# DeepDeadNet

Benjamin Wakefield
Northeastern University
wakefield.b@northeastern.edu

## ABSTRACT

DeepDeadNet is a web application built with Python and Django, and leverages the Musika and Demucs libraries. It allows users to generate new music in the style of the Grateful Dead, and isolate or mix each instrument to their preference. The goal of the project is to provide an accessible graphical user interface to people who are interested in AI music generation, but are unfamiliar with command-line interfaces. Additionally, this project combines multiple neural networks into a single workflow, in order to demonstrate how they can be applied for a practical use case. Broadly, this project aims to explore the current potential of music generation and alteration by artificial intelligence.
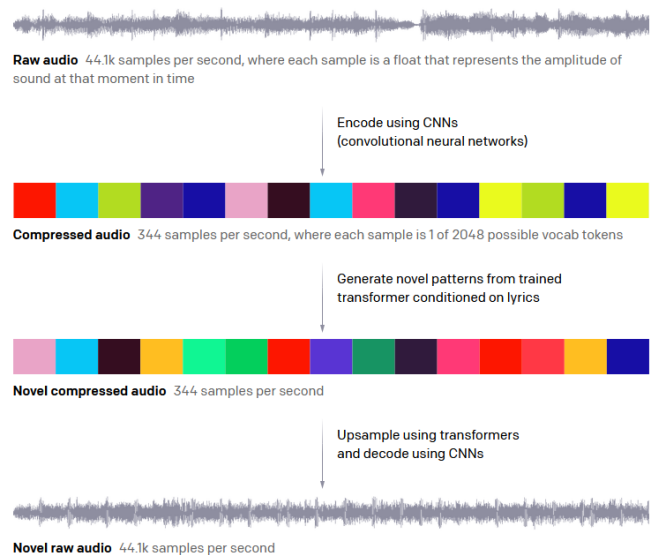
## 1. INTRODUCTION



Recent advancements in AI have made it possible to teach machines to recognize patterns in music and generate new compositions based on those patterns. By training an AI model on the Grateful Dead's music, we can create new pieces that capture the essence of the band's unique style. This project has the potential to push the boundaries of what AI can do in the realm of creative expression and could pave the way for new forms of music production. Additionally, it provides a unique listening experience for the Deadheads who just can't get enough!

## 2. BACKGROUND

Jukebox, MuseNet, Riffusion, and Musika are all AI music generation models that use deep neural networks. However, they use different types of neural network architectures and techniques to generate music, and have varying levels of usability.

## 2.1 Jukebox



**Raw audio** 44.1k samples per second, where each sample is a float that represents the amplitude of sound at that moment in time

Encode using CNNs (convolutional neural networks)

**Compressed audio** 344 samples per second, where each sample is 1 of 2048 possible vocab tokens

Generate novel patterns from trained transformer conditioned on lyrics

**Novel compressed audio** 344 samples per second

Upsample using transformers and decode using CNNs
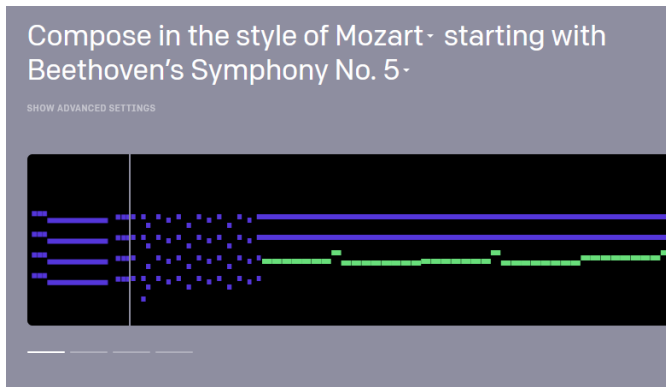
**Novel raw audio** 44.1k samples per second

OpenAI Jukebox[1] is an AI system developed by OpenAI that uses a combination of deep neural networks and natural language processing to generate music in a variety of styles and genres. The system is based on a large corpus of music that has been preprocessed and organized into a hierarchy of musical concepts, such as scales, chords, and melody lines. When a user inputs a musical prompt, such as a short melody or a genre, the system generates a sequence of musical events that follow the input and are consistent with the hierarchical structure of the music corpus.

OpenAI Jukebox is trained on a dataset of over a million songs, and its deep neural networks have billions of parameters. The system is capable of generating original compositions in a wide range of styles and genres, from classical and jazz to pop and hip-hop. The generated music can be customized by adjusting parameters such as the tempo, key, and instrumentation.

Although the results from Jukebox are stunning, the time required both to train a model and generate audio is excessive. It takes hundreds of hours to train a model, and generation requires a few hours for each second of audio, causing it to not be practical for this project.
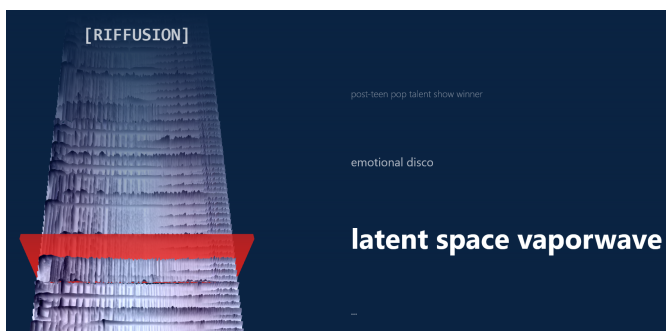
## 2.2 MuseNet



MuseNet[2] is an AI system developed by OpenAI that uses deep neural networks to generate music in a wide range of genres and styles. The system is based on a generative model architecture that uses a combination of hierarchical latent variables and autoregressive models to generate long sequences of musical events. MuseNet is trained on a massive dataset of MIDI files that cover a broad range of musical genres and styles, including classical, jazz, pop, and rock. The system can generate original compositions of varying lengths and complexity, and users can specify parameters such as the genre, tempo, and instrumentation to customize the output.

One of the key features of MuseNet is its ability to generate music that is both diverse and coherent. The system is able to learn and incorporate complex musical structures and patterns from the training data, enabling it to generate music that is musically consistent and coherent, while also being diverse and unpredictable.

Due to the fact that MuseNet is designed to process MIDI datasets, it would not be practical for this project. Although there is a limited amount of MIDI transpositions for Grateful Dead songs, the amount of raw audio available far outweighs this. However, a future project might involve training a MuseNet model on the guitar solos of Jerry Garcia (the lead guitarist of the Grateful Dead) by manually transcribing them to MIDI.

## 2.3 Riffusion



Riffusion[4] is an open-source music generation system developed by the Riffusion research group at the University of Montreal. The system uses a combination of deep neural networks and evolutionary algorithms to generate novel and creative musical sequences, such as melodies, chords, and rhythms.
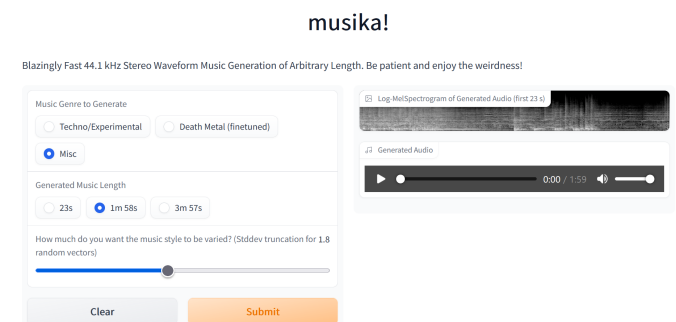
The system is based on a generative model architecture that uses deep recurrent neural networks (RNNs) to learn the statistical patterns and relationships in a large corpus of MIDI data. The system also employs an evolutionary algorithm to guide the generation process, allowing the model to explore and discover new and creative musical sequences that are not present in the training data.

One of the unique features of Riffusion is its ability to generate music that is both diverse and coherent. The system can generate complex and intricate musical structures and patterns, while also maintaining a consistent musical style and coherence. The system is also highly customizable, allowing users to specify various parameters such as the style, genre, and complexity of the music.

Since Riffusion performs audio generation quickly, it is a reasonable choice for this project. However, there are currently no published methods for reliably training a Riffusion model, which presents a roadblock. Additionally, the fidelity of the audio produced by Riffusion is very low, since it relies on converting compressed images of spectrograms to audio. Nevertheless, as neural networks designed for processing images improve, and the resolution of the images is increased, Riffusion could yield very promising results.

## 2.4 Musika



Musika [3] is a music generation system that can be trained on a single consumer GPU and allows for much faster than real-time generation of music of any length on a consumer CPU. The system uses a two-step process to generate music: first, it learns a compact representation of spectrogram magnitudes and phases with adversarial autoencoders, and then it trains a Generative Adversarial Network (GAN) on this representation for a particular music domain. A latent coordinate system enables generating arbitrarily long sequences of excerpts in parallel, while a global context vector allows the music to remain stylistically coherent through time. Musika is designed to be user-controllable, allowing for novel ways of composing or performing music. The system is evaluated quantitatively to assess the quality of the generated samples and is showcased for user control in piano and techno music generation.

Due to the incredibly fast speed of music generation, Musika was chosen as a foundation for DeepDeadNet. It also satisfies the requirement of being able to 'easily' train a neural network on a dataset of raw audio. Compared to the other techniques, Musika is the most suitable for the goal of this project. Although OpenAI's Jukebox yields much better results, it can take multiple hours to generate a single second of audio, while Musika is able to generate a few minutes of audio in a few seconds.
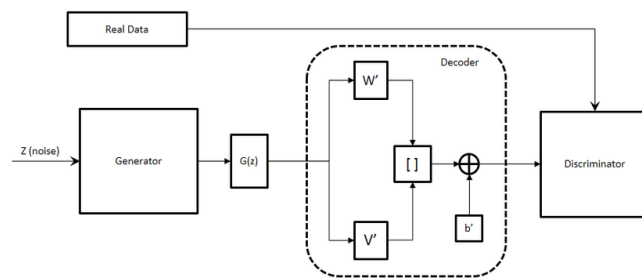
## 3. PROJECT DESCRIPTION

## 3.1 Data Consolidation

In order to train a Musika model, a huge quantity of input data is required. Thankfully, the live performances of the Grateful Dead have been archived extensively, and can be downloaded in bulk. However, the recordings are not all high-quality, and needed to be cleaned up. First, each individual audio file was normalized to a maximum peak of 0 decibels. This was done in order to prevent unnatural fluctuations in volume in the generated audio. Additionally, each track commonly contained audience noise at the beginning and end. To prevent this from bleeding into the models, each track was batch processed in order to remove the first and last ten seconds.
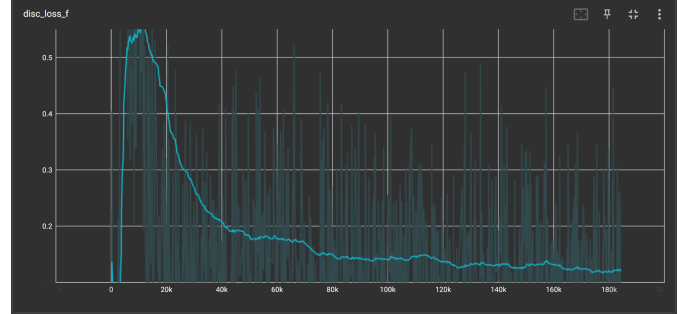


The consolidation process yielded a dataset containing 68 GB of raw audio, which includes 771 songs across 35 different performances. In total, this resulted in 99 hours of audio that the models could be trained on, which was sufficient. However, before the learning could begin, the raw audio needed to be encoded into arrays of latent vectors using the NumPy library. Essentially, this compresses the amplitude values of the waveforms into latent vectors that the neural network can understand. The figure above shows the raw view of one of these arrays.

## 3.2 Training the Musika Models



Musika utilizes the Generative Adversarial Network (GAN) architecture, which is a machine learning technique that is well suited to audio. GAN models are trained using a two-part process, where a generator network and a discriminator network are trained simultaneously in an adversarial manner. The generator takes as input a random vector and produces fake samples, while the discriminator takes as input both real and fake samples and tries to distinguish between them. During training, the discriminator tries to correctly identify the real samples from the fake ones, while the generator tries to produce samples that are indistinguishable from the real ones. This creates a feedback loop where the generator learns to produce more realistic samples as the discriminator becomes better at distinguishing them. The training process continues until the generator produces samples that are similar enough to the real samples, or until a certain number of iterations or loss threshold is reached. GAN training can be challenging, as it involves finding the right balance between the generator and discriminator networks to prevent the generator from producing easily detectable fake samples, while still allowing it to generate diverse and realistic samples.



Multiple models were trained using the Musika library, with varying datasets and training parameters. The main model was trained on the previously mentioned dataset of Grateful Dead live recordings, spanning from 1968 to 1990, with 99 hours of audio in total. This model was trained from scratch (not continuing from an existing checkpoint), for 500k iterations at a low learning rate of 0.0005. The results are consistent with the example models provided in the Musika library. Pictured above is a graph of the discriminator loss values during training, which represent how closely the generated audio matches the provided dataset, with lower being more closely.
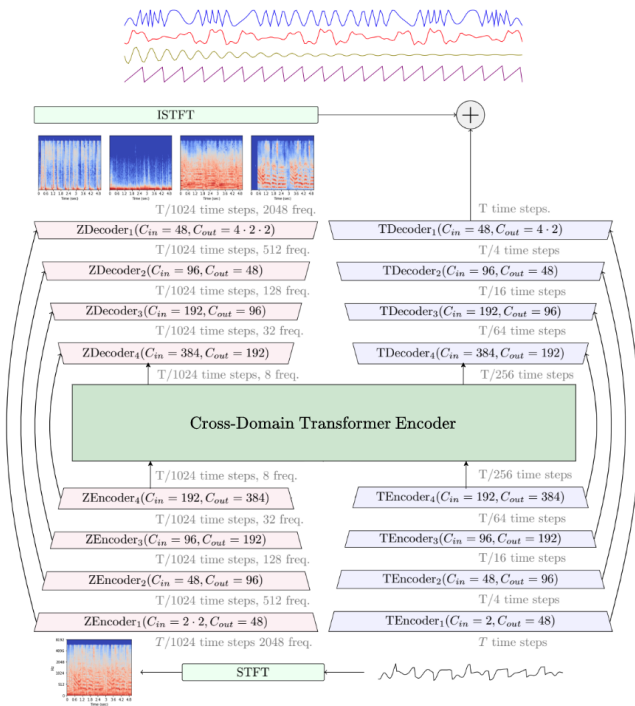


Additional models were finetuned from the Musika base model on specific albums, which contain about 2-3 hours of audio each. These albums include Cornell '77, Veneta '72, Winterland '73, and Nassau '90. Each of the finetuned
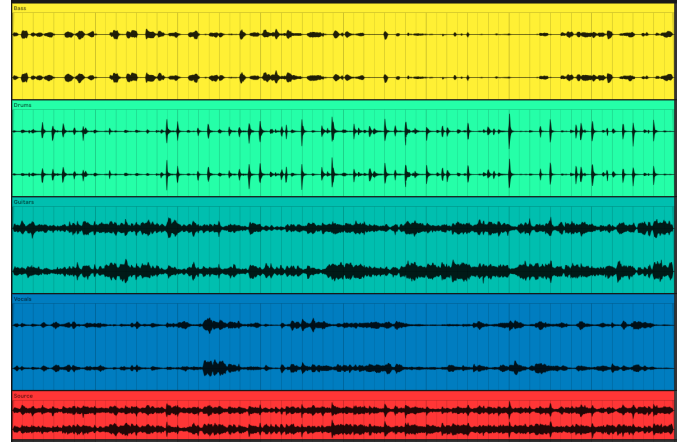
models was trained for only 100k iterations, at a variable learning rate that changed based on the loss values from the discriminator. These models also yielded results in line with the example models in the Musika library, and started producing listenable audio much faster than the main model.

## 3.3 Demucs Source Separation

Separating individual instruments is a surprisingly complex problem in data analysis. Once individual tracks have been combined into a single mix, it is almost impossible for a computer to perfectly separate them out again, because some data is lost during the mixing process. A common analogy for this process is if you were to make a smoothie, give it to someone else who did not watch you make it, and then ask them what fruits they think you used. This is simple for one or two ingredients, but more than that would be very difficult.



Demucs is a deep learning-based method for music source separation that can extract individual audio tracks from a mixed audio signal. It uses a neural network architecture based on a type of recurrent neural network called a Transformer, which has been shown to be effective for processing sequential data such as audio signals. The network is trained to predict the magnitude and phase of the individual sources from the mixed signal. Once the network is trained, it can be used to extract individual sources from new mixed signals in real-time. Demucs works by iteratively refining the predicted sources by passing them through the network multiple times. This allows the network to gradually separate the individual sources from the mixed signal. Demucs has been shown to outperform other state-of-the-art methods for music source separation, and can be used for a variety of applications, including music transcription, remixing, and analysis. Pictured above is a diagram of the Demucs architecture.



Originally, Spleeter was used in order to process the generated audio from the Musika models. However, the results contained a large quantity of artifacts. The Spleeter model struggled with differentiating between vocals and accompaniment, due to the less than ideal quality of the generated audio. Much better results were seen with Demucs, which is used in the final web application. Overall, Demucs is able to separate the bass, vocals, and percussion at much higher fidelity than Spleeter. Pictured above is the resultant stems from running Demucs on a track generated by the main Musika model.

## 3.4 DeepDeadNet



Putting all of these processes together results in DeepDeadNet, an application that can generate music in the style of the Grateful Dead in real-time. Users can select from the main model, or show-specific models, and generate audio of a specified length in that style. They can then adjust the volume of each instrument to their liking. Since the two models that serve as the basis for application (Musika and Demucs) have very low computational requirements, this application can run on a server with minimal costs. Most of the computational resources are allocated to the training process, which has already been completed, the results of which are stored in the models.

## 4. FUTURE WORK

DeepDeadNet was built with scalability in mind, and can be repurposed for any audio dataset. As music generation

and processing models improve over time, they can be seamlessly updated in application, and it will become more useable. Additonally, since the processing occurs in real-time, an endless Grateful Dead radio could be created that leverages DeepDeadNet, allowing users to enjoy new music in the style of the Grateful Dead forever.

It is also worth mentioning that the audio produced by the Musika models is not high-fidelity, and it is rare to find a result that is truly listenable. However, the results could be sampled by producers for use in a larger project, as the spectral quality of the generated audio is quite unique. Additionally, due to the nature of the encoders used for Musika, the audio must be compressed heavily, which results in inevitable low fidelity. However, a trainable encoder is currently being developed, which will help alleviate this issue.

## 5. CONCLUSIONS

This project showed me how difficult it can be to train a neural network. Although it's tempting to allow the training simply run its course, manual adjustments at critical moments are often necessary to prevent the network from overfitting or generally spiralling out of control. There is also a certain element of randomness to training a neural network, which requires immense patience to overcome. On the other hand, this project served as a catalyst for my exploration of various neural network architectures, including those unrelated to audio. More specifically, I've been exploring Stable Diffusion for artwork generation (the basis for Riffusion), as well as the application of ChatGPT's API in video games for realistic character interactions. Moving forward, I will continue to seek opportunities in the world of AI, and am very excited to see how it will grow.

Overall, this project aimed to explore the current potential of music generation and enhancement by artificial intelligence. DeepDeadNet proves that current machine learning models are able to distinguish between not only different styles of music, but also different instruments, all in real-time. The AI revolution has only just begun, yielding amazing results in art and text generation, and it seems that music is following suit. But at this stage, nothing beats the real thing.

## 6. ACKNOWLEDGMENTS

## 7. REFERENCES

[1] P. Dhariwal, C. Jun, C. Payne, J. W. Kim, A. Radford, and I. Sutskever. Jukebox: A generative model for music, 2020.

[2] C. Hawthorne, E. Elsen, J. Song, D. Roberts, I. Simon, C. Raffel, J. Engel, and S. Oore. Neural generation of multi-track music with hierarchical variational autoencoders. In *Proceedings of the 34th Conference on Neural Information Processing Systems (NeurIPS)*, 2019.

[3] M. Pasini and J. Schlüter. Musika! fast infinite waveform music generation, 2022.

[4] Z. Rasheed, S. Oore, J. Arsenault, A. Bergeron, M. McKinney, and G. Tzanetakis. Riffusion: Creative ai for music generation. In *Proceedings of the 21st International Society for Music Information Retrieval Conference (ISMIR)*, 2020.