

PAT 498/598 (Winter 2025)

# Music & AI

## **Lecture 18: Music Search & Recommendation**

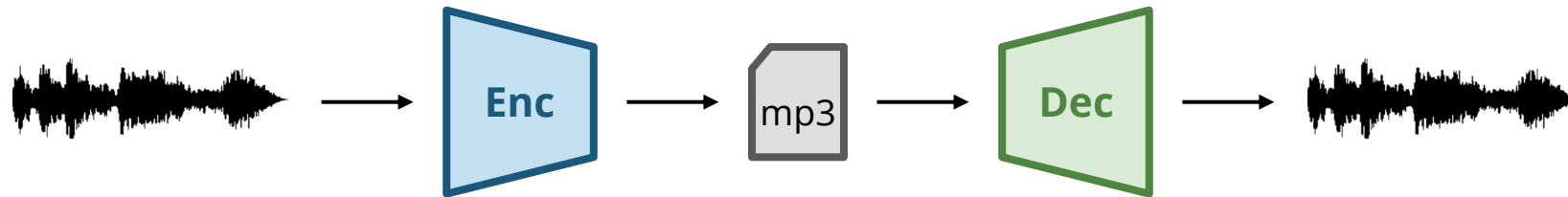
Instructor: Hao-Wen Dong



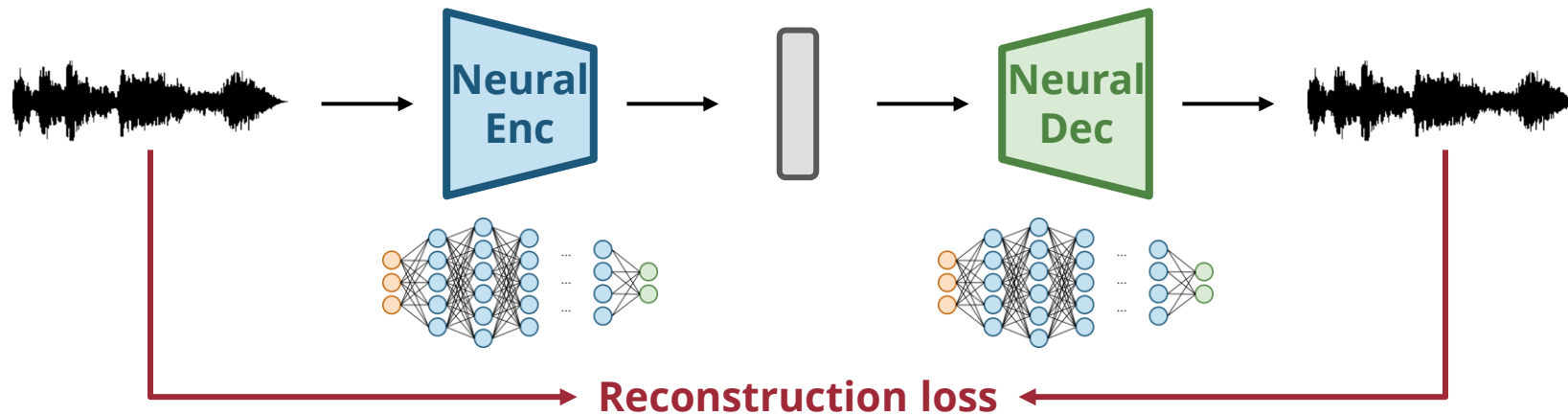
SCHOOL OF MUSIC, THEATRE & DANCE  
PERFORMING ARTS TECHNOLOGY  
UNIVERSITY OF MICHIGAN

# (Recap) Neural Codec

## Traditional Codec

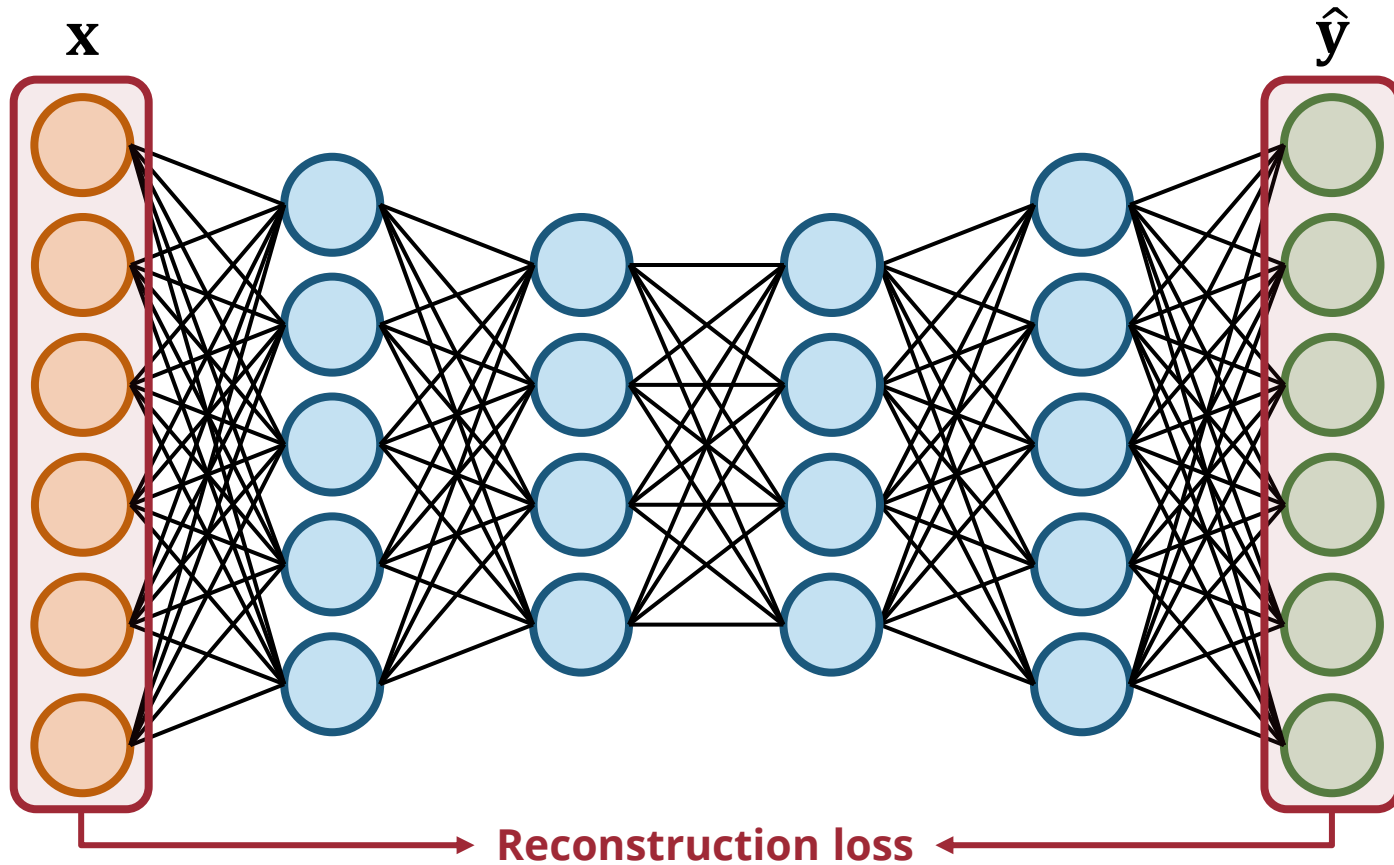


## Neural Codec

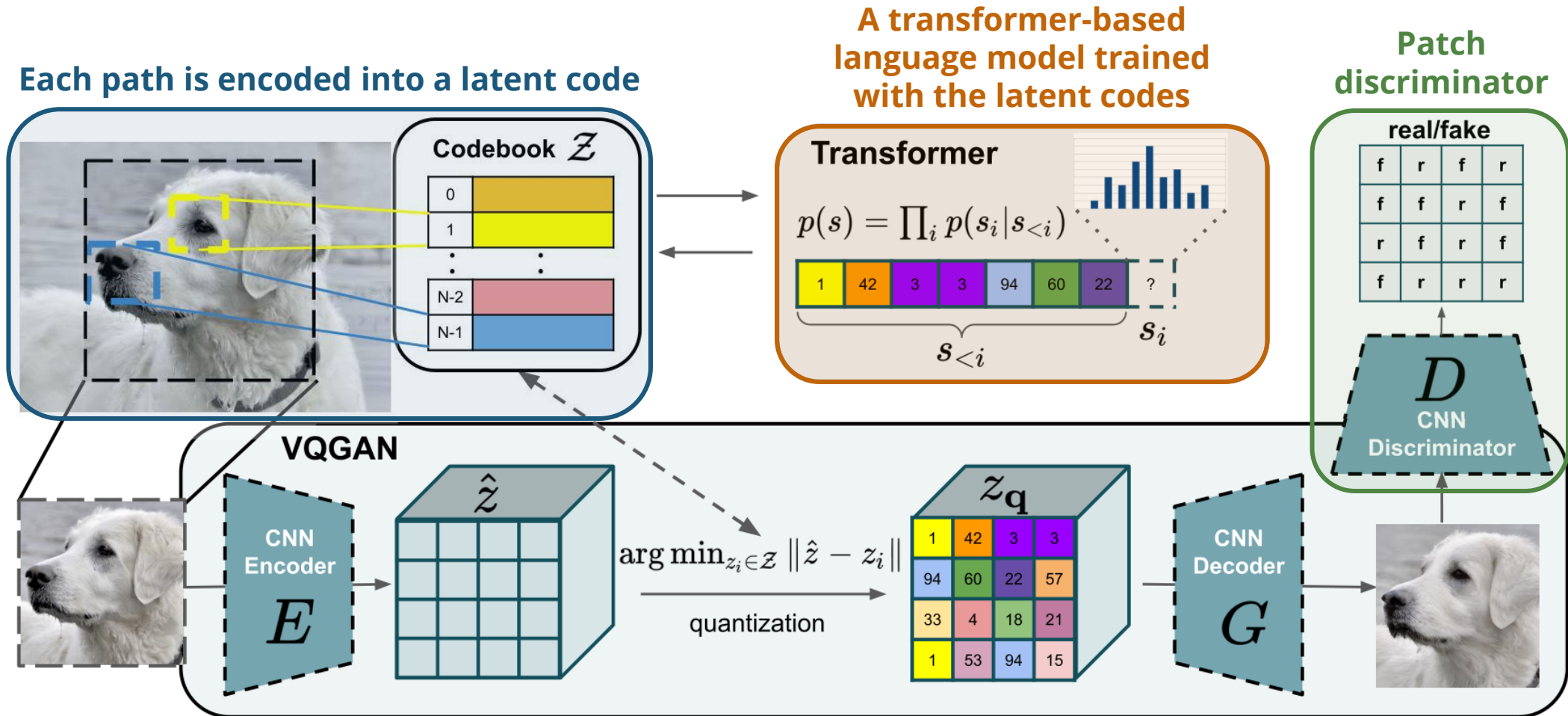


# (Recap) Autoencoders

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



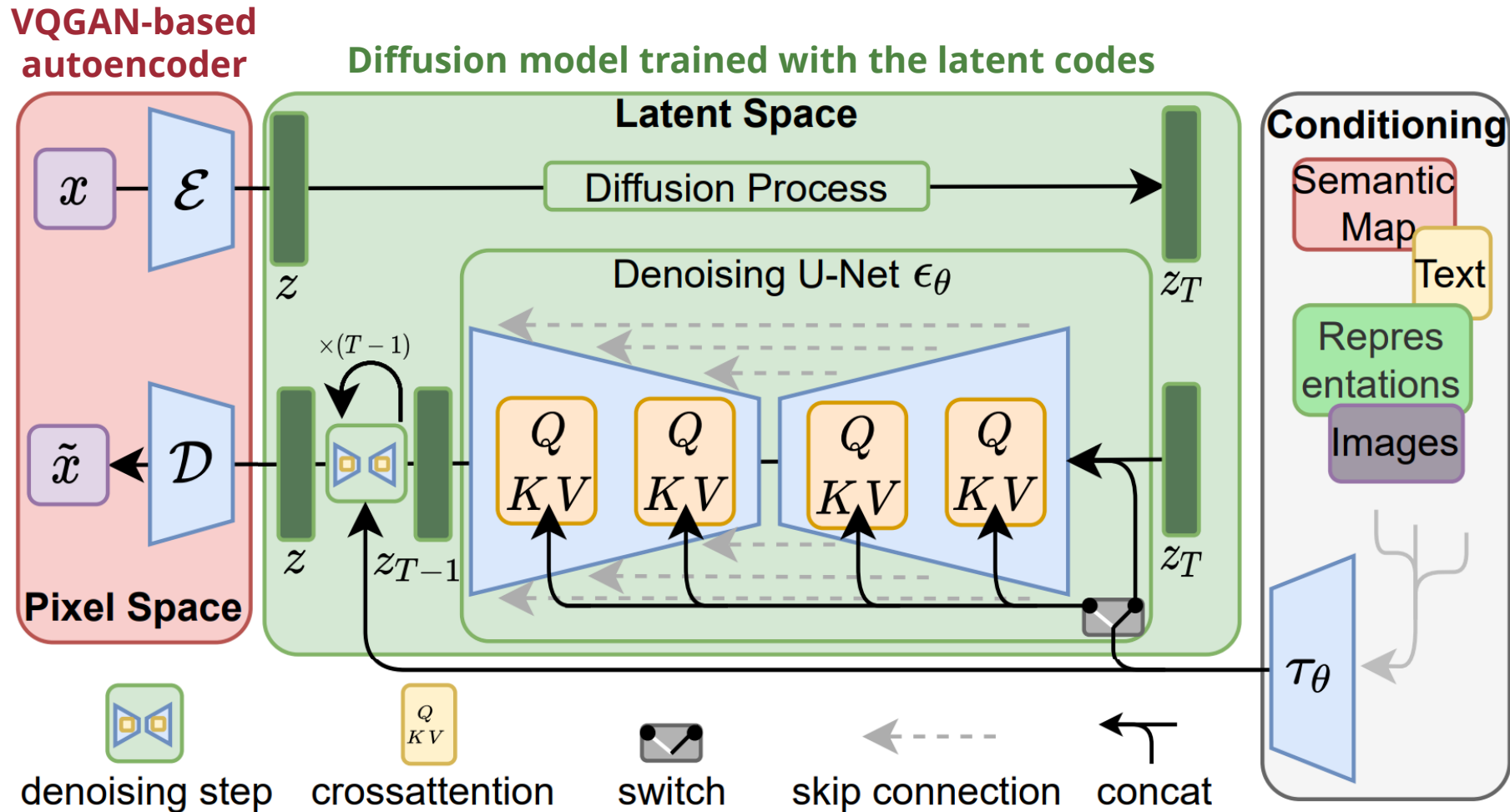
# (Recap) VQGAN



(Source: Esser et al., 2021)

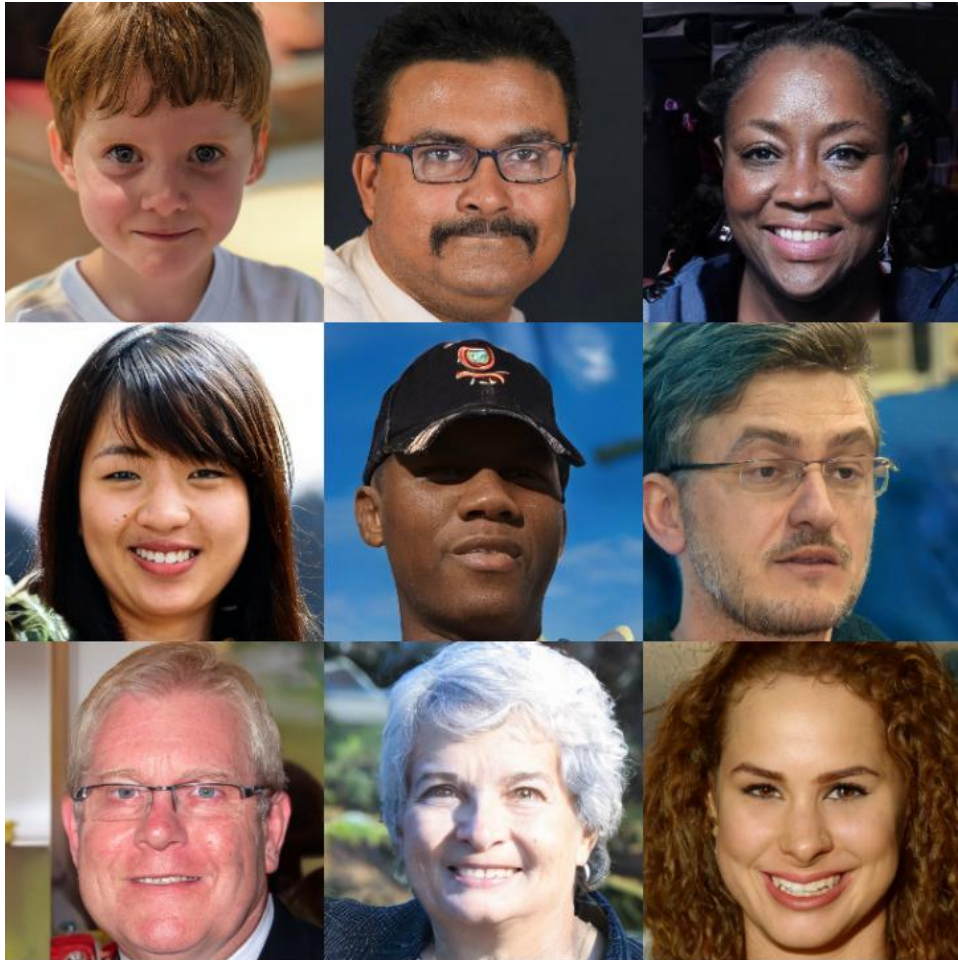
**A VQGAN is a VQVAE equipped with adversarial loss**

# (Recap) Latent Diffusion Models (LDMs)



(Source: Rombach et al., 2022)

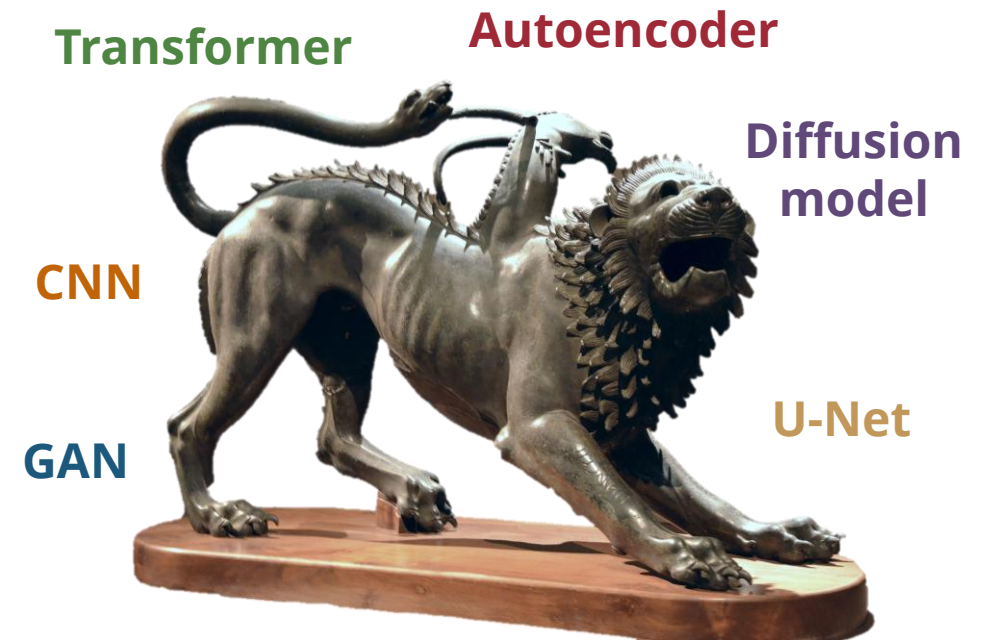
# (Recap) LDMs: Examples



(Source: Rombach et al., 2022)

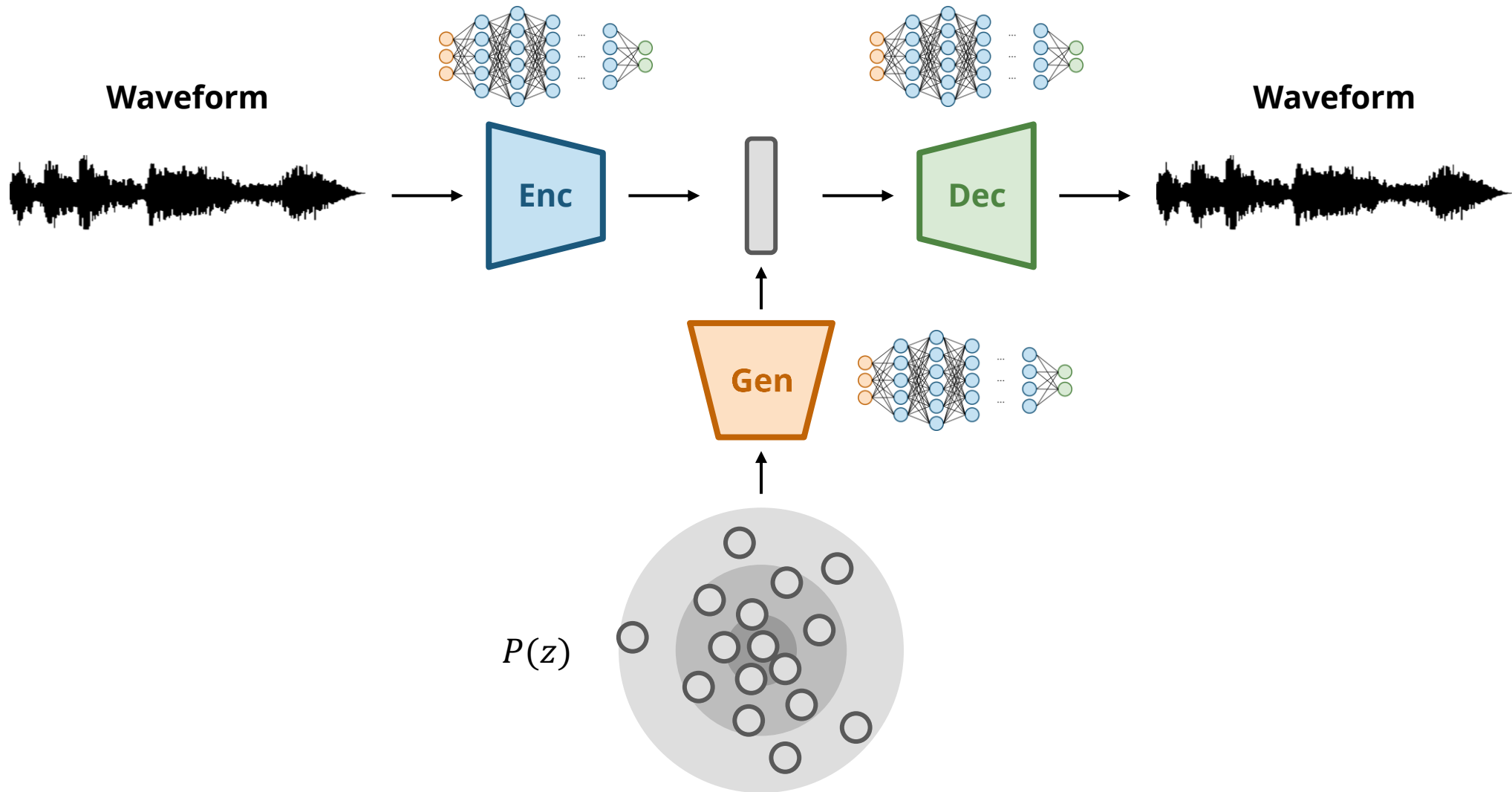
# (Recap) Latent Diffusion Model is a Chimera

- **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)

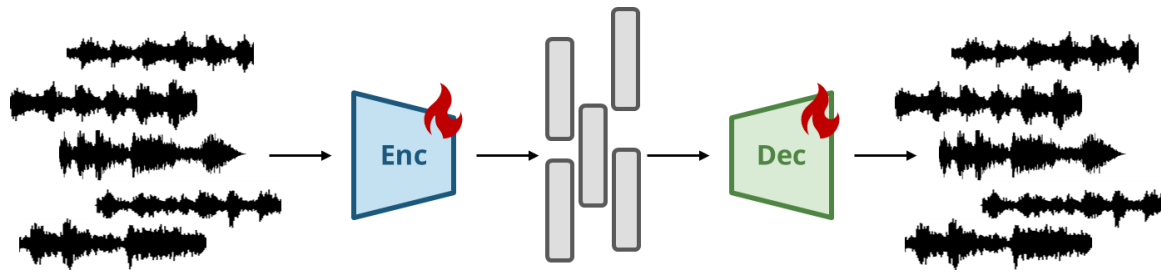
# (Recap) Latent-based Audio Synthesis



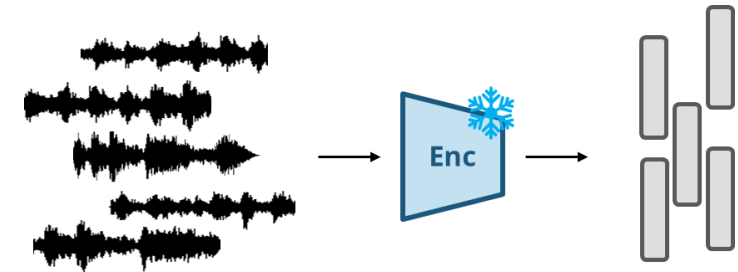


# (Recap) Pipeline

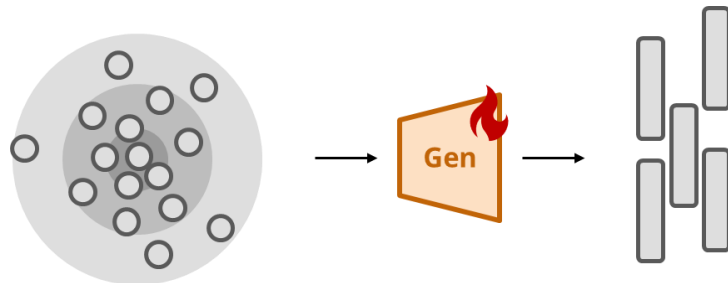
Step 1: Train an Autoencoder



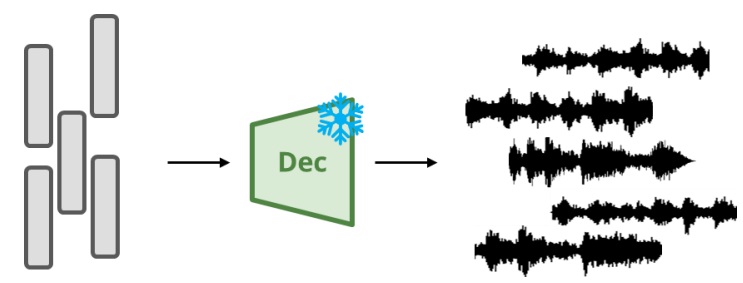
Step 2: Compute the Latent Vectors



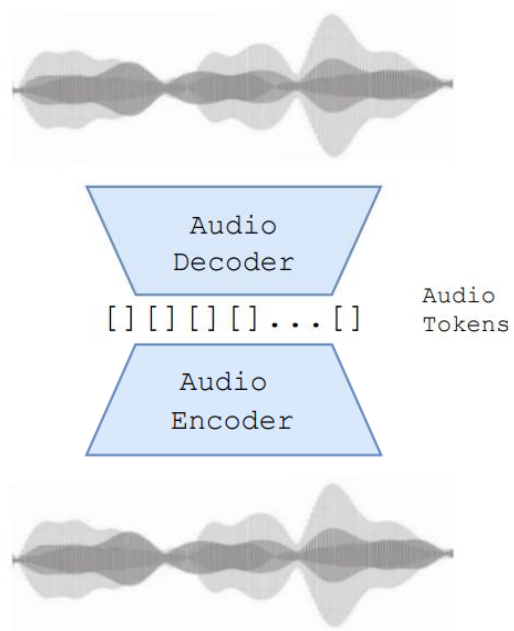
Step 3: Train a Latent Generative Model



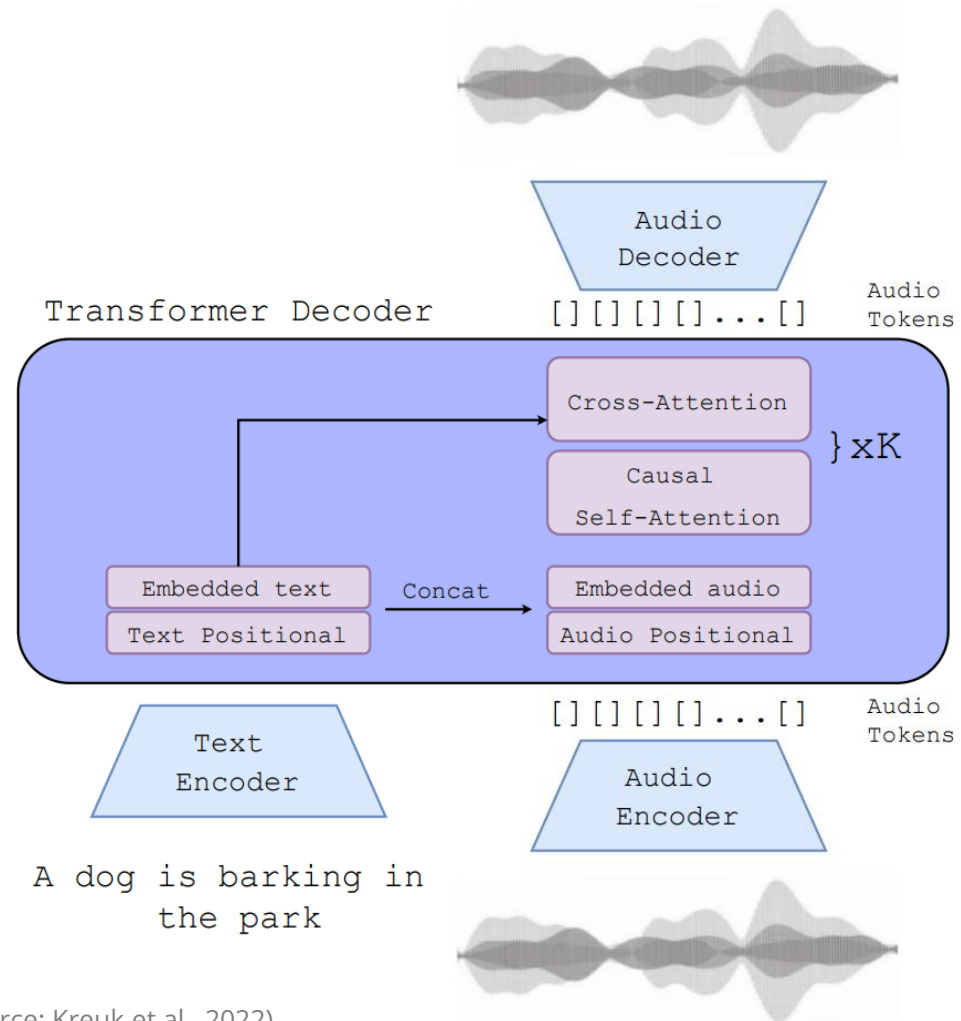
Step 4: Decode the Latent Vectors



# (Recap) AudioGen (Kreuk et al., 2023)



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



(Source: Kreuk et al., 2022)

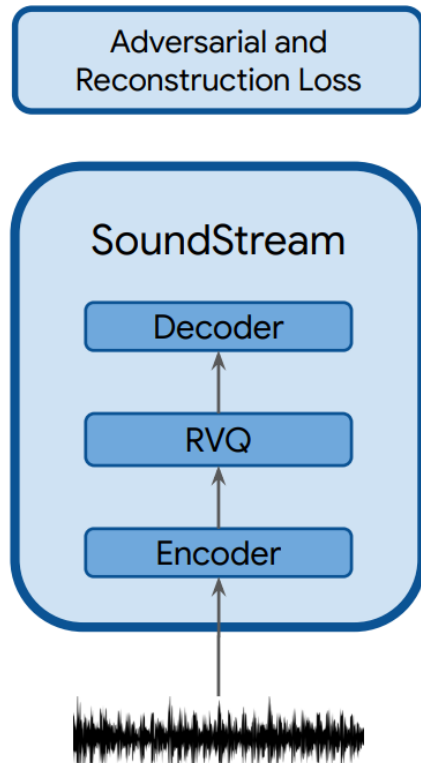
## (Recap) MusicGen (Copet et al., 2023)

- 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/)

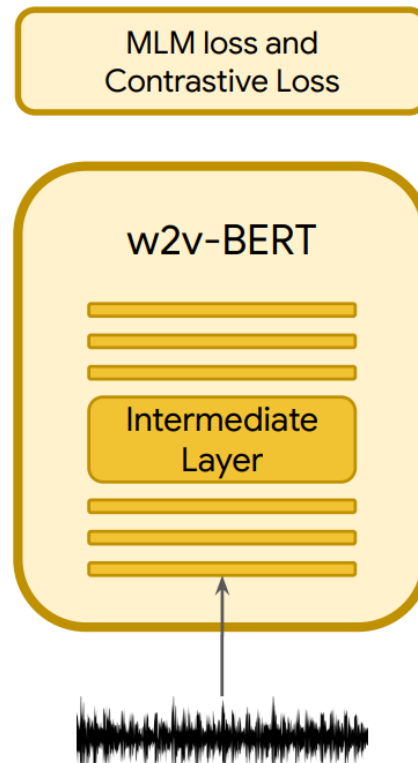
# (Recap) MusicLM (Agostinelli et al., 2023)

## Audio autoencoder

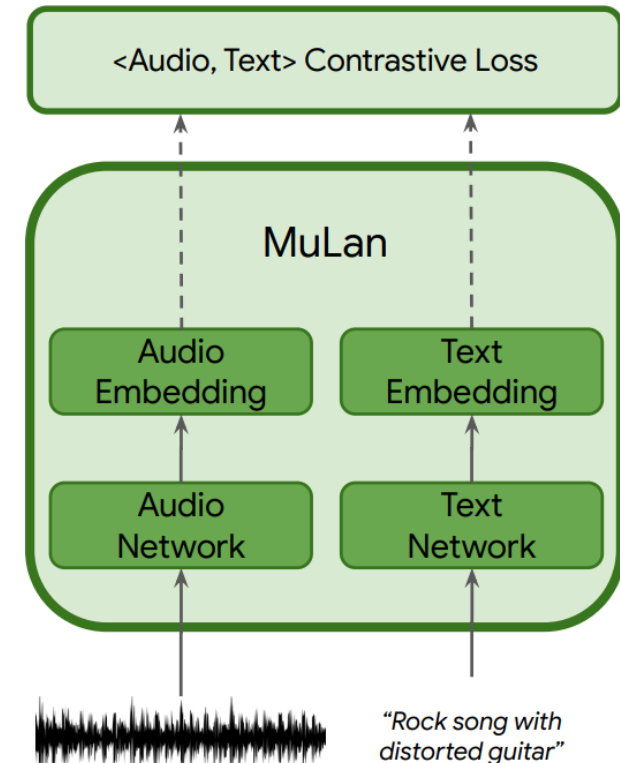


**106k songs, 8.2k hours**

## Semantic representation



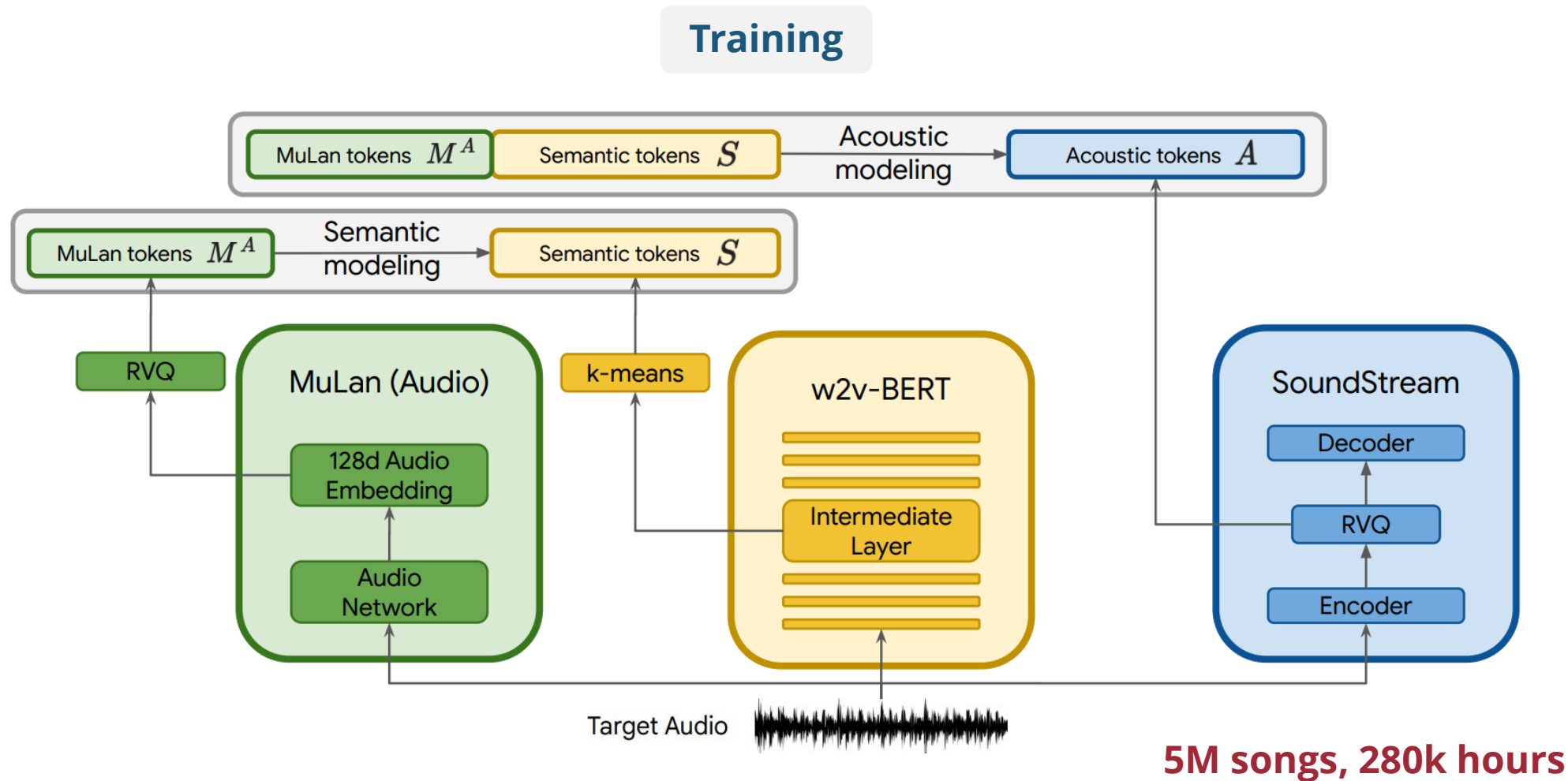
## Text-music correspondence



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

(Source: Agostinelli et al., 2022)

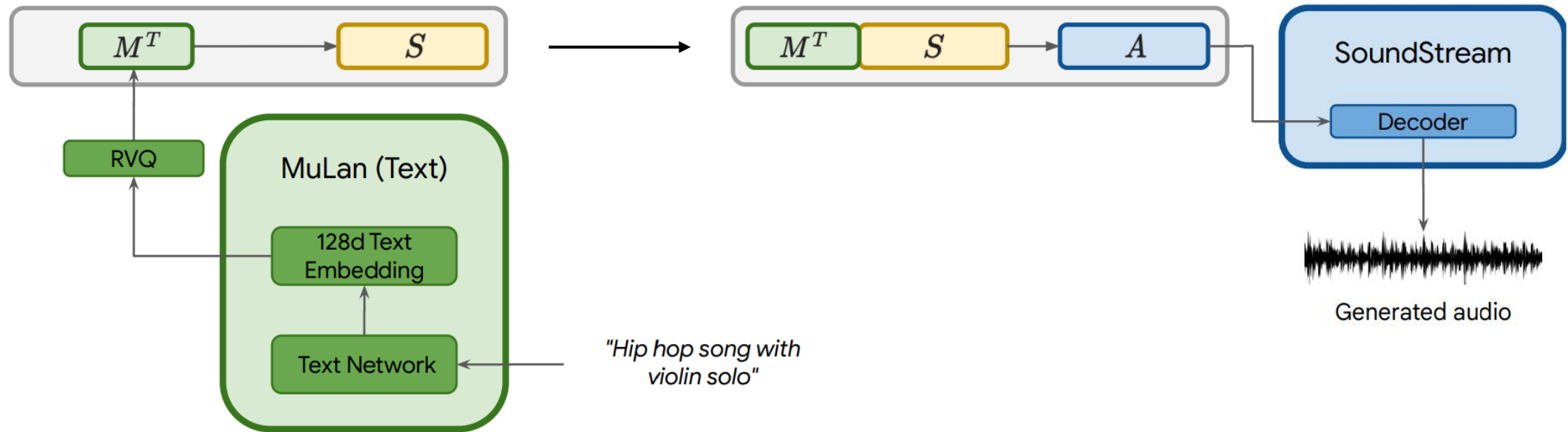
# (Recap) MusicLM (Agostinelli et al., 2023)



(Source: Agostinelli et al., 2022)

# (Recap) MusicLM (Agostinelli et al., 2023)

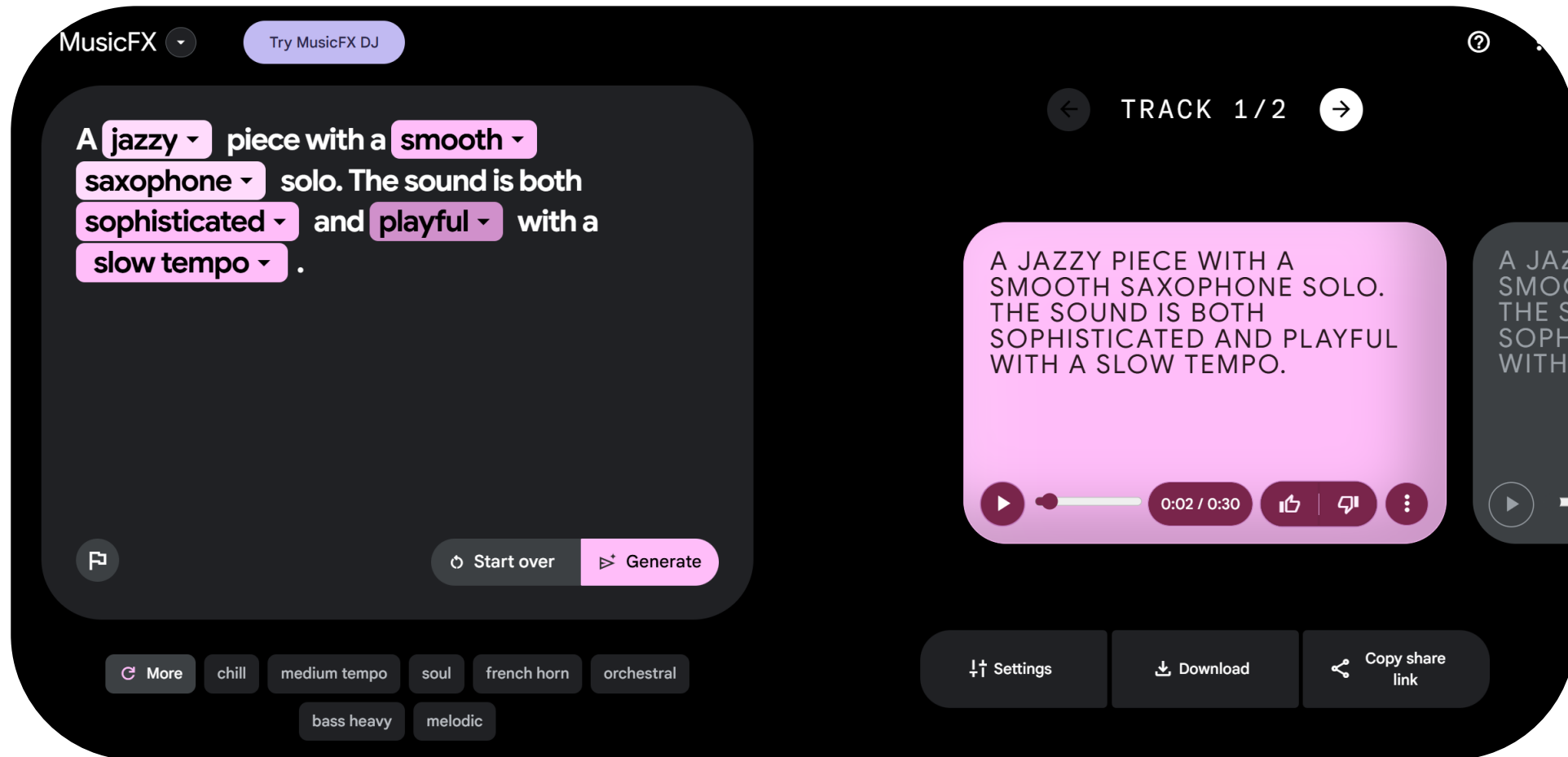
## Inference



(Source: Agostinelli et al., 2022)

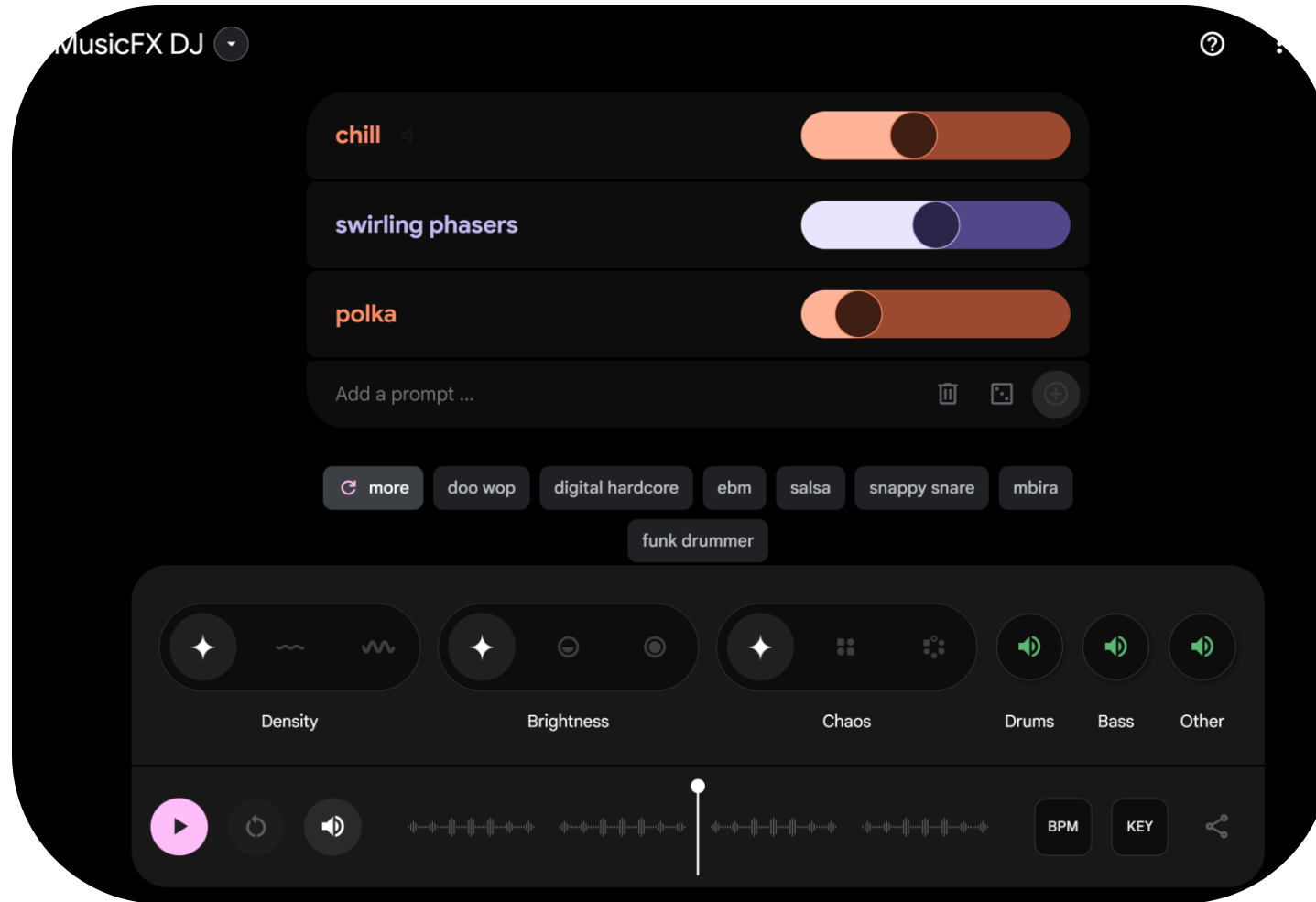
[google-research.github.io/seanet/musiclm/examples/](https://google-research.github.io/seanet/musiclm/examples/)

# (Recap) Music FX (2024)



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

# (Recap) Music FX DJ (2024)



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

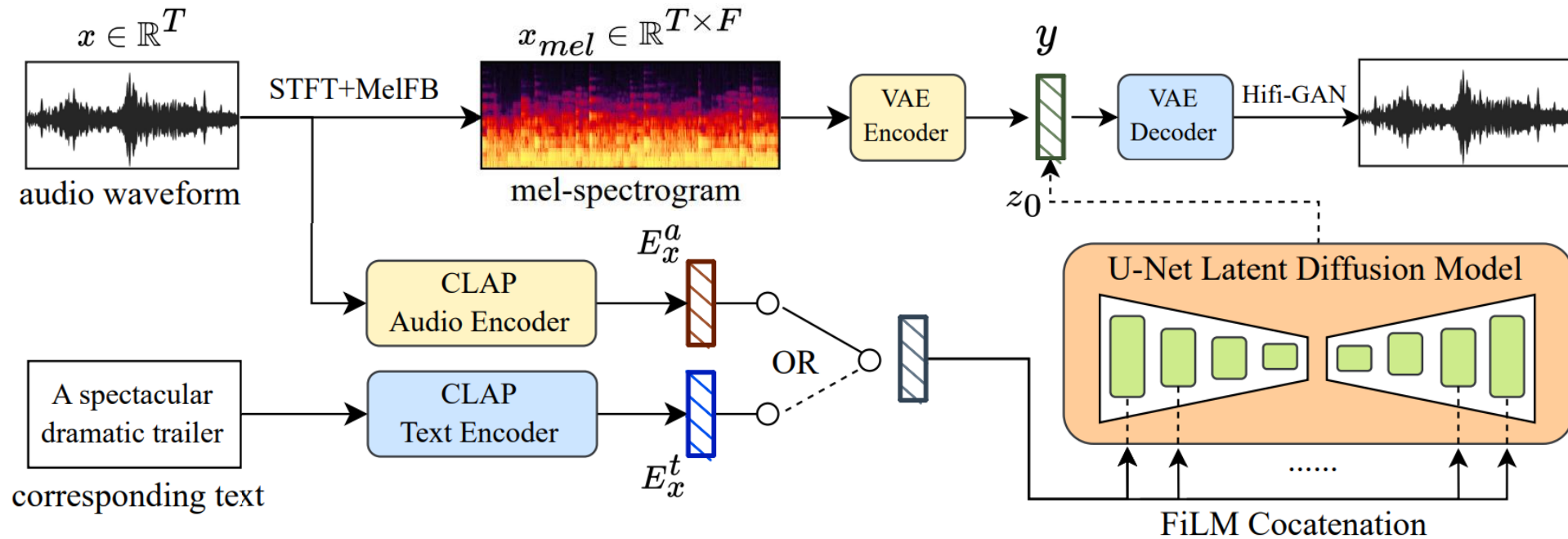


# (Recap) Music FX DJ (2024)



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

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



(Source: Ke et al., 2023)

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

(Recap) Example: **MusicLDM** (Chen et al., 2023)



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

# (Recap) unloop (Garcia et al., 2023)

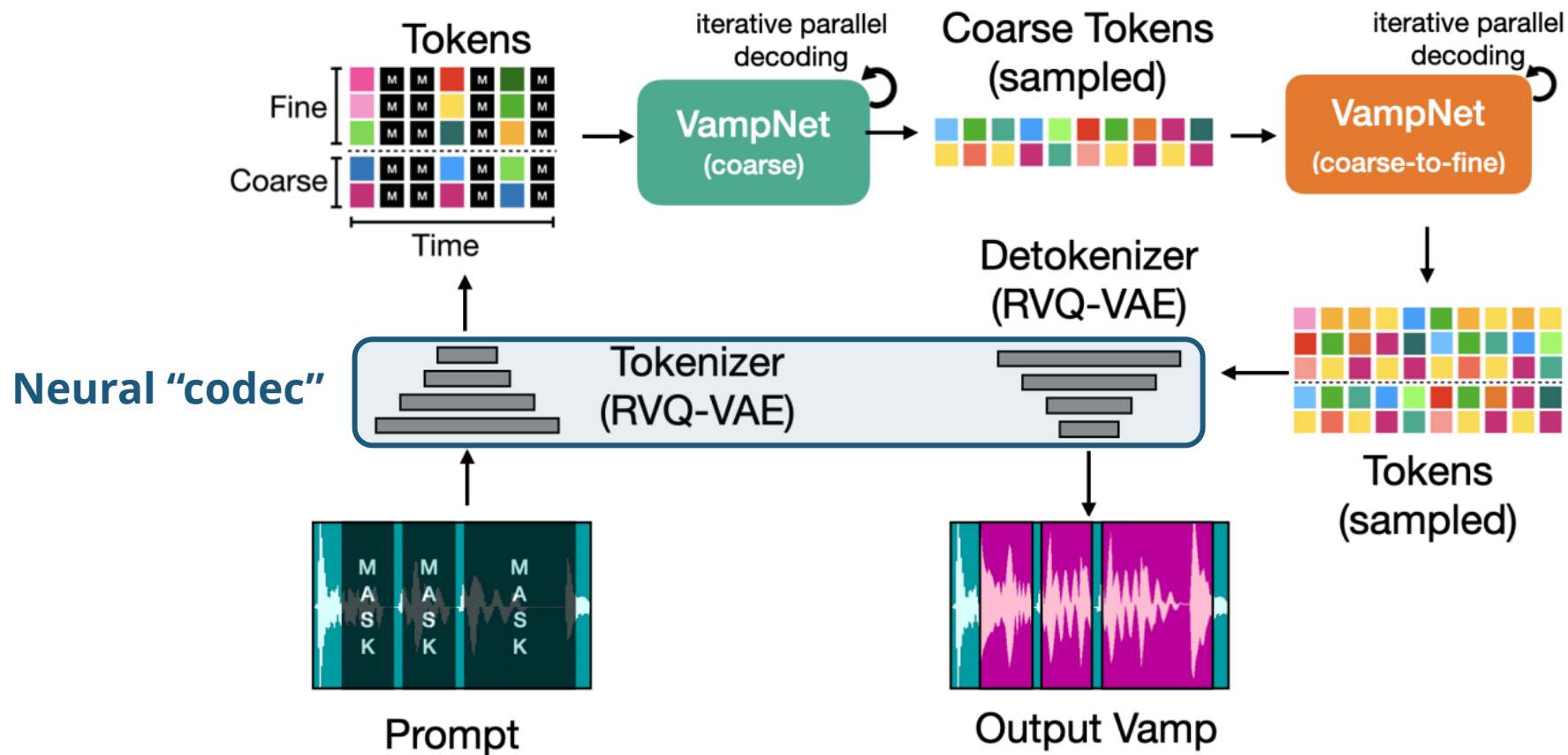


[youtu.be/yzBI8Vcjd2s](https://youtu.be/yzBI8Vcjd2s)

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

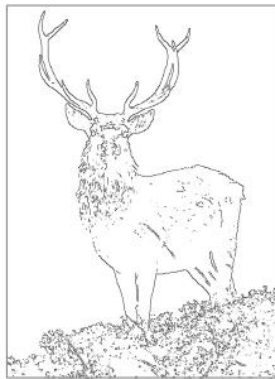


# (Recap) VampNet (Garcia et al., 2023)



(Source: Garcia et al., 2023)

# (Recap) ControlNet (Zhang et al., 2023)



Input Canny edge



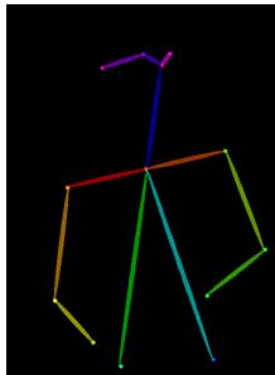
Default



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



“..., quaint city Galic”



Input human pose



Default



“chef in kitchen”

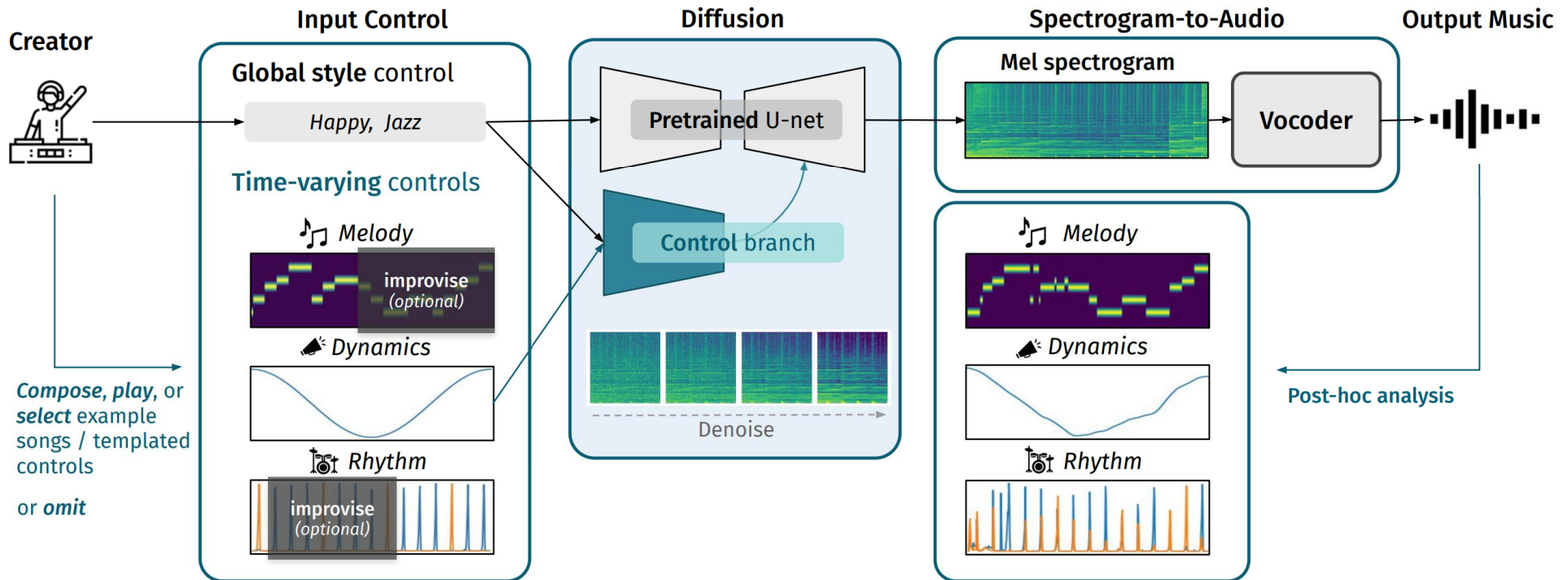


“Lincoln statue”

(Source: Zhang et al., 2023)

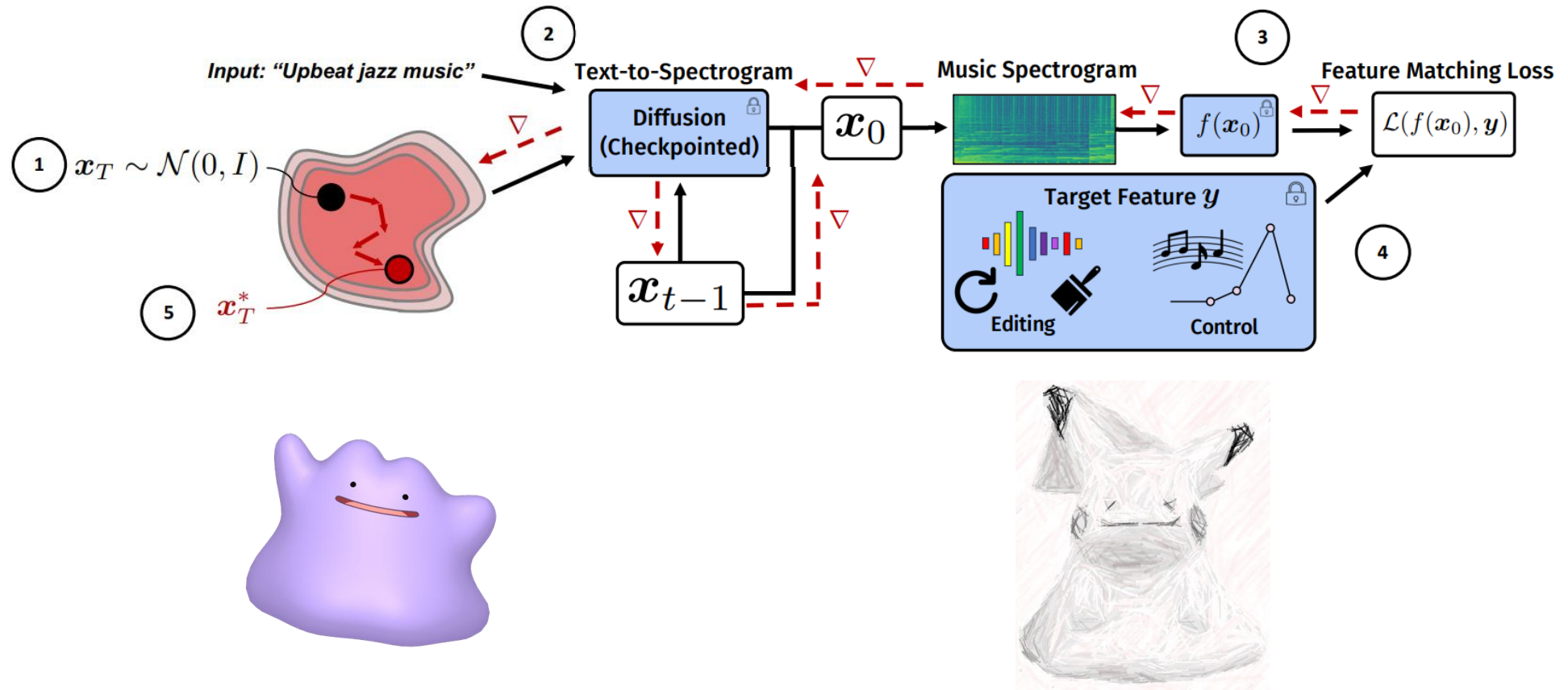
Can we **add controls** to a trained text-to-image diffusion model?

# (Recap) Music ControlNet (Wu et al., 2024)



(Source: Wu et al., 2024)

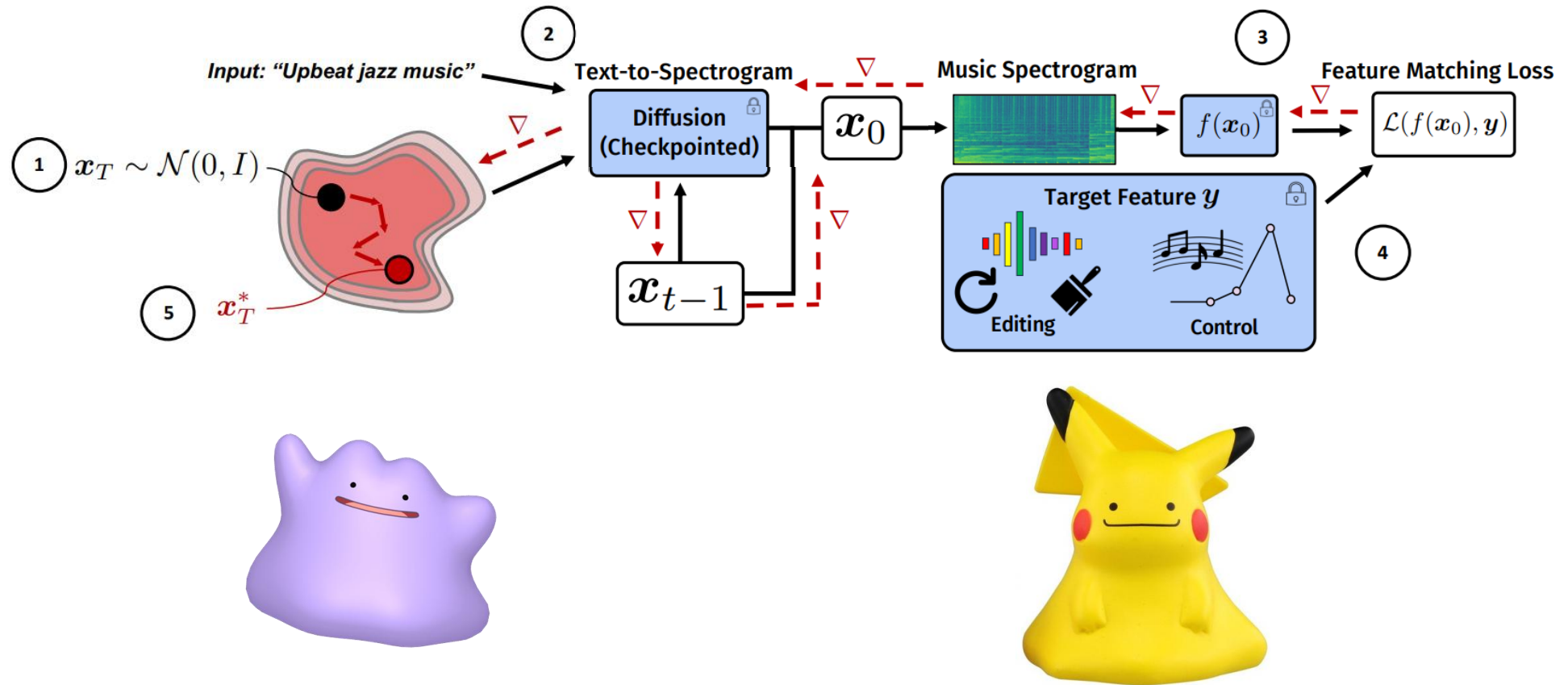
# (Recap) DITTO (Novack et al., 2024)



(Source: Novack et al., 2024)



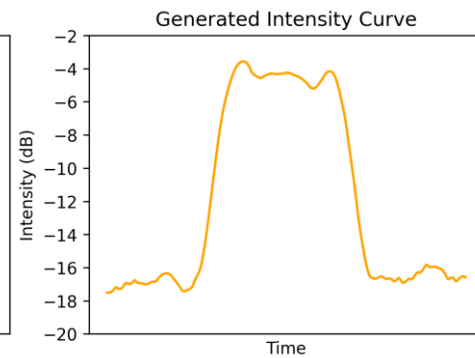
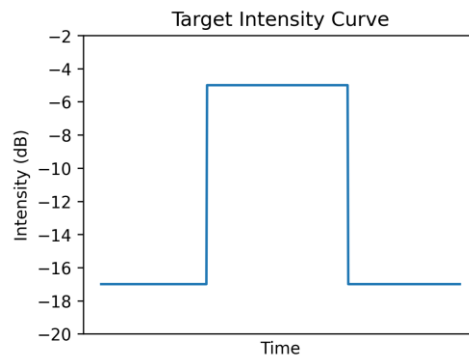
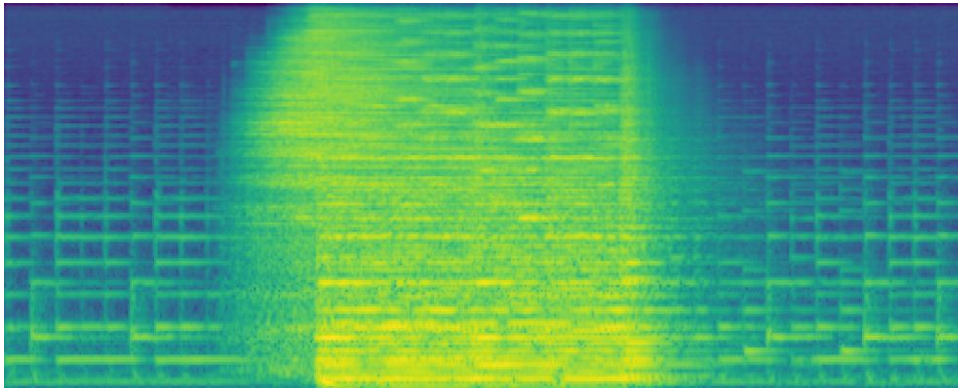
# (Recap) DITTO (Novack et al., 2024)



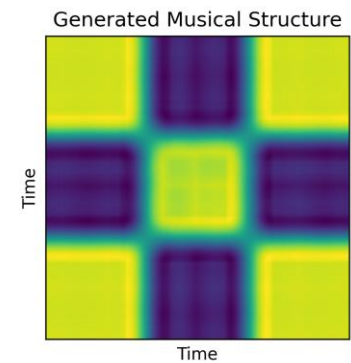
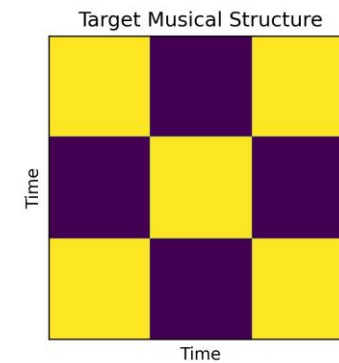
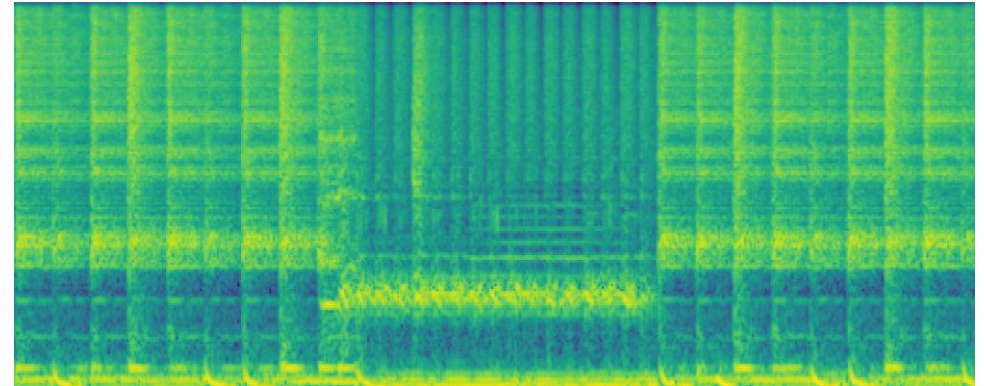
(Source: Novack et al., 2024)

# (Recap) DITTO (Novack et al., 2024)

## Intensity control



## Structure control



(Source: Novack et al., 2024)

# Music Fingerprinting

# Shazam & Siri

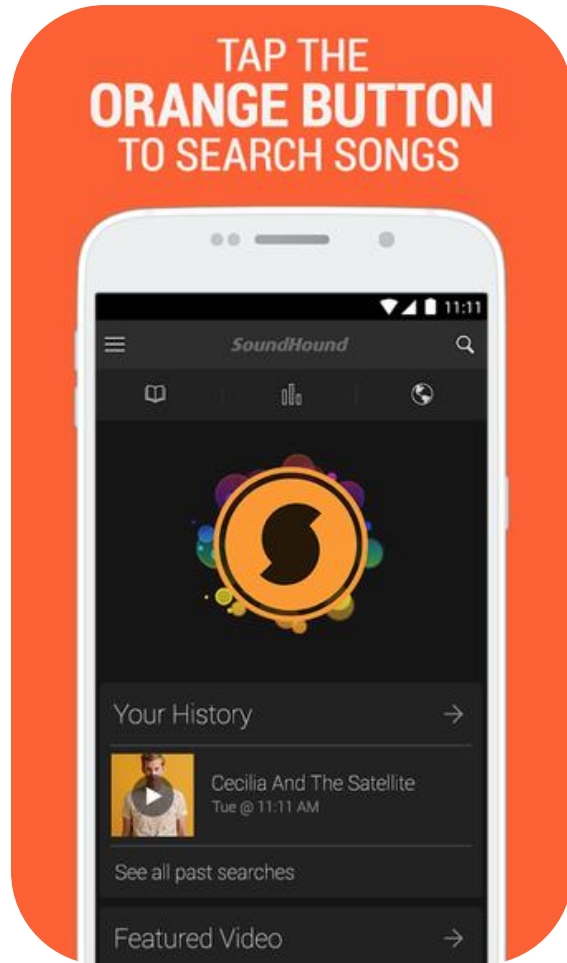


(Source: Shazam User Guide)

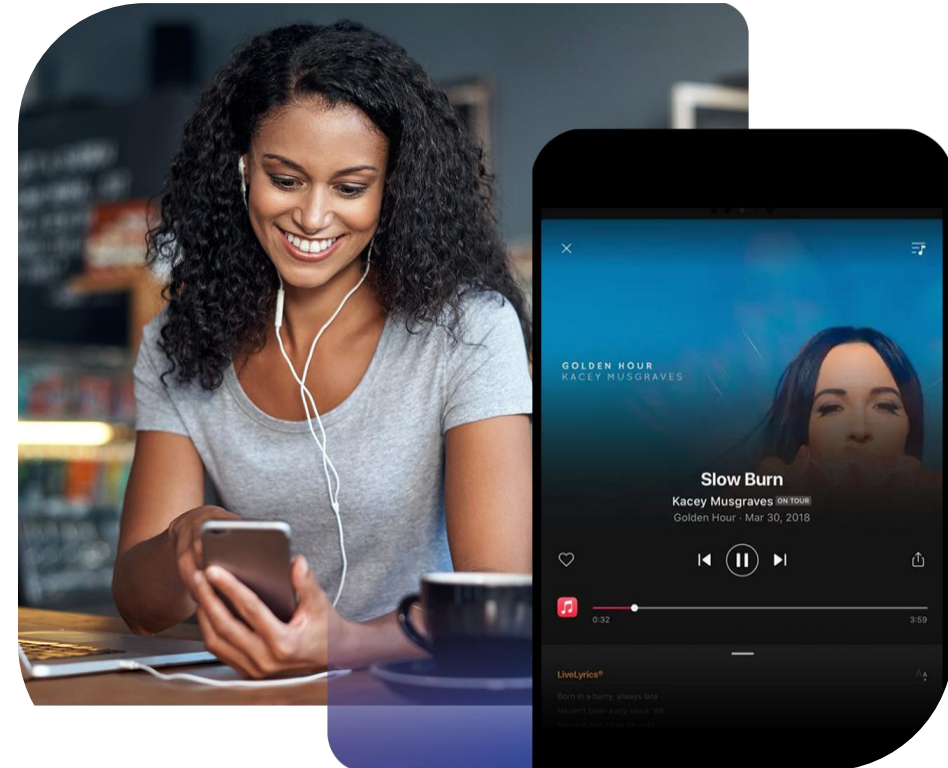


(Source: OSXDaily)

# SoundHound

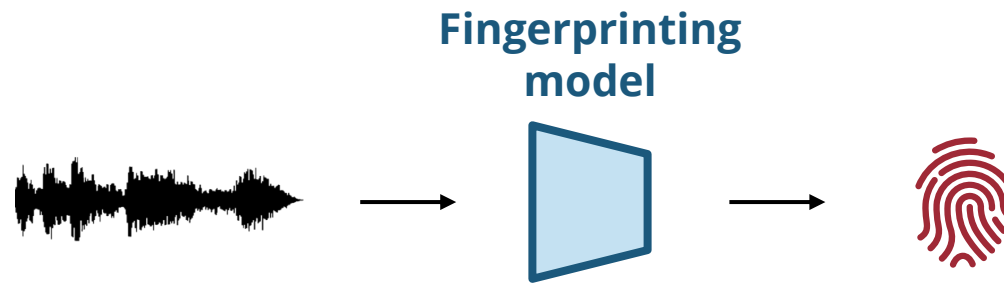


(Source: CNet)

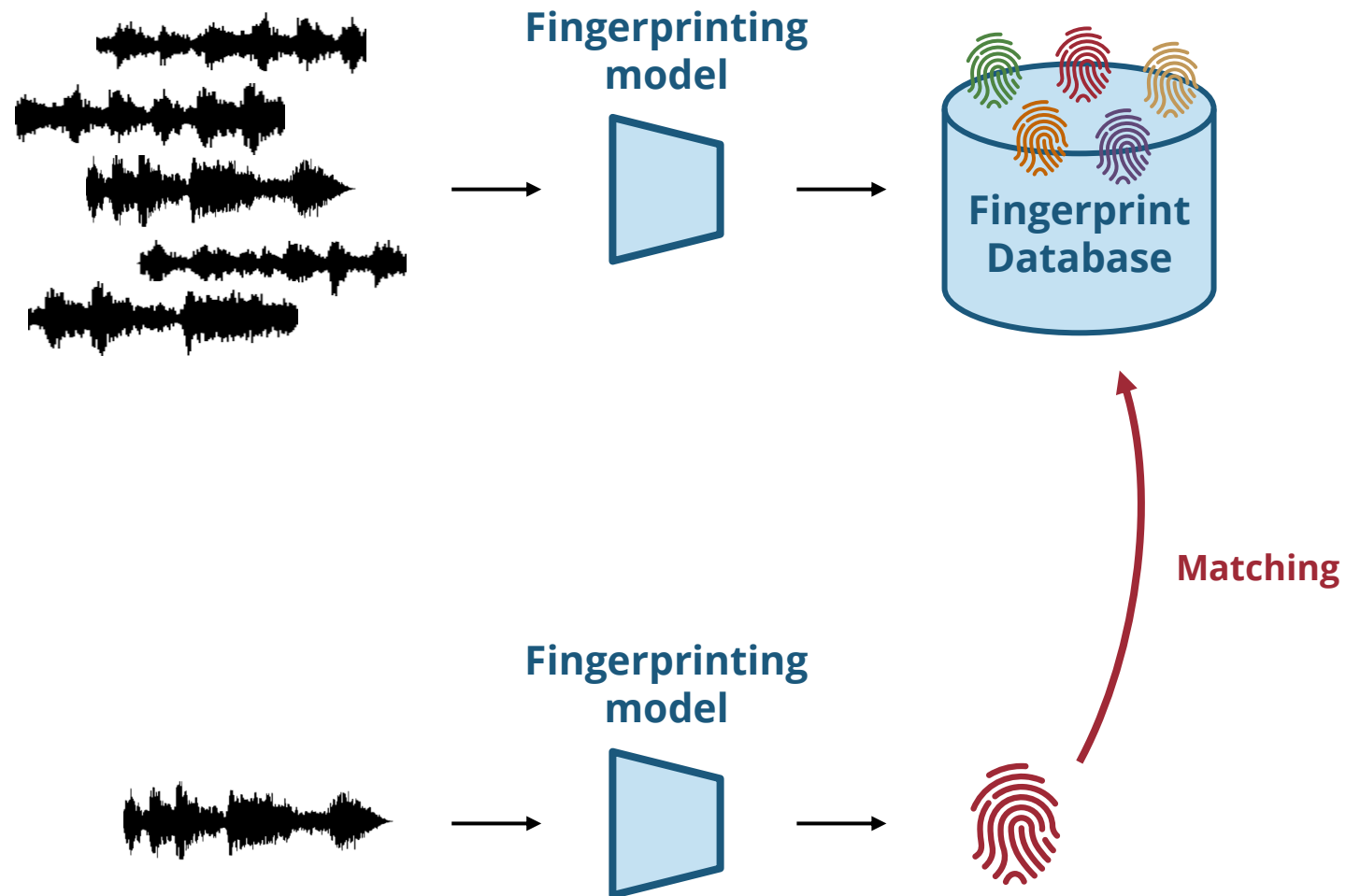


(Source: SoundHound)

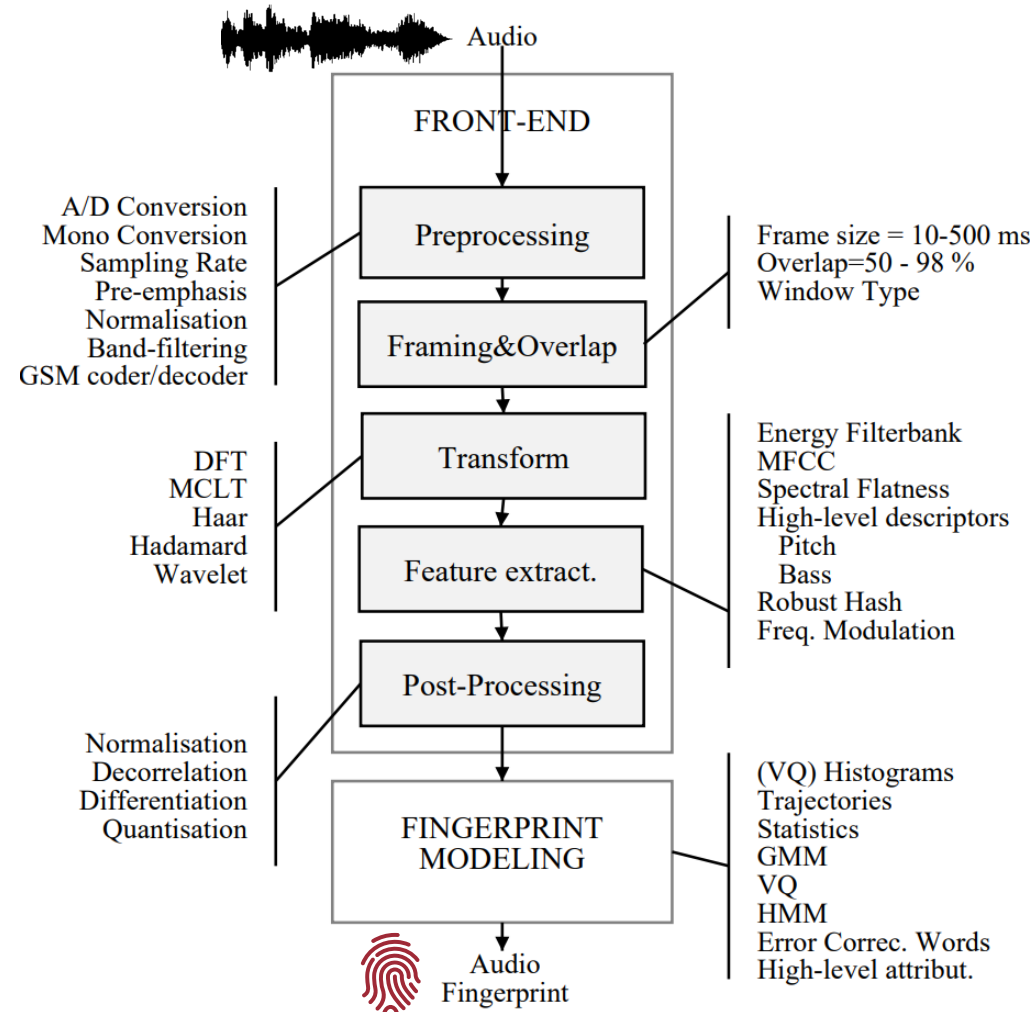
# Audio Fingerprinting



# Audio Fingerprinting for Audio Identification



# Audio Fingerprinting

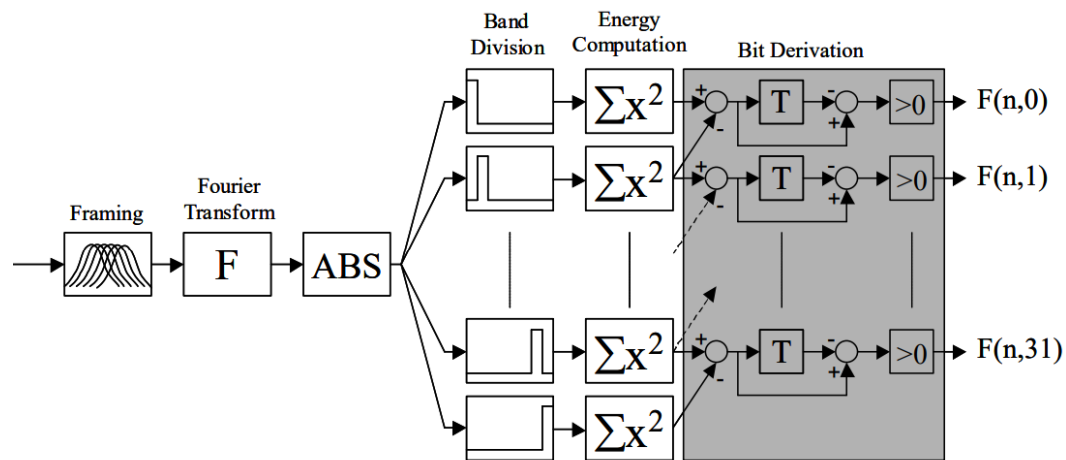


(Source: Cano et al., 2002)

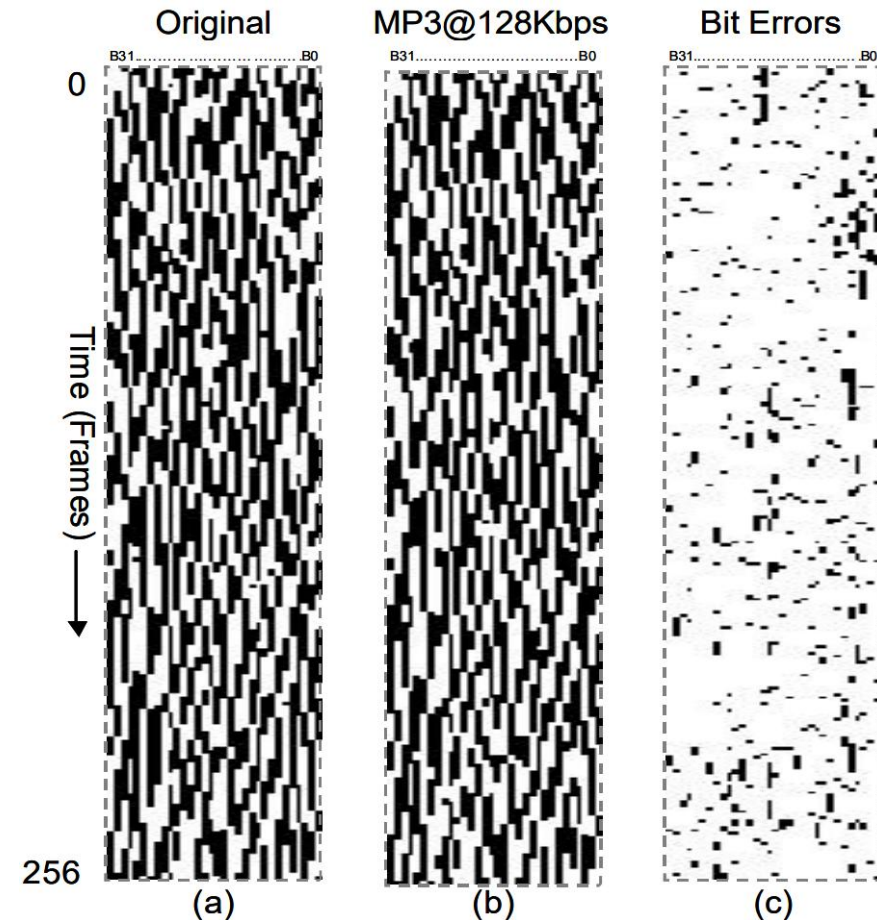


# Audio Fingerprinting: Examples (Haitsma et al., 2002)

## Fingerprint extraction

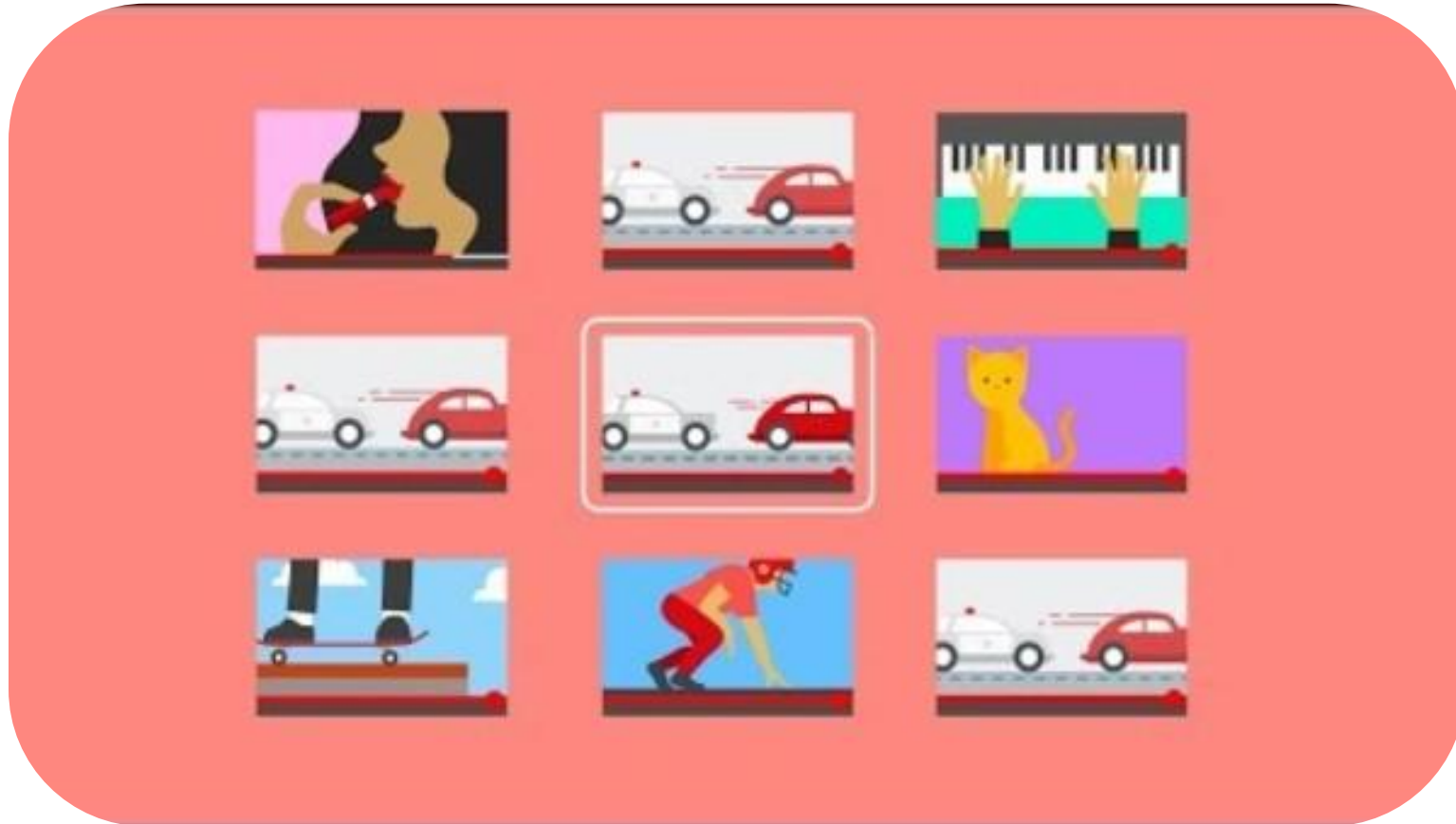


## Example fingerprint blocks



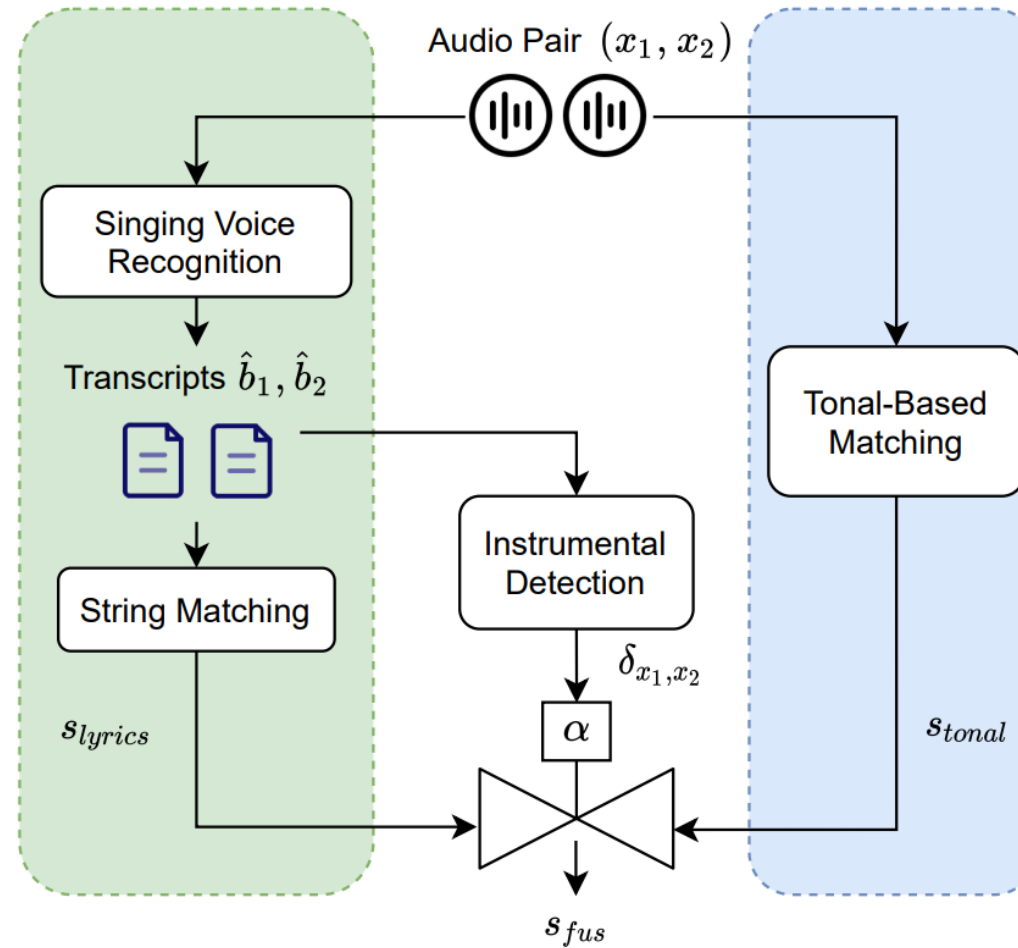
(Source: Haitsma et al., 2002)

# YouTube's Content ID



[youtu.be/9g2U12SsRns](https://youtu.be/9g2U12SsRns)

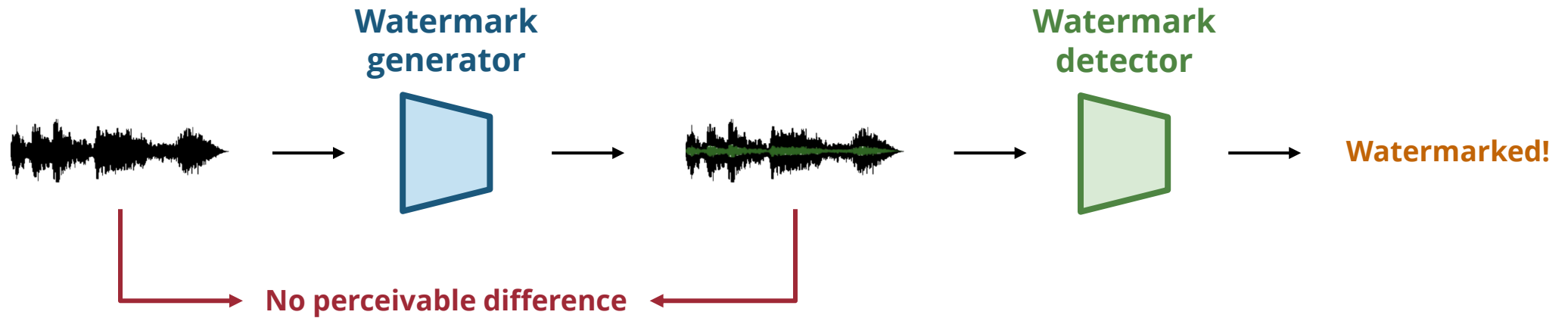
# Cover Song Identification



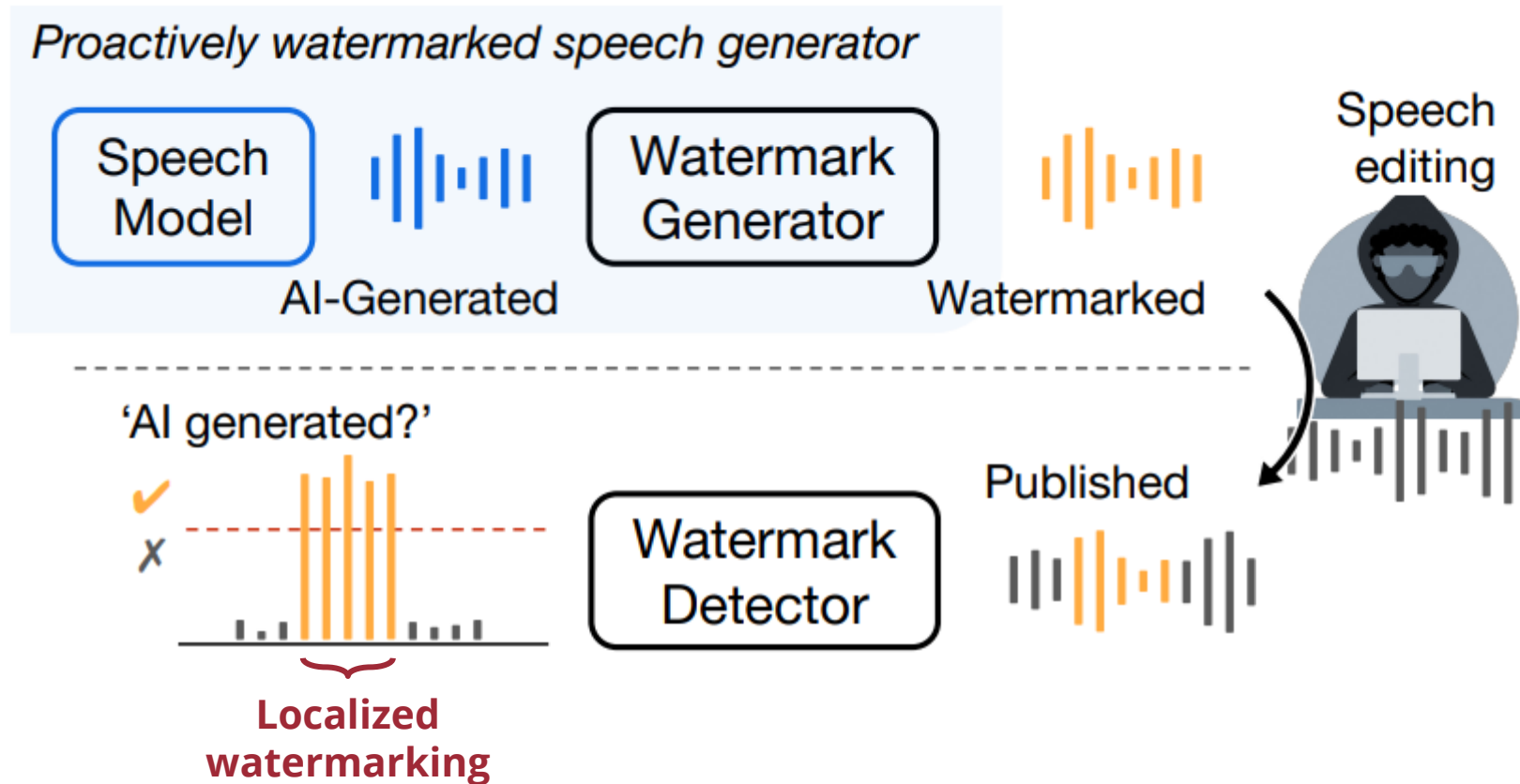
(Source: Vaglio et al., 2021)

# Music Watermarking

# Audio Watermarking

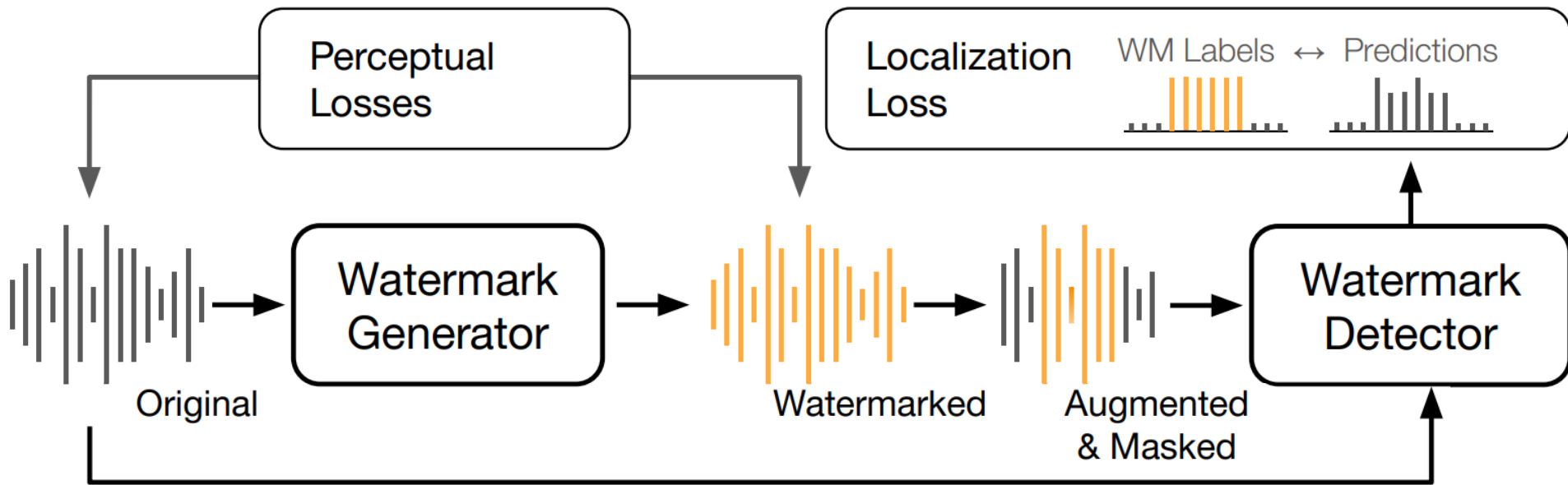


# Audio Watermarking Against Generated Audio



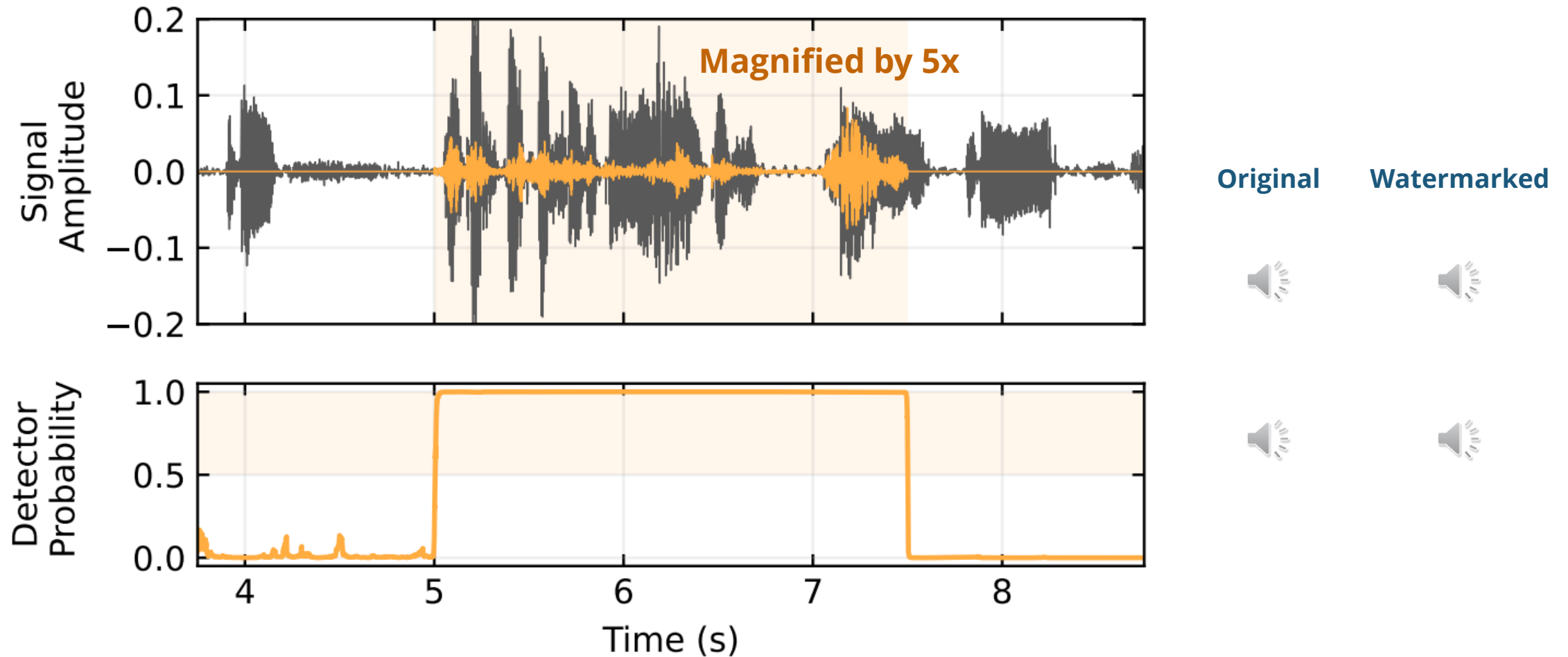
(Source: Roman et al., 2024)

# AudioSeal (Roman et al., 2024)



(Source: Roman et al., 2024)

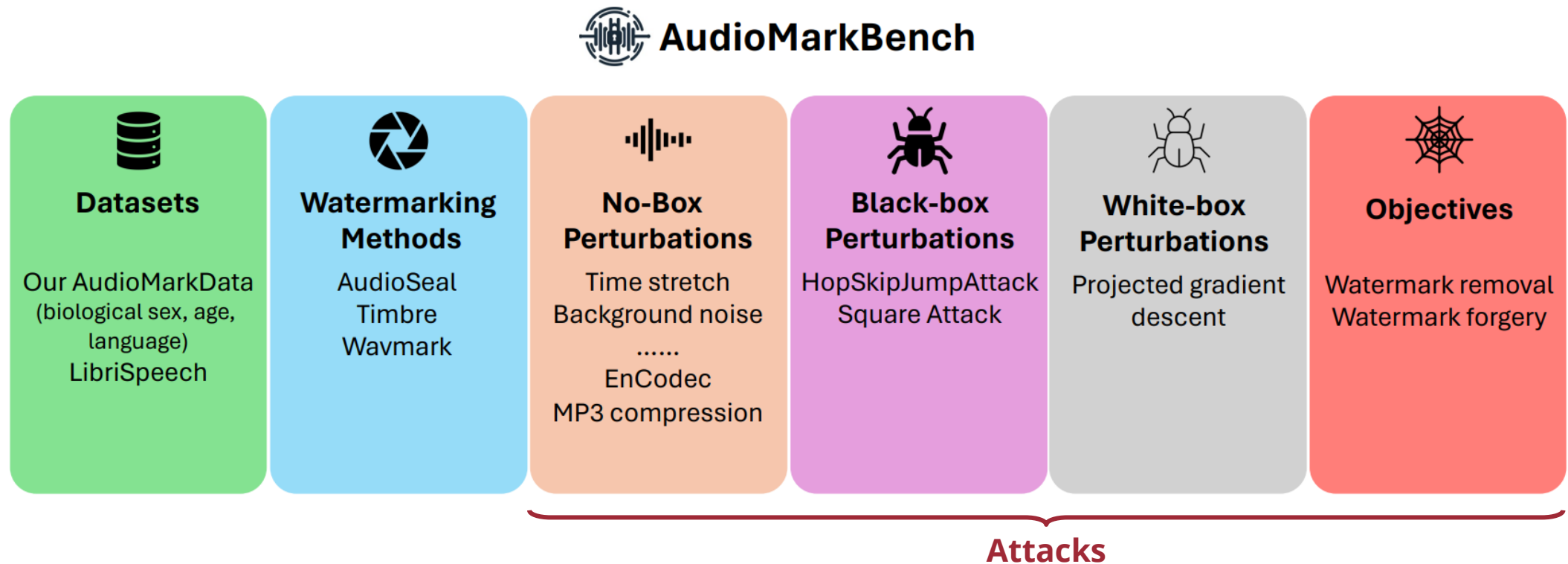
# AudioSeal (Roman et al., 2024)



(Source: Roman et al., 2024)



# AudioMarkBench (Liu et al., 2024)



(Source: Liu et al., 2024)

# MusicFX's SynthID (Google)

## About MusicFX

MusicFX is an experimental technology that allows you to generate your own music. Certain queries that mention specific artists or include vocals will not be generated.

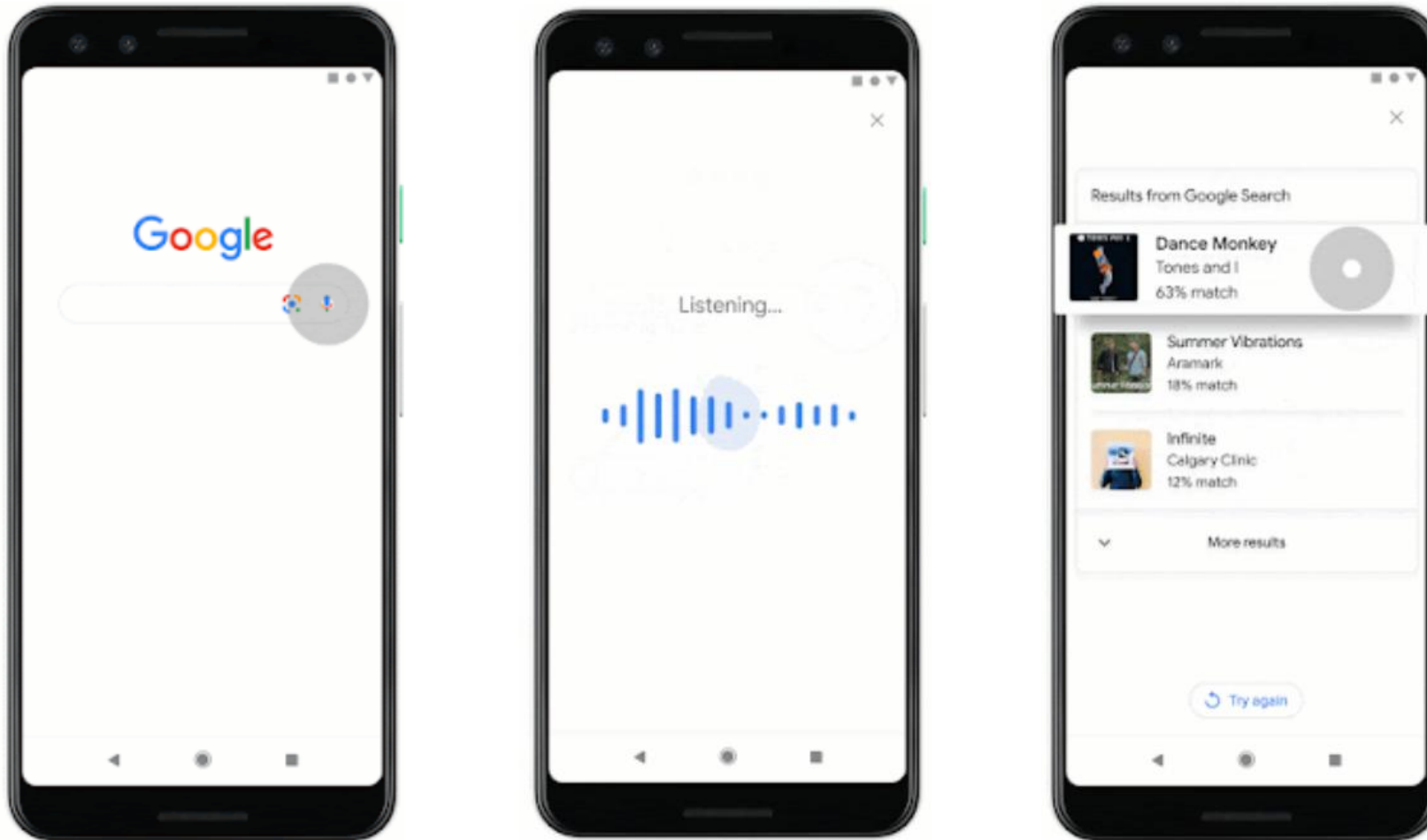
MusicFX is powered by Google's [MusicLM](#) and uses Google DeepMind's novel watermarking technology, [SynthID](#) to embed a digital watermark in the outputs.

We need your help to improve AI for everybody. Generated audio and prompt suggestions are experimental. You can [report](#) content under our policies or applicable laws, or give feedback by clicking the flag icon so we can improve AI responsibly together.

Got it

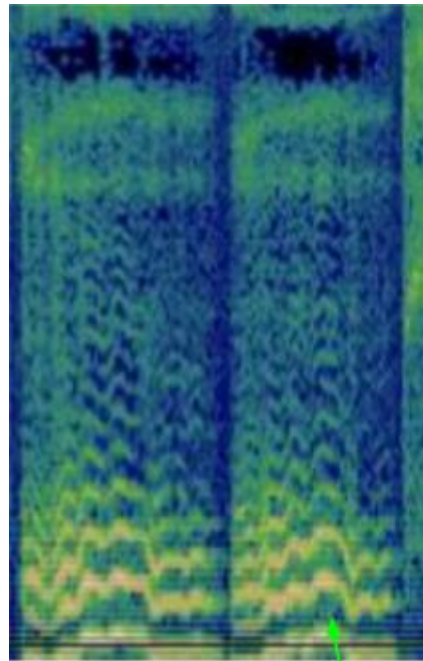
# Query by Humming

# Hum to Search (Google)

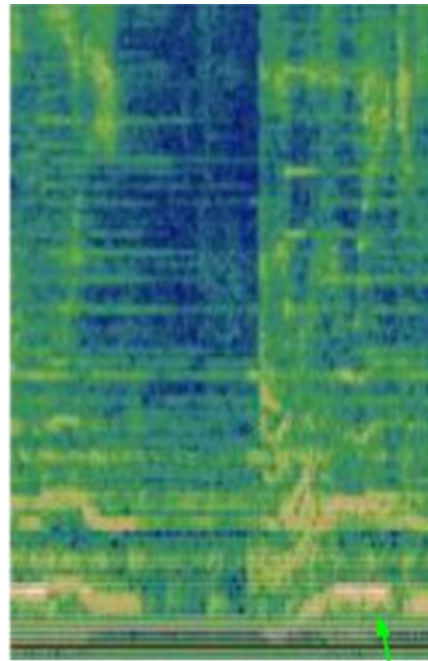


(Source: Google Research Blog)

# Hum to Search (Google)



Humming



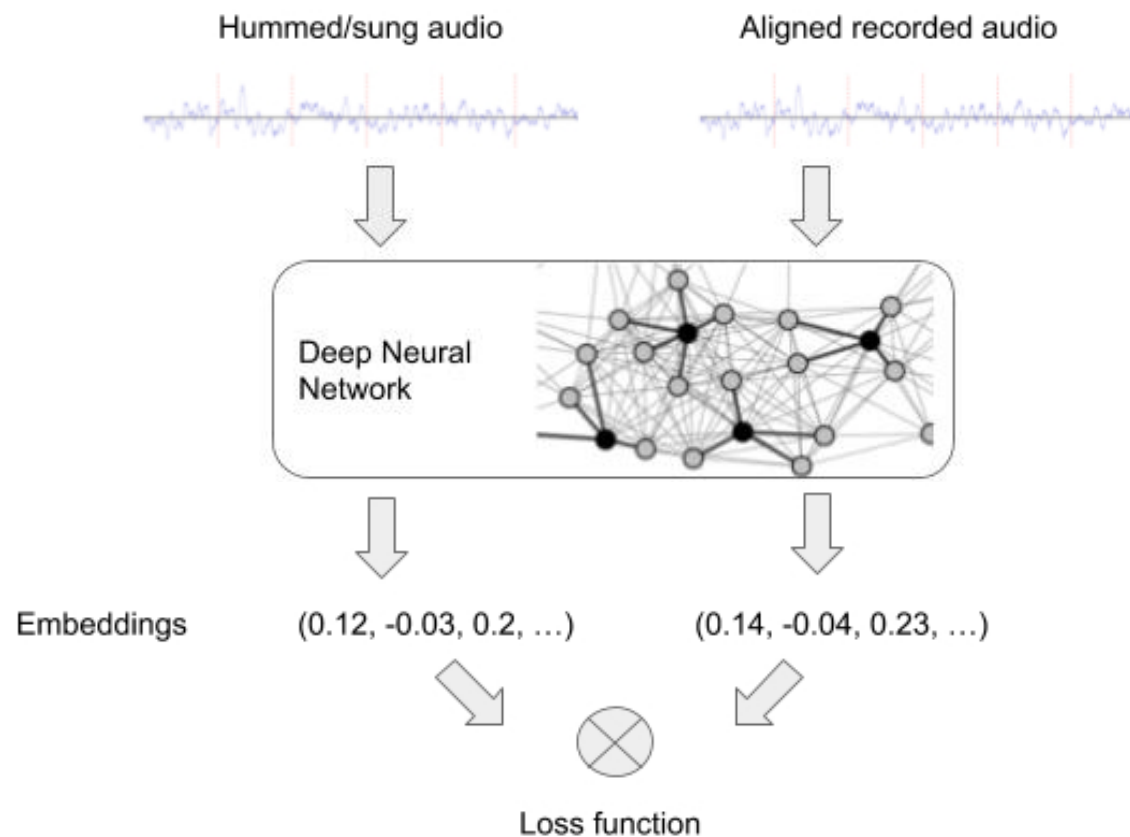
Studio Recording



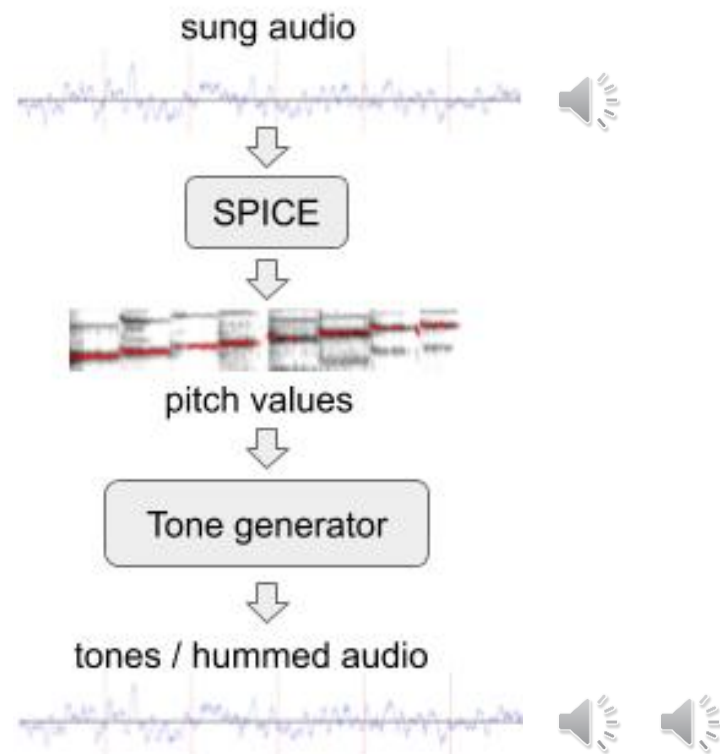
(Source: Google Research Blog)

# Hum to Search (Google)

## Audio encoder



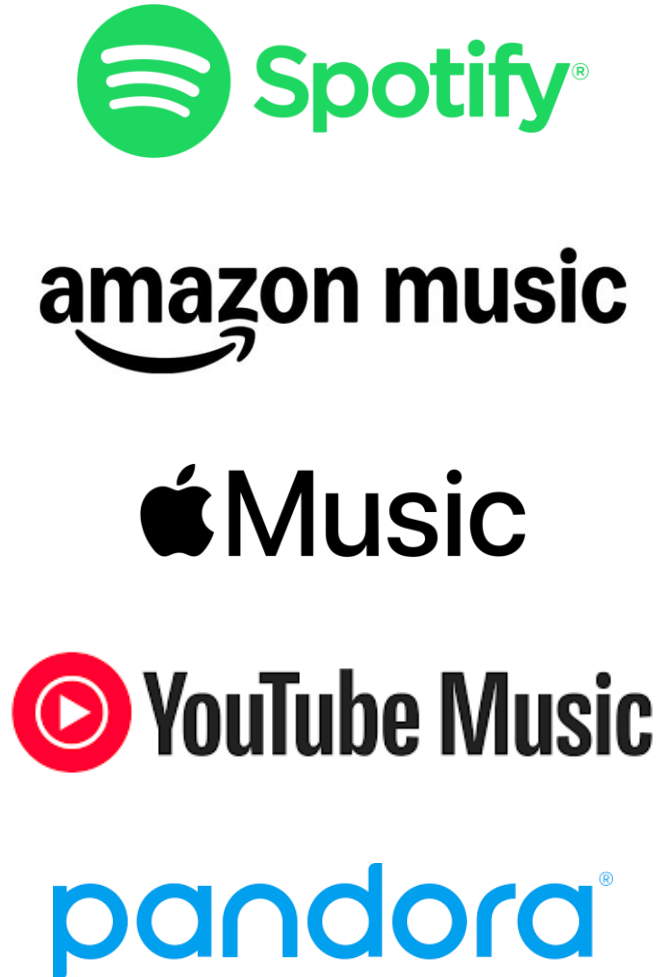
## Data augmentation



(Source: Google Research Blog)

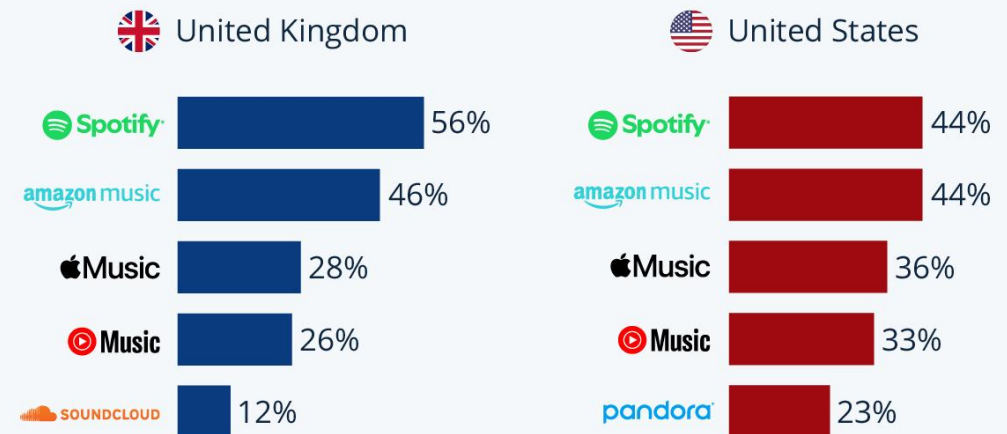
# Music Recommendation

# Music Recommendation



## The Most Loved Digital Audio Streaming Platforms

Share of respondents who have paid for audio downloads or streaming services from the following platforms\*



\* in the 12 months prior to the survey  
2,362 (UK)/4,944 (USA) respondents (18-64 y/o) surveyed Jul. 2023-Jun. 2024  
Source: Statista Consumer Insights



statista

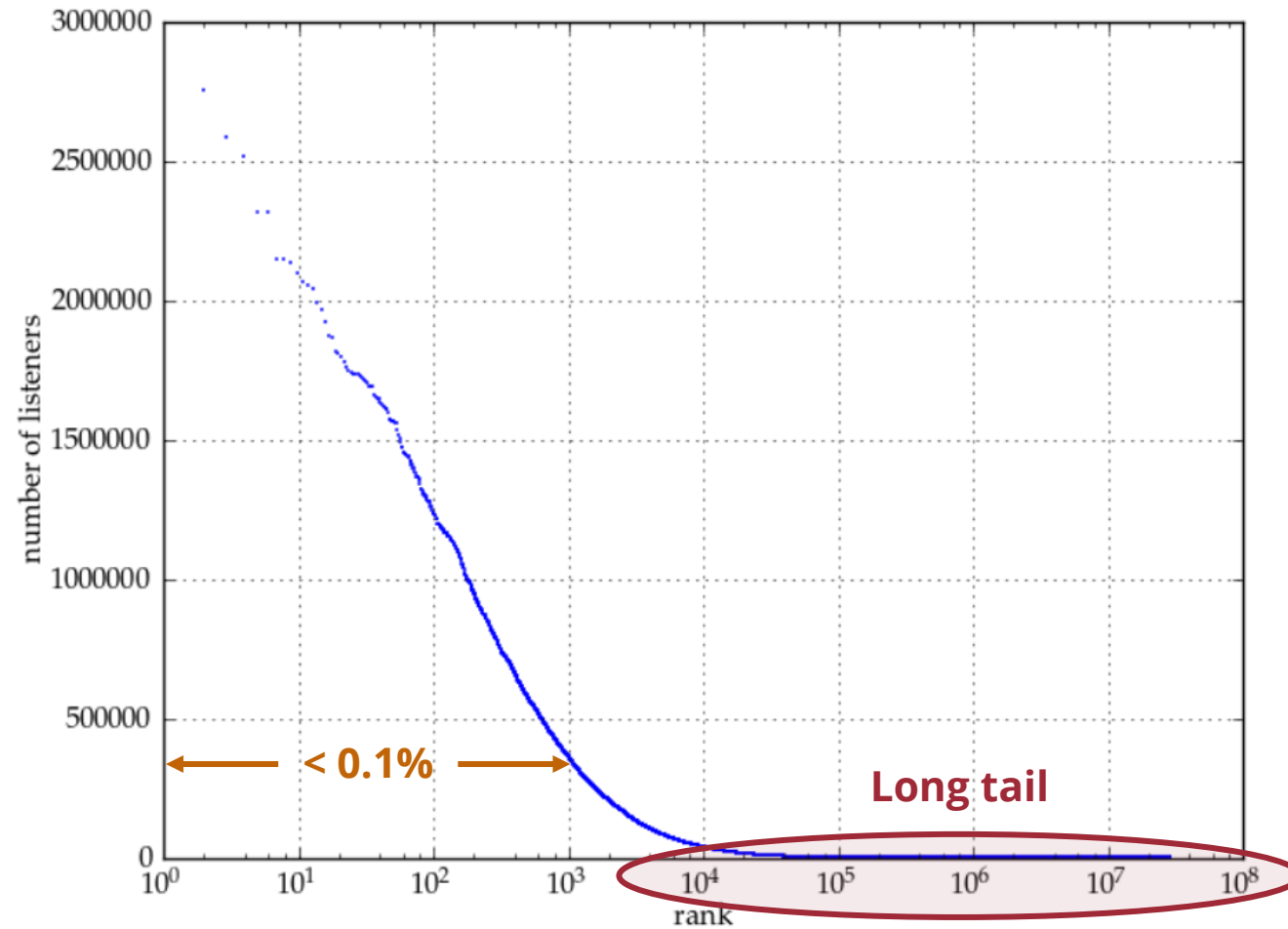


# Music Recommendation

What to play next?


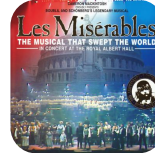





# The Long Tail Problem

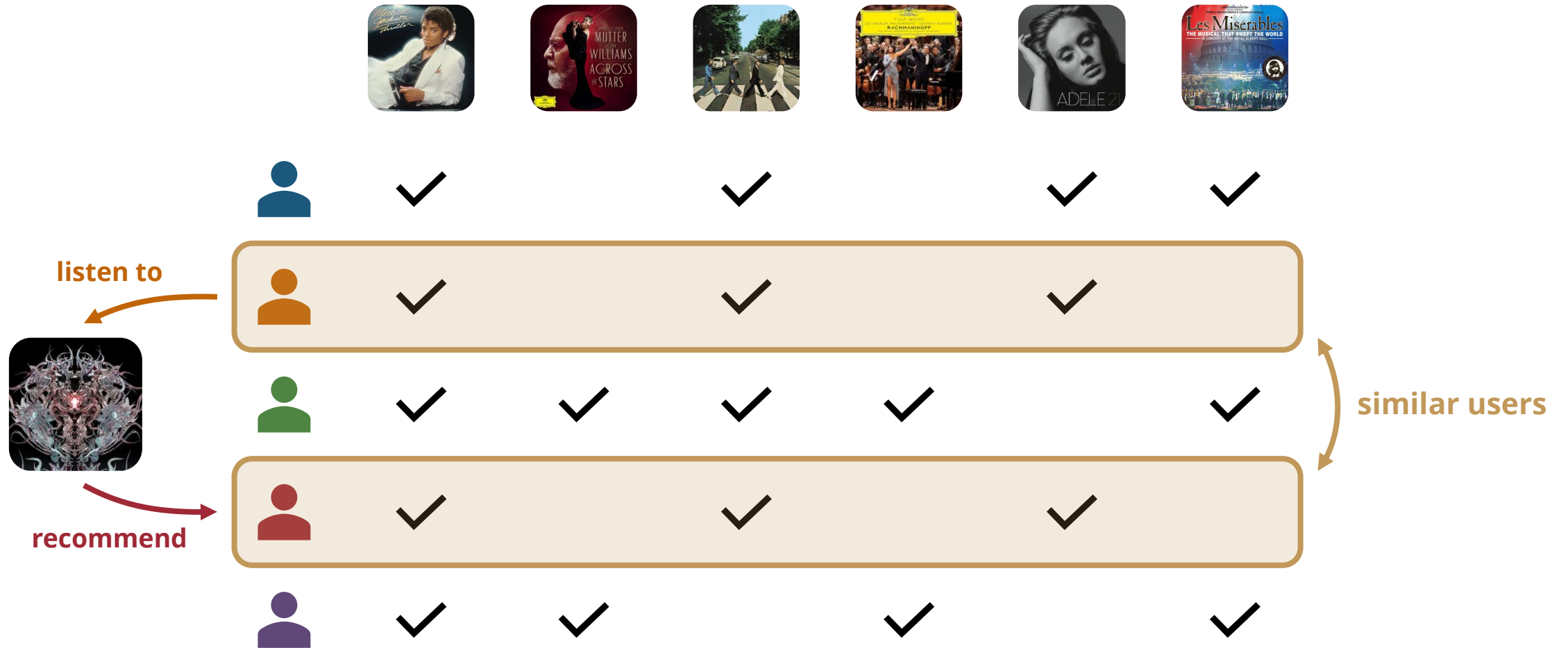


(Source: Levy & Bosteels, 2010)




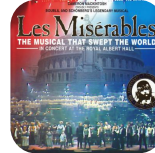



# Collaborative Filtering

						
	✓		✓		✓	✓
	✓		✓		✓	
	✓	✓	✓	✓		✓
	✓		✓		✓	
	✓	✓		✓		✓

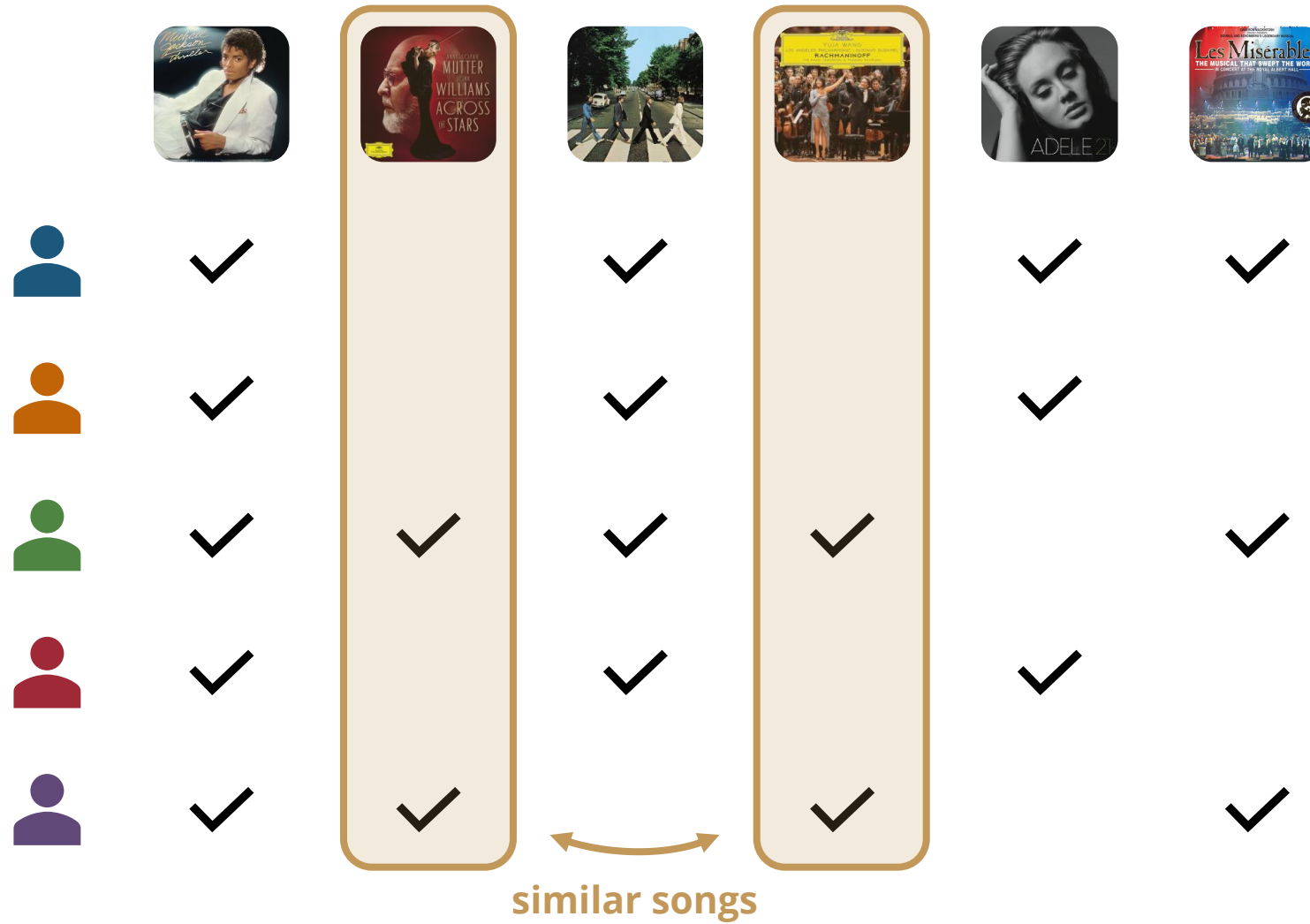
# User-based Collaborative Filtering



# Collaborative Filtering

						
	✓		✓		✓	✓
	✓		✓		✓	
	✓	✓	✓	✓		✓
	✓		✓		✓	
	✓	✓		✓		✓

# Item-based Collaborative Filtering



# Item-based Collaborative Filtering



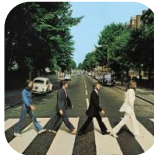










# Item-based Collaborative Filtering

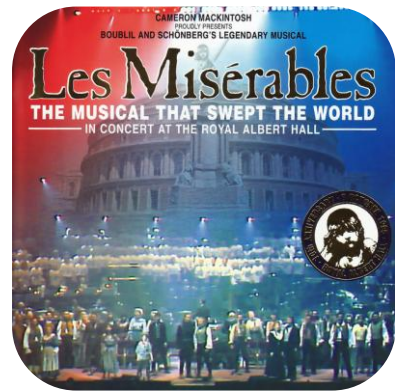




# Cold Star Problem

						New item 
	✓		✓		✓	?
	✓		✓		✓	?
	✓	✓	✓	✓		?
	✓		✓		✓	?
	✓	✓		✓		?

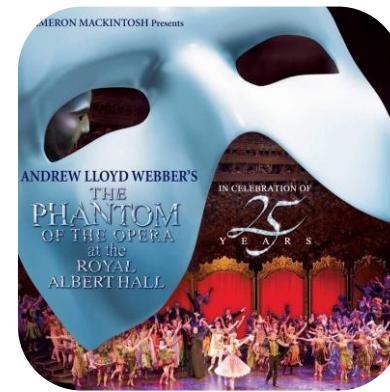
# Content-based Filtering



**Musical**

**Live concert version**

**Concert at Royal Albert Hall**



**Musical**

**Live concert version**

**Concert at Royal Albert Hall**



**similar songs**

# Cold Star Problem

	✓		✓		✓	✓
	✓		✓		✓	
	✓	✓	✓	✓		✓
	✓		✓		✓	
New user		?	?	?	?	?

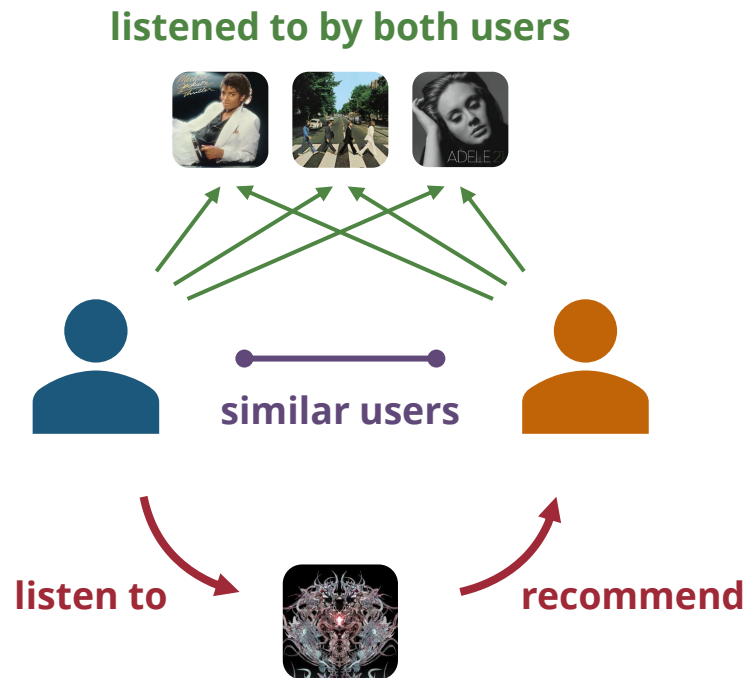
# User Profile Modelling

<b>Data type</b>	<b>Example</b>
Demographic	Age, marital status, gender etc.
Geographic	Location, city, country etc.
Psychographic	<i>Stable</i> : interests, lifestyle, personality etc. <i>Fluid</i> : mood, attitude, opinions etc.

(Source: Song et al., 2012)

# Collaborative Filtering vs Content-based Filtering

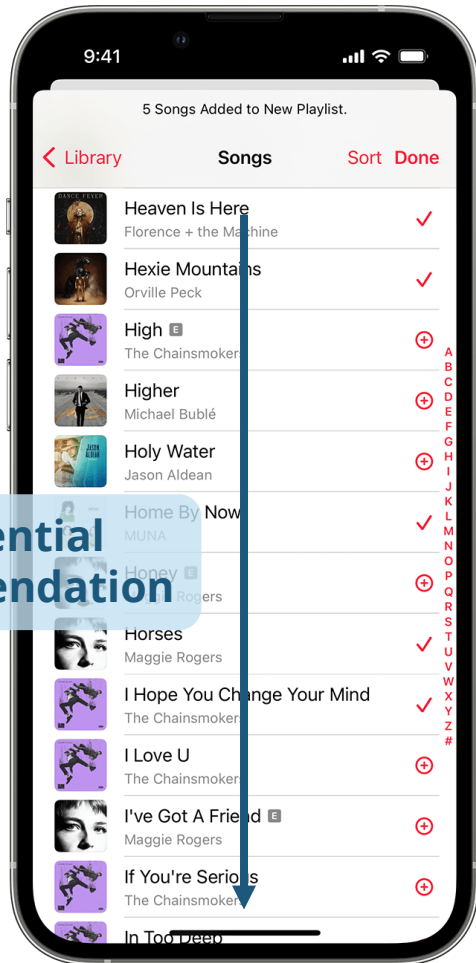
## Collaborative filtering



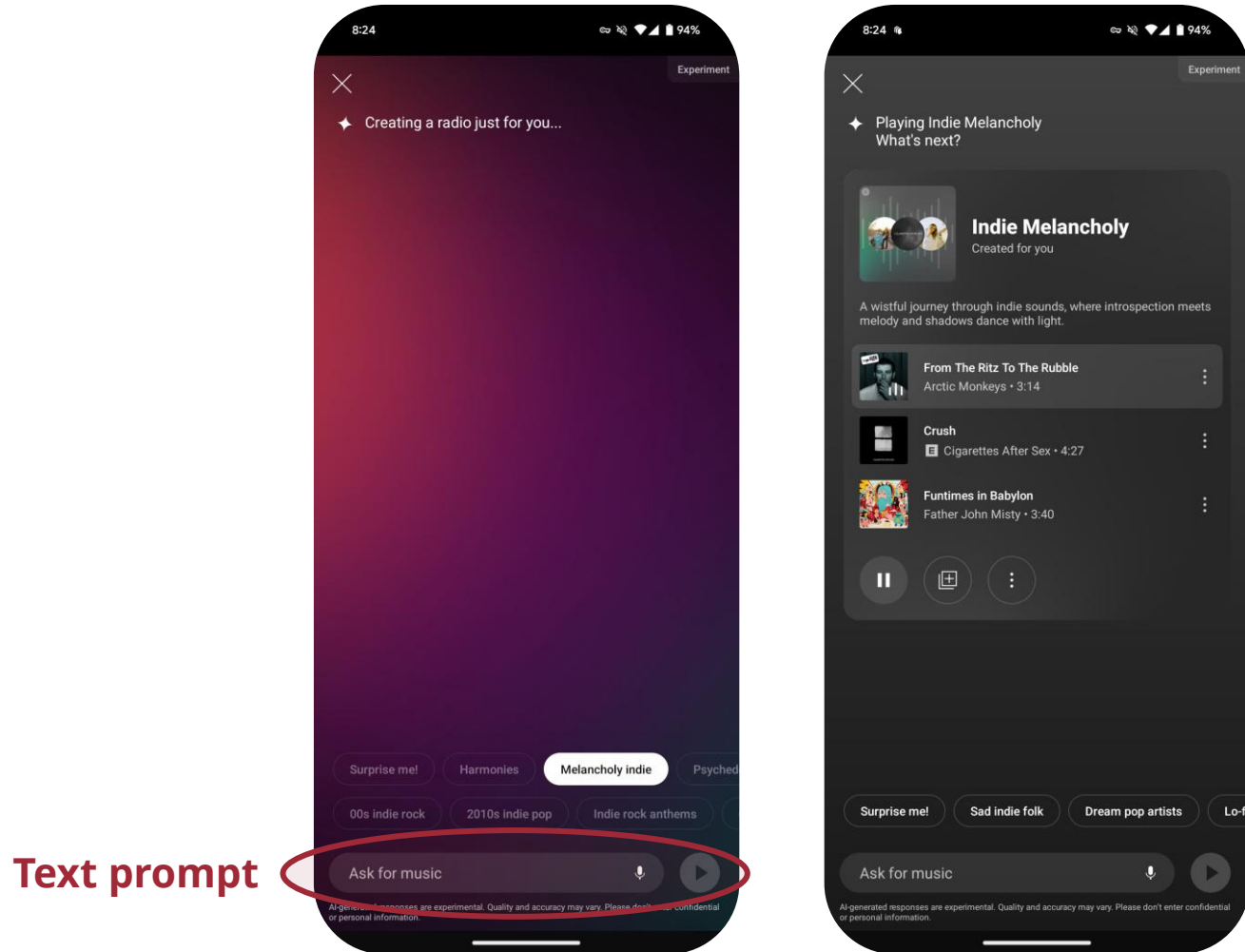
## Content-based filtering



# Music Playlist Generation



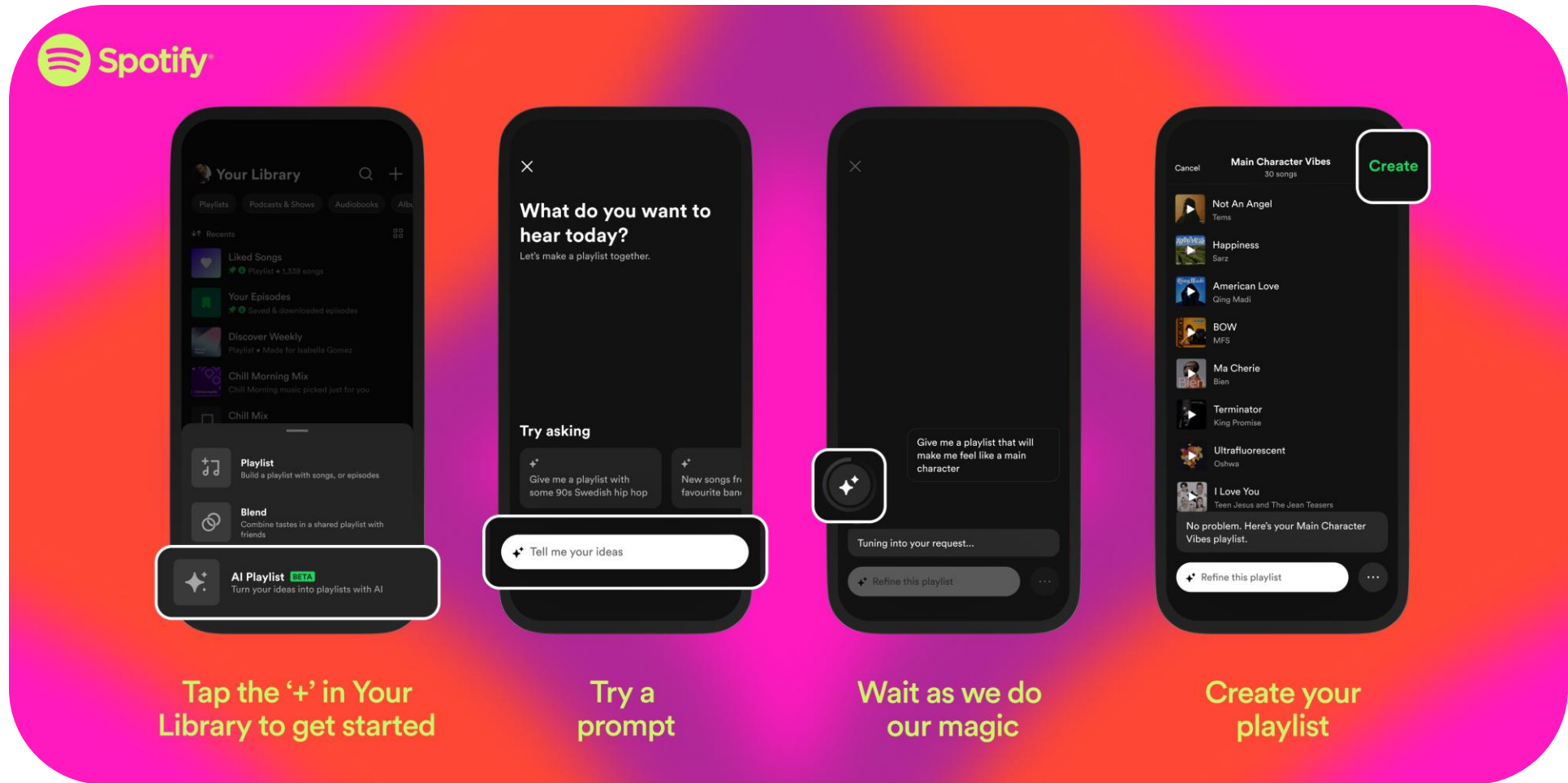
# Ask Music (YouTube Music)



Text prompt

(Source: Android Police)

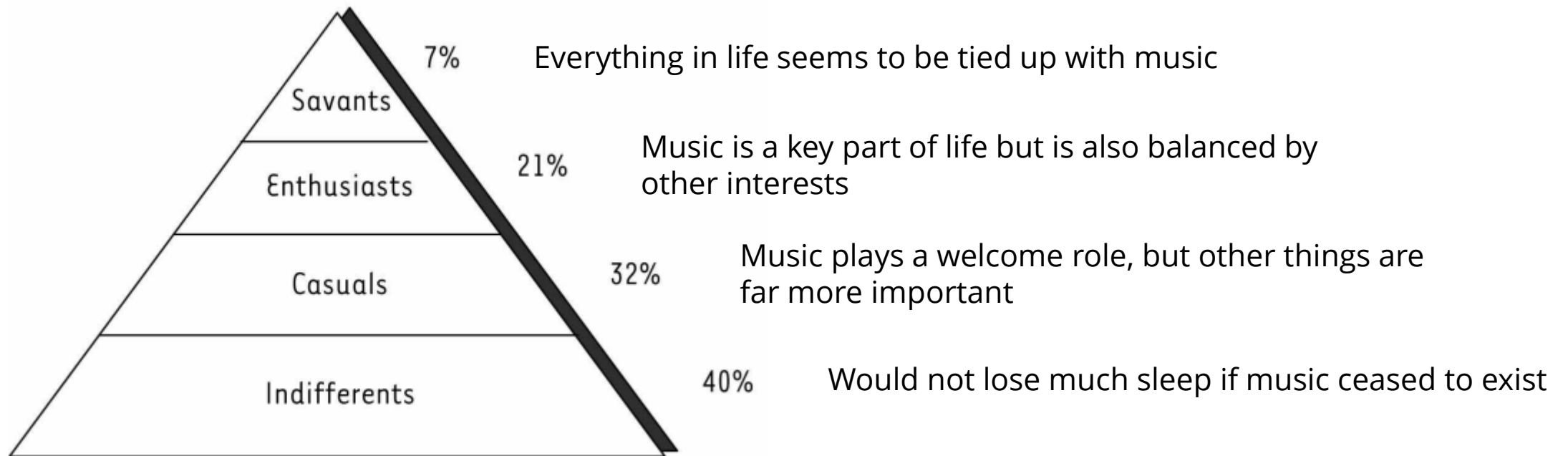
# AI Playlist (Spotify)



(Source: Spotify)



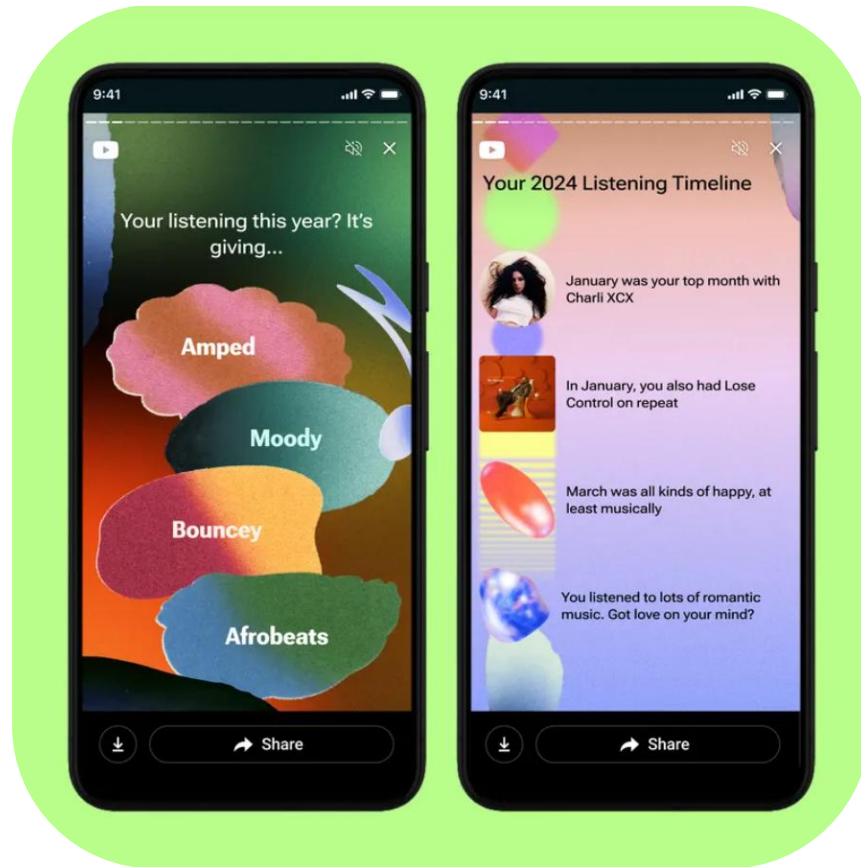
# User Listening Experience Modelling



(Source: Jennings, 2007)

# Listening Behavior Analysis

## YouTube's Music Recap



(Source: YouTube)

## Spotify's Listening Personality



(Source: Spotify)