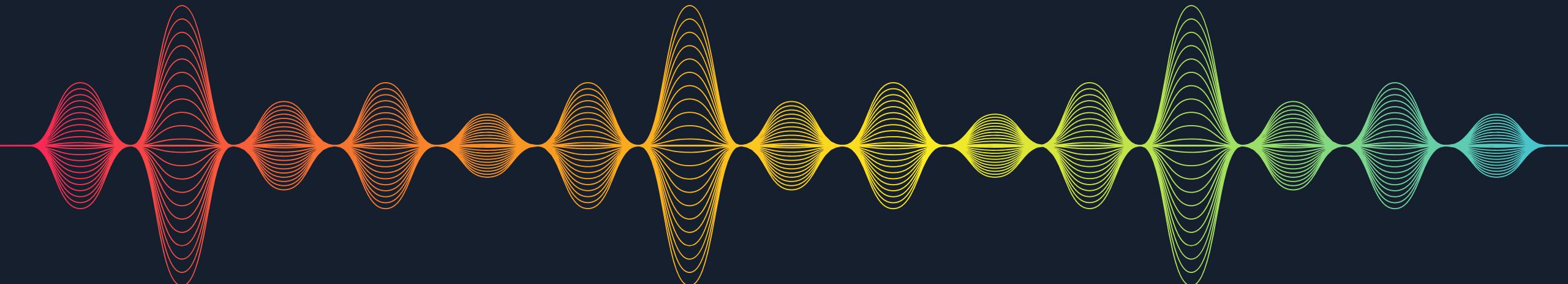


# Musical Style Transfer in Sparse Audio Representation Domains

Flor Sanders

*Sparse and Low-Dimensional Models for High-Dimensional Data*



“

**Without music,  
life would be a mistake.**

— Friedrich Nietzsche

# Table of Contents

1

Musical Style Transfer

2

Music &  
Signal Processing

3

Autoencoders

4

Experiments & Results

5

State of the Art

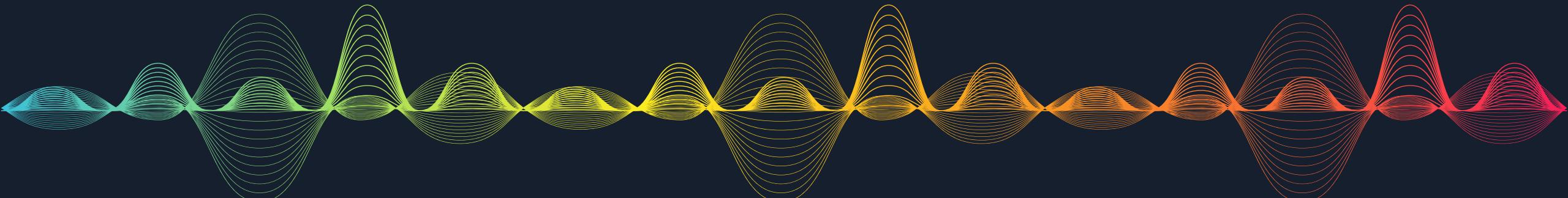
6

Conclusion



# 1

# Musical Style Transfer



# Musical Style Transfer



## ASSUMPTION

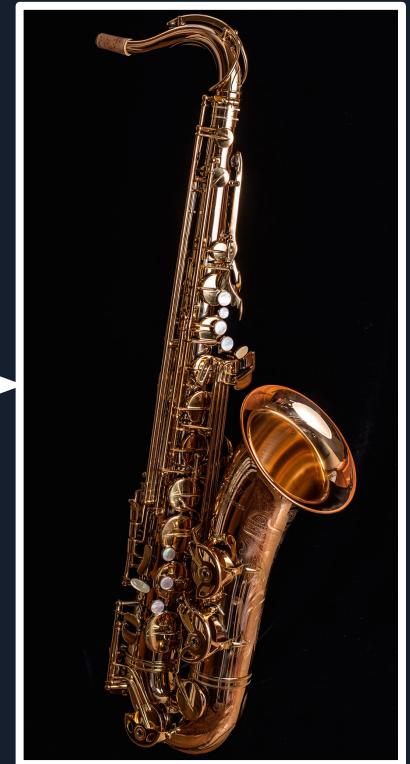
**Music** = content & style

**Content** = domain-variant

**Style** = domain-invariant

## TASK

Modify the style while preserving the content



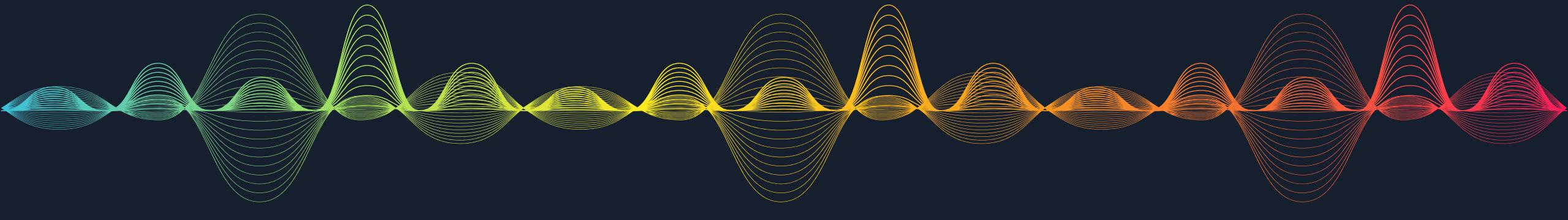
## ISSUE

These concepts are only vaguely defined



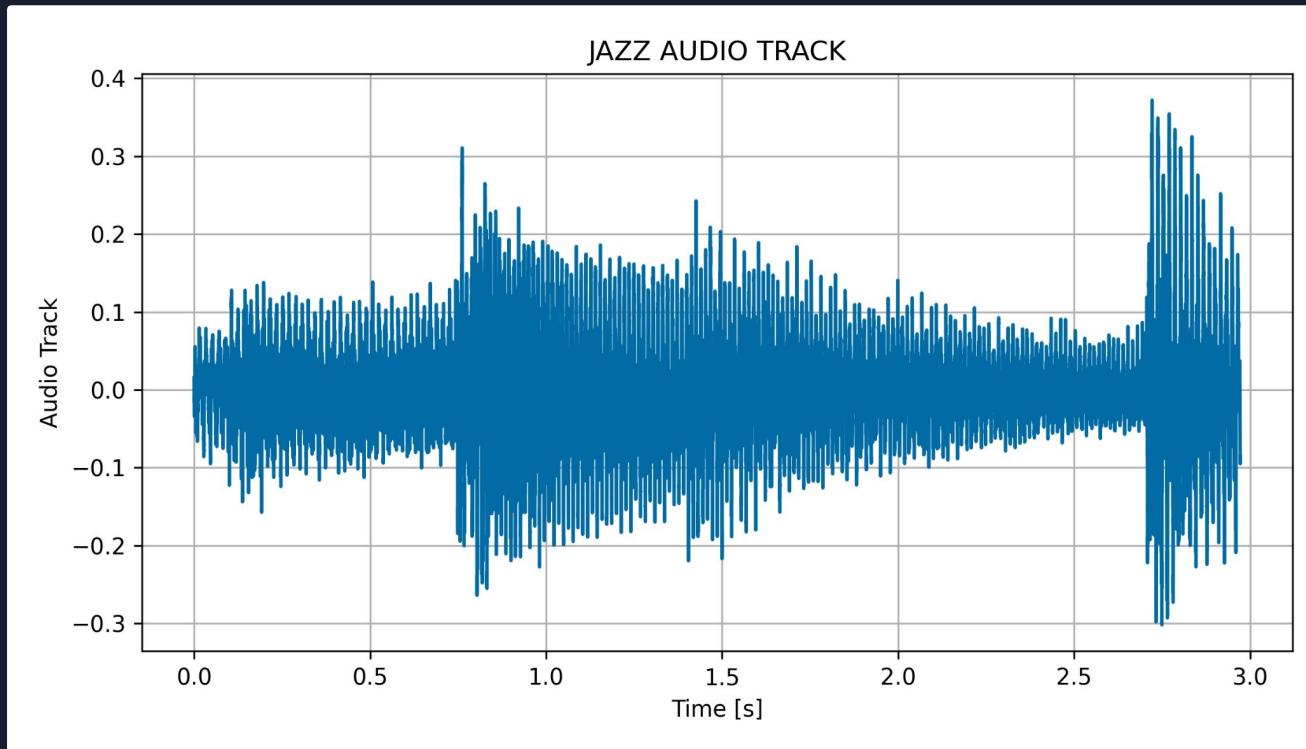
# 2

# Music & Signal Processing





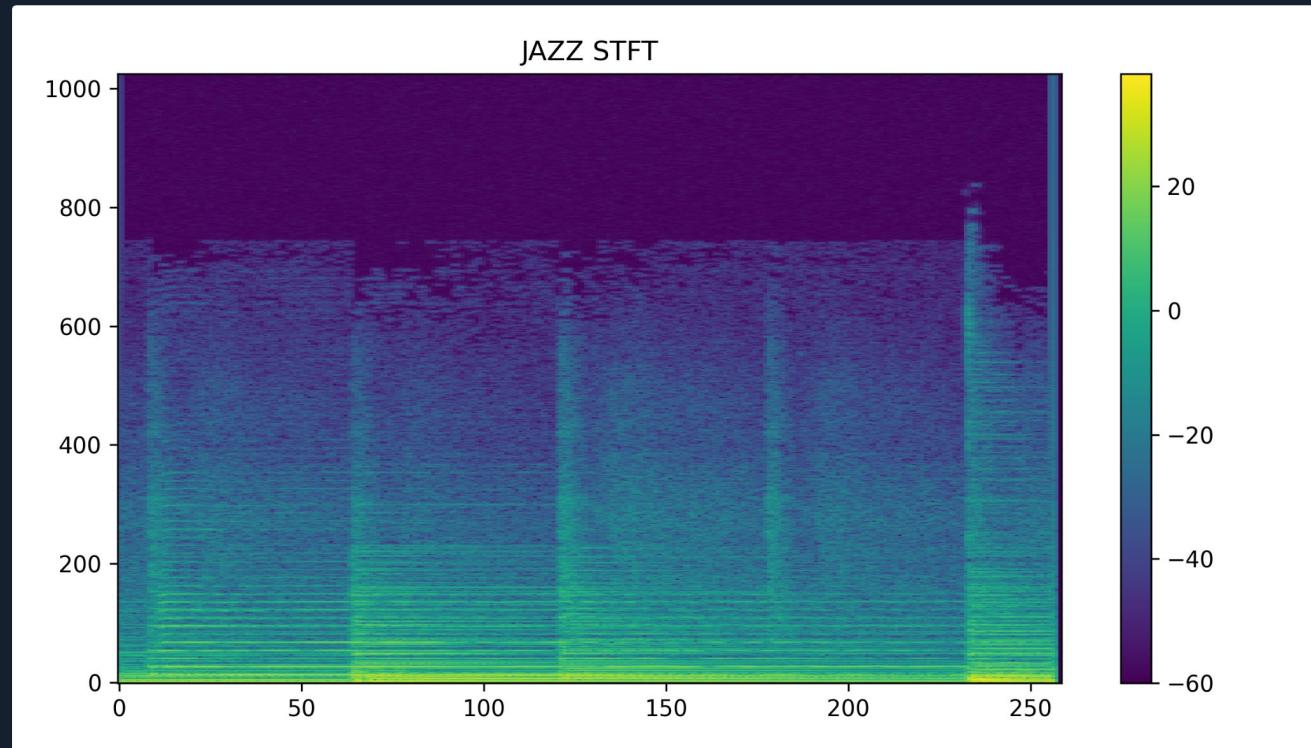
# Music as a Signal



Where is the structure?



# STFT Spectrogram Feature Space

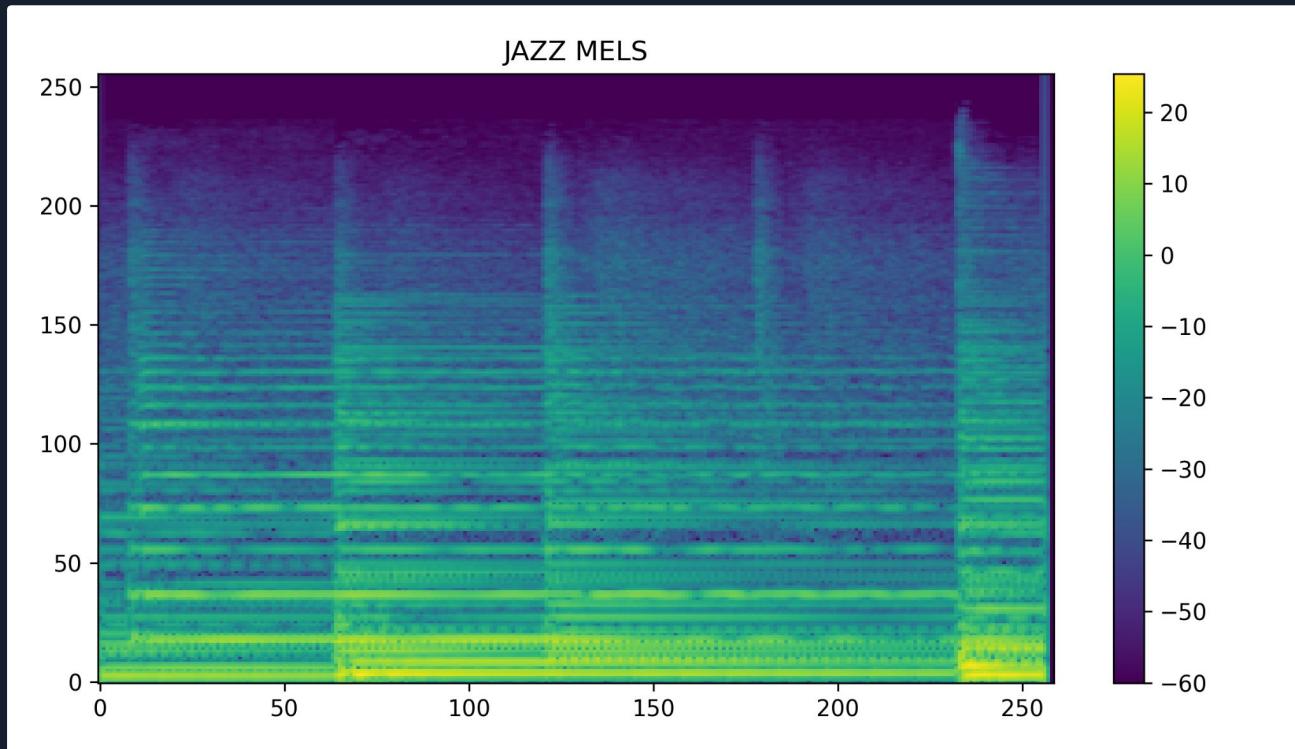


Poor use of space...



$$m = 2595 \log_{10} \left( 1 + \frac{f}{700} \right)$$

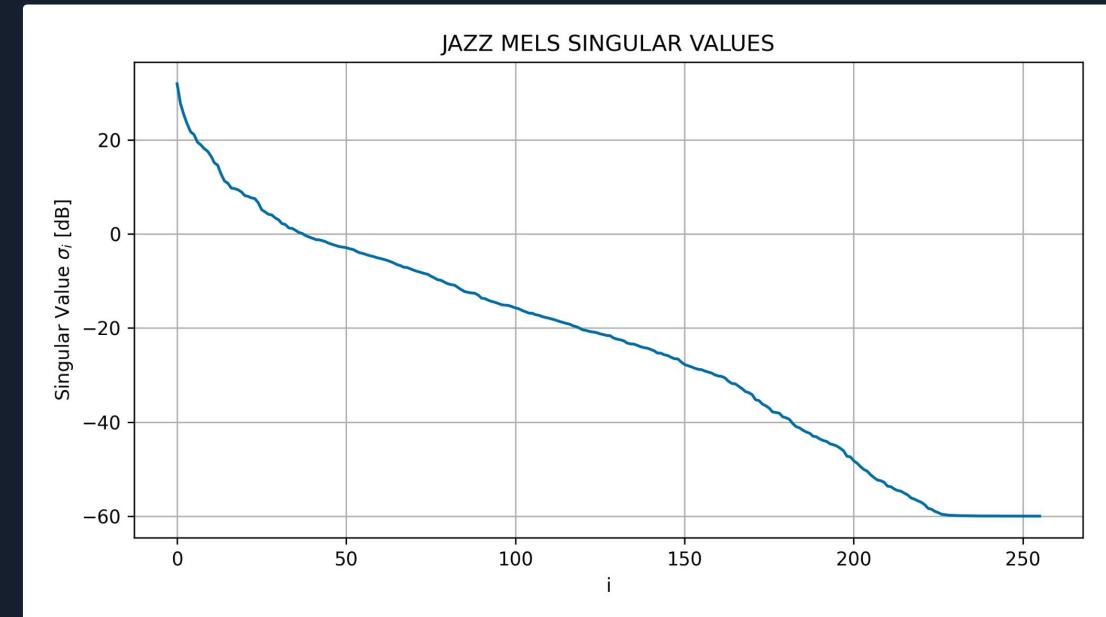
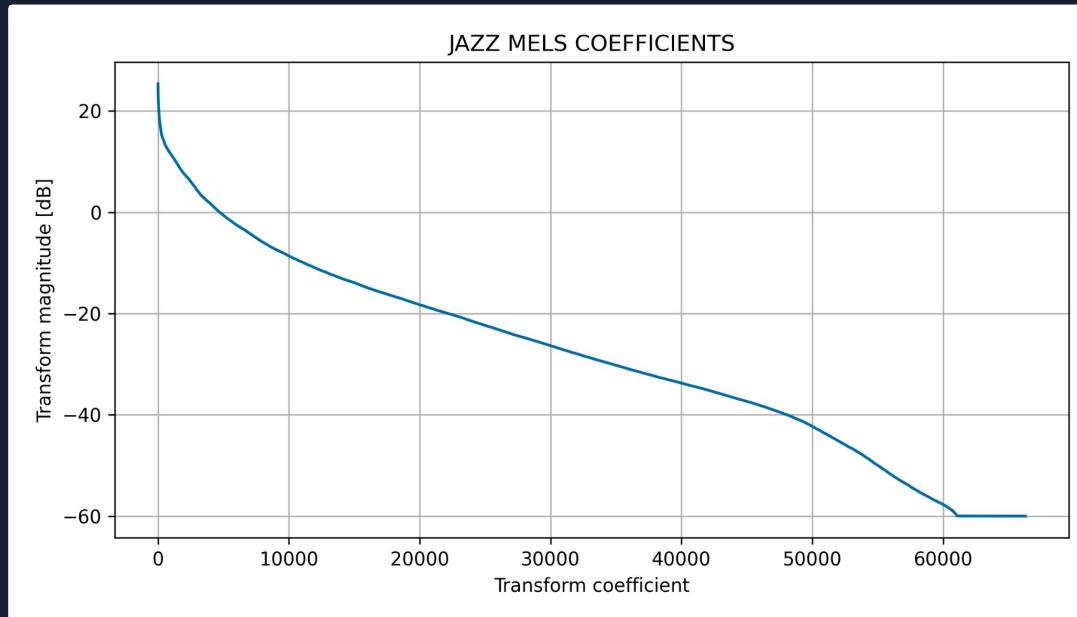
# Mel Spectrogram Feature Space



Equidistant frequencies for human hearing



# Mel Spectrogram Feature Space



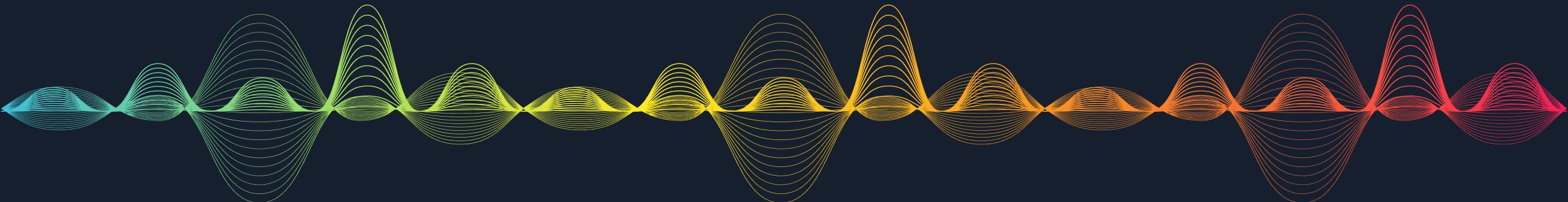
Sparse



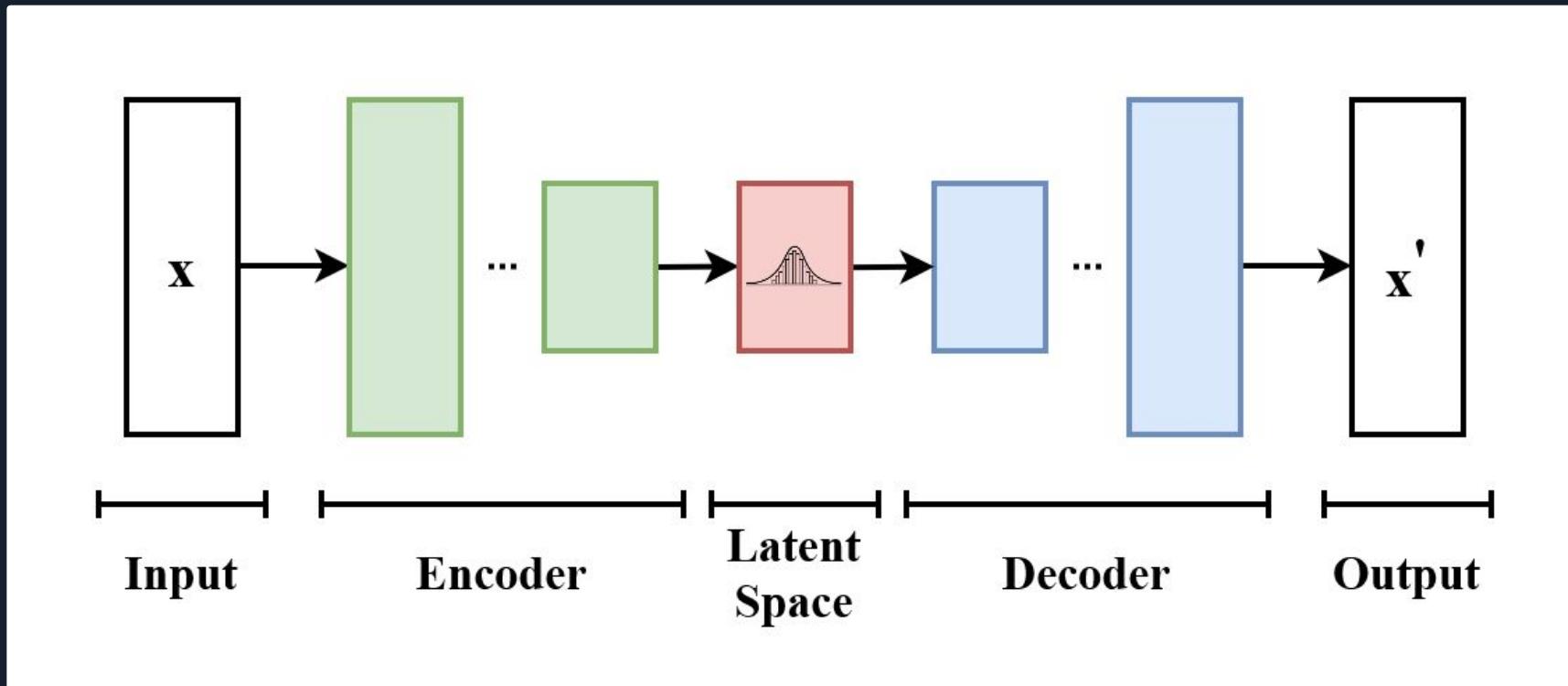
Low Rank

# 3

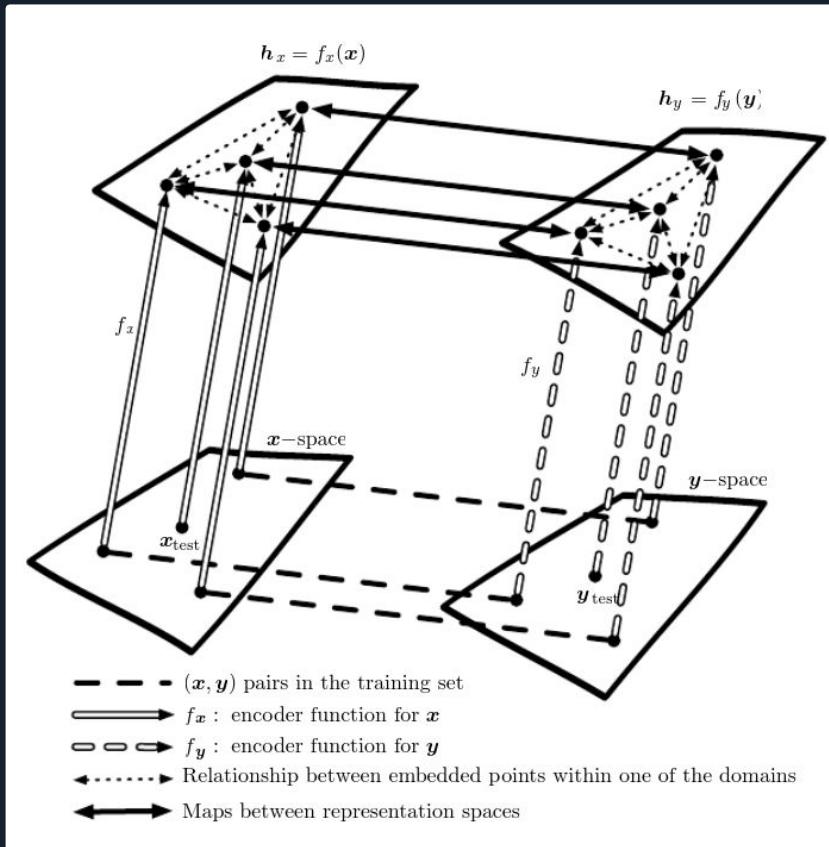
# Autoencoders



# Autoencoders

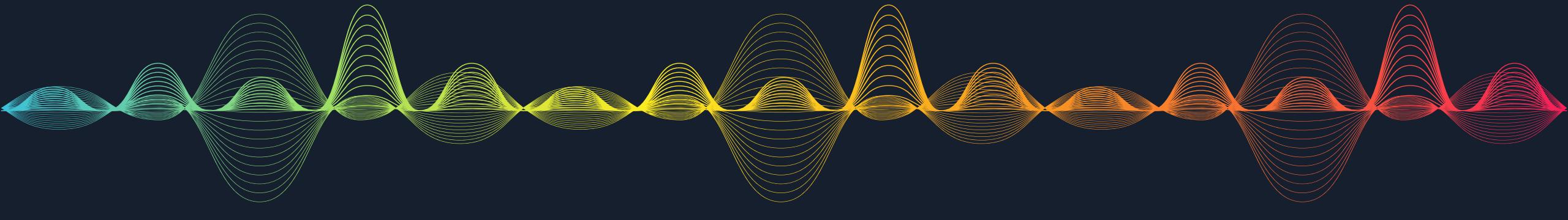


# Representation Learning

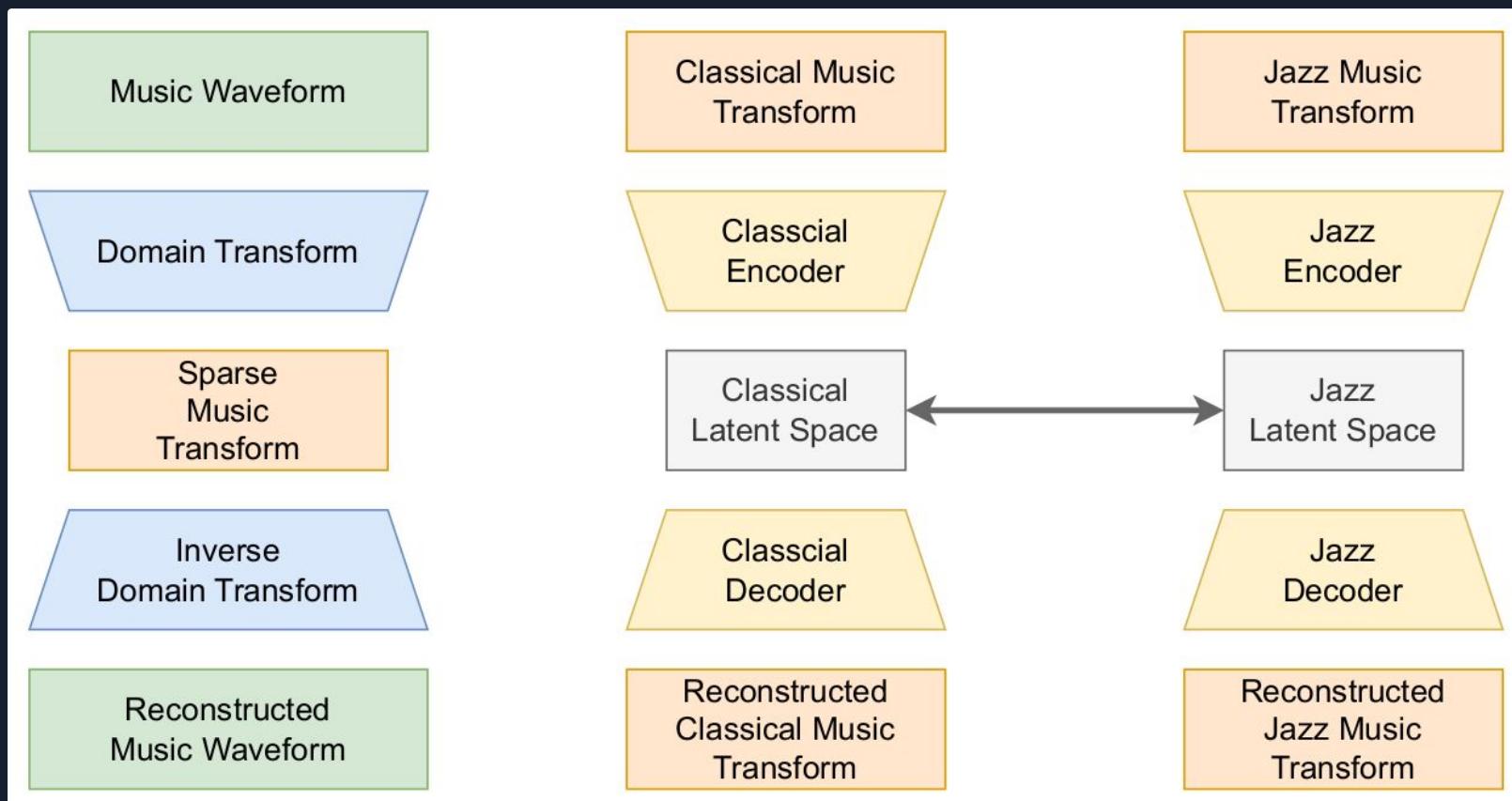


# 4

## Experiments & Results



# The Idea



# The Data

## 2 MUSICNET

Start	End	Instrument	Note	Measure	Beat	Note Value
45.29	45.49	Violin	G5	21	3	Eighth
48.99	50.13	Cello	A#3	24	2	Dotted Half
82.91	83.12	Viola	C5	51	2.5	Eighth

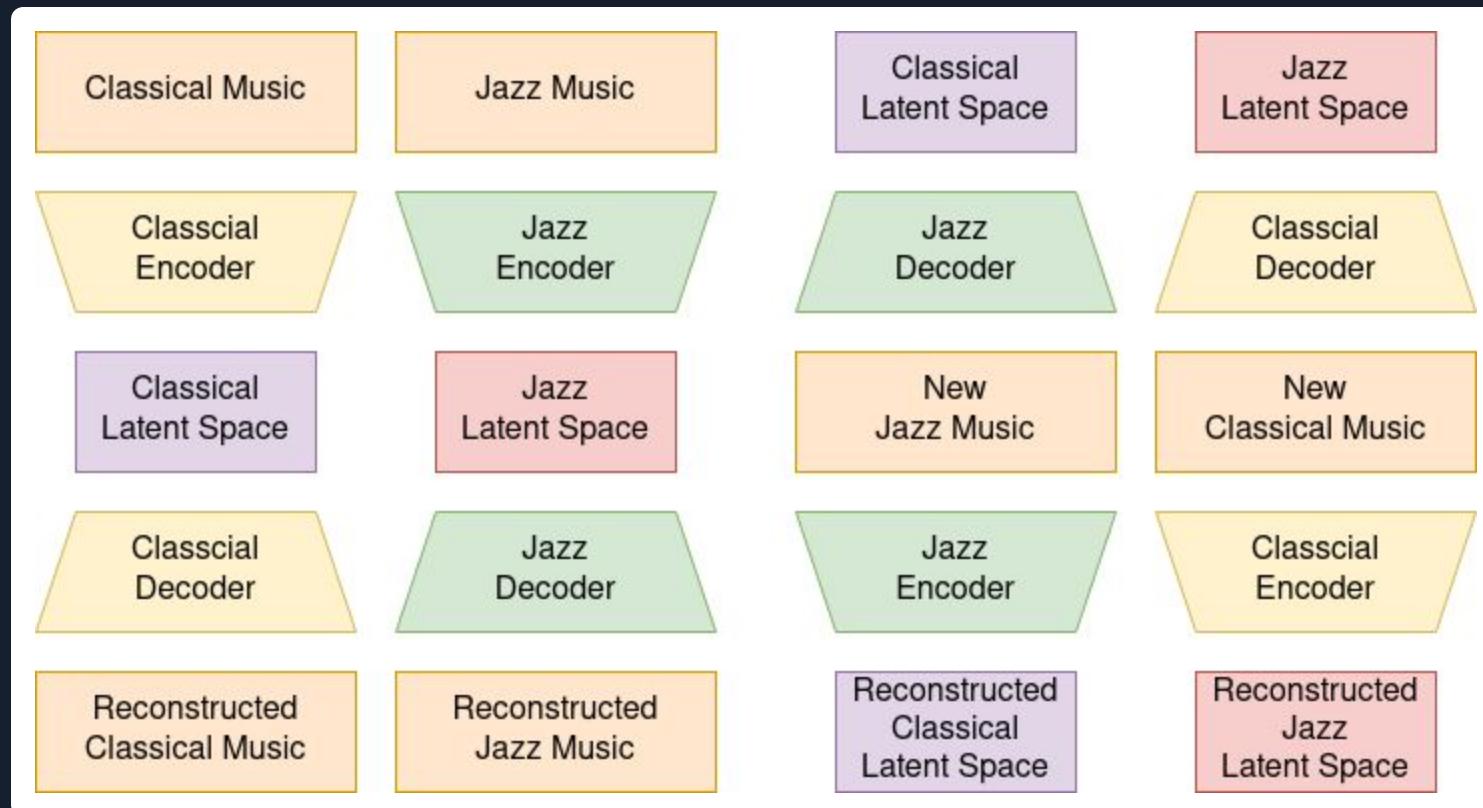


Classical Music

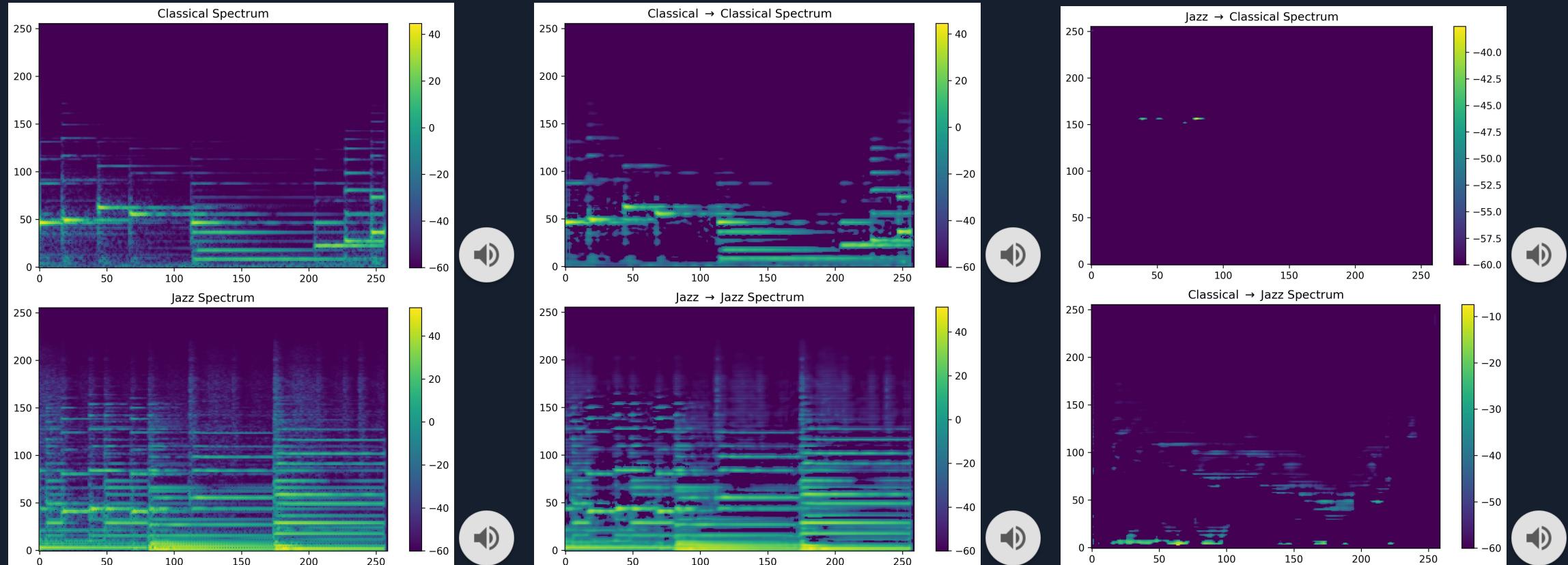


Jazz Music

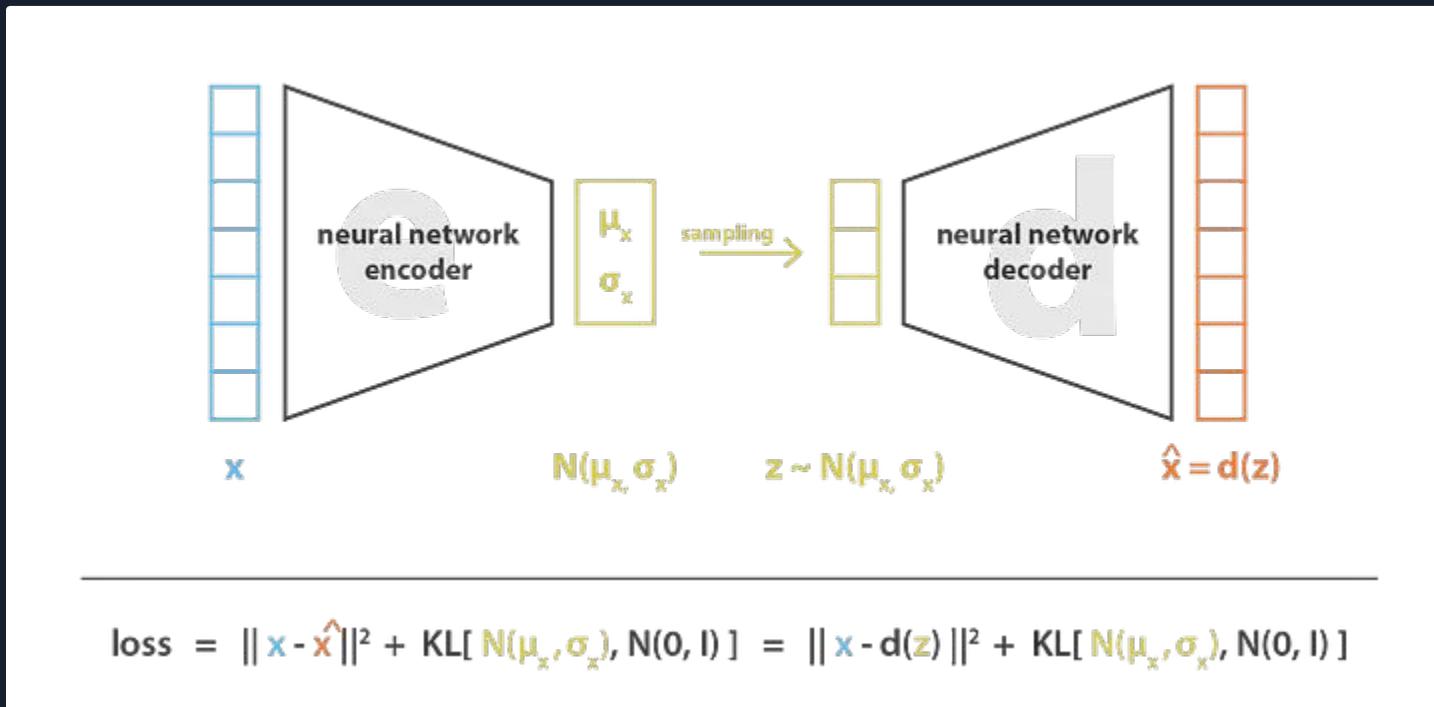
# 1 - Naive Twin Autoencoders



# 1 - Naive Twin Autoencoders



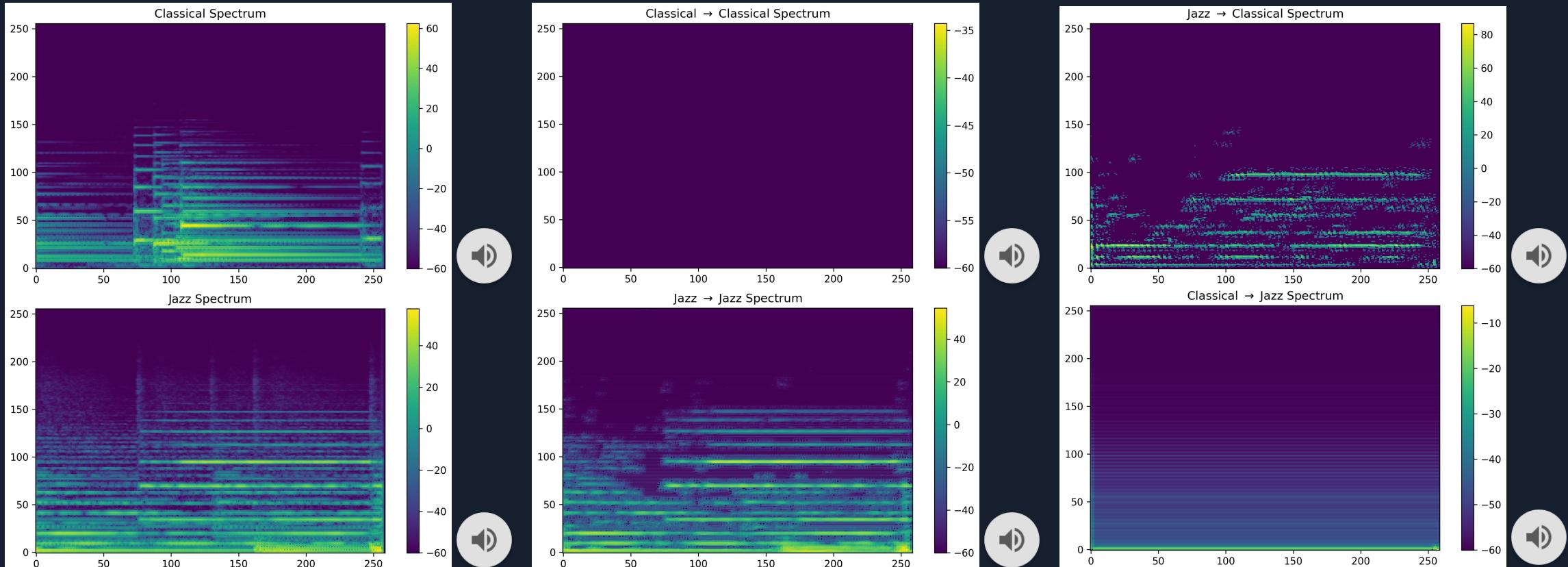
# 2 - Variational Twin Autoencoders



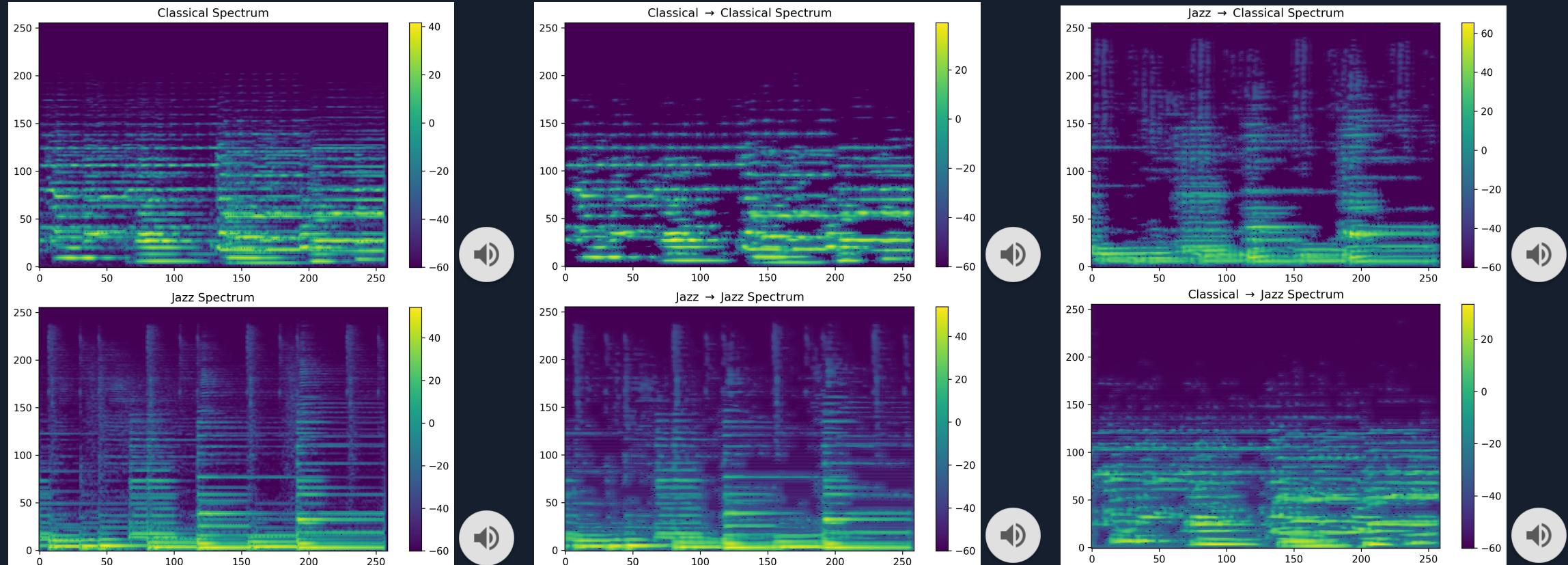
$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right),$$



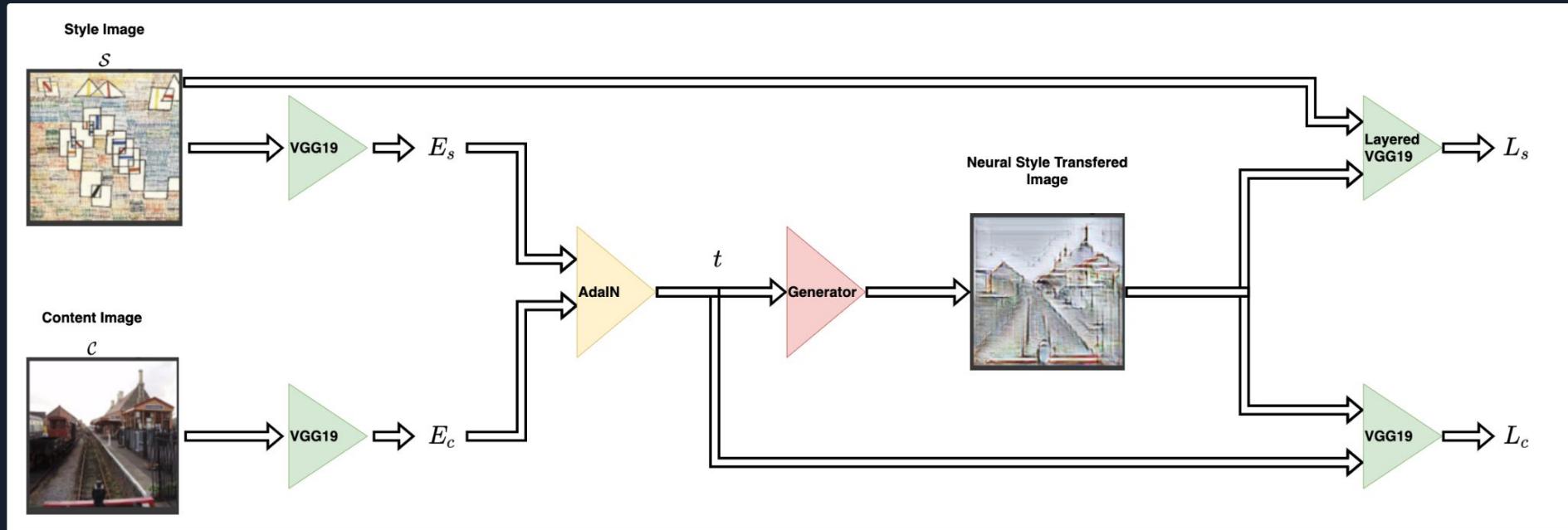
# 2 - Variational Twin Autoencoders (5 Epochs)



# 2 - Variational Twin Autoencoders



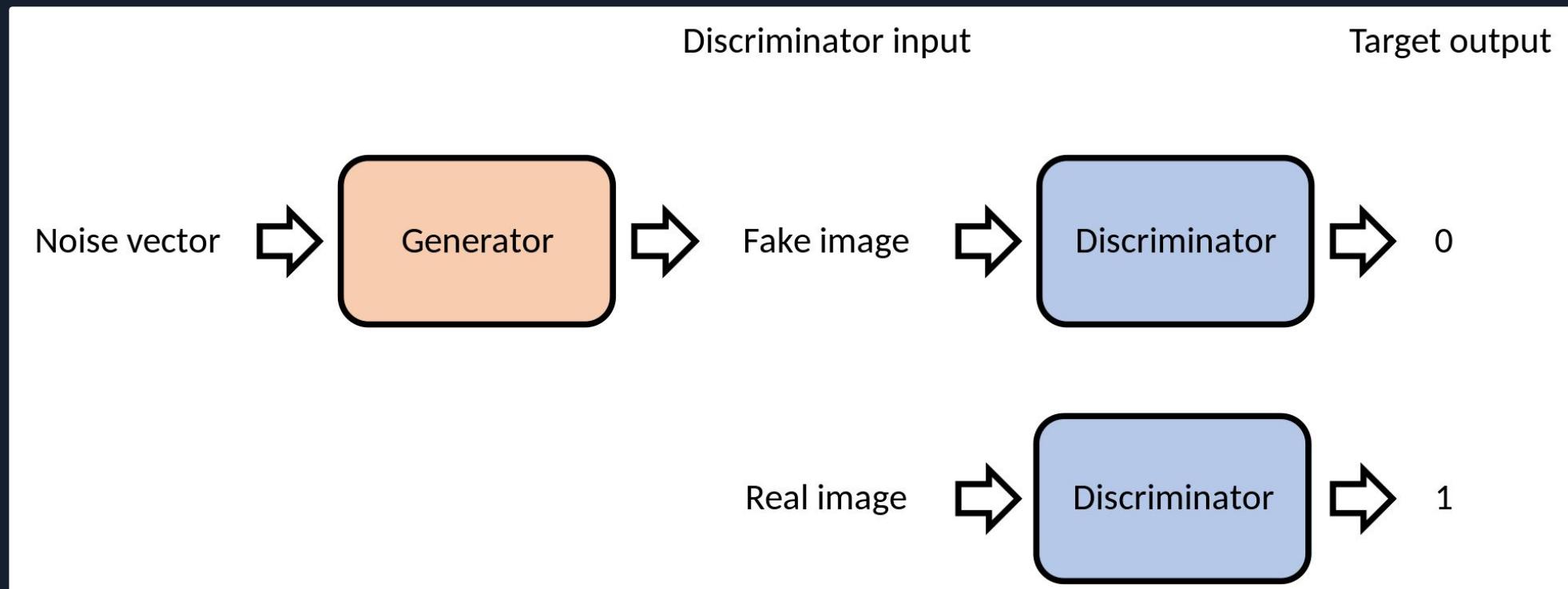
# 3 - GAN + AdaIN Twin Autoencoders



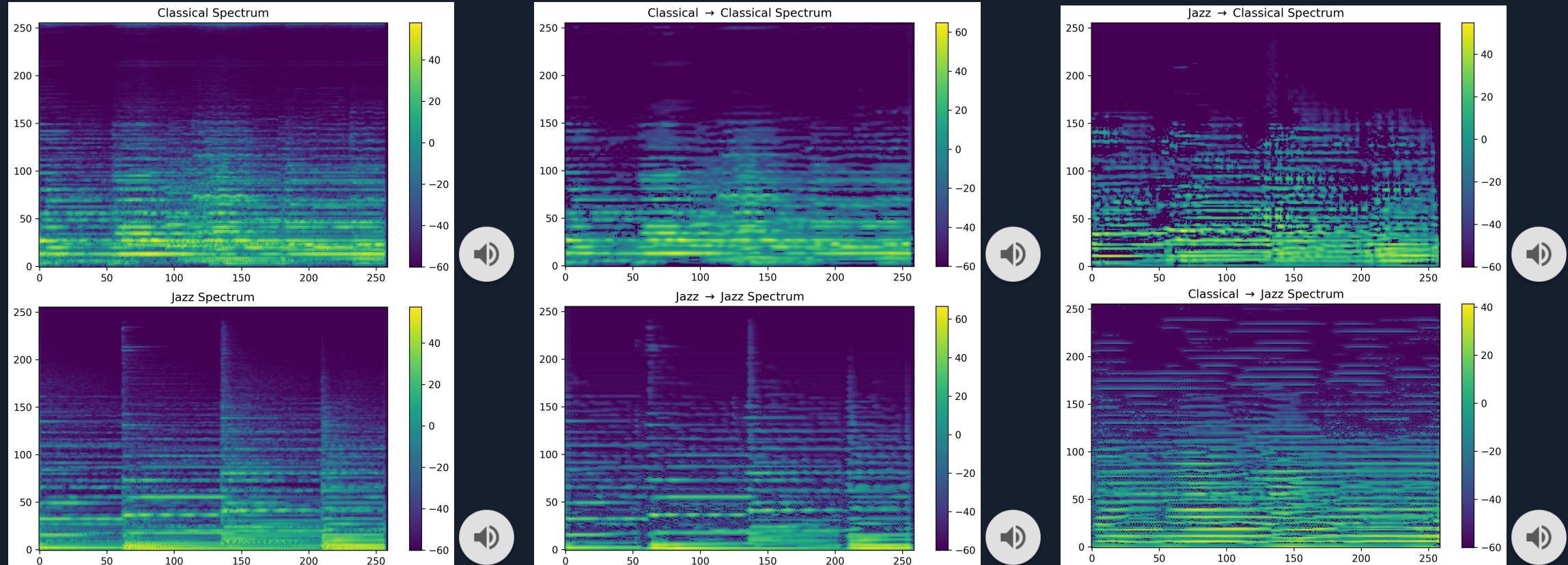
$$AdaIn(x, y) = \sigma(y)\left(\frac{x - \mu(x)}{\sigma(x)}\right) + \mu(y)$$



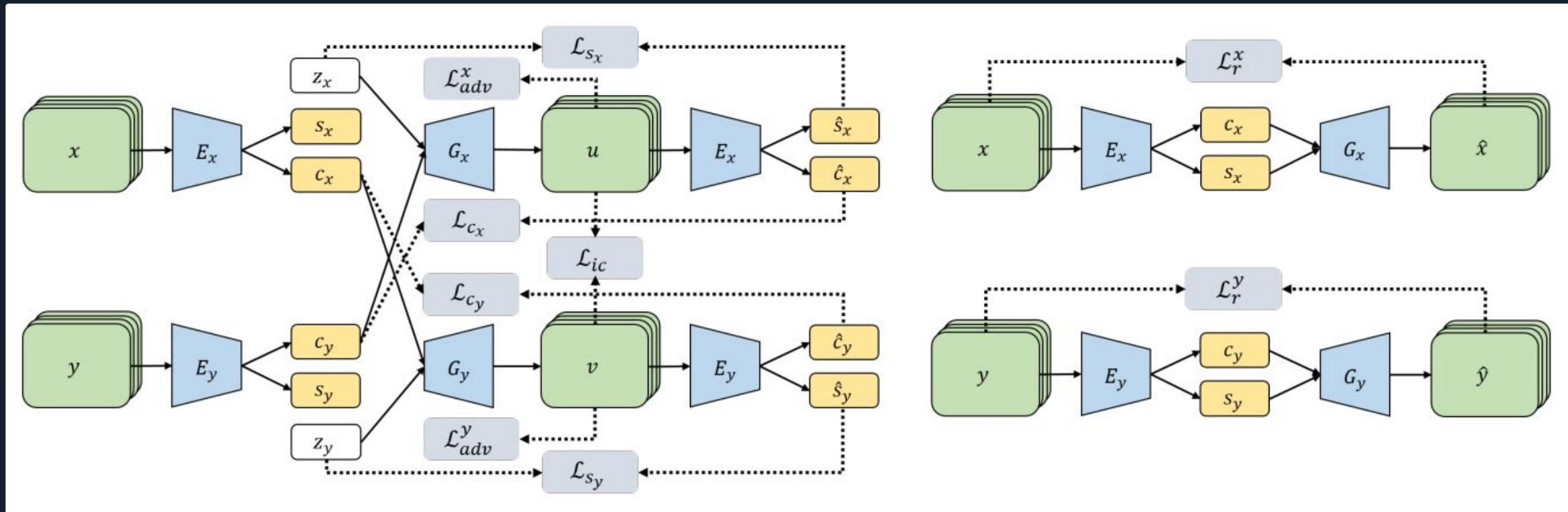
# 3 - GAN + AdaIN Twin Autoencoders



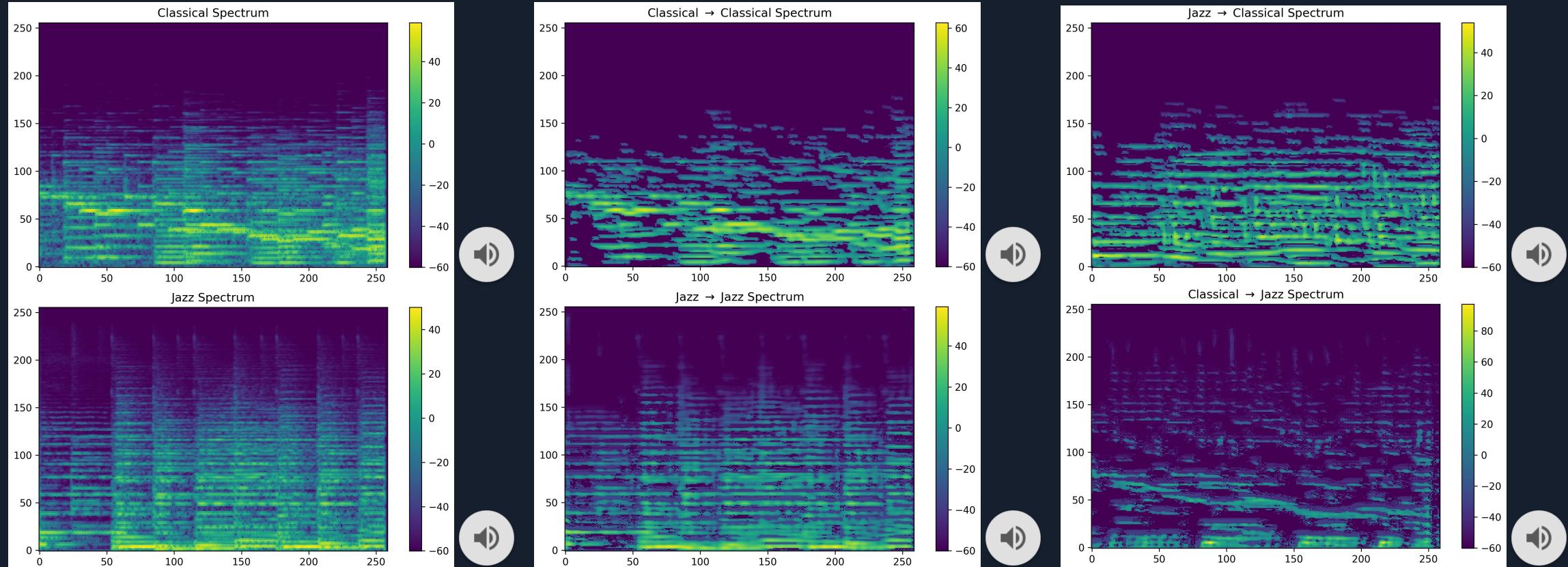
# 3 - GAN + AdaIN Twin Autoencoders



# 4 - Variational + GAN + Adain Twin Autoencoders

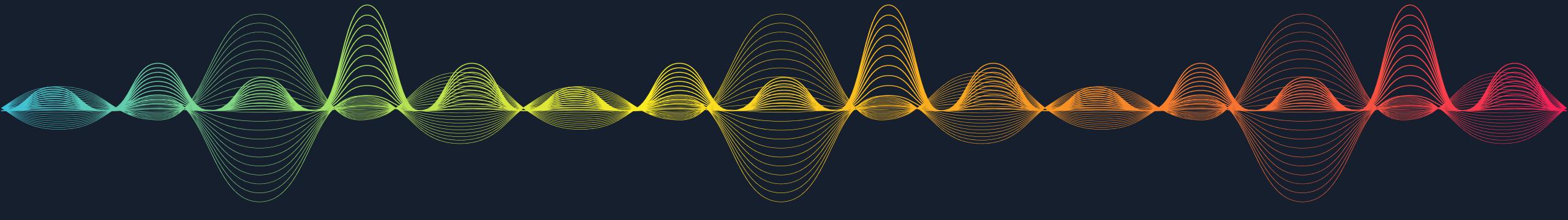


# 4 - Variational + GAN + Adain Twin Autoencoders

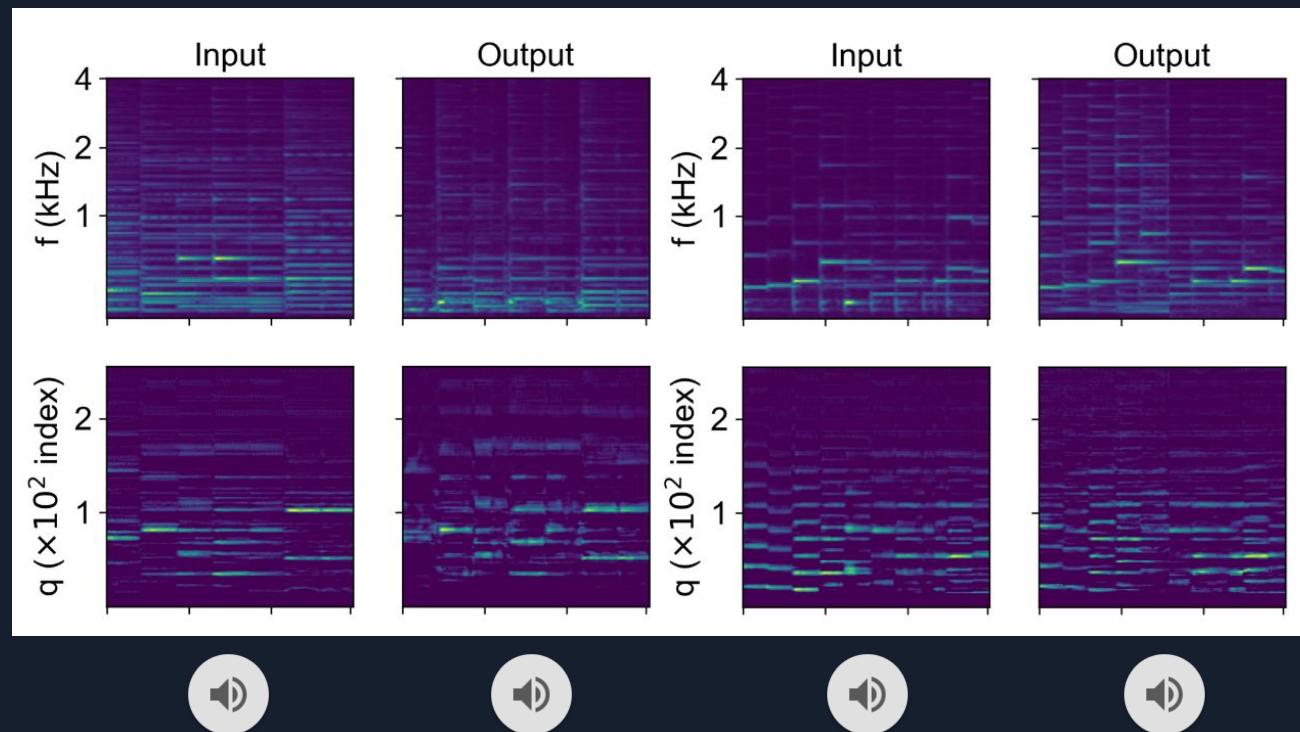


# 5

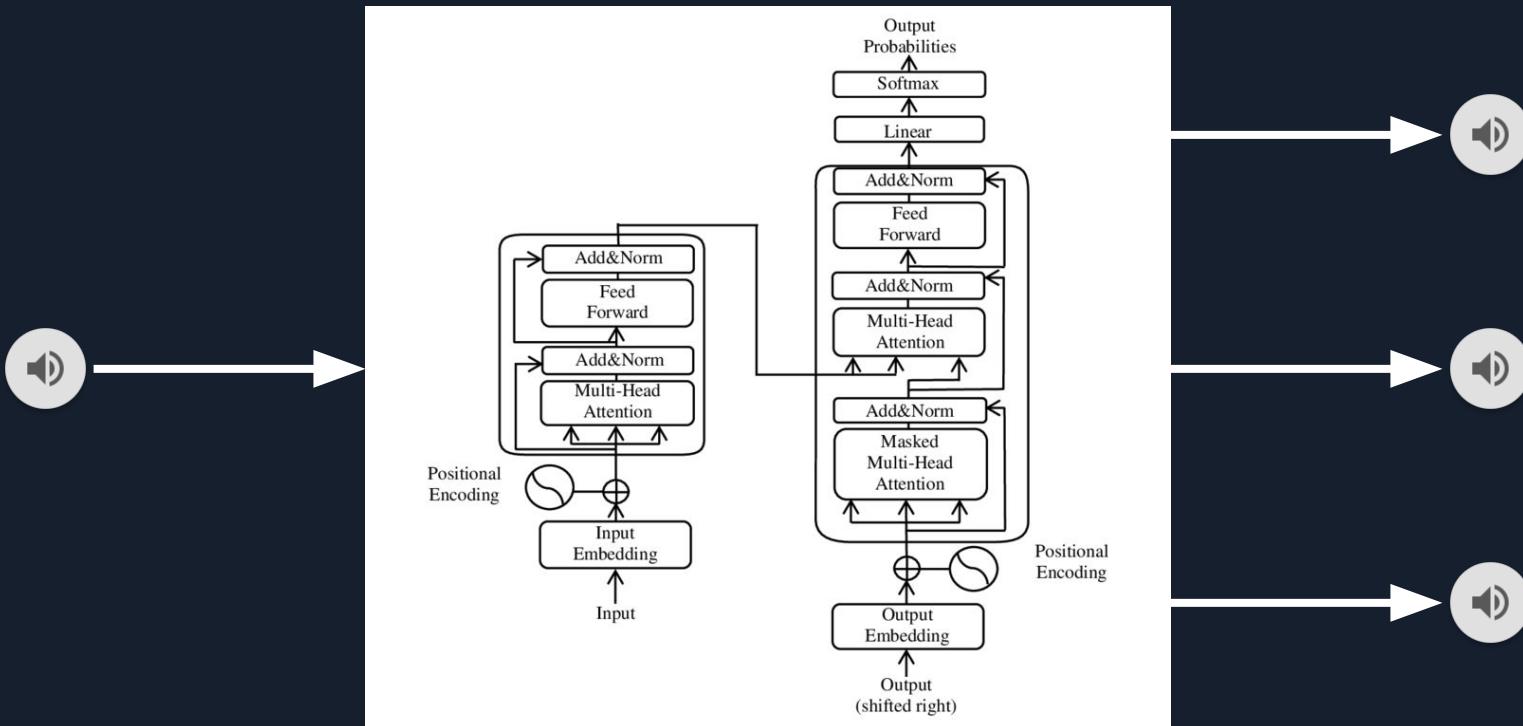
# State of the Art



# (2018) Play as You Like: Timbre-enhanced Multi-modal Music Style Transfer

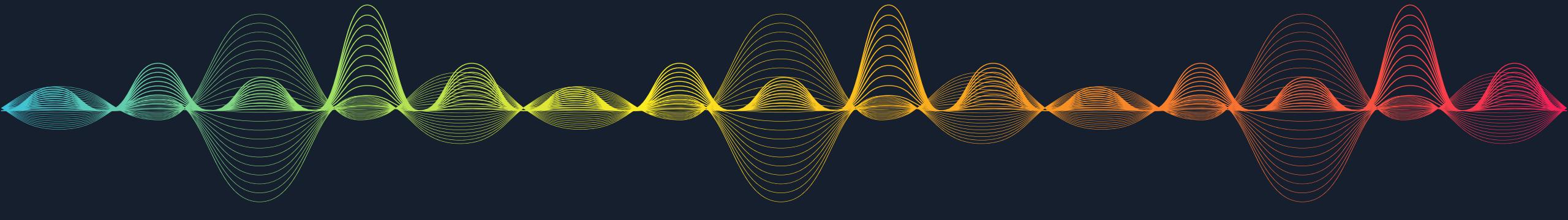


# (2023) Meta MusicGen: Simple and Controllable Music Generation

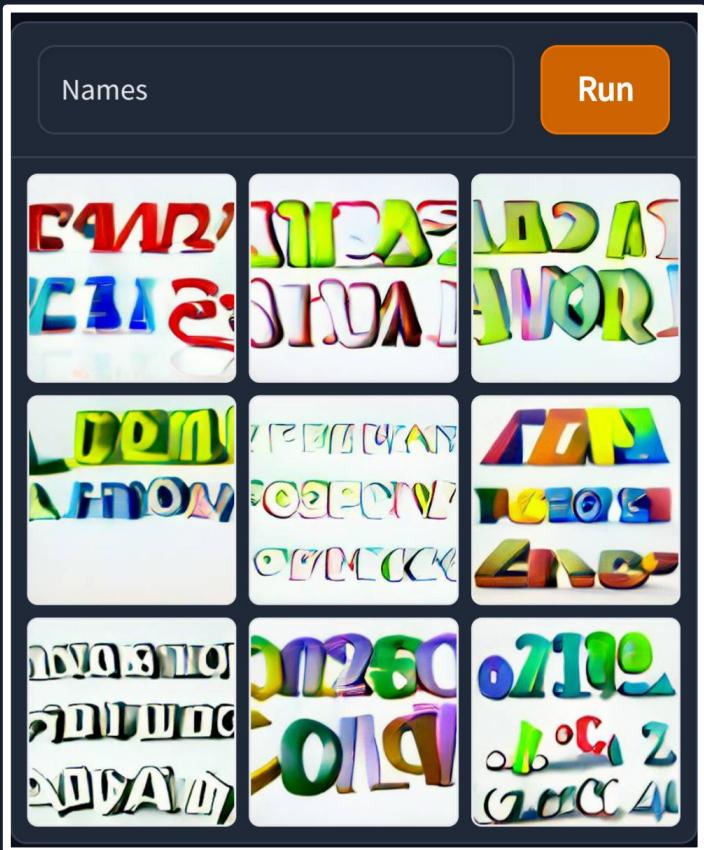


# 6

# Conclusion

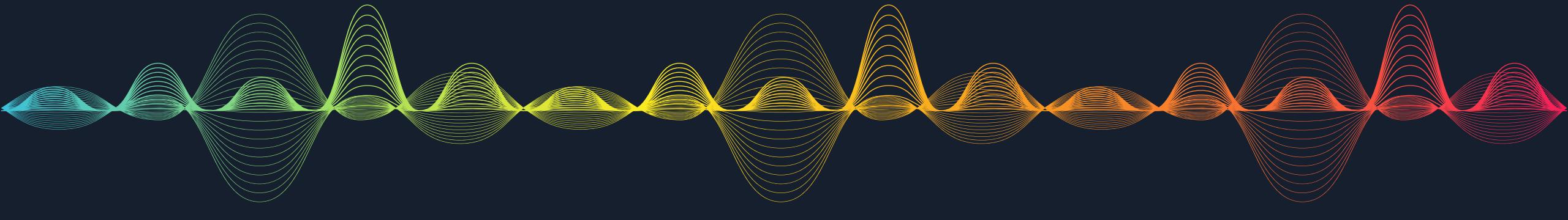


# Style Transfer / Generative AI is a Difficult Problem

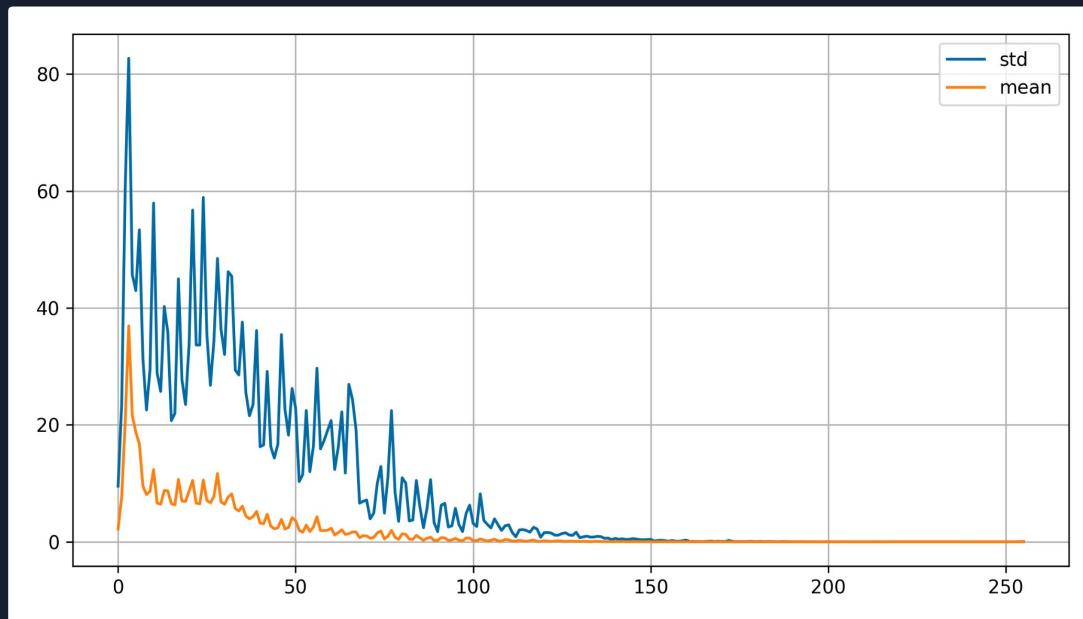


# 7

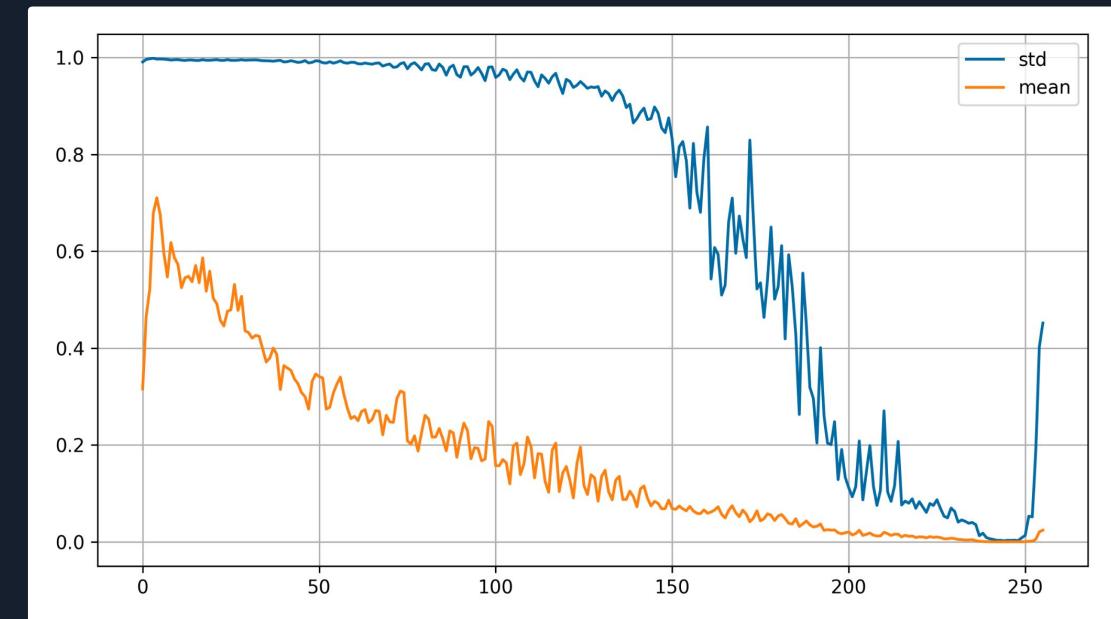
# Appendices



# Mel Spectrogram Normalization



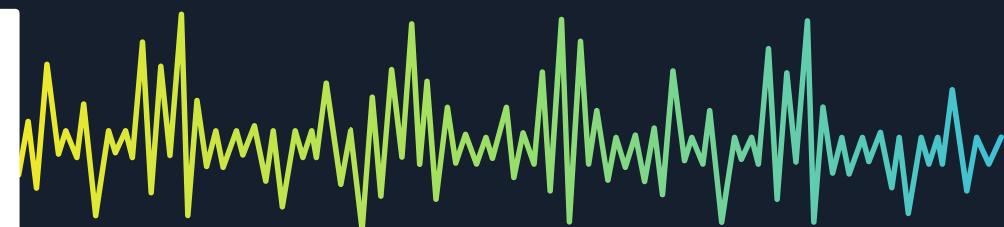
Before Normalization



After Normalization



$$\frac{\log_{10} (1 + X/\sigma_X)}{\sigma_{\log_{10}(f)}}$$



# Audio Reconstruction

---

**Algorithm 1** Griffin-Lim algorithm (GLA)

---

Fix the initial phase  $\angle c_0$

**Initialize**  $c_0 = s \cdot e^{\cdot i \angle c_0}$

**Iterate** for  $n = 1, 2, \dots$

$$c_n = P_{\mathcal{C}_1}(P_{\mathcal{C}_2}(c_{n-1}))$$

**Until convergence**

$$x^* = \mathbf{G}^\dagger \mathbf{c}_n$$

---

$$P_{\mathcal{C}_1}(c) = \mathbf{G}\mathbf{G}^\dagger c$$

$$P_{\mathcal{C}_2}(c) = s \cdot e^{\cdot i \angle c}.$$

