



Carnegie Mellon University

Deep Music Features for Music Recommendation Systems Proposal

Nicholas Magal

Problem Statement



Recommender Systems

Recommender Systems

Content Based Filtering vs Collaborative Filtering

Content Based Filtering

Recommend items based on similar liked items

Advantages

- No data needed from other users
- Can recommend niche items that other users are not interested in

Disadvantages

- Need to hand craft features
- Limited Recommendations

Collaborative Filtering

Recommend items based on similarities between users and liked items

Advantages

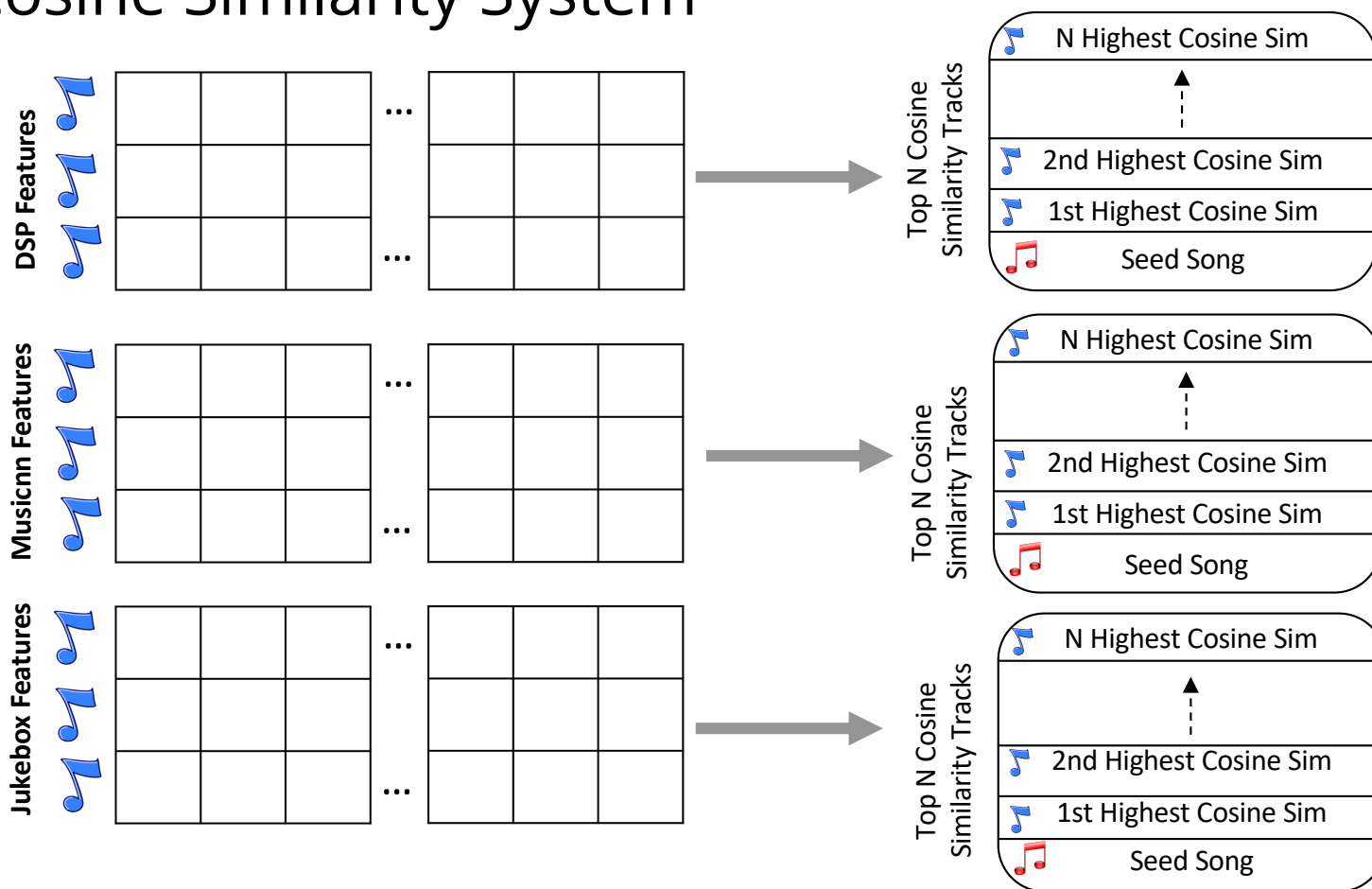
- Incorporates more information
- No domain knowledge needed

Disadvantages

- Cold start problem

Proposed Approach

Cosine Similarity System



Dataset

GTZAN

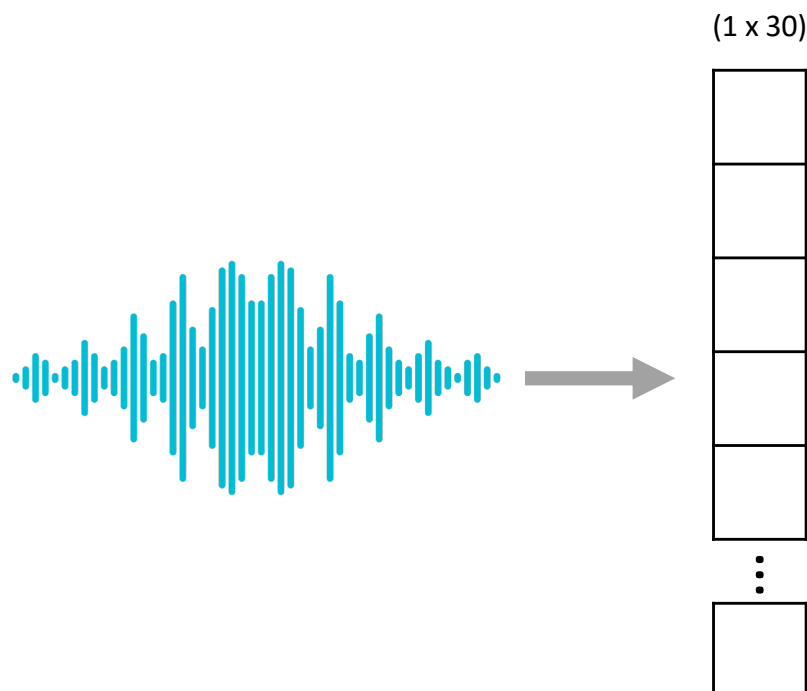
- 1,000 half minute music excerpts
- 10 different genre labels
 - Blues, classical, country, disco, hiphop, jazz, metal, pop, reggae, rock

Feature Representations

DSP Features

Timbral Texture Features

- Mean and variance of:
 - Spectral Centroid
 - Brightness
 - Spectral roll-off
 - Frequency concentration
 - Zero crossing rate
 - Noisiness
- 12 MFCCs



Deep Learning Features

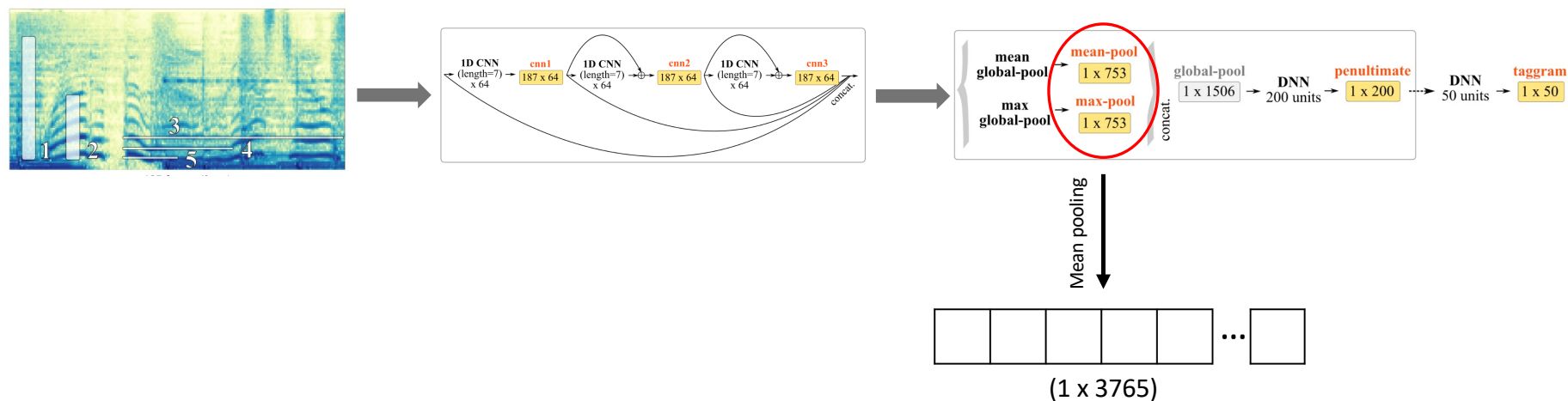
Musicnn Features

High Level Overview

- Musically motivated convolutional neural network
- Pretrained on Million Song Dataset for auto-tagging



Musicnn Features Architecture



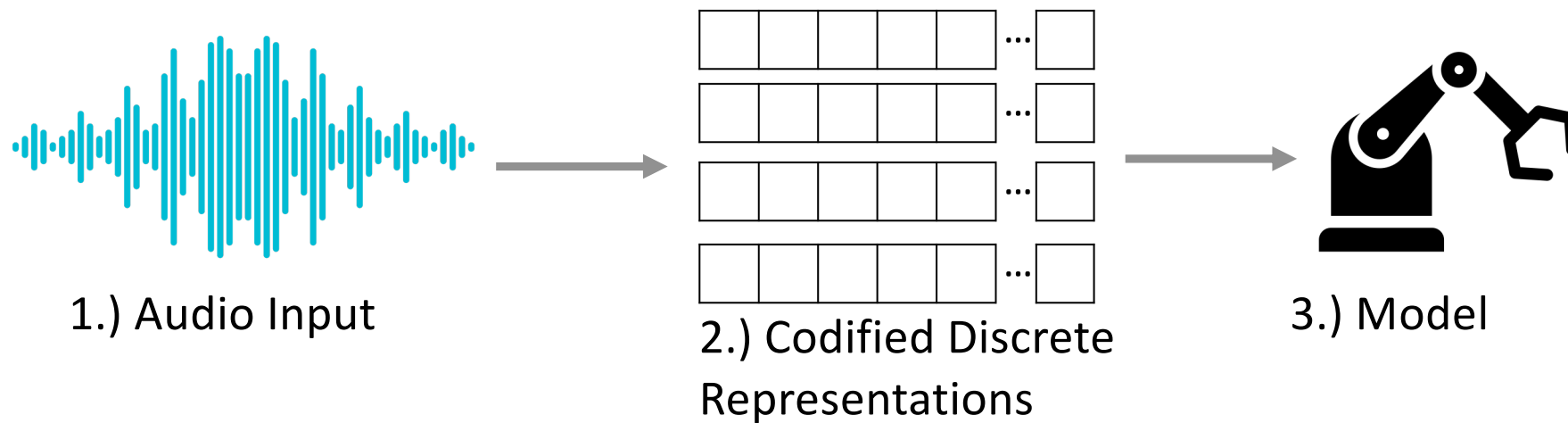
Jukebox Features

High Level Overview

- Generative Model based on Hierarchical VQ-VAE and Transformer architecture
- Pretrained on 1.2 million songs scraped by OpenAI
- Performs 30% on average better than other embeddings on genre detection, music emotion classification, auto-tagging, and key detection

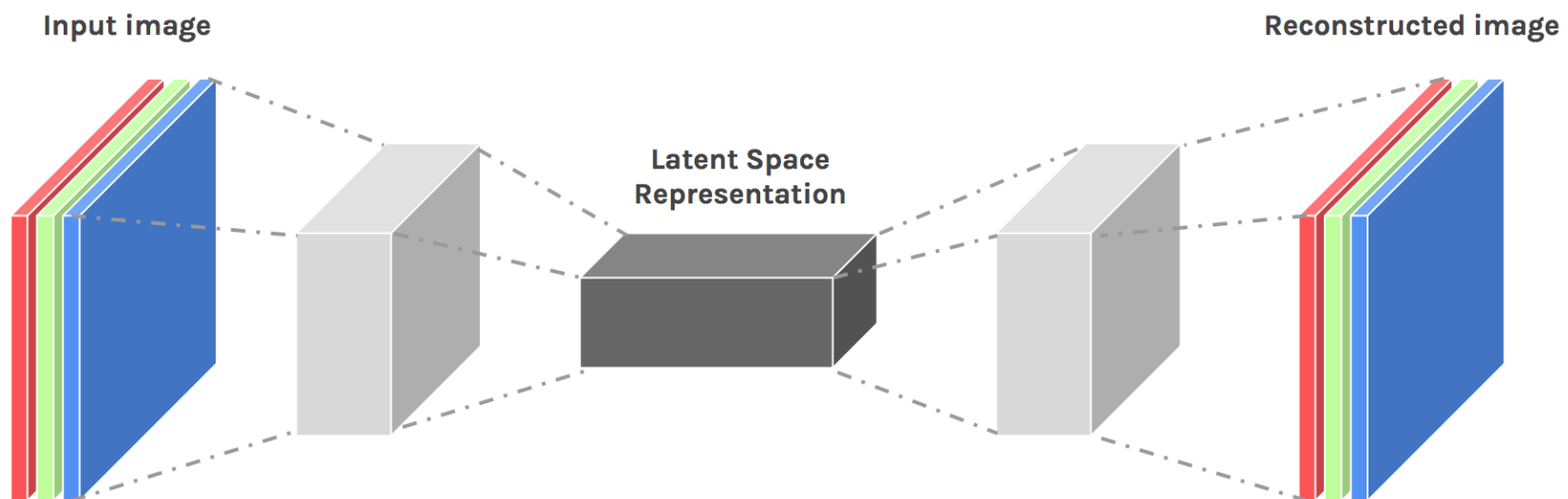
Jukebox Features

High Level Overview



Quick Detour

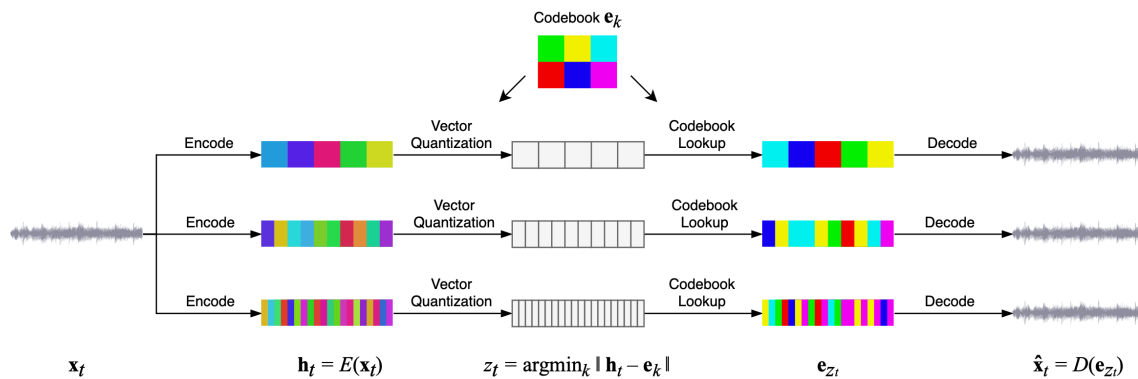
Autoencoders



Jukebox Features

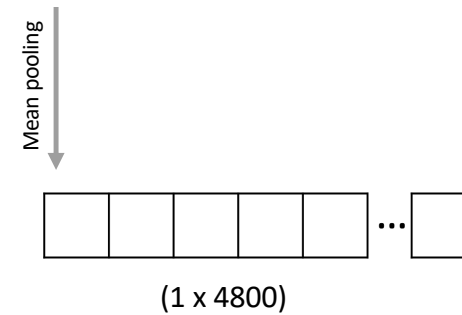
Model Details

1.) Generate Codified Discrete Representations



2.) Model Codified Discrete Representations

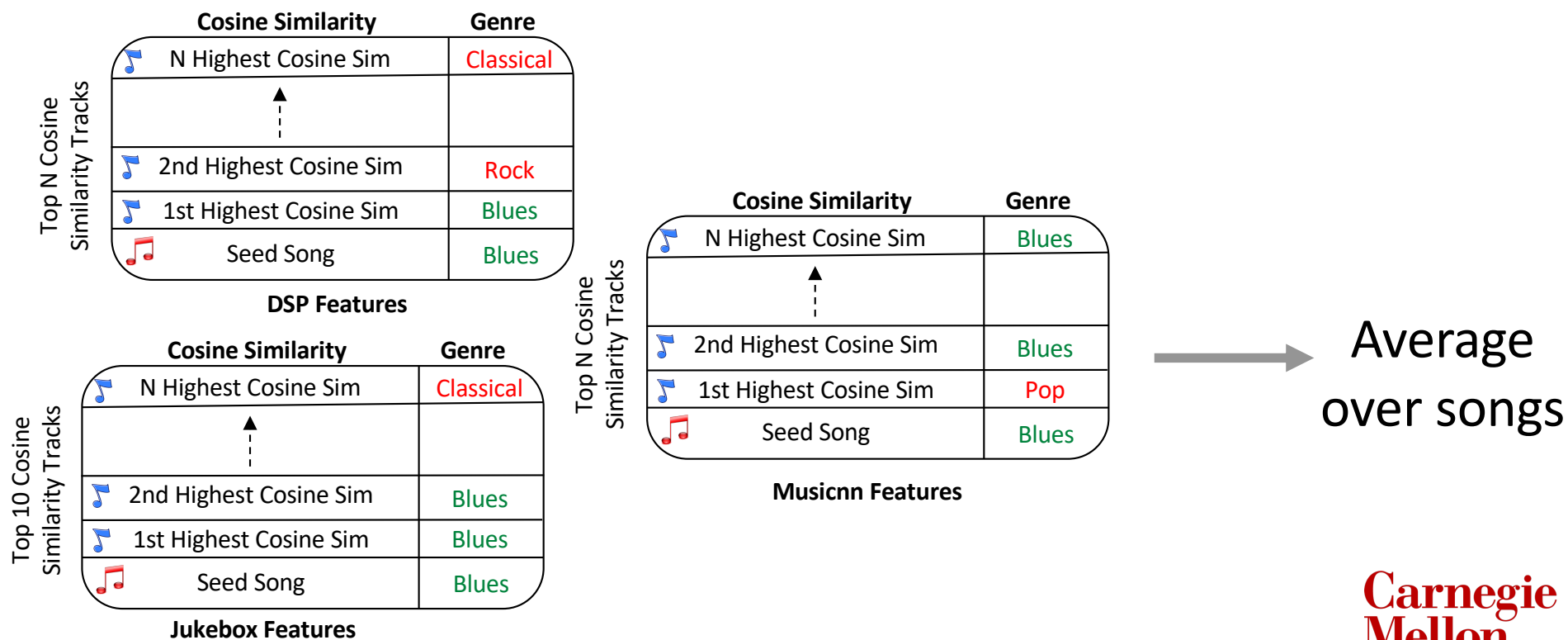
$$p(\mathbf{z}) = p(\mathbf{z}^{\text{top}}, \mathbf{z}^{\text{middle}}, \mathbf{z}^{\text{bottom}})$$



Results

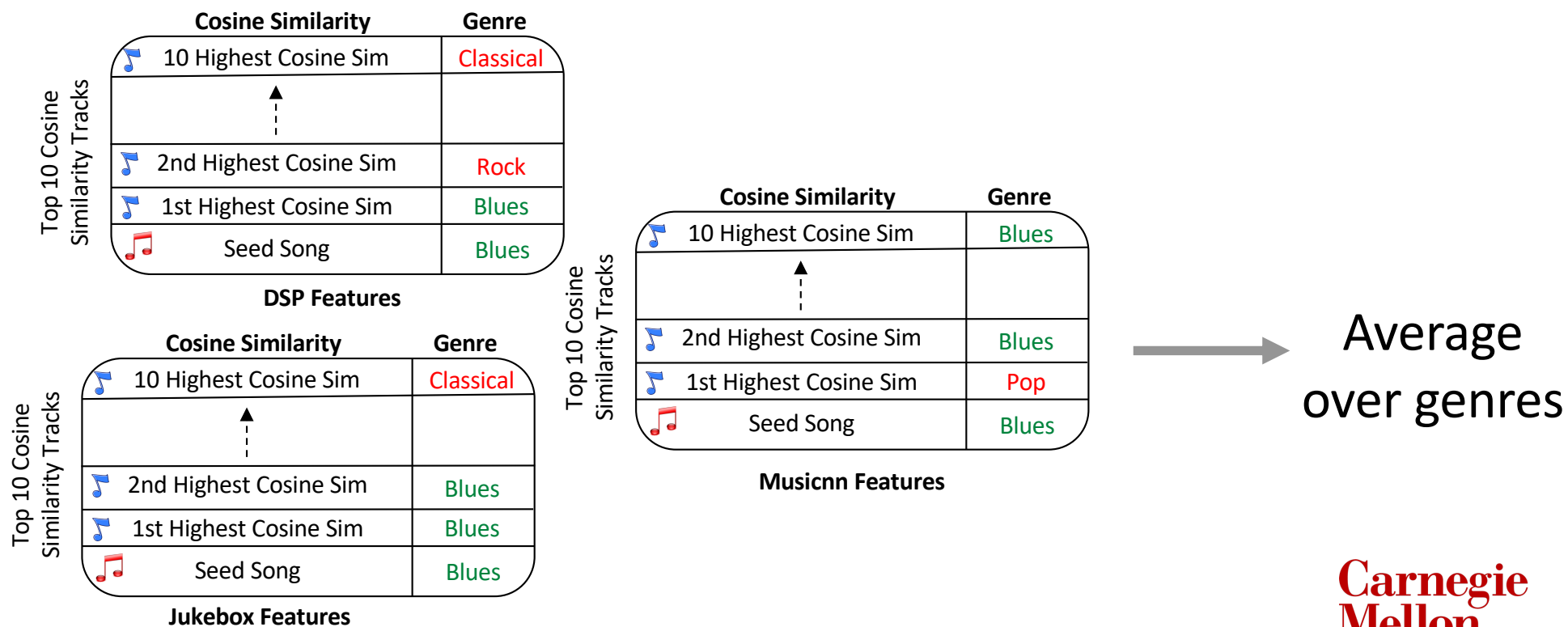
Quantitative Evaluation

Average Matching Genre of Top N Tracks

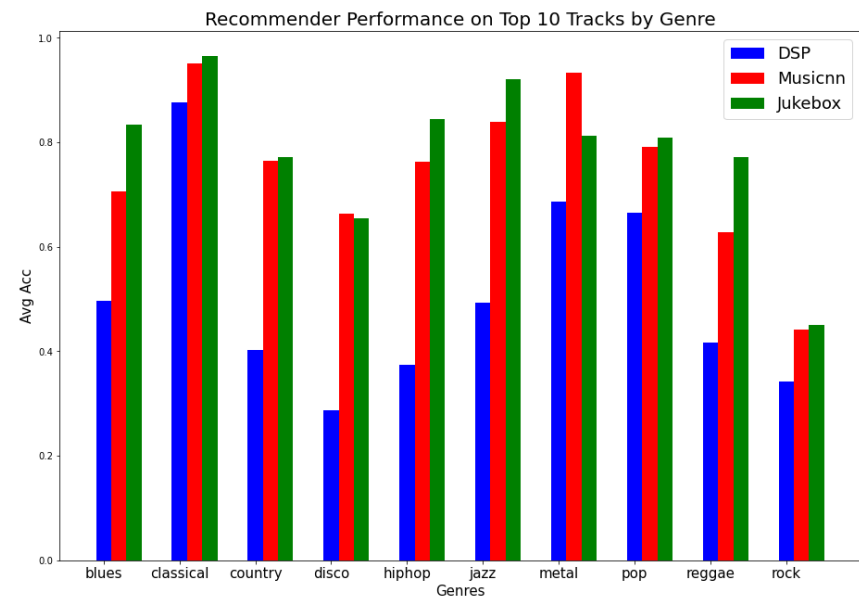
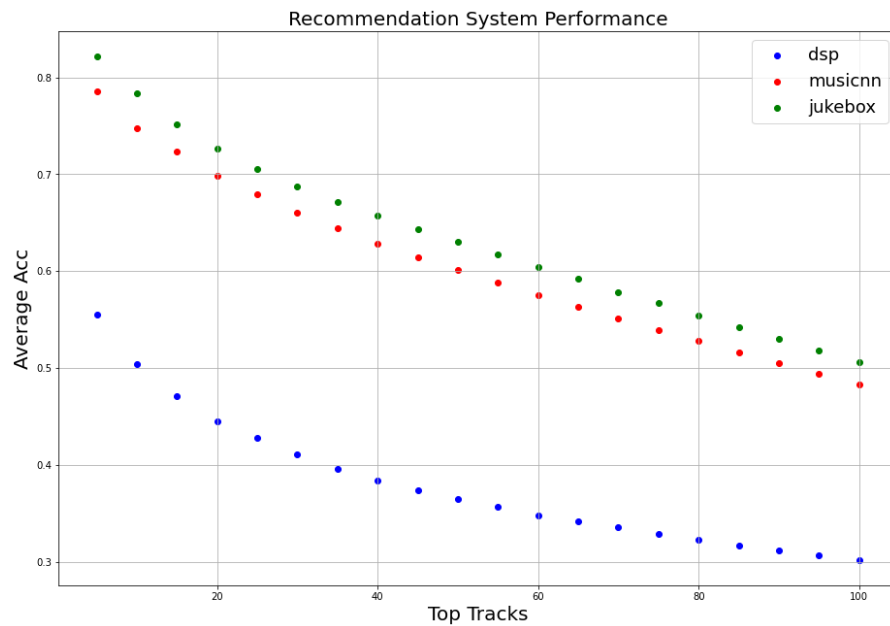


Quantitative Evaluation

Average Matching Genre of Top 10 Tracks

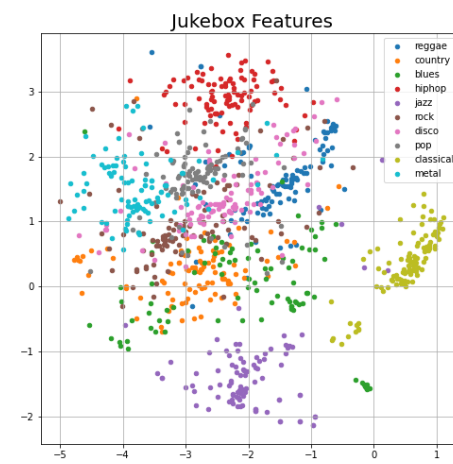
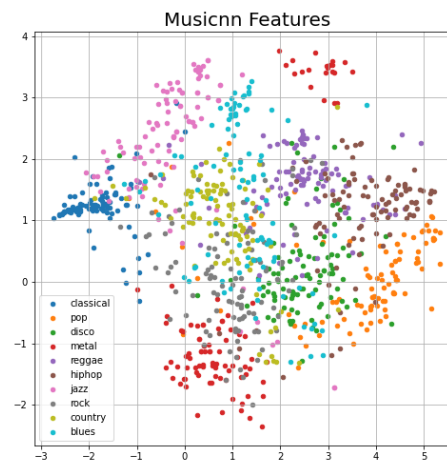
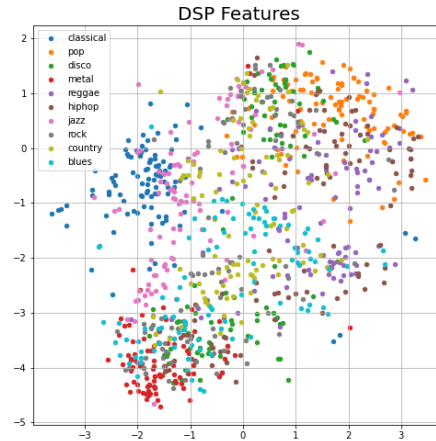


Results



Results

TSNE Plots



Proposed Auditory Evaluation



Online Survey:

1. 3 prechosen seed songs
2. From each feature set, return the top N most similar song
3. User rates the recommendation on a scale of 1-10 based on lyrics, rhythm, melody, and overall traits of each track.

Thank You!

Citations

[1]: Google Developers, “Content-based filtering.”

<https://developers.google.com/machinelearning/recommendation/content-based/basics>, Jul 2018

[2]: Google, “Collaborative filtering.”

<https://developers.google.com/machinelearning/recommendation/collaborative/basics>, Jul 2018.

[3]: Sturm, Bob L. "An analysis of the GTZAN music genre dataset." *Proceedings of the second international ACM workshop on Music information retrieval with user-centered and multimodal strategies*. 2012.

[4]: G. Tzanetakis and P. Cook, “Musical genre classification of audio signals,” IEEE Transactions on speech and audio processing, vol. 10, no. 5, pp. 293–302, 2002.

[5]: J. Pons and X. Serra, “Musicnn: pre-trained convolutional neural networks for music audio tagging,” 2019.

[6]: Castellon, Rodrigo, Chris Donahue, and Percy Liang. "Codified audio language modeling learns useful representations for music information retrieval." *arXiv preprint arXiv:2107.05677* (2021).