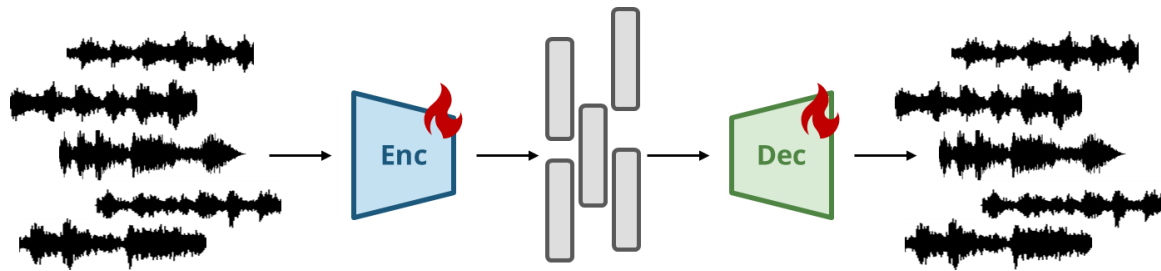
[youtu.be/yzBI8Vcjd2s](https://youtu.be/yzBI8Vcjd2s)

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

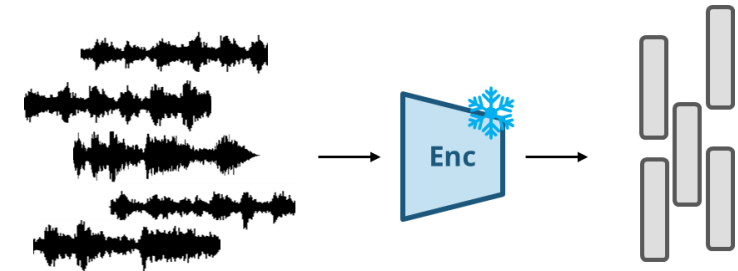


# Pipeline

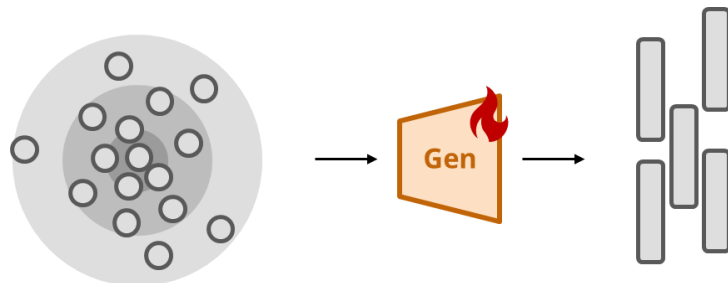
Step 1: Train an Autoencoder



Step 2: Compute the Latent Vectors



Step 3: Train a Latent Generative Model



Step 4: Decode the Latent Vectors

