

# ARE WE THERE YET? A BRIEF SURVEY OF MUSIC EMOTION PREDICTION DATASETS, MODELS AND OUTSTANDING CHALLENGES

Jaeyong Kang<sup>1</sup>

Dorien Herremans<sup>1</sup>

<sup>1</sup> Singapore University of Technology and Design, Singapore

jaeyong\_kang@sutd.edu.sg, dorien\_herremans@sutd.edu.sg

## ABSTRACT

Deep learning models for music have advanced drastically in the last few years. But how good are machine learning models at capturing emotion these days and what challenges are researchers facing? In this paper, we provide a comprehensive overview of the available music-emotion datasets and discuss evaluation standards as well as competitions in the field. We also provide a brief overview of various types of music emotion prediction models that have been built over the years, offering insights into the diverse approaches within the field. Through this examination, we highlight the challenges that persist in accurately capturing emotion in music. Recognizing the dynamic nature of this field, we have complemented our findings with an accompanying GitHub repository. This repository contains a comprehensive list of music emotion datasets and recent predictive models<sup>1</sup>.

## 1. INTRODUCTION

Music has long been revered for its profound ability to evoke and convey emotions, transcending cultural, linguistic, and geographical barriers. Researchers from various fields have been captivated by the intricate interplay between music and human emotions for decades. Classic works such as Meyer’s paper [1] and pioneering studies conducted by scholars such as Seashore [2] and Hevner [3] in the early 20th century laid the groundwork for understanding the emotional impact of music. However, understanding and quantifying this poses a significant challenge due to its multifaceted nature [4]. In this paper, we provide an overview of the current state-of-the-art along with challenges and directions for future work.

With the advent of technology and data-driven methodologies, new avenues have opened up for exploring the complex relationship between music and emotions, such

as deep learning models. Despite these advancements, accurate predictions of emotions from music remains elusive. This is due to a number of reasons, including the subjective, personal nature of emotional perception, as well as bias and limitations in the current datasets, and challenges in benchmarking. We discuss these and other challenges at length in Section 5.

Despite these challenges, the potential applications of Music Emotion Recognition (MER) systems are vast and varied. In the healthcare domain, we may find personalized music recommendation systems that guide a listener to different emotional states [5], or we may even see MER systems used to build large-scale datasets on which we can train generative music AI systems that can be controlled to generate music with specific emotions [6, 7]. In the same line of thought, [8] developed a Brain-Computer Interface system that can provide musical feedback about the listener’s current emotional state and subsequently influence this state. Additionally, MER systems may help facilitate emotional analysis during music composition, interactive experiences in media and entertainment, as well as inform market research [9]. The implications of understanding music and emotions extend across diverse domains. In general, following the idea of positive psychology [10], understanding music emotions can be used to enhance the user experience when designing various systems.

In this paper, we do not aim to provide a comprehensive overview of MER models, instead, we focus on discussing datasets, evaluation approaches, and identifying key challenges as well as future directions. We only briefly touch upon some of the more recent MER models. For a more comprehensive overview of MER machine learning models, the reader is referred to [9, 11–14].

In the next section, we provide an extensive overview of the available emotion-annotated music datasets. This is followed by a discussion of the evaluation practices in the field (Section 3). After that, we describe a selected list of recent models and approaches in Section 4. Finally, Section 5 dives into the remaining challenges and future direction for the field of MER, followed by a general conclusion.

<sup>1</sup> <https://github.com/AMAAI-Lab/awesome-MER/>



**Table 1.** Overview of Music Emotion Datasets (Sorted by Year).

Dataset	Year	# of instances	Length	Type	Categorical	Dimensional	Static/ Dynamic	Perceived/ Induced
MoodsMIREX [15]	2007	269	30s	MP3	5 labels	-	Static	Perceived
CAL500 [16]	2007	500	full	MP3	174 labels	-	Static	Perceived
Yang-Dim [17]	2008	195	25s	WAV	-	Russell	Static	Perceived
MoodSwings [18]	2008	240	15s	MP3	-	Russell	Dynamic	Perceived
NTWICM [19]	2010	2,648	full	MP3	-	Russell	Static	Perceived
Soundtracks [20]	2011	470	15s-1m	MP3	6 labels	3 dimensions	Static	Perceived
DEAP [21]	2012	120	60s	YouTube id	-	Russell	Static	Induced
Panda et al.’s dataset [22]	2013	903	30s	MP3, MIDI	21 labels	-	Static	Perceived
Solymani et al.’s dataset [23]	2013	1000	45s	MP3	-	Russell	Both	Perceived
Emotify [24]	2016	400	60s	MP3	GEMS	-	Static	Induced
Moodo [25]	2016	200	15s	WAV	-	Russell	Static	Perceived
CH818 [26]	2017	818	30s	MP3	-	Russell	Static	Perceived
4Q-emotion [27]	2018	900	30s	MP3	Quadrants	-	Static	Perceived
MediaEval DEAM [28]	2018	2,058	45s	MP3	-	Russell	Both	Perceived
PMemo [29]	2018	794	full	MP3	-	Russell	Both	Induced
RAVDESS [30]	2018	1,012	full	MP3, MP4	5 labels	-	Static	Perceived
DMDD [31]	2018	18,644	full	Audio, Lyrics	-	Russell	Static	Perceived
MTG-Jamendo [32]	2019	18,486	full	MP3	56 labels	-	Static	Perceived
VGMIDI [33]	2019	200	full	MIDI	-	Russell	Dynamic	Perceived
Turkish Music Emotion [34]	2019	400	30s	MP3	4 labels	-	Static	Perceived
EMOPIA [35]	2021	1,087	30s-40s	Audio, MIDI	Quadrants	-	Static	Perceived
MER500 [36]	2020	494	10s	WAV	5 labels	-	Static	Perceived
Music4all [37]	2020	109,269	30s	WAV	-	3 dimensions	Static	Perceived
CCMED-WCMED [38]	2020	800	8-20s	WAV	-	Russell	Static	Perceived
MuSe [39]	2021	90,001	full	Audio	-	Russell (V-A-D)	Static	Perceived
HKU956 [40]	2022	956	full	MP3	-	Russell	Static	Induced
MERP [41]	2022	54	full	WAV	-	Russell	Both	Perceived
MuVi [42]	2022	81	full	YouTube id	GEMS	Russell	Both	Perceived
YM2413-MDB [43]	2022	699	full	WAV, MIDI	35 labels	-	Static	Perceived
MusAV [44]	2022	2,092	full	WAV	-	Russell	Static	Perceived
EmoMV [45]	2023	5,986	30s	WAV	6 labels	-	Static	Perceived
Indonesian Song [46]	2023	476	full	WAV	3 labels	-	Static	Perceived
TROMPA-MER [47]	2023	1,161	30s	WAV	11 labels	-	Static	Perceived
EMMA [48]	2024	364	30s-60s	WAV	GEMS	-	Static	Perceived
SiTunes [49]	2024	300	full	WAV	-	Russell	Static	Induced

## 2. DATASETS

Table 1 presents an extensive overview of emotion-annotated music datasets. These datasets vary in size, annotation granularity, and focus on either perceived or induced emotions. Before delving deeper into the datasets, it is essential to understand the emotion representations commonly referenced in the field.

**Emotion representations** One of the earliest emotion representation models in music is Hevner’s affective ring [50], developed in 1936. Based on extensive experimental studies, Hevner’s model categorizes music emotions into eight fundamental categories: dignified, sad, dreamy, serene, graceful, happy, exciting, and vigorous. In current-day MER research, we see that Russell’s Circumplex Model of Affect [51] is widely used. This model characterizes emotions along two dimensions: valence and arousal. Valence represents the degree of positive or negative emotion, while arousal reflects the intensity of emotion, ranging from passive to activated states. Russell’s original model contained a third dimension: dominance [51]. This dimension is typically omitted as it can be hard to annotate, although some researchers have argued to re-include it [52]. Russell’s model is a *dimensional* model, as

it consists of continuous values along multiple dimensions (valence/arousal). Hevner’s model, on the other hand, is *categorical* as it consists of discrete emotion labels.

Thayer’s two-dimensional model [53] focuses on energetic arousal and tense arousal as the primary dimensions of emotion. Thayer suggests that valence can be inferred from the combination of energetic and tense arousal levels.

Categorical models include the Geneva Emotional Music Scales (GEMS) [24]. This model was specifically designed for music-induced emotions and consists of 45 emotion tags grouped into nine categories, including amazement, solemnity, tenderness, nostalgia, calmness, power, joyful activation, tension, and sadness. The datasets listed in Table 1 use a variety of emotion representation models, with Russell’s model being the most popular dimensional representation. Some datasets do not use a specific emotion model, for instance, the MTG-Jamendo dataset consists of freely assigned tags by the listeners, resulting in a diverse and comprehensive set of 56 tags, spanning from ‘melancholic’ to ‘upbeat’. A subset of this dataset is used for the ‘Emotion and Theme Recognition in Music’ Task of MediaEval. In Table 1, the number of tags used to represent the emotions are listed in the column ‘Categorical’.

**Static versus dynamic** Regardless of which emotion representation model is being used, we notice two fundamentally different approaches: static versus dynamic annotations. In a static setting, the listeners indicate the emotion for the entire song or fragment. In a dynamic annotation setting, the listener continuously indicates the emotion throughout the song or fragment. For instance, the MTG-Jamendo dataset [32] offers categorical tags for each full-length song. The annotation is static, but multiple tags are allowed per song. The MoodSwings dataset [18], on the other hand, offers dynamic annotations of valence/arousal for every second of the musical fragments. Finally, some datasets offer both, for instance, MERP [41] provides both static GEM labels for the entire song, as well as dynamic valence/arousal ratings for every 1s of a song.

**Induced versus perceived** Emotion labels in the datasets can represent perceived or induced emotion. Perceived emotions are emotions that listeners consciously recognize while listening to the music. Induced emotions, on the other hand, are emotions that listeners experience as a result of listening to the music, but are not necessarily consciously recognized. Most of the datasets in Table 1 use perceived emotion labels, which are easier to annotate. For induced emotion labels, researchers have used psycho-physiological measurements ranging from electromyogram (EMG), volume pulse (BVP), electrocardiograms (ECG), skin conductance, respiration rate, heart rate, to electroencephalograms (EEG). There have been many studies in psychology that use such biosensors to explore how music can influence our emotions [54–58]. Since these studies include medical data, the datasets are not often public. However, there are a few datasets with music and its induced emotions. Firstly the DEAP dataset [21] includes EEG, facial video recordings, as well as peripheral physiological signals that were recorded while they watched music videos. In addition, they collected perceived emotion ratings in terms of arousal, valence, like/dislike, dominance and familiarity. Second, the HKU956 dataset [40] records five kinds of physiological signals (i.e., heart rate, electrodermal activity, blood volume pulse, inter-beat interval, and skin temperature) of participants as they listen to music, along with reported emotions in the arousal and valence dimensions. Finally, SiTunes [49] includes physiological signals (i.e., heart rate, activity intensity, activity step, and activity type) measured both before and after music listening, alongside environmental data (i.e., time of day, weather information, and location) recorded during users’ daily lives. The MER models that can predict *induced* emotions have various medical applications, such as curating playlists to guide patients to different emotional states [5]. An interesting study by [59] concludes that music tends to induce the emotion that is perceived, enabling researchers to use perceived emotion datasets for developing emotion-inducing models.

**Modalities** Music comes in many formats, the most common one being ‘audio’. Audio files can either be raw waveforms or compressed .mp3 files. However, we should not neglect the MIDI format, a popular format still of-

ten used by music producers, composers and performers. Whereas most datasets contain audio files, only a handful focus on MIDI: 1) the VGMIDI dataset [33], which contains continuous valence/arousal ratings for MIDI files of piano arrangements for video game soundtracks; 2) the Panda et al.’s dataset [22] (used for the MIREX competition), which provides a diverse collection of audio clips, lyrics, and aligned MIDI files, contains labels for five emotion clusters derived from a cluster analysis of online tags; and 3) the EMOPIA dataset [35], which contains paired piano music audio with MIDI that has emotion annotations of the four high/low valence/arousal quadrants.

Many sources can elicit emotion. For instance, video, or lyrics may also affect our perceived or induced emotions. Some datasets offer alternative multimedia streams such as emotion-rated music videos from a variety of genres, including pop, rock, classical, and jazz as in the DEAP dataset [21]; or text of lyrics with the annotated music in the DMDD dataset [31]. Finally, the MuVi dataset [42] even offers isolated modality ratings for music videos. In this study, the raters were presented with either the music video, the music alone, or the muted video. The final dataset contains ratings for each of these modalities separately as well as together. This allows [42] to build a model on pure isolated modalities, which proved to be more accurate than a traditional model.

**Dataset size** Compared to affective datasets available in other domains (e.g., the Sentiment140 dataset [60] in NLP, which contains 1.6 million tweets annotated as positive or negative), the size of the available datasets is still very limited, with the largest dataset containing 109,269 instances. In addition to the number of instances in datasets, we also notice a difference in the length of the instances. A number of datasets (e.g. MERP [41], VGMIDI [33], and HKU956 [40]) offer ratings for full-length songs. In datasets with dynamic ratings throughout the song, this may provide a means for researchers to analyse how our emotions evolve throughout a song. Other datasets focus on short fragments, often ranging from 30 seconds to 1-minute fragments (e.g. DEAM [28], EMOPIA [35], and EMMA [48]), but even as short as 10s as is the case for the MER500 dataset [36].

In sum, there is a variety of emotion-annotated datasets available as shown in Table 1. They differ in terms of emotion representation models, as well as music format and perceived versus induced emotion labels. The challenges and future research directions related to datasets are discussed in detail in Section 5.

### 3. EVALUATION PROTOCOLS

Evaluation metrics play a crucial role in assessing the performance of MER systems. Depending on which type of emotion ratings are being predicted (dimensional versus categorical), the evaluation metrics change as we are dealing with a regression or a classification task respectively. Commonly used evaluation metrics for categorical MER systems include accuracy, precision, Area under the Curve (AUC), and confusion matrices. In the case of regression

MER models, we see metrics such as precision, recall, F1-score, mean squared error (MSE),  $R^2$ , and Pearson correlation coefficient [61]. These metrics provide insights into the effectiveness of MER models when it comes to emotions represented by dimensional models.

Evaluating MER systems goes beyond just defining a common metric. Due to the inherent differences between datasets, a comparison across datasets is often not possible. For an in-depth discussion on this, the reader is referred to Section 5. Within one dataset, however, it is possible to establish benchmarks and compare the performance of different models, provided that the train/test split is shared. There are some initiatives to facilitate the comparison between models, such as competitions, as well as individual papers that offer clear data splits and metrics [62–66].

Competitions and challenges provide valuable platforms for evaluating and comparing the performance of MER systems. These initiatives often involve standardized datasets, evaluation protocols, and metrics, enabling researchers to benchmark their algorithms against state-of-the-art methods. Existing benchmarking initiatives and competitions in MER include the ‘Audio K-POP Mood Classification’ task in MIREX (Music Information Retrieval Evaluation eXchange) (last organized in 2019)<sup>2</sup>, the ‘Emotion in Music’ task in MediaEval (last organized in 2015)<sup>3</sup>, and the ‘Emotion and Theme Recognition in Music Using Jamendo’ in MediaEval (last organized in 2021)<sup>4</sup>.

Given the noisiness of labels in music-emotion datasets [41], one could also argue to include listening tests to verify the accuracy of the predicted labels, or even to perform extrinsic evaluation through a downstream task. The latter could be implementing MER in an emotion-inducing system, e.g. by performing playlist recommendations based on predicted emotions from music in the dataset, and subsequently evaluating if the listeners’ emotion indeed reaches the target emotion.

#### 4. MODELS AND APPROACHES

In this section, we briefly touch upon some existing MER models. This is by no means a complete overview, for this, the reader is referred to other survey papers [9, 11–14]. The aim of this section is merely to point out the current state-of-the-art and various approaches over the last few years.

Some of the earliest attempts at music emotion prediction involved rule-based approaches and hierarchical frameworks. For instance, Feng et al. [67] used Computational Media Aesthetics (CMA) to analyze tempo and articulation, mapping them into four mood categories: happiness, anger, sadness, and fear. They achieved a total precision of 67% and a total recall of 66%. Lu et al. [68] developed a hierarchical framework for automatically extracting music emotion from acoustic data. They em-

ployed music intensity to represent the energy dimension of Thayer’s model while using timbre and rhythm to capture the stress dimension. They achieved an average accuracy of mood detection of up to 86.3%. These early models laid the groundwork for subsequent advancements in MER, paving the way for the adoption of more sophisticated techniques, including deep learning approaches.

In the last few years, various deep-learning MER models have been developed. Many of these leverage common techniques for time-series data, including Recursive Neural Networks (RNNs) such as Long-Short Time Memory networks (LSTMs) [42, 69], and more recently, Transformer architectures [46, 70]. Trying to identify the current state-of-the-art MER model is a near-impossible task, due to reasons discussed in the previous section, for instance, the difference in labels across datasets. Looking at a recent competition, the 2021 Emotion and Theme Recognition in Music Using MTG-Jamendo dataset, we see that the best-performing model [66] reaches an PR-AUC-macro of 0.150872. This model uses a Convolutional neural networks (CNN) architecture.

A distinction between the various models can be made based on the input modality. Some models are exclusively based on MIDI such as the multi-task architecture proposed by [71], and as such use a token-based representation. Given the limited size of MIDI datasets, the accuracy of such models is limited, e.g., the model by [71] reaches 67.56% accuracy when predicting between four emotion classes. Most of the existing MER models are based on audio, and hence they take as input raw audio. This is often converted into Mel-spectrograms which are then processed through convolutional neural networks (CNNs) [72], or the audio could be directly processed through a WaveNet architecture, which is a type of temporal CNN [73]. In recent years, audio embeddings pretrained on large-scale datasets have become available, providing researchers with tools to enable transfer learning. For instance, [74] incorporated a VGGish network with Transformers, achieving mean squared error (MSE) and Pearson correlation coefficient (PCC) values of 0.117 and 0.655 for arousal, respectively, and 0.170 and 0.575 for valence, respectively.

Alternatives to these approaches include directly extracting spectral features (e.g., MFCCs, spectral centroids) with libraries such as OpenSmile, which has a configuration file specifically for the emotion recognition task [75]. [62] uses this approach and achieves  $R^2$  scores of 0.378 and 0.638 for predicting dynamic valence and arousal values, respectively, on the MediaEval dataset.

Finally, musically meaningful features may be included, such as Rhythmic features (e.g., tempo, beat histogram), or note features (e.g., pitch), as implemented by Shi et al. [76] and Panda et al. [27], respectively. Shi et al. achieved a precision of 92.8% for 4 emotion categories (calm, sad, pleasant, and excited), while Panda et al. achieved an F1-score of 76.0% for 4 emotion categories (Quadrants).

Other input modalities may include text (lyrics), or video. In the case of the former, some models extract the

<sup>2</sup>[https://www.music-ir.org/mirex/wiki/2019:Audio\\_K-POP\\_Mood\\_Classification/](https://www.music-ir.org/mirex/wiki/2019:Audio_K-POP_Mood_Classification/)

<sup>3</sup><http://www.multimediaeval.org/mediaeval2015/emotioninmusic2015/>

<sup>4</sup><https://multimediaeval.github.io/2021-Emotion-and-Theme-Recognition-in-Music-Task/>

sentiment from the lyrics using Natural Language Processing (NLP) tools [77], or they use the entire lyrics with an embedding model [70]. These features are then combined with audio-based features or analyzed independently to predict the emotional content of music. Including the lyrics does not always improve the sentiment prediction, particularly in predicting arousal [77], however, [31] did manage to increase the performance of a MER model by using a word2vec [78] embedding trained on 1.6 million lyrics. Finally, models that include video modalities typically use pretrained networks to capture image and video features, and thus improve model performance. For instance, [45] uses ResNet-50 [79] and FlowNetS [80], respectively.

We have only provided a glimpse into the existing MER models in this section. From the performance of the various models, we see, however, that much improvement can still be made. We discuss some of the remaining challenges in the next section.

## 5. CHALLENGES AND FUTURE DIRECTIONS

Whereas the first publications on music and emotion surfaced in the 1930s [50], we have not yet reached the stage where models can reach human-like performance (e.g. > 95% accuracy). The field of MER still faces several challenges and various opportunities for future exploration remain, ranging from overcoming data limitations to the integration of emerging technologies. Understanding and addressing these challenges is crucial for advancing the field and unlocking its full potential.

**Dataset limitations** One of the primary challenges in MER is the scarcity of large, diverse, copyright-cleared, emotion-annotated datasets. Limited datasets hinder the development and evaluation of robust MER models, leading to potential biases and generalization issues. We have recently seen advances in this area with the release of larger datasets such as MTG-Jamendo [32], Music4all [37], MuSe [39], which contain 18k, 109k, and 90k instances respectively. Whereas MTG-Jamendo and MuSe are available under a Creative Commons licence, however, Music4all contains copyrighted tracks. In addition, when exploring datasets, we notice that they are often skewed towards one particular genre. For instance, the DEAM dataset [28] consists mostly of rock and electronic music genres, while the VGMIDI dataset is focused solely on video game soundtracks.

To deal with the current dataset size limitations, techniques such as unsupervised pretrained may be helpful. With this technique, latent representations are first learned using unlabeled datasets. This evolution has led to some available large-scale audio encoders such as the Variational Auto Encoder used in [81], and the AST Audio Spectrogram Transformer presented in [82], as well as various recently developed neural audio encoders (e.g. Descript Audio Codec [83]). These novel pretrained representations may help deal with the limiting size of emotion-annotated datasets.

**Subjective labels** The subjective and variable nature

of emotion perception also poses a significant challenge. Emotions are inherently complex and subjective, varying across individuals, cultures, and contexts [41]. For instance, researchers have observed significant dependence between the number of years of musical training [41, 84], gender [42], familiarity with songs [42], culture [41, 85], genre preference [41], and even age [41]. The latter study makes an argument for including profile information in the emotion prediction model to personalize the predictions. However, not many datasets include this type of information other than MERP.

In addition, many datasets have a cultural bias, meaning that they are annotated by people from the same culture, often simply because of the locality of the experiment or language constraints. However, it has been established that people from different countries or cultural backgrounds have different perceptions of emotion for the same music fragment [86]. For instance, [41] have raters from both the US as well as India, and see a significant difference in their annotations. HKU956 [40] go even further and include the rater’s responses to a personality test: ‘Ten Item Personality Measure’. In future work, one could develop datasets annotated by annotators from different cultural backgrounds, that include this information about the raters. The personalized nature of emotion ratings can cause a low inter-rater reliability, a metric of agreement between raters often calculated using Cronbach’s Alpha [87].

**Noisy labels** When creating datasets, an additional challenge arises: it can be hard to identify the emotion perceived from music, especially when working with valence and arousal models. In fact, Russel’s model [51] included a third dimension: dominance. This dimension is typically omitted because of the ambiguity in annotation. In general, categorical models, although less precise, are often easier to annotate [59]. This personal variability and ambiguity in emotion labels introduces noise in dataset labels, causing low inter-rater reliability, and making it challenging to train accurate and reliable MER models.

Machine learning models may also offer useful techniques to deal with noisy labels, as discussed in the survey by [88]. These could include Noisy Graph Cleaning (NGC) [89], Joint Training with Co-Regularization (JoCoR) [90], and Robust Curriculum Learning (RoCL) [91].

**Annotation interfaces** When creating a static dataset with static annotations, some standard tools can be used, including PsyToolkit [92]. However, to create any larger-scale dataset, the annotations are often done through crowdsourcing services such as Amazon Mechanical Turk<sup>5</sup> (e.g. for MERP [41], DEAM [28], and Solyman’s dataset [23]). These services offer access to a large ‘army’ of annotators, which may come at the cost of accuracy. There are, however, a number of techniques that can be used to filter out some of this annotation noise. This includes limiting the annotations to Master raters (raters with a known track record that often work at a premium price) [41], by including qualification tasks to assess participants’ understanding of the dimensional model [23], or

<sup>5</sup> <https://www.mturk.com>

by using multiple ground truth questions that are shown to all raters [41]. Finally, inter-rater reliability can be used to filter out low-quality annotations [41]. When doing this, it is important to keep in mind personal characteristics, which may cause different people to rate music differently. Hence, inter-rater reliability should ideally be calculated by taking into account the rater’s profile features.

An additional difficulty is that the amount of available interfaces for music emotion annotation is very limited. Especially when it comes to time-continuous annotations of valence and arousal. [41] released their dynamic annotation interface<sup>6</sup> that hooks into Amazon mTurk. Kim et al. [18] also introduced, an annotation interface called MoodSwings, designed to record dynamic emotion labels.

**Benchmarking** The ImageNet competition has been instrumental in establishing a clear performance benchmark among computer vision models [93]. While there have been similar competitions for music emotion prediction (see Section 3), none of these competitions ran in the last three years, indicating the lack of a current benchmark for MER systems. To establish benchmarks on individual datasets, we have to revert to individual model papers, such as [23, 28]. It is unfortunately not always clear which train/test split these systems use, making direct comparisons hard. One way to facilitate an easy overview would be to leverage Leaderboard features on popular websites such as Papers With Code<sup>7</sup> and HuggingFace<sup>8</sup>.

Finally, it is also hard to make comparisons across datasets, as many of them use different emotion representations. While some of these differences may be bridged by using synonyms, it is currently not a common practice to compare across datasets. However, if we are able to merge different emotion representations and train across multiple datasets, then we can start leveraging larger-scale models. Bridging between datasets that use continuous representations and categorical models is the least evident. [94] have identified a mapping between arousal/valence and categorical labels. This has been used by [6] to facilitate emotion-controlled lead sheet generation, and it may offer a way to merge continuous and categorical datasets into a larger-scale dataset for music emotion.

**MIDI** The forgotten format in MER. Emotion originates from many different aspects of the music, including the tonal tension [7], instrumentation and timbre [95], production quality [96], expressiveness of the performance [97], harmony [98]. The symbolic MIDI format only captures parts of these [12], which may explain the lack of models for predicting emotion from MIDI. Another reason may simply be the lack of emotion-annotated datasets (three in total). Even though MIDI is an abstraction of music, it is still a widely used format by music producers and performers and warrants its own models for emotion prediction. Even more, such datasets may enable generative music systems (which are typically trained on MIDI) to be controlled by emotion [6, 99].

<sup>6</sup><https://github.com/dorienh/MERP>

<sup>7</sup><https://paperswithcode.com/sota>

<sup>8</sup><https://huggingface.co/docs/competitions/en/leaderboard>

**Multimodal predictions** Our sensory input is multimodal. As such, some of the emotion-annotated datasets enable us to look at multiple modalities. For instance, DEAP [21] offers biofeedback data, i.e. EEG recordings, and frontal face videos from participants. Similarly, HKU956 [40] offers physiological signals including heart rate, electrodermal activity, blood volume pulse, inter-beat interval, and skin temperature. SiTunes [49], on the other hand, provides physiological and environmental situation recordings collected via smart wristband devices.

These biological data can serve as the induced emotion labels, e.g. EEG signals can be translated into human emotions [100]. Increased datasets with different types of physiological signals enable researchers to focus on creating models for induced versus perceived emotion detection. This offers new avenues to use biofeedback in music emotion mediation applications through smart devices.

Sometimes other emotion-inducing modalities are present, such as video, or lyrics. In this case, it is important to consider the influence of each of these modalities. In the case of video and music, Chua et al. [42] have studied the influence of exposing participants not only to the music but also the muted video, as well as the music videos. They found that the music modality explains most of the variance in arousal values, and both music and video modalities explain the variance in valence values. Phuong et al. [74] explore the influence of using only audio features and only video features to predict emotions from movie fragments. They found that the prediction is most accurate when both modalities are used. However, when predicting from a single modality, the audio model is most accurate.

**Real-time** Many of the currently existing models are not implemented as an easy-to-use library, nor are they quick to run. They often require a GPU and may take several minutes to run. There are use cases, however, for real-time emotion recognition systems, as they would be able to integrate into therapeutic emotion detection systems [5], mood guidance playlist systems, as well as more commercial systems such as advertisement targeting systems.

In sum, the challenges mentioned above provide direct opportunities and future directions to further advance the exciting field of MER.

## 6. CONCLUSION

MER is a promising field with various practical applications. With the rise of large-language models, we have seen impressive performance in various tasks. The field of MER, however, still seems to be lagging. Given the importance of large training datasets to facilitate the training of LLMs, we provide a comprehensive overview and discussion of the existing datasets for MER.

We also explore current state-of-the-art models and dive into evaluation methods such as metrics as well as competitions, leaderboards, and benchmarks within the MER field. With this knowledge, we discuss the current challenges of the MER field at length and provide concrete future directions and emerging trends such as real-time systems and multimodal prediction systems.

In closing, this survey serves as a valuable resource for the MER community, offering insights into the current state-of-the-art, as well as a discussion of challenges and inspiration for future directions.

## 7. ACKNOWLEDGEMENTS

This work has received SEED funding from SUTD TL under grant number RTDS S 22 14 04 1.

## 8. REFERENCES

- [1] M. Leonard, "Emotion and meaning in music," *Chicago: University of Chicago*, 1956.
- [2] C. E. Seashore, "The psychology of music," *Music Educators Journal*, vol. 23, no. 4, pp. 30–33, 1937.
- [3] K. Hevner, "The affective character of the major and minor modes in music," *The American Journal of Psychology*, vol. 47, no. 1, pp. 103–118, 1935.
- [4] P. N. Juslin and D. Västfjäll, "Emotional responses to music: The need to consider underlying mechanisms," *Behavioral and brain sciences*, vol. 31, no. 5, pp. 559–575, 2008.
- [5] K. R. Agres, R. S. Schaefer, A. Volk, S. van Hooren, A. Holzapfel, S. Dalla Bella, M. Müller, M. De Witte, D. Herremans, R. Ramirez Melendez *et al.*, "Music, computing, and health: a roadmap for the current and future roles of music technology for health care and well-being," *Music & Science*, vol. 4, p. 2059204321997709, 2021.
- [6] D. Makris, K. R. Agres, and D. Herremans, "Generating lead sheets with affect: A novel conditional seq2seq framework," in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2021, pp. 1–8.
- [7] D. Herremans and E. Chew, "Morpheus: generating structured music with constrained patterns and tension," *IEEE Transactions on Affective Computing*, vol. 10, no. 4, pp. 510–523, 2017.
- [8] K. R. Agres, A. Dash, and P. Chua, "Affectmachine-classical: a novel system for generating affective classical music," *Frontiers in Psychology*, vol. 14, p. 1158172, 2023.
- [9] Y.-H. Yang and H. H. Chen, "Machine recognition of music emotion: A review," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 3, no. 3, pp. 1–30, 2012.
- [10] M. E. Seligman and M. Csikszentmihalyi, *Positive psychology: An introduction*. American Psychological Association, 2000, vol. 55, no. 1.
- [11] D. Han, Y. Kong, J. Han, and G. Wang, "A survey of music emotion recognition," *Frontiers of Computer Science*, vol. 16, no. 6, p. 166335, 2022.
- [12] Y. E. Kim, E. M. Schmidt, R. Migneco, B. G. Morton, P. Richardson, J. Scott, J. A. Speck, and D. Turnbull, "Music emotion recognition: A state of the art review," in *Proc. ismir*, vol. 86, 2010, pp. 937–952.
- [13] X. Yang, Y. Dong, and J. Li, "Review of data features-based music emotion recognition methods," *Multimedia systems*, vol. 24, pp. 365–389, 2018.
- [14] M. Barthelet, G. Fazekas, and M. Sandler, "Music emotion recognition: From content-to context-based models," in *From Sounds to Music and Emotions: 9th International Symposium, CMMR 2012, London, UK, June 19-22, 2012, Revised Selected Papers 9*. Springer, 2013, pp. 228–252.
- [15] X. Hu and J. S. Downie, "Exploring mood metadata: Relationships with genre, artist and usage metadata." in *ISMIR*, 2007, pp. 67–72.
- [16] D. Turnbull, L. Barrington, D. Torres, and G. Lanckriet, "Towards musical query-by-semantic-description using the cal500 data set," in *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, 2007, pp. 439–446.
- [17] Y.-H. Yang, Y.-C. Lin, Y.-F. Su, and H. H. Chen, "A regression approach to music emotion recognition," *IEEE Transactions on audio, speech, and language processing*, vol. 16, no. 2, pp. 448–457, 2008.
- [18] Y. E. Kim, E. M. Schmidt, and L. Emelle, "Moodswings: A collaborative game for music mood label collection." in *Ismir*, vol. 8, 2008, pp. 231–236.
- [19] B. Schuller, J. Dorfner, and G. Rigoll, "Determination of nonprototypical valence and arousal in popular music: features and performances," *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2010, pp. 1–19, 2010.
- [20] T. Eerola and J. K. Vuoskoski, "A comparison of the discrete and dimensional models of emotion in music," *Psychology of Music*, vol. 39, no. 1, pp. 18–49, 2011.
- [21] S. Koelstra, C. Muhl, M. Soleymani, J.-S. Lee, A. Yazdani, T. Ebrahimi, T. Pun, A. Nijholt, and I. Patras, "Deap: A database for emotion analysis; using physiological signals," *IEEE transactions on affective computing*, vol. 3, no. 1, pp. 18–31, 2011.
- [22] R. E. S. Panda, R. Malheiro, B. Rocha, A. P. Oliveira, and R. P. Paiva, "Multi-modal music emotion recognition: A new dataset, methodology and comparative analysis," in *10th International symposium on computer music multidisciplinary research (CMMR 2013)*, 2013, pp. 570–582.
- [23] M. Soleymani, M. N. Caro, E. M. Schmidt, C.-Y. Sha, and Y.-H. Yang, "1000 songs for emotional analysis of music," in *Proceedings of the 2nd ACM international workshop on Crowdsourcing for multimedia*, 2013, pp. 1–6.

- [24] M. Zentner, D. Grandjean, and K. R. Scherer, “Emotions evoked by the sound of music: characterization, classification, and measurement.” *Emotion*, vol. 8, no. 4, p. 494, 2008.
- [25] M. Pesek, G. Strle, A. Kavčič, and M. Marolt, “The moodo dataset: Integrating user context with emotional and color perception of music for affective music information retrieval,” *Journal of New Music Research*, vol. 46, no. 3, pp. 246–260, 2017.
- [26] X. Hu and Y.-H. Yang, “The mood of chinese pop music: Representation and recognition,” *Journal of the Association for Information Science and Technology*, vol. 68, no. 8, pp. 1899–1910, 2017.
- [27] R. Panda, R. Malheiro, and R. P. Paiva, “Musical texture and expressivity features for music emotion recognition,” in *19th International Society for Music Information Retrieval Conference (ISMIR 2018)*, 2018, pp. 383–391.
- [28] A. Aljanaki, Y.-H. Yang, and M. Soleymani, “Developing a benchmark for emotional analysis of music,” *PloS one*, vol. 12, no. 3, p. e0173392, 2017.
- [29] K. Zhang, H. Zhang, S. Li, C. Yang, and L. Sun, “The pmemo dataset for music emotion recognition,” in *Proceedings of the 2018 acm on international conference on multimedia retrieval*, 2018, pp. 135–142.
- [30] S. R. Livingstone and F. A. Russo, “The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english,” *PloS one*, vol. 13, no. 5, p. e0196391, 2018.
- [31] R. Delbouys, R. Hennequin, F. Piccoli, J. Royo-Letelier, and M. Moussallam, “Music mood detection based on audio and lyrics with deep neural net,” *arXiv preprint arXiv:1809.07276*, 2018.
- [32] D. Bogdanov, M. Won, P. Tovstogan, A. Porter, and X. Serra, “The mtg-jamendo dataset for automatic music tagging.” *ICML*, 2019.
- [33] L. N. Ferreira and J. Whitehead, “Learning to generate music with sentiment,” *arXiv preprint arXiv:2103.06125*, 2021.
- [34] M. B. Er and I. B. Aydilek, “Music emotion recognition by using chroma spectrogram and deep visual features,” *International Journal of Computational Intelligence Systems*, vol. 12, no. 2, pp. 1622–1634, 2019.
- [35] H.-T. Hung, J. Ching, S. Doh, N. Kim, J. Nam, and Y.-H. Yang, “Emopia: A multi-modal pop piano dataset for emotion recognition and emotion-based music generation,” *arXiv preprint arXiv:2108.01374*, 2021.
- [36] M. Velankar, “Mer500 — music emotion recognition,” <https://www.kaggle.com/>, Jun. 2020, accessed March 10, 2024.
- [37] I. A. P. Santana, F. Pinhelli, J. Donini, L. Catharin, R. B. Mangolin, V. D. Feltrim, M. A. Domingues *et al.*, “Music4all: A new music database and its applications,” in *2020 International Conference on Systems, Signals and Image Processing (IWSSIP)*. IEEE, 2020, pp. 399–404.
- [38] J. Fan, Y.-H. Yang, K. Dong, and P. Pasquier, “A comparative study of western and chinese classical music based on soundscape models,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 521–525.
- [39] C. Akiki and M. Burghardt, “Muse: The musical sentiment dataset,” *Journal of Open Humanities Data*, vol. 7, 2021.
- [40] X. Hu, F. Li, and R. Liu, “Detecting music-induced emotion based on acoustic analysis and physiological sensing: A multimodal approach,” *Applied Sciences*, vol. 12, no. 18, p. 9354, 2022.
- [41] E. Y. Koh, K. W. Cheuk, K. Y. Heung, K. R. Agres, and D. Herremans, “Merp: a music dataset with emotion ratings and raters’ profile information,” *Sensors*, vol. 23, no. 1, p. 382, 2022.
- [42] P. Chua, D. Makris, D. Herremans, G. Roig, and K. Agres, “Predicting emotion from music videos: exploring the relative contribution of visual and auditory information to affective responses,” *arXiv preprint arXiv:2202.10453*, 2022.
- [43] E. Choi, Y. Chung, S. Lee, J. Jeon, T. Kwon, and J. Nam, “Ym2413-mdb: A multi-instrumental fm video game music dataset with emotion annotations,” *arXiv preprint arXiv:2211.07131*, 2022.
- [44] D. Bogdanov, X. Lizarraga Seijas, P. Alonso-Jiménez, and X. Serra, “Musav: A dataset of relative arousal-valence annotations for validation of audio models,” 2022.
- [45] H. T. P. Thao, G. Roig, and D. Herremans, “Emomv: Affective music-video correspondence learning datasets for classification and retrieval,” *Information Fusion*, vol. 91, pp. 64–79, 2023.
- [46] A. S. Sams and A. Zahra, “Multimodal music emotion recognition in indonesian songs based on cnn-lstm, xlnet transformers,” *Bulletin of Electrical Engineering and Informatics*, vol. 12, no. 1, pp. 355–364, 2023.
- [47] J. S. Gómez-Cañón, N. Gutiérrez-Páez, L. Porcaro, A. Porter, E. Cano, P. Herrera-Boyer, A. Gkiokas, P. Santos, D. Hernández-Leo, C. Karreman *et al.*, “Trompa-mer: an open dataset for personalized music emotion recognition,” *Journal of Intelligent Information Systems*, vol. 60, no. 2, pp. 549–570, 2023.



- [48] H. Strauss, J. Vigl, P.-O. Jacobsen, M. Bayer, F. Talamini, W. Vigl, E. Zangerle, and M. Zentner, “The emotion-to-music mapping atlas (emma): A systematically organized online database of emotionally evocative music excerpts,” *Behavior Research Methods*, pp. 1–18, 2024.
- [49] V. Grigorev, J. Li, W. Ma, Z. He, M. Zhang, Y. Liu, M. Yan, and J. Zhang, “Situnes: A situational music recommendation dataset with physiological and psychological signals,” in *Proceedings of the 2024 ACM SIGIR Conference on Human Information Interaction and Retrieval*, 2024, pp. 417–421.
- [50] K. Hevner, “Experimental studies of the elements of expression in music,” *The American journal of psychology*, vol. 48, no. 2, pp. 246–268, 1936.
- [51] J. A. Russell, “A circumplex model of affect,” *Journal of personality and social psychology*, vol. 39, no. 6, p. 1161, 1980.
- [52] I. Bakker, T. Van Der Voordt, P. Vink, and J. De Boon, “Pleasure, arousal, dominance: Mehrabian and russell revisited,” *Current Psychology*, vol. 33, pp. 405–421, 2014.
- [53] R. E. Thayer, *The biopsychology of mood and arousal*. Oxford University Press, 1990.
- [54] K. Trochidis, D. Sears, D.-L. Tr an, and S. McAdams, “Psychophysiological measures of emotional response to romantic orchestral music and their musical and acoustic correlates,” in *From Sounds to Music and Emotions: 9th International Symposium, CMMR 2012, London, UK, June 19-22, 2012, Revised Selected Papers 9*. Springer, 2013, pp. 44–57.
- [55] V. N. Salimpoor, M. Benovoy, K. Larcher, A. Dagher, and R. J. Zatorre, “Anatomically distinct dopamine release during anticipation and experience of peak emotion to music,” *Nature neuroscience*, vol. 14, no. 2, pp. 257–262, 2011.
- [56] J. Jaimovich, N. Coghlan, and R. B. Knapp, “Emotion in motion: A study of music and affective response,” in *From Sounds to Music and Emotions: 9th International Symposium, CMMR 2012, London, UK, June 19-22, 2012, Revised Selected Papers 9*. Springer, 2013, pp. 19–43.
- [57] R. A. McFarland, “Relationship of skin temperature changes to the emotions accompanying music,” *Biofeedback and Self-regulation*, vol. 10, pp. 255–267, 1985.
- [58] J. Kim and E. Andr e, “Emotion recognition based on physiological changes in music listening,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 30, no. 12, pp. 2067–2083, 2008.
- [59] Y. Song, S. Dixon, M. T. Pearce, and A. R. Halpern, “Perceived and induced emotion responses to popular music: Categorical and dimensional models,” *Music Perception: An Interdisciplinary Journal*, vol. 33, no. 4, pp. 472–492, 2016.
- [60] A. Go, R. Bhayani, and L. Huang, “Twitter sentiment classification using distant supervision,” *CS224N project report, Stanford*, vol. 1, no. 12, p. 2009, 2009.
- [61] I. Cohen, Y. Huang, J. Chen, J. Benesty, J. Benesty, J. Chen, Y. Huang, and I. Cohen, “Pearson correlation coefficient,” *Noise reduction in speech processing*, pp. 1–4, 2009.
- [62] K. W. Cheuk, Y.-J. Luo, B. Balamurali, G. Roig, and D. Herremans, “Regression-based music emotion prediction using triplet neural networks,” in *2020 international joint conference on neural networks (ijcnn)*. IEEE, 2020, pp. 1–7.
- [63] S. T. Rajamani, K. Rajamani, and B. Schuller, “Emotion and theme recognition in music using attention-based methods,” 2020.
- [64] M. Mayerl, M. V otter, A. Peintner, G. Specht, and E. Zangerle, “Recognizing song mood and theme: Clustering-based ensembles,” in *MediaEval*, 2021.
- [65] H. H. Tan, “Semi-supervised music emotion recognition using noisy student training and harmonic pitch class profiles,” *arXiv preprint arXiv:2112.00702*, 2021.
- [66] V. Bour, “Frequency dependent convolutions for music tagging,” in *MediaEval*, 2021.
- [67] Y. Feng, Y. Zhuang, and Y. Pan, “Music information retrieval by detecting mood via computational media aesthetics,” in *Proceedings IEEE/WIC international conference on web intelligence (WI 2003)*. IEEE, 2003, pp. 235–241.
- [68] L. Lu, D. Liu, and H.-J. Zhang, “Automatic mood detection and tracking of music audio signals,” *IEEE Transactions on audio, speech, and language processing*, vol. 14, no. 1, pp. 5–18, 2005.
- [69] J. Grekow, “Music emotion recognition using recurrent neural networks and pretrained models,” *Journal of Intelligent Information Systems*, vol. 57, no. 3, pp. 531–546, 2021.
- [70] S. A. Suresh Kumar and R. Rajan, “Transformer-based automatic music mood classification using multi-modal framework,” *Journal of Computer Science & Technology*, vol. 23, 2023.
- [71] J. Qiu, C. Chen, and T. Zhang, “A novel multi-task learning method for symbolic music emotion recognition,” *arXiv preprint arXiv:2201.05782*, 2022.
- [72] R. Sarkar, S. Choudhury, S. Dutta, A. Roy, and S. K. Saha, “Recognition of emotion in music based on deep convolutional neural network,” *Multimedia Tools and Applications*, vol. 79, no. 1, pp. 765–783, 2020.

- [73] T.-N. Do, M.-T. Nguyen, H.-D. Nguyen, M.-T. Tran, and X.-N. Cao, "Hcmus at mediaeval 2020: Emotion classification using wavenet feature with specaugment and efficientnet." in *MediaEval*, 2020.
- [74] H. T. P. Thao, B. Balamurali, D. Herremans, and G. Roig, "Attendaffectnet: Self-attention based networks for predicting affective responses from movies," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 8719–8726.
- [75] F. Eyben, M. Wöllmer, and B. Schuller, "Opensmile: the munich versatile and fast open-source audio feature extractor," in *Proceedings of the 18th ACM international conference on Multimedia*, 2010, pp. 1459–1462.
- [76] Y.-Y. Shi, X. Zhu, H.-G. Kim, and K.-W. Eom, "A tempo feature via modulation spectrum analysis and its application to music emotion classification," in *2006 IEEE International Conference on Multimedia and Expo*. IEEE, 2006, pp. 1085–1088.
- [77] T. Krols, Y. Nikolova, and N. Oldenburg, "Modality in music: Predicting emotion in music from high-level audio features and lyrics," *arXiv preprint arXiv:2302.13321*, 2023.
- [78] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [79] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [80] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, and T. Brox, "Flownet: Learning optical flow with convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2758–2766.
- [81] J. Melechovsky, Z. Guo, D. Ghosal, N. Majumder, D. Herremans, and S. Poria, "Mustango: Toward controllable text-to-music generation," *NAACL*, 2024.
- [82] Y. Gong, Y.-A. Chung, and J. Glass, "Ast: Audio spectrogram transformer," *arXiv preprint arXiv:2104.01778*, 2021.
- [83] R. Kumar, P. Seetharaman, A. Luebs, I. Kumar, and K. Kumar, "High-fidelity audio compression with improved rvqgan," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [84] C. F. Lima and S. L. Castro, "Emotion recognition in music changes across the adult life span," *Cognition and Emotion*, vol. 25, no. 4, pp. 585–598, 2011.
- [85] X. Wang, Y. Wei, L. Heng, and S. McAdams, "A cross-cultural analysis of the influence of timbre on affect perception in western classical music and chinese music traditions," *Frontiers in Psychology*, vol. 12, p. 732865, 2021.
- [86] H. Lee, F. Hoeger, M. Schoenwiesner, M. Park, and N. Jacoby, "Cross-cultural mood perception in pop songs and its alignment with mood detection algorithms," *arXiv preprint arXiv:2108.00768*, 2021.
- [87] J. M. Bland and D. G. Altman, "Statistics notes: Cronbach's alpha," *Bmj*, vol. 314, no. 7080, p. 572, 1997.
- [88] H. Song, M. Kim, D. Park, Y. Shin, and J.-G. Lee, "Learning from noisy labels with deep neural networks: A survey," *IEEE transactions on neural networks and learning systems*, 2022.
- [89] Z.-F. Wu, T. Wei, J. Jiang, C. Mao, M. Tang, and Y.-F. Li, "Ngc: A unified framework for learning with open-world noisy data," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 62–71.
- [90] H. Wei, L. Feng, X. Chen, and B. An, "Combating noisy labels by agreement: A joint training method with co-regularization," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 13 726–13 735.
- [91] T. Zhou, S. Wang, and J. Bilmes, "Robust curriculum learning: from clean label detection to noisy label self-correction," in *International Conference on Learning Representations*, 2020.
- [92] G. Stoet, "Psytoolkit: A novel web-based method for running online questionnaires and reaction-time experiments," *Teaching of Psychology*, vol. 44, no. 1, pp. 24–31, 2017.
- [93] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [94] G. Paltoglou and M. Thelwall, "Seeing stars of valence and arousal in blog posts," *IEEE Transactions on Affective Computing*, vol. 4, no. 1, pp. 116–123, 2012.
- [95] S. McAdams and B. L. Giordano, "The perception of musical timbre," 2014.
- [96] D. Ronan, J. D. Reiss, and H. Gunes, "An empirical approach to the relationship between emotion and music production quality," *arXiv preprint arXiv:1803.11154*, 2018.
- [97] P. N. Juslin and R. Timmers, "Expression and communication of emotion in music performance," *Handbook of music and emotion: Theory, research, applications*, pp. 453–489, 2010.

- [98] M. M. Farbood, "A quantitative, parametric model of musical tension," Ph.D. dissertation, Massachusetts Institute of Technology, 2006.
- [99] H. H. Tan and D. Herremans, "Music fadernets: Controllable music generation based on high-level features via low-level feature modelling," *Proc. of ISMIR*, 2020.
- [100] M. M. Rahman, A. K. Sarkar, M. A. Hossain, M. S. Hossain, M. R. Islam, M. B. Hossain, J. M. Quinn, and M. A. Moni, "Recognition of human emotions using eeg signals: A review," *Computers in Biology and Medicine*, vol. 136, p. 104696, 2021.