

Lyrics Matter: Exploiting the Power of Learnt Representations for Music Popularity Prediction

Anonymous ACL submission

Abstract

Accurately predicting the popularity of a music is a critical challenge in the music industry given the potential benefits to artists, producers and streaming platforms. Historically, research on music success was focused on factors such as audio features and extrinsic metadata (e.g., artist demographics, listener trends), or advancing prediction model architecture. This paper addresses the under-explored area of exploiting lyrical content to predict music popularity. We present a novel automated pipeline that uses LLMs to extract mathematical representations from lyrics, capturing their semantic and syntactic structure, while preserving sequential information. These features are then integrated into a novel multimodal architecture, HitMusicLyricNet, combining audio, lyrics, and social metadata for predicting popularity score. Our method outperforms the available baseline in end-to-end deep learning architecture for music popularity prediction on the SpotGenTrack (SPD) dataset. We achieve an overall 9% and 20% improvement in prediction model performance metrics MAE and MSE respectively. We confirm that the improvements result from the introduction of our lyrics feature engineering pipeline (LyricsAENet) in our model architecture, HitMusicLyricNet.

1 Introduction

In 2023, the global recorded music market generated \$28.6 billion¹ in revenues. With the advent of social media and streaming services, it has become increasingly difficult to define mathematical metrics for music success. Music popularity prediction could help the music industry and artists in maximizing future success of a newly released song.

Research in music popularity prediction has been driven by the advancements in machine learning with researchers applying classical ML approaches

to predict popularity using acoustic features, and further with the growth of social networks, information about music consumers' tastes capturing consumer response and their evolving music preferences. Advancements in deep learning further sharpen the prediction model capability of capturing and learning complex patterns of evolving music taste, and researchers have worked on incorporating multiple modalities such as audio, lyrics and social metadata to predict song success (Zangerle et al., 2019) (Martín-Gutiérrez et al., 2020). In all these works, the popularity score is typically defined as the time the song remains on the Billboard Top charts, and the evaluation metrics used include MAE, MSE, R^2 for regression, and accuracy, precision, recall, and F1 for classification. Recent developments in large language models have led to further research in music-related fields such as recommendation systems, sentiment/emotion analysis, data augmentation, understanding and composing song lyrics, using song lyrics text as the data source (Rossetto et al., 2023), (Sable et al., 2024), (Ma et al., 2024), (Ding et al., 2024). Music Popularity Prediction research has still not fully exploited the power of lyrics in the models, while recent research have shown lyrics contributing significantly to song popularity (Yu et al., 2023). Through our work, we address the gap in the existing literature with the following chief contributions:

- **LyricsAENet:** A novel automated lyric feature extraction pipeline that uses LLMs to encode music lyrics into rich, learned representations.
- **HitMusicLyricNet:** An end to end multimodal deep learning architecture which predicts the popularity score in range (1,100) and outperformed current baseline by 9% and 20% in MAE and MSE metrics respectively .

The next section reviews related work. This

¹IFPI Report'23

is followed by a discussion of our methods, the dataset and our experiments.

2 Related Work

Music Popularity Prediction. Studied as a classification or regression problem in a supervised learning fashion, where a model learns to predict either binary class labels (hit or no-hit) or generate a continuous popularity score. These predictions are based on a comprehensive analysis of a song’s features and associated social factors. Song popularity is primarily measured using charts like Billboard², and UK Singles Charts³, which rank songs based on sales, radio airplay, and streaming activity. Researchers determine success metrics based on these rankings. Other measures include economic metrics like merchandise sales, and engagement metrics such as interactions on social media and streaming services (Seufitelli et al., 2023).

Traditional research focused on using various machine learning techniques, including Logistic Regression, Decision Trees, Support Vector Machines (SVM), Bayesian Networks, Naive Bayes, Random Forest Ensemble, XGBoost, and K-Nearest Neighbors (KNN). These approaches advanced further to neural networks and deep learning techniques, building much stronger predictive models. A significant number of studies (Bischoff et al., 2009), (Dorien Herremans and Sørensen, 2014), (Zangerle et al., 2019), (Silva et al., 2022) focused on using acoustic characteristics of songs along with metadata that includes factors such as social influences. Other works such as (Dhanaraj and Logan, 2005), (Singhi and Brown, 2015b), (Martín-Gutiérrez et al., 2020) also emphasized the importance of song lyrics in determining song success using handcrafted text-based features that captured sentiment, emotions, and the syntactic structure of lyrics. These studies were often limited by their capabilities to capture central expressions of the song’s lyrics.

Multiple datasets have been released to drive research further and quench the thirst of data-heavy deep learning models. This includes Million Song Dataset⁴, SpotGenTrack⁵, and AcousticBrainz⁶ sourced from different platforms like Spotify, Bill-

board, Genius⁷, Youtube, and others. These datasets comprise a wide range of features, from low-level features like Mel-Frequency Cepstral Coefficients (MFCCs), lyrics text, and temporal features to high-level audio features such as danceability and loudness. Additionally, they include metadata on artists, albums, genres, demographics, and other relevant information.

Learned Representations of Lyrics. Lyrics form an integral part of music and carry a deep emotional meaning, which can strongly influence how listeners feel—sometimes even more than the song’s acoustic features alone (Singhi and Brown, 2015a). Yet, lyrics have often been overlooked as compared to acoustic attributes and social metrics of songs (Seufitelli et al., 2023). Earlier studies used methods like Probabilistic Latent Semantic Analysis (PLSA) (Hofmann, 1999) to capture the semantic content of lyrics, which helped researchers understand their role in defining a “hit” song (Dhanaraj and Logan, 2005). Later work moved beyond basic semantic analysis, focusing on more detailed features. For instance, (Hirjee and Brown, 2010) and (Singhi and Brown, 2014) relied on various rhyme and syllable characteristics to predict hit songs using only their lyrics, while other researchers applied Latent Dirichlet Allocation (LDA) (Blei et al., 2003) to discover thematic topics within lyrics (Ren et al., 2016).

Progress of deep learning techniques advanced the use of multimodal approaches that combine lyrics with audio and metadata, using stylometric analysis to extract lyric text features (Martín-Gutiérrez et al., 2020). Sentiment analysis also emerged as a way to glean emotional insights from lyrics when predicting popularity (Raza and Nanath, 2020). More recent research has turned to learned lyric representations, such as embeddings (Kamal et al., 2021) (McVicar et al., 2022), which offer a more robust way to capture lyrical meaning. (Barman et al., 2019) demonstrated that these distributed representations can effectively predict both genre and popularity, reducing the need for handcrafted features. Datasets such as Music4All-Onion (Moscati et al., 2022) provide lyric embeddings that make it easier to study how lyrical content relates to a song’s success. Finally, a recent study found that a song’s lyrical uniqueness has a significant contribution towards its popularity (Yu et al., 2023), using TF-IDF for lyric vector

²Billboard Hot 100

³Official Singles Chart Top 100

⁴Million Song Dataset

⁵SpotGenTrack

⁶AcousticBrainz

⁷Genius.com

representation; however, this approach inherently lacks the capacity to capture deeper sequential and contextual nuances, emphasizing the growing importance of learning robust, richer representations of lyrics to better understand what makes certain songs resonate with audiences.

To the best of our knowledge, there are limitations in existing literature for efficient automated lyrics feature extraction that are expressive and capture the underlying complexity of song lyrics. Thus, we have built a novel pipeline to exploit the power of Large Language Models. It has the potential to provide rich lyric representations that encapsulate both semantic and syntactic understanding, while preserving the sequential structure of the lyrics.

3 Methodology

In this section, we provide the theoretical foundation of our approach. We begin by defining the problem of music popularity prediction in mathematical equations. This is followed by explaining the baseline approach and its implementation, including details of the dataset. Finally, we present a formal description of our proposed architecture.

3.1 Problem Formulation

Given a song S , its features are represented in a multi-dimensional space $X \in \mathbb{R}^d$, which comprises three key modalities: audio waveform $w \in \mathbb{R}^k$, lyrical text $l \in \mathbb{R}^m$, and metadata attributes $m \in \mathbb{R}^p$, where $d = k + m + p$ represents the total dimensionality of our feature space. Our primary objective is to extract meaningful features from the song lyrics to effectively encode each song into a unique vector representation. Next, the prediction task is formulated as learning a mapping function $f : X \rightarrow Y$, where we minimize the expected prediction error: $\mathbb{E}[(f(X) - Y)^2]$ across the training distribution. Here, $Y \in \mathbb{R}$ represents the continuous popularity score.

3.2 Baseline Methodology

We trained *HitMusicNet*, a multimodal end-to-end Deep Learning architecture as proposed by (Martín-Gutiérrez et al., 2020) and validated the results using the SpotGenTrack Popularity Dataset (SPD). The model outputs a popularity score between 1 and 100, using audio features, text features, and metadata containing artist and demographic information as inputs. A complete description of the feature set used is provided in Table 1.

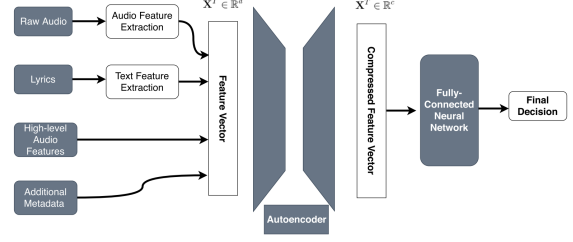


Figure 1: Diagram of the HitMusicNet pipeline outlining the principal functionalities and data components. Image src (Martín-Gutiérrez et al., 2020).

Feature Type	Features
Text Features	Sentence count, Avg words, Word count, Avg syllables/word, Sentence similarity, Vocabulary wealth
High-Level Audio	Danceability, Energy, Key, Loudness, Mode, Speechiness, Acousticness, Instrumentalness, Liveness, Valence, Tempo, Duration, Time Signature
Low-Level Audio	Mel-spectrogram, MFCCs, Tonnetz, Chromagram, Spectral Contrast, Centroid, Bandwidth, Zero-Crossing Rate
Meta-Data Features	Artist followers, Artist popularity, Available markets

Table 1: Summary of features used in the HitMusicNet architecture (Martín-Gutiérrez et al., 2020).

HitMusicNet architecture as shown in Fig 1, employs an autoencoder for feature compression through two encoder layers with dimensions $d/2$ and $d/3$, followed by a bottleneck layer of $d/5$. Each layer uses ReLU activation, and the output layer employs a sigmoid activation for reconstruction. The autoencoder was trained using the Adam optimizer and an MSE loss function. The compressed features are then passed through a fully connected neural network with four layers, where the number of neurons in each layer is scaled by factors $\alpha = 1$, $\beta = 1/2$, and $\gamma = 1/4$. The model is trained using an 80%-20% train-test split with stratified cross-validation (SCV) using $k = 5$. These settings helped us in effectively replicating the baseline results on the SPD dataset.

3.3 Dataset

The SpotGenTrack Popularity Dataset (SPD) proposed by (Martín-Gutiérrez et al., 2020) and used in this research contains 101,939 tracks, 56,129 artists, and 75,511 albums sourced using Spotify and Genius APIs. The data was gathered from 26 countries where Spotify is available, including the top 50 playlists per category for each country.

Popularity scores for tracks range between 1 and 100 and are provided by Spotify based on internal metrics. The scores follow a Gaussian distribution with $\mu = 40.02$ and a standard deviation of $\sigma = 16.79$. The dataset contains low-level audio features extracted using audio waveform, text features extracted using stylometric analysis of lyrics. High-level audio features and metadata are sourced from Spotify. The lyrics in the SPD dataset had to be cleaned and pre-processed to align with the objectives of this research. We inspected long tails of lyrics length distribution and observed that extremely short or long entries typically contained irrelevant content such as random numbers, out-of-context text, or placeholder text. Based on this analysis, we retained songs with lyrics lengths between 100 and 7000 characters. Furthermore, we filtered the dataset to include only English lyrics, which comprised approximately 60% of the total data. These steps resulted in a clean dataset comprising 51,319 tracks, 30,024 unique artists, and 39,371 unique albums. The resulting popularity distribution, as shown in Fig 2, had $\mu = 41.11$ and a standard deviation of $\sigma = 17.51$, retaining original data characteristics.

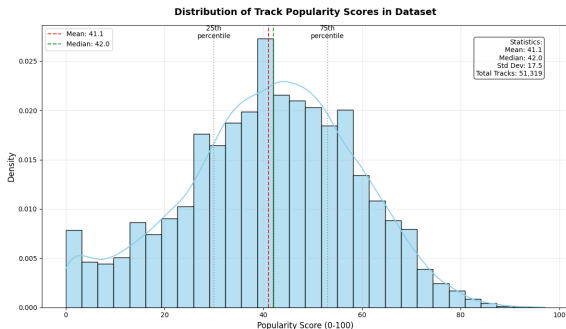


Figure 2: Popularity Distribution in cleaned SpotGen-Track (SPD) with $\mu = 41.11$ and a standard deviation of $\sigma = 17.51$.

3.4 HitMusicLyricNet

This section details our proposed HitMusicLyricNet, an end-to-end multimodal deep learning architecture built upon the foundation of HitMusicNet. HitMusicLyricNet comprises of three key components: AudioAENet, LyricsAENet, and MusicFuseNet. AudioAENet compresses the low-level audio features. LyricsAENet compresses the lyric embeddings into a fixed-size representation using an Autoencoder, thereby encoding information while reducing noise. MusicFuseNet then combines these compressed audio and lyric representa-

tions with metadata and high-level audio features as described in Table 1.

In the HitMusicNet architecture, a single auto-encoder compressed the combined feature vector of audio, lyrics, and metadata. We hypothesize that this can lead to information loss, particularly for the less abundant lyrics and metadata features. We believe that lyrics and metadata features should be fed directly into the popularity prediction network to retain their predictive power for song popularity. Furthermore, our reasoning behind the new approach of introducing distinct Autoencoders for audio and lyrics is based on the bipolar and directional nature of lyrics embeddings, requiring a different architecture for compression (Balazy et al., 2021).

3.4.1 AudioAENet

The Autoencoder used for compression has a similar architecture to that of MusicAENet, but takes in only low-level audio features as described in Table 1 for compression. For input dimension $d = 209$, it gradually compresses the data to dimension $d/2$, $d/3$, and $d/5$. The output layer employs a sigmoid activation for reconstruction, whereas all remaining layers use ReLU activation functions. The model is trained using the Adam optimizer with a MSE loss function, achieving a loss value in the range of $1e-5$, indicating negligible loss in compression.

3.4.2 LyricsAENet

LyricsAENet implements a tied-weights autoencoder architecture (Li and Nguyen, 2019) designed to reduce parameter size and risk of overfitting. Compressing lyric embeddings is susceptible to overfitting due to high dimensionality. The encoder follows a progressive compression with the following dimensions ($d/2$, $d/4$, $d/8$), followed by bottleneck layers ($d/12$ or $d/16$). The decoder mirrors the structure in reverse order, utilizing the transpose of the encoder weight. The progressive dimensional reduction is designed to minimize reconstruction losses in compressed embeddings extracted out of language models and LLMs such as BERT (Devlin et al., 2019), LLaMA 3 Herd (Grattafiori et al., 2024), and OpenAI’s embedding models⁸.

We use Scaled Exponential Linear Unit (SELU) (Klambauer et al., 2017) as the activation function for its self-normalizing characteristics and the ability to handle the bipolar nature of embeddings. Comparative analyses include alternate activation

⁸Open AI text Embedding Model

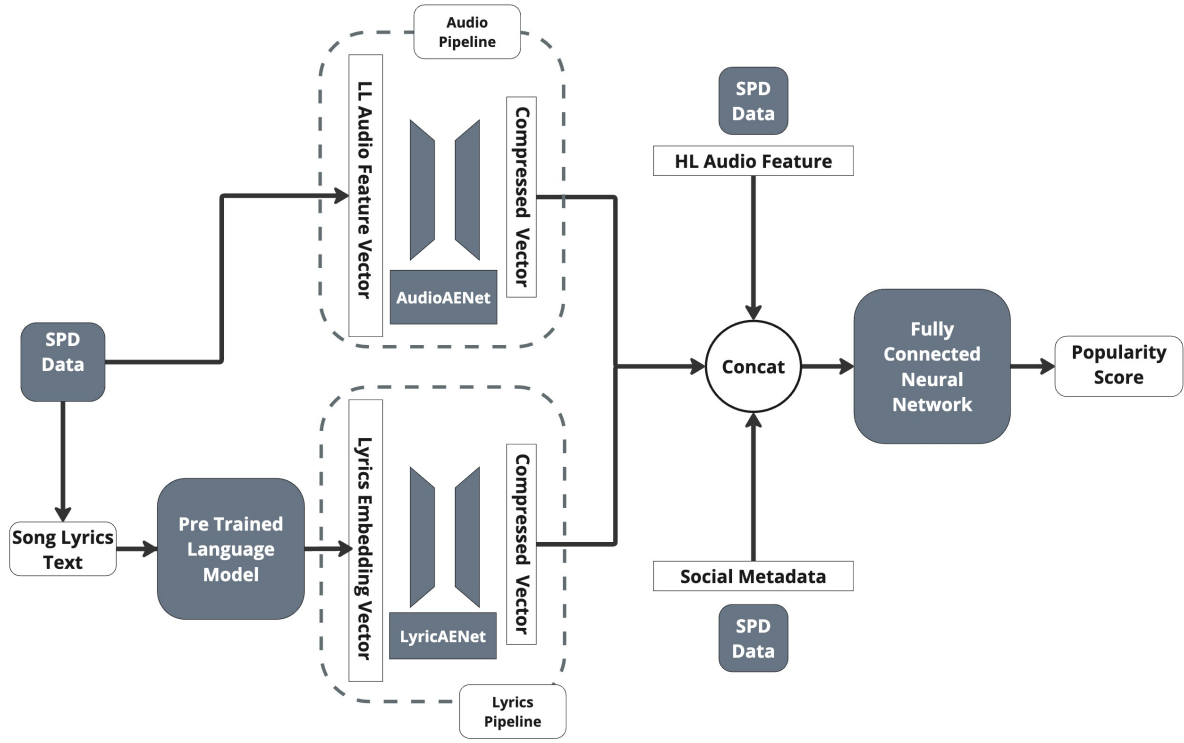


Figure 3: Block-Schema of the *HitMusicLyricNet* architecture comprising of two Autoencoders and a Fully Connected NN predicting popularity score. 'HL' stands for high-level and 'LL' stands for low-level.

functions such as the Sigmoid Linear Unit (SiLU) (Elfwing et al., 2018) and the Gaussian Error Linear Unit (GELU) (Hendrycks and Gimpel, 2016). LyricsAENet was trained using the Adam optimizer with a MSE loss function, achieving loss values of approximately $1e-5$. To further refine the training process, we incorporated a directional loss function inspired by (Bałazy et al., 2021) to preserve the directional characteristics of the embeddings during compression. This combined loss function is defined as:

$$L(Y, \bar{Y}) = \alpha_1 \cdot \text{MSE}(Y, \bar{Y}) + \alpha_2 \cdot \text{CD}(Y, \bar{Y}), \quad (1)$$

where $\text{MSE}(Y, \bar{Y})$ represents the Mean Squared Error. $\text{CD}(Y, \bar{Y})$ denotes the Cosine Distance, which captures the angular similarity between the vectors Y and \bar{Y} . The constants α_1 and α_2 control the relative importance of the two loss terms.

3.4.3 MusicFuseNet

We employ a similar architecture configuration as MusicPopNet by (Martín-Gutiérrez et al., 2020) for our MusicFuseNet. It uses a concatenation of compressed audio feature vectors from AudioAENet, compressed lyrics embeddings vectors from LyricAENet, high-level audio features and metadata as mentioned in Table 1. The output of this neural

net is a popularity score in the range $[0, 1]$. The architecture consists of a fully connected network with scaling parameters of $(1, 1/2, 1/3)$ and ReLU activation functions, followed by a Sigmoid activation in the final layer, as empirically validated by (Martín-Gutiérrez et al., 2020). To train the model, we used the Adam optimizer with an MSE loss function and applied dropout regularization to mitigate overfitting.

4 Experiments and Results

Using the **Code**⁹ to implement HitMusicNet and selecting the configuration details described in Section 3.2, we trained HitMusicNet on the SPD dataset with an 80-20 split. To replicate the results obtained by (Martín-Gutiérrez et al., 2020), we employed Stratified Cross-Validation (SCV) with $k=5$ folds and used MAE and MSE as performance metrics. As Table 4 shows, we achieved similar results on all performance metrics, validating our training and testing strategy. Further, removing the lyrics text features proposed by (Martín-Gutiérrez et al., 2020) did not degrade the metrics, so we dropped those features for further experiments.

To train HitMusicLyricNet, we extracted lyric

⁹Github: HitMusicNet

LyricsAENet Config	MAE (Train)	MAE (Val)	MAE (Test)
<i>SELU, MSE</i>	0.0769	0.0746	0.0775
<i>SILU, MSE</i>	0.0736	0.0731	0.0790
<i>GELU, MSE</i>	0.0740	0.0731	0.0792
<i>SELU, Dir.</i>	0.0741	0.0740	0.0799

Table 2: Results of training and testing HitMusicLyricNet on cleaned SPD data with various LyricAENet configurations (activation function, loss function), using BERT Large embeddings throughout. ‘Dir’ indicates directional loss 1.

Embeddings Model	MAE (Train)	MAE (Val)	MAE (Test)
<i>BERT large</i>	0.0793	0.0784	0.0786
<i>Llama 3.1 8B</i>	0.0774	0.0759	0.0795
<i>Llama 3.2 1B</i>	0.0775	0.0754	0.0800
<i>Llama 3.2 3B</i>	0.0781	0.0766	0.0798
<i>OpenAI Small</i>	0.0746	0.0738	0.0788
<i>OpenAI Large</i>	0.0761	0.0743	0.0772

Table 3: Results of training and testing HitMusicLyricNet on cleaned SPD data with different lyric embeddings sent to LyricAENet (Selu activation, MSE loss).

Model Config	Dataset Config	MSE (Train)	MSE (Val)	MAE (Train)	MAE (Val)	MAE (Test)
<i>HitMusicNet</i>	SPD	0.0116	0.0115	0.0836	0.0851	0.0862
<i>HitMusicNet w/o lyrics</i>	SPD	0.0114	0.0116	0.0843	0.0859	0.0870
<i>HitMusicLyricNet</i>	*SPD	0.0095	0.0091	0.0761	0.0743	0.0772
<i>HitMusicLyricNet w/o lyrics</i>	*SPD	0.0109	0.0113	0.0818	0.0841	0.0852

Table 4: Performance comparisons with the baseline (HitMusicNet) on SPD and SPD* data respectively, where SPD* denotes cleaned SPD data. Here, we report the best results from Table 3.

embeddings from language models. For open-source models (BERT, Llama), we downloaded vanilla weights from Hugging Face¹⁰ and loaded its vanilla configuration. We used Nvidia A100 GPU for compute requirements. After tokenizing lyrics, we forward-passed them through each model, extracted the last hidden-layer states, and applied max/mean pooling to obtain fixed-size vectors for our Autoencoder. Specifically for BERT, we considered mean pooling and concat (max + CLS token). To get embeddings from OpenAI text models, we used the API endpoint, costing \$3 for the small model and \$6 for the large. We then studied LLM model architecture and its training corpus effects on music popularity prediction with BERT, BERT Large, Llama 3.1 8B, Llama 3.2 1B, Llama 3.2 3B, and OpenAI text embeddings (small and large).

After extracting these embeddings, we examined different activation layers (Selu, Silu, Gelu) for embedding compression using LyricsAENet and introduced a directional loss function with $\alpha_1 = 0.5$ and $\alpha_2 = \frac{0.5}{5}$ as suggested by (Bałazy et al., 2021), alongside our standard MSE loss for LyricsAENet, to see their impact on the HitMusicLyricNet performance metric MAE. As reported in Table 2, using SELU with the MSE loss function in LyricsAENet yielded the least MAE error while training HitMusicLyricNet on popularity prediction. Directional loss produced comparable metrics but not enough

improvement to be included further. Other activation functions performed closely, but for simplicity and observing 1–2% randomness error, we proceeded with SELU and MSE.

Next, we compressed embeddings for different LLM models. While we experimented with two variants of BERT (small and large) and considered mean embeddings and concat (max + CLS token) embeddings, here we only report results for BERT large with mean embeddings, as it yielded the best results as seen in Table 3. All the Llama variants had very close performance metrics, whereas the OpenAI large text embedding model surpassed all of them. We attribute these small differences ($\sim 2\%$ variation) in HitMusicLyricNet’s performance to variations in each model’s training data and architecture, since none was specifically trained for our downstream task, leading to large differences in rich embedding representation.

Hence with HitMusicLyricNet used the OpenAI large text embeddings and the SELU activation with MSE loss function in lyricsAENet. Overall, we achieved close to a 9% improvement compared to the SOTA architecture, despite training on 40% less data. Dropping the lyrics feature pipeline and retraining and testing HitMusicLyricNet led to performance metrics comparable to that of HitMusicNet, validating the effectiveness of our proposed lyric feature extraction pipeline using LLMs and the overall advancements in music popularity prediction pipeline.

¹⁰Hugging Face

5 Conclusion and Future Work

The work presented in this paper showcases the power of leveraging lyrics to predict the popularity of a song, with the help of LLMs with capabilities of capturing the deeper meaning of sentences using embeddings. We believe that advancements in music-aware language models will lead to more explainable and expressive lyric features based on domain-specific knowledge. This research presented a novel architecture, HitMusicLyricNet, along with an ablation study. HitMusicLyricNet beats the SOTA by 9% by incorporating lyric embeddings and improving upon the SOTA architecture. With advancements in compression techniques and multimodal learning architecture, we believe accuracy and commercial use can be improved. Furthermore, with advancements in audio representation learning using neural audio codecs, richer music audio representations can be scoped into the study. Current research aggregates features over an entire song. However, contemporary phenomena of virality suggest that local features within different musical segments need to be studied deeply and cannot be ignored given the micro-content consumption driven by platforms like Instagram and SnapChat.

6 Limitation

Our research results are potentially limited to the music genres represented in our dataset and may not generalize across genres, demographics, and cultural contexts. Some limitations arise as a result of the choice of dataset used in our study, SpotGenTrack. The findings are highly dependent on the quality and size of the SpotGenTrack dataset. The dataset has been cleaned to filter out lyrics that are not in the English language. Though this reduced the size of the raw dataset by 40%, it limits the model's ability to be generalized across different languages and associated cultural contexts. The use of LLMs such as BERT and Llama 3 in our model will lead to a lack of domain-specific context, as horizontal LLMs are not typically trained or fine-tuned on music-focused data. While adequate measures have been made to address the risk of overfitting, the risk cannot be completely eliminated due to the high dimensionality of the data. The lyric embedding vectors are flowing downstream and are used to predict the popularity of a song. Finally, since we are assessing the quality of lyric embeddings using the performance metrics

of downstream tasks (music popularity prediction), this requires a further examination to evaluate the intrinsic qualities of lyric embeddings vector in capturing rich representation. We are limited by the explainability of our lyrics feature vector which can be further improved by using fine-tuned LLMs for music lyric analysis and explainable feature extraction.

References

- Klaudia Bałazy, Mohammadreza Banaei, Rémi Lebret, Jacek Tabor, and Karl Aberer. 2021. [Direction is what you need: Improving word embedding compression in large language models](#). In *Proceedings of the 6th Workshop on Representation Learning for NLP (RepL4NLP-2021)*, pages 322–330, Online. Association for Computational Linguistics.
- Manash Pratim Barman, Kavish Dahekar, Abhinav Anshuman, and Amit Awekar. 2019. [It's only words and words are all i have](#). In *Advances in Information Retrieval: 41st European Conference on IR Research, ECIR 2019, Cologne, Germany, April 14–18, 2019, Proceedings, Part II*, page 30–36, Berlin, Heidelberg. Springer-Verlag.
- Kerstin Bischoff, Claudiu S. Firan, Mihai Georgescu, Wolfgang Nejdl, and Raluca Paiu. 2009. Social knowledge-driven music hit prediction. In *Advanced Data Mining and Applications*, pages 43–54, Berlin, Heidelberg. Springer Berlin Heidelberg.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. [Latent dirichlet allocation](#). *Journal of Machine Learning Research*, 3:993–1022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ruth Dhanaraj and Beth Logan. 2005. [Automatic prediction of hit songs](#). In *International Society for Music Information Retrieval Conference*.
- Shuangrui Ding, Zihan Liu, Xiaoyi Dong, Pan Zhang, Rui Qian, Conghui He, Dahua Lin, and Jiaqi Wang. 2024. Songcomposer: A large language model for lyric and melody composition in song generation. *arXiv preprint arXiv:2402.17645*.
- David Martens Dorien Herremans and Kenneth Sörensen. 2014. [Dance hit song prediction](#). *Journal of New Music Research*, 43(3):291–302.
- Stefan Elfving, Eiji Uchibe, and Kenji Doya. 2018. Sigmoid-weighted linear units for neural network

546 function approximation in reinforcement learning.
 547 *Neural networks*, 107:3–11.

548 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
 549 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
 550 Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-
 551 ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh
 552 Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-
 553 tra, Archie Sravankumar, Artem Korenev, Arthur
 554 Hinsvark, Arun Rao, Aston Zhang, Aurelien Ro-
 555 driguez, Austen Gregerson, Ava Spataru, Baptiste
 556 Roziere, Bethany Biron, Binh Tang, Bobbie Chern,
 557 Charlotte Caucheteux, Chaya Nayak, Chloe Bi,
 558 Chris Marra, Chris McConnell, Christian Keller,
 559 Christophe Touret, Chunyang Wu, Corinne Wong,
 560 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-
 561 lonsius, Daniel Song, Danielle Pintz, Danny Livshits,
 562 Danny Wyatt, David Esiobu, Dhruv Choudhary,
 563 Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
 564 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,
 565 Elina Lobanova, Emily Dinan, Eric Michael Smith,
 566 Filip Radenovic, Francisco Guzmán, Frank Zhang,
 567 Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis An-
 568 derson, Govind Thattai, Graeme Nail, Gregoire Mi-
 569 alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,
 570 Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
 571 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-
 572 han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
 573 Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,
 574 Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
 575 Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,
 576 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang,
 577 Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,
 578 Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-
 579 teng Jia, Kalyan Vasuden Alwala, Karthik Prasad,
 580 Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth
 581 Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,
 582 Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal
 583 Lakhota, Lauren Rantala-Yearly, Laurens van der
 584 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,
 585 Louis Martin, Lovish Madaan, Lubo Malo, Lukas
 586 Blecher, Lukas Landzaat, Luke de Oliveira, Madeline
 587 Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar
 588 Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew
 589 Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-
 590 badur, Mike Lewis, Min Si, Mitesh Kumar Singh,
 591 Mona Hassan, Naman Goyal, Narjes Torabi, Niko-
 592 lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,
 593 Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick
 594 Alrassy, Pengchuan Zhang, Pengwei Li, Petar Va-
 595 sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,
 596 Praveen Krishnan, Punit Singh Koura, Puxin Xu,
 597 Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj
 598 Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,
 599 Robert Stojnic, Roberta Raileanu, Rohan Maheswari,
 600 Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-
 601 nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan
 602 Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-
 603 hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-
 604 hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-
 605 ran Narang, Sharath Raparthy, Sheng Shen, Shengye
 606 Wan, Shruti Bhosale, Shun Zhang, Simon Van-
 607 denhende, Soumya Batra, Spencer Whitman, Sten
 608 Sootla, Stephane Collot, Suchin Gururangan, Syd-

ney Borodinsky, Tamar Herman, Tara Fowler, Tarek
 Sheasha, Thomas Georgiou, Thomas Scialom, Tobias
 Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal
 Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh
 Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-
 ginie Do, Vish Vogeti, Vitor Albiero, Vladan Petro-
 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-
 ney Meers, Xavier Martinet, Xiaodong Wang, Xi-
 aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-
 feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-
 schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen,
 Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,
 Zacharie DelPierre Coudert, Zheng Yan, Zhengxing
 Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-
 vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,
 Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,
 Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei
 Baevski, Allie Feinstein, Amanda Kallet, Amit San-
 gani, Amos Teo, Anam Yunus, Andrei Lupu, And-
 res Alvarado, Andrew Caples, Andrew Gu, Andrew
 Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-
 dani, Annie Dong, Annie Franco, Anuj Goyal, Apar-
 ajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,
 Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-
 dan, Beau James, Ben Maurer, Benjamin Leonhardi,
 Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi
 Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-
 cock, Bram Wasti, Brandon Spence, Brani Stojkovic,
 Brian Gamido, Britt Montalvo, Carl Parker, Carly
 Burton, Catalina Mejia, Ce Liu, Changhan Wang,
 Changkyu Kim, Chao Zhou, Chester Hu, Ching-
 Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-
 ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,
 Daniel Kreymer, Daniel Li, David Adkins, David
 Xu, Davide Testuggine, Delia David, Devi Parikh,
 Diana Liskovich, Didem Foss, Dingkan Wang, Duc
 Le, Dustin Holland, Edward Dowling, Eissa Jamil,
 Elaine Montgomery, Eleonora Presani, Emily Hahn,
 Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-
 ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,
 Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat
 Ozgenel, Francesco Caggioni, Frank Kanayet, Frank
 Seide, Gabriela Medina Florez, Gabriella Schwarz,
 Gada Badeer, Georgia Sweet, Gil Halpern, Grant
 Herman, Grigory Sizov, Guangyi, Zhang, Guna
 Lakshminarayanan, Hakan Inan, Hamid Shojanaz-
 eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun
 Habeeb, Harrison Rudolph, Helen Suk, Henry As-
 pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim
 Damla, Igor Molybog, Igor Tufanov, Ilias Leontiadis,
 Irina-Elena Veliche, Itai Gat, Jake Weissman, James
 Geboski, James Kohli, Janice Lam, Japhet Asher,
 Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-
 nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy
 Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
 Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-
 Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,
 Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-
 delwal, Katayoun Zand, Kathy Matosich, Kaushik
 Veeraraghavan, Kelly Michelena, Keqian Li, Ki-
 ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle
 Huang, Lailin Chen, Lakshya Garg, Lavender A,
 Leandro Silva, Lee Bell, Lei Zhang, Liangpeng
 Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-

673	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	J. Kamal, P. Priya, M. R. Anala, and G. R. Smitha. 2021.	735
674	Martynas Mankus, Matan Hasson, Matthew Lennie,	A classification based approach to the prediction of	736
675	Matthias Reso, Maxim Groshev, Maxim Naumov,	song popularity. In <i>2021 International Conference</i>	737
676	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	<i>on Innovative Computing, Intelligent Communication</i>	738
677	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	<i>and Smart Electrical Systems (ICSES)</i> , pages 1–5.	739
678	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	IEEE.	740
679	Mike Macey, Mike Wang, Miquel Jubert Hermoso,		
680	Mo Metanat, Mohammad Rastegari, Munish Bansal,	Günter Klambauer, Thomas Unterthiner, Andreas Mayr,	741
681	Nandhini Santhanam, Natascha Parks, Natasha	and Sepp Hochreiter. 2017. Self-normalizing neural	742
682	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	networks. <i>Advances in neural information process-</i>	743
683	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	<i>ing systems</i> , 30.	744
684	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,		
685	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	Ping Li and Phan-Minh Nguyen. 2019. On random deep	745
686	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-	weight-tied autoencoders: Exact asymptotic analysis,	746
687	dro Rittner, Philip Bontrager, Pierre Roux, Piotr	phase transitions, and implications to training. In <i>In-</i>	747
688	Dollar, Polina Zvyagina, Prashant Ratanchandani,	<i>ternational Conference on Learning Representations.</i>	748
689	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel		
690	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	Yinghao Ma, Anders Øland, Anton Ragni, Bleiz Mac-	749
691	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	Sen Del Sette, Charalampos Saitis, Chris Donahue,	750
692	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	Chenghua Lin, Christos Plachouras, Emmanouil	751
693	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,	Benetos, Elona Shatri, et al. 2024. Foundation	752
694	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	models for music: A survey. <i>arXiv preprint</i>	753
695	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	<i>arXiv:2408.14340.</i>	754
696	Satadru Pan, Saurabh Mahajan, Saurabh Verma,		
697	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	David Martín-Gutiérrez, Gustavo Hernández Peñaloza,	755
698	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	Alberto Belmonte-Hernández, and Federico Álvarez	756
699	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	García. 2020. A multimodal end-to-end deep learn-	757
700	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	ing architecture for music popularity prediction.	758
701	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	<i>IEEE Access</i> , 8:39361–39374.	759
702	Stephanie Max, Stephen Chen, Steve Kehoe, Steve		
703	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	Matt McVicar, Bruno Di Giorgi, Baris Dundar, and	760
704	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	Matthias Mauch. 2022. Lyric document embeddings	761
705	Subramanian, Sy Choudhury, Sydney Goldman, Tal	for music tagging.	762
706	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,		
707	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	Marta Moscati, Emilia Parada-Cabaleiro, Yashar Deld-	763
708	Matthews, Timothy Chou, Tzook Shaked, Varun	joo, Eva Zangerle, and Markus Schedl. 2022.	764
709	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	Music4all-onion – a large-scale multi-faceted	765
710	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	content-centric music recommendation dataset. In	766
711	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	<i>Proceedings of the 31st ACM International Con-</i>	767
712	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	<i>ference on Information & Knowledge Management,</i>	768
713	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	CIKM '22, page 4339–4343, New York, NY, USA.	769
714	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo	Association for Computing Machinery.	770
715	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,		
716	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	Agha Haider Raza and Krishnadas Nanath. 2020. Pre-	771
717	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	dicting a hit song with machine learning: Is there an	772
718	Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary	apriori secret formula? In <i>2020 International Con-</i>	773
719	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	<i>ference on Data Science, Artificial Intelligence, and</i>	774
720	Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd	<i>Business Analytics (DATABIA)</i> , pages 111–116.	775
721	of models. Preprint, arXiv:2407.21783.		
722	Dan Hendrycks and Kevin Gimpel. 2016. Gaus-	Jing Ren, Jialie Shen, and Robert J. Kauffman. 2016.	776
723	sian error linear units (gelus). <i>arXiv preprint</i>	What makes a music track popular in online social	777
724	<i>arXiv:1606.08415.</i>	networks? In <i>Proceedings of the 25th International</i>	778
725	Hussein Hirjee and Daniel G. Brown. 2010. Rhyme ana-	<i>Conference Companion on World Wide Web, WWW</i>	779
726	lyzer: An analysis tool for rap lyrics. In <i>International</i>	'16 Companion, page 95–96, Republic and Canton of	780
727	<i>Society for Music Information Retrieval Conference</i>	Geneva, CHE. International World Wide Web Con-	781
728	<i>(ISMIR)</i> . ISMIR. Late-Breaking Demo.	ferences Steering Committee.	782
729	Thomas Hofmann. 1999. Probabilistic latent semantic	Federico Rossetto, Jeffrey Dalton, and Roderick	783
730	indexing. In <i>Proceedings of the 22nd Annual Inter-</i>	Murray-Smith. 2023. Generating multimodal aug-	784
731	<i>national ACM SIGIR Conference on Research and</i>	mentations with llms from song metadata for music	785
732	<i>Development in Information Retrieval</i> , SIGIR '99,	information retrieval. In <i>Proceedings of the 1st Work-</i>	786
733	page 50–57, New York, NY, USA. Association for	<i>shop on Large Generative Models Meet Multimodal</i>	787
734	Computing Machinery.	<i>Applications, LGM3A '23</i> , page 51–59, New York,	788
		NY, USA. Association for Computing Machinery.	789

- Prof. R.Y. Sable, Aqsa Sayyed, Baliraje Kalyane, Kosheen Sadhu, and Prathamesh Ghatole. 2024. [Enhancing music mood recognition with llms and audio signal processing: A multimodal approach](#). *International Journal for Research in Applied Science and Engineering Technology*.
- Danilo B. Seufitelli, Gabriel P. Oliveira, Mariana O. Silva, Clarisse Scofield, and Mirella M. Moro. 2023. [Hit song science: a comprehensive survey and research directions](#). *Journal of New Music Research*, 52:41 – 72.
- Mariana O. Silva, Gabriel P. Oliveira, Danilo B. Seufitelli, Anisio Lacerda, and Mirella M. Moro. 2022. [Collaboration as a driving factor for hit song classification](#). In *Proceedings of the Brazilian Symposium on Multimedia and the Web, WebMedia '22*, page 66–74, New York, NY, USA. Association for Computing Machinery.
- Abhishek Singhi and Daniel G. Brown. 2014. [Hit song detection using lyric features alone](#). In *International Society for Music Information Retrieval Conference (ISMIR): Late-Breaking Demo*, Waterloo, Canada. University of Waterloo, Cheriton School of Computer Science, ISMIR. Late-Breaking Demo.
- Abhishek Singhi and Daniel G. Brown. 2015a. [Can song lyrics predict hits?](#) In *Proceedings of the 11th International Symposium on Computer Music Multidisciplinary Research (CMMR)*, pages 457–471.
- Anurag Singhi and David G. Brown. 2015b. [Can song lyrics predict hits](#). In *International Symposium on Computer Music Multidisciplinary Research*, pages 457–471. The Laboratory of Mechanics and Acoustics.
- Yulin Yu, Pui Yin Cheung, Yong-Yeol Ahn, and Paramveer S. Dhillon. 2023. [Unique in what sense? heterogeneous relationships between multiple types of uniqueness and popularity in music](#). *Proceedings of the International AAAI Conference on Web and Social Media*, 17(1):914–925.
- Eva Zangerle, Michael Vötter, Ramona Huber, and Yi-Hsuan Yang. 2019. [Hit song prediction: Leveraging low- and high-level audio features](#). In *International Society for Music Information Retrieval Conference*.