

# MUSAV: A DATASET OF RELATIVE AROUSAL-VALENCE ANNOTATIONS FOR VALIDATION OF AUDIO MODELS

Dmitry Bogdanov Xavier Lizarraga-Seijas Pablo Alonso-Jiménez Xavier Serra  
Music Technology Group, Universitat Pompeu Fabra, Spain

dmitry.bogdanov@upf.edu, xavier.lizarraga@upf.edu, pablo.alonso@upf.edu

## ABSTRACT

We present MusAV, a new public benchmark dataset for comparative validation of arousal and valence (AV) regression models for audio-based music emotion recognition. To gather the ground truth, we rely on relative judgments instead of absolute values to simplify the manual annotation process and improve its consistency. We build MusAV by gathering comparative annotations of arousal and valence on pairs of tracks, using track audio previews and metadata from the Spotify API. The resulting dataset contains 2,092 track previews covering 1,404 genres, with pairwise relative AV judgments by 20 annotators and various subsets of the ground truth based on different levels of annotation agreement. We demonstrate the use of the dataset in an example study evaluating nine models for AV regression that we train based on state-of-the-art audio embeddings and three existing datasets of absolute AV annotations. The results on MusAV offer a view of the performance of the models complementary to the metrics obtained during training and provide insights into the impact of the considered datasets and embeddings on the generalization abilities of the models.

## 1. INTRODUCTION

Audio-based music emotion recognition is a popular task in music information retrieval (MIR) that has recently gained more presence in the context of industrial applications. It is relevant for building systems for navigation of music collections, music search, exploration, and recommendation, and diverse applications that can benefit from MIR, such as audio branding or music therapy.

There are two types of approaches to emotion recognition in MIR following research in music psychology and affective computing [1–3]. The categorical approach considers different discrete categories of emotions (or moods<sup>1</sup>) or their clusters [4, 5] separately. It relies on

<sup>1</sup> Even though some researchers distinguish the terms “emotion” and “mood”, with moods being longer-term perceptions of musical input, we

taxonomies of descriptive mood tags and is frequently addressed by research on music classification and auto-tagging [6–9]. In contrast, the dimensional approach proposes representations on a continuous scale for several dimensions [10], allowing for a direct comparison of different moods, which is convenient for many applications.

The dimensional approach is based on existing research in music psychology which proposes two-dimensional or three-dimensional representations, including *arousal* (energy and stimulation), *valence* (pleasantness and positivity), and *dominance* (potency and control) [11] or, alternatively, *depth* [12] or *tension* [13], with many representations inheriting from the circumplex model of emotion by Russell [14]. In general, the 2D arousal/valence (AV) representation is a common model widely adopted in affective computing in different domains including music.

Various MIR researchers have worked on building datasets of AV annotations of music and training machine learning models for their automatic regression from audio [10, 15–21]. These datasets have been created with different methodologies, music collections, and participating annotators. However, there is no common benchmark dataset that could be conveniently used to compare models proposed by researchers and trained on different datasets. Existing studies report model performances using dataset splits without validation of the trained models on external datasets, which has been found to be very informative in other music classification tasks [22–24], providing insights on the generalization abilities and preventing overoptimistic conclusions.

In this paper, we propose to establish a common dataset for complementary evaluation on external data. We describe our methodology for building such a dataset, taking into account music genre diversity, and using it for evaluation of AV regression models. In contrast to many previous studies, we use comparisons between pairs of songs as ground truth instead of absolute values (coordinates) within the 2D AV space to make our annotations easier to gather and potentially more reliable, as suggested by previous works on relative emotion annotation in music and other domains [25–28]. In addition, we apply loudness normalization to avoid bias in arousal annotations [29], which has not been considered in previous datasets. The proposed validation of AV emotion recognition models provides a complementary view on their performance giving an opportunity to estimate generalization capabilities

will use both terms interchangeably in this paper for simplicity.



of the models on a common ground truth.

Following this methodology, we build a dataset based on audio previews and metadata available via Spotify API and evaluate a selection of AV regression models based on state-of-the-art audio embeddings. We analyze annotation agreement, propose strategies for building refined subsets of the dataset with different levels of consistency of annotations, and discuss the performance of the AV models.

## 2. RELATED WORK

Music emotion recognition is challenging because of the biases in cultural background, generation, genre, and personality [2, 30–33]. Nevertheless, this task gathered research attention in MIR from early on [34,35] given potential applications. Table 1 summarizes public AV datasets previously used in research. They all contain music audio excerpts and crowdsourced explicit arousal/valence annotations (absolute values that characterize each track or comparisons of track pairs) except for the MuSe dataset.

### 2.1 AV datasets with absolute values

There is a considerable variety of approaches to AV regression [10, 36–38] based in different audio features, including MediaEval campaigns in 2013-2015 [39–42]. We highlight the three most commonly used datasets containing absolute value annotations:

- **EmoMusic** [17] has been presented for the MediaEval 2013 Emotion in Music Task [39]. It contains 744 full audio tracks as well as 45-second excerpts. All audio is sourced from Free Music Archive (FMA). The excerpts are manually annotated with AV values characterizing the overall feel of the tracks as well as dynamic AV values at different rates, additionally summarized over the segment length (mean and stdev).
- **DEAM** [16] contains 1,802 audio excerpts (58 full-length songs and 1,744 excerpts of 45 seconds). The audio comes from several sources including FMA, Jamendo, and the MedleyDB dataset. The dataset similarly contains the overall AV values and dynamic values at a one per second rate and their summary (mean and stdev). This dataset has been derived from EmoMusic and used for the MediaEval 2013-2015 Emotion in Music Task [42] and in more recent studies [43, 44].
- **MuSe** [20] provides track-level valence, arousal and dominance values derived from social tags associated with music tracks on Last.fm<sup>2</sup> by using a dictionary of emotional ratings of words [45]. The dataset includes annotations for 90,408 songs, however the audio is not directly available. Instead, the tracks are identified by metadata, including Spotify IDs for 61,630 tracks. Nevertheless, only 41,021 30-second audio previews are currently accessible via Spotify API.<sup>3</sup>

Importantly, the EmoMusic and DEAM datasets are limited in coverage and they do not represent a large va-

Dataset	# tracks	Type	Source
EmoMusic [17]	744 ft/exc	abs	MTurk
DEAM [16]	1,802 ft/exc	abs	MTurk
MuSe [20]	41,021 exc	abs	Last.fm tags
MER-TAFFC [37]	900 exc	quad	manual
CCMED-WCMED [46]	800 exc	rel	CrowdFlower
EMusic [26]	140 exc	rel	CrowdFlower
MusAV	2,092 exc	rel	manual

**Table 1.** Public music datasets for AV regression and the proposed MusAV dataset. ft: full tracks, exc: excerpts, abs: ranged absolute values, quad: quadrants, rel: relative annotations.

riety of music available on commercial digital music platforms. The MuSe dataset has a significantly larger size and coverage, including 835 genres, achieved by sampling Last.fm using a diverse set of mood labels. Yet, its downside is that it is possibly noisy due to the tags-to-AV mapping. As a compromise, Panda et al. [37] propose to infer AV annotations from AllMusic<sup>4</sup> emotion tags, but they only use them to create a balanced annotation pool that is then manually validated. The resulting dataset (MER-TAFFC) contains 900 30-second track previews annotated by four AV emotion quadrants. In our work, we also follow an automated music preselection approach and prioritize large genre coverage while keeping the annotations manual. We can then compare AV regression models trained on EmoMusic, DEAM, and MuSe using our new dataset.

### 2.2 Relative annotations

Some researchers in affective computing highlighted the disadvantages of rating-based emotion annotation by absolute values and propose to use relative annotations [27,28]. In MIR this has been considered in few studies. Yang and Chen [25] propose to gather relative AV annotations and employ learning-to-rank algorithms to train models predicting absolute AV values. They discuss the limitations of absolute value rating-based annotations and show that relative annotations are significantly easier and have more within-subject and between-subject reliability. Their annotation experiment involved a corpus of 1,240 pop songs (30-second segments) and 99 annotators with an average of 4.3 annotators per song. However, they considered relative annotations only for valence and the audio for the dataset is not publicly available.

The idea of relative AV annotations has been further explored by Fan et al. for the case of experimental music [26] (the EMusic dataset) and classical music [46] (CCMED-WCMED). For the former dataset, they crowdsource pairwise track AV comparisons for 140 track segments from 9 genres by 823 annotators gathering up to three annotators per pair. The latter contains 800 track segments from Western and Chinese classical music with pairwise comparisons by 989 annotators. In addition, a similar study proposes relative ground truth for soundscape emotion recognition [47].

<sup>2</sup> <https://www.last.fm/>

<sup>3</sup> As of May 13, 2022.

<sup>4</sup> <https://allmusic.com>

### 3. ANNOTATION APPROACH

We follow a methodology that allows to avoid some of the limitations related to subjective annotation of arousal and valence with absolute values and instead consider comparative annotations on pairs of music tracks [25, 26]. Our main motivations are the following:

- Many practical applications are concerned not with the absolute AV values of songs, but with rankings songs according to these values. In such cases, relative annotations are convenient to validate ranking performance.
- For annotators, pairwise comparisons can be easier to understand and make decisions. They might require less effort than annotating with ranges of continuous values or Likert scales, which have more cognitive load [25, 48].
- There is evidence that relative annotations have higher between-subject and within-subject agreement [25, 27].
- There are simple strategies to refine annotations according to the agreement between different annotators.
- Model evaluation can be a simple comparison of the ground-truth ordering of two tracks with the ordering according to the predicted AV values for all track pairs.
- Depending on how the models are trained, different models might predict AV values within different value ranges, which should be accounted for when comparing the performance metrics such as RMSE. Relative comparisons on pairs of tracks allows using simple common metrics that are compatible with all models.

To collect the dataset ground truth, we need an annotation tool able to reproduce pairs of music excerpts and collect the user’s input. Such a tool should address potential sources of bias and simplify and speed up the annotation process. We define several requisites to accomplish this:

- We are interested in relative judgements about a pair of tracks (A and B) that can be formulated as the following question: “Which song has more *music property X*”. The following choices are considered: *A*, *B*, or *same*. To minimize potential biases toward any of the choices, none of them should be selected by default.
- It may be difficult to maintain consistency when answering non-factual questions. Therefore, we consider that the interface should display multiple pairs in the same *page* to give the user an opportunity to improve coherence of their annotations before submitting the page.
- Loudness has a large impact in music perception. Higher loudness correlates to higher perceived arousal [29]. To minimize this effect, the annotation tool should include loudness normalization. Previous studies that proposed interfaces for AV annotations [16, 17, 25, 26] ignored this issue and did not include any normalization.

### 4. THE MUSAV DATASET

We created our dataset following the described methodology, using Spotify API as a source for music audio previews and genre metadata. Using Spotify allows us to access a wide range of music, while the 30-second previews

it provides are sufficiently long to capture an overall perceived emotion with AV annotations.

#### 4.1 Preparing the annotation pool

We collected a list of 5,716 genres from *everynoise.com*<sup>5</sup> which corresponds to the genre taxonomy of the Spotify API<sup>6</sup> and contains broad genres as well as specific sub-genres. We then used the API’s *Search* method to select random tracks for each genre. We generated multiple queries for each genre (using the genre tag, a wildcard search string starting with a random character, and a random market) and picked a random track from the list of the returned results for each query. This method allowed us to diversify music coverage and avoid popularity bias, downloading 17,574 track previews for 4,386 genres (up to 15 tracks per genre). All audio previews are 30-second long MP3 files with a 96 kbit/s bitrate. Each preview has a corresponding metadata file obtained with the API’s *Get Track* method, including artist and album metadata and various audio analysis features.

We then analyzed the loudness of the audio previews to discard tracks with atypical levels, computing the integrated loudness in LUFS [49] of each track with the *Essentia* audio analysis library [50]. Based on the distribution of the obtained values, we kept tracks within the range between -20 and -5 LUFS, which represents the range of healthy loudness levels for the majority of mastered music. As a result, we reduced our pool to 15,979 tracks by 3,630 genres.

We organized the tracks into triplets with three pairwise comparisons each, allowing for additional inconsistency checks according to gathered relations within each triplet. We randomly assigned the tracks to two types of triplets: *genre-triplets* with all tracks sharing the same genre (one triplet per genre) and *global-triplets* containing tracks from various genres (the remaining tracks) to account for a use-case of distinguishing emotions within the same genre. All resulting 5,326 triplets contain unique tracks. We randomly split all generated triplets into annotation chunks, each one containing 100 triplets with 80% being global-triplets and 20% genre-triplets.

#### 4.2 Annotation tool and process

We implemented a custom tool according to the requirements in Section 3. For each pairwise comparison (a pair), we used the *wavsurfer-js*<sup>7</sup> player to display a navigable representation of the waveforms. To prevent loudness bias, we normalize the songs to a common level of -20 LUFS. We computed the normalization factors from the LUFS values precomputed in the dataset preparation step, and converted them to linear gain units as expected by *wavsurfer-js*.

Our interface formulates two questions: “Which song has more arousal?” and “Which song has more va-

<sup>5</sup> <https://everynoise.com>

<sup>6</sup> <https://developer.spotify.com/documentation/web-api/reference>

<sup>7</sup> <https://wavsurfer-js.org/>

lence?”. For both questions, the choices are *A*, *B*, or *same* arousal/valence. The tool shows a configurable number of pairs on each screen page (6 by default) and a submit button that stores the answers from the current page and renders the next one. We present the annotator with multiple pairs on the same page, as this can facilitate double-checking decisions made across pairs to minimize inconsistencies. Additionally, it simplifies navigation of the annotation interface, reducing the amount of necessary mouse movements and clicks. Finally, the tool has a page counter and displays each pair’s ID to facilitate reporting any possible issues. The annotator outputs one JSON file per pair containing the answers to the two questions.

Figure 1 depicts the annotation tool. The source code for the annotation tool is publicly available online.<sup>8</sup> The tool is distributed as a Docker web application.

Due to the limitations on effort and availability of our annotators, we have proceeded with 7 annotation chunks which account to 2,100 tracks assigned to 700 triplets with 2,100 pairwise track comparisons. Overall, we gathered annotations from 20 participants, including authors’ colleagues and students, with a background in music and technology. Each chunk was presented to three different annotators. Every annotator was given a single chunk, with an exception of one annotator who worked with two chunks. All annotators were instructed about the meaning of arousal and valence beforehand following their common definition [16, 42] and were asked to focus on perceived emotion [51]. Participants were aware of the subjectivity of the task and we encouraged to provide their subjective opinion. In total, we gathered 6,255 comparative arousal and valence judgments on pairs of tracks after discarding 15 pairs that the annotators reported having non-music tracks (speech) and duplicated tracks. These annotations involve 2,092 track previews by 1,404 genres.

### 4.3 Annotation agreement and consistency

By having multiple people annotate the same chunks of audio, we can measure the agreement between annotators. Computing ordinal Krippendorff’s alpha, we obtained values of 0.48 for arousal and 0.39 for valence, which indicates a fair to moderate level of agreement, which is consistent with previous studies [26, 42, 46].

For building our ground truth for arousal and valence, we defined two types of agreement for pairwise comparisons of tracks by three different annotators. If all three annotators agreed on a pair of tracks with the same answer (*A-A-A*, *B-B-B*, or *same-same-same*) the annotations for this pair were considered to be in *full agreement*. If only two annotators agreed, we checked whether the third annotator was in a soft (e.g., *A-A-same*, *same-same-A*) or hard (*A-A-B* or *B-B-A*) disagreement. In the case of the former, we considered the annotations for the pair to be in *majority agreement*. Table 2 presents the agreement statistics for all gathered annotations.

In addition, we checked whether pairwise comparisons contradict each other within triplets (that is, whether they

Agreement	Arousal		Valence	
	# pairs	%	# pairs	%
FA+MA	1,448	69.4	1,341	64.3
FA	975	46.8	810	38.8
FA+MA, CT	738	35.4	606	29.1
FA, CT	519	24.9	381	18.3

**Table 2.** Number and percentage of annotated track pairs with different levels of annotator agreement and consistency for arousal and valence. FA+MA: pairs with full or majority agreement. FA: pairs with full agreement. CT: only pairs belonging to consistent triplets.

are geometrically inconsistent). For example, for three tracks *X*, *Y*, and *Z* forming a triplet, if  $X > Y$  and  $Y > Z$ , but  $X \leq Z$ , such triplet and all its pairs are considered inconsistent. We considered triplets as consistent only if all constituent pairs had full or majority agreement in the annotations and no contradictions have been found.

As a result, we generated different subsets of the annotations, with 69.4% and 64.3% of track pairs having at least some level of agreement and 24.9% and 18.3% passing the most strict conditions (pairs with full agreement, belonging to consistent triplets) for arousal and valence, accordingly.

Finally, as we had two types of triplets, we checked the effect of genre on the agreement rate: 67% and 61% of pairs in global-triplets had either full or majority agreement compared to 76% and 75% in the case of genre-triplets for arousal and valence, accordingly. This observation revealed that it was slightly easier to reach agreement on pairs of tracks coming from the same genre than from different genres.

### 4.4 Dataset contents

We provide the following contents as part of the dataset, available online:<sup>9</sup>

- Metadata for the entire annotation pool.<sup>10</sup> Each triplet is identified by a triplet ID and contains track Spotify IDs, triplet type (global-triplet or genre-triplet) and genre information.
- Split of the annotation pool into annotation chunks.<sup>11</sup>
- Raw comparative arousal and valence annotations on track pairs by anonymized annotators.
- Processed ground-truth annotations with different levels of agreement and consistency (full and major agreement with/without triplet consistency).
- Track audio previews and metadata gathered from the Spotify API for the annotated chunks.<sup>12</sup>
- Dataset metadata statistics (e.g., genre distribution).
- Scripts to reproduce the creation of the dataset.

<sup>9</sup> <https://mtg.github.io/musav-dataset>

<sup>10</sup> All annotation metadata is licensed under CC BY-NC-SA 4.0.

<sup>11</sup> It is possible to expand the dataset by annotating more chunks.

<sup>12</sup> Available under request for non-commercial scientific research purposes only. Any publication of results based on this data must cite Spotify API as the source of the data.

<sup>8</sup> <https://github.com/MTG/musav-annotator>

### Task: arousal\_and\_valence

You are on page #1/75

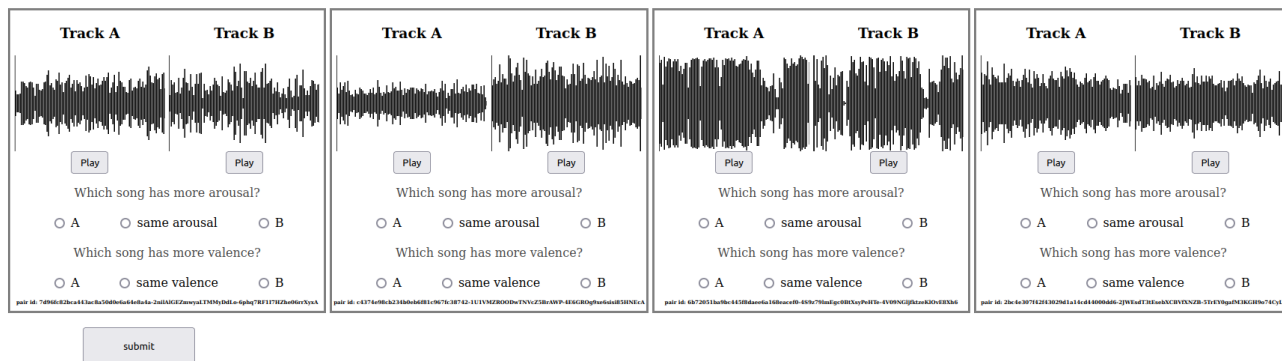


Figure 1. Screenshot of the annotation tool.

## 5. EXPERIMENTS

We demonstrate the dataset in use on the example of evaluating AV regression models based on audio embeddings.

### 5.1 Models

We created AV regression models based on three types of audio embeddings:

- **MusiCNN-MSD** (*musicnn*) is a music auto-tagging CNN with filter shapes motivated by music domain [52] trained on a subset of the Million Song Dataset. Its embedding layer has 200 units.
- **VGGish** (*vggish*) is a VGG architecture with an embedding layer of 128 units trained on audio from YouTube videos mapped to general purpose audio labels derived from their metadata [53]. The embeddings from this model were previously used for AV regression in combination with support vector regression [46].
- **EffNet-Discogs** (*effnet*) is an EfficientNet architecture trained to predict the music styles tags from Discogs.<sup>13</sup> The model produces 1280-dimensional embeddings and it is publicly available as part of Essentia models [54].<sup>14</sup>

The embeddings are extracted on short one-, two-, and three-second audio excerpts for *vggish*, *effnet*, and *musicnn* models, accordingly. They are then used by the downstream regression models that we train. These models provide arousal-valence inference for each embedding vector, with a variable batch size with batch normalization. As Figure 2 depicts, we use a fully connected layer with a linear activation function, preceded by batch normalization and dropout. We also apply L1-L2 and L2 regularizers in the fully connected layer and dropout as regularization methods.

For training, we used three different datasets: DEAM, EmoMusic and MuSe, which we selected based on their music coverage as more appropriate for general use, with the goal to incorporate the resulting models as part of Es-



Figure 2. Arousal-valence backend model architecture.

sentia [24]. Each of them provides different audio excerpts (with 30-45 second duration) and arousal-valence values characterizing the overall emotion of each track. We extracted the embeddings with the pretrained models and used them as features. We followed a standard data splitting and loading strategy used in previous music classification publications [24]. To generate a train/test split stratified in terms of AV quadrants, we use Z-score normalization. However, we did not normalize the training data. The AV value range in all the datasets is from 1 to 9.

Our models operate on short audio chunks and allow us to generate sequences of AV predictions over time with a new prediction every 1-3 seconds, depending on the receptive field of the embedding model used. Therefore, we compute the average of predictions on chunks to estimate the overall arousal/valence of a track.

In Table 3, we report Root Mean Square Error (RMSE) and Coefficient of Determination  $R^2$  ( $R^2$ ) commonly used for evaluation of regression models [17, 55, 56], obtained

	Arousal		Valence	
	$R^2$	RMSE	$R^2$	RMSE
deam-effnet	0.404	0.913	0.335	0.909
deam-musicnn	<b>0.417</b>	<b>0.894</b>	<b>0.400</b>	<b>0.818</b>
deam-vggish	0.396	0.963	0.344	0.963
emomusic-effnet	0.420	1.090	0.375	<b>0.948</b>
emomusic-musicnn	<b>0.451</b>	<b>1.030</b>	0.363	0.966
emomusic-vggish	0.429	1.037	<b>0.376</b>	0.973
muse-effnet	<b>0.143</b>	<b>1.148</b>	<b>0.089</b>	<b>1.581</b>
muse-musicnn	0.141	1.320	0.085	2.509
muse-vggish	<b>0.143</b>	<b>1.148</b>	0.085	1.584

Table 3. Evaluation metrics for the AV regression models on the held-out (testing) sets of the corresponding training datasets. The best values for each dataset are in bold.

<sup>13</sup> <https://blog.discogs.com/en/genres-and-styles>

<sup>14</sup> <https://essentia.upf.edu/models.html>

# track pairs	Arousal				Valence			
	FA+MA	FA	FA+MA, CT	FA, CT	FA+MA	FA	FA+MA, CT	FA, CT
	1413	950	716	502	1310	787	588	368
deam-effnet	72.28	75.44	72.60	74.84	61.59	63.38	63.91	65.51
deam-musicnn	78.81	81.04	76.92	78.41	59.75	61.98	62.33	62.90
deam-vggish	78.40	82.14	79.33	81.55	62.32	64.86	66.47	67.83
emomusic-effnet	82.57	86.55	84.75	87.61	71.29	75.41	73.77	78.55
emomusic-musicnn	85.61	89.21	84.78	87.63	<b>74.80</b>	<b>78.76</b>	76.53	80.29
emomusic-vggish	<b>86.42</b>	<b>90.30</b>	<b>86.86</b>	<b>89.73</b>	70.81	77.03	74.51	<b>81.16</b>
muse-effnet	59.92	60.99	59.00	62.11	62.14	63.78	61.54	64.35
muse-musicnn	63.96	66.59	64.84	68.55	67.72	70.77	69.03	71.01
muse-vggish	66.34	69.03	64.63	68.00	62.27	66.35	62.50	68.22
Spotify API	83.31	86.67	83.17	85.95	73.44	74.59	<b>77.51</b>	77.68

**Table 4.** External validation results on the proposed MusAV dataset (the percentage of track pairs with a correctly predicted ordering). FA+MA: pairs with full or majority agreement. FA: pairs with full agreement. CT: only pairs belonging to consistent triplets. The highest values are marked in bold. The top three AV models are marked in gray.

on the datasets used for training the models.

### 5.2 External validation on MusAV

We used our new dataset to validate the performance of the models. To this end, we assessed whether the ground truth ordering of track pairs coincided with the ordering according to the AV values predicted by the models. In addition, we also evaluated arousal and valence estimations provided by Spotify API and computed from audio as an additional reference.<sup>15</sup> This reference possibly represents a common state of the art in industrial systems. We ensured that our external validation set is independent of the datasets used for training the models: EmoMusic and DEAM contain non-commercial music unavailable on Spotify, while the intersection with MuSe, for which we also used Spotify track previews, includes 24 tracks that we filtered out for our evaluation.

For simplicity, we discarded all ground-truth annotations marking two songs as equivalent (13% and 15% of the ground-truth pairs with full or majority agreement in the case of arousal and valence, respectively). Thus, we focused only on examples with clear difference in arousal or valence. Table 4 presents the accuracy of the models in terms of the percentage of track pairs with correct ordering. We report the results on different subsets of the AV ground truth to demonstrate various evaluation possibilities.

### 5.3 Discussion

Having a new common ground truth for all models, our external validation shows the impact of the training dataset and embeddings.

Remarkably, models trained on the EmoMusic dataset perform the best for both arousal and valence regression. This is surprising, given that this dataset is smaller and less diverse than DEAM, which was derived from EmoMusic. On the other side, models based on MuSe have the worst performance in the case of arousal. Even though MuSe is the largest dataset in terms of size and coverage of commercially-available music, it appears to be too noisy

<sup>15</sup> We consider the "energy" descriptor in the Spotify API as arousal.

to be able to train efficient models for arousal. Notably, it is the only dataset out of three relying on user-generated tags instead of explicit AV annotations, and the employed process for mapping tags to AV values might be inherently noisier.

Second, given a dataset, the choice of embedding model also matters. For example, the *effnet* embeddings trained on a large music style dataset appear to be inefficient for emotion recognition. The models based on them are consistently worst in the case of arousal (with all three datasets used for training) and valence (with EmoMusic and MuSe used for training). In turn, our validation reveals high performance of the *vggish* embeddings in many cases, possibly due to their generalization ability which was previously evidenced in literature [24]. This observation contradicts the results obtained in the respective held-out sets, where it did not have a remarkable performance overall.

Finally, in our validation, some of the considered AV models have performance competitive with an industrial reference. Still, all of the considered models only achieve up to 90% accuracy for arousal and 81% for valence. This is in line with evidence that predicting valence (as well as its annotation) is generally considered more complex than arousal [17, 42].

## 6. CONCLUSIONS

We present a new public dataset of relative AV annotations for validation of audio-based AV models. To build it, we employ a methodology that maximizes coverage in terms of genres to gather our annotation pool and allows to assess consistency of annotations on triplets of songs. The dataset is based on audio previews from Spotify API which allows validating performance on diverse types of commercially-available music. As an example, we train and evaluate AV regression models based on three common AV datasets and three types of pretrained audio embeddings and show how such a benchmarking can provide valuable complementary information about model performances. The resulting pretrained models are publicly available as part of Essentia models.

## 7. ACKNOWLEDGEMENTS

This research was carried out under the project Musical AI - PID2019-111403GB-I00/AEI/10.13039/501100011033, funded by the Spanish Ministerio de Ciencia e Innovación and the Agencia Estatal de Investigación. We also thank Juan Sebastián Gómez Cañón for his suggestions and all participating annotators.

## 8. REFERENCES

- [1] 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 *International Society for Music Information Retrieval Conference (ISMIR 2010)*, 2010.
- [2] Y.-H. Yang and H. H. Chen, *Music emotion recognition*. CRC Press, 2011.
- [3] J. Grekow, *From content-based music emotion recognition to emotion maps of musical pieces*. Springer International Publishing, 2018.
- [4] X. Hu, M. Bay, and J. S. Downie, "Creating a simplified music mood classification ground-truth set," in *International Conference on Music Information Retrieval (ISMIR 2007)*, 2007.
- [5] C. Laurier, M. Sordo, J. Serra, and P. Herrera, "Music mood representations from social tags," in *International Society for Music Information Retrieval Conference (ISMIR 2009)*, 2009.
- [6] X. Downie, C. Laurier, and M. Ehmann, "The 2007 MIREX audio mood classification task: Lessons learned," in *International Symposium on Music Information Retrieval (ISMIR 2008)*, 2008.
- [7] C. Laurier and P. Herrera, "Audio music mood classification using support vector machine," in *International Conference on Music Information Retrieval (ISMIR 2007)*, 2007.
- [8] K. Bischoff, C. S. Firan, R. Paiu, W. Nejdl, C. Laurier, and M. Sordo, "Music mood and theme classification—a hybrid approach," in *International Society for Music Information Retrieval Conference (ISMIR 2009)*, 2009.
- [9] D. Bogdanov, A. Porter, P. Tovstogan, and M. Won, "MediaEval 2019: Emotion and theme recognition in music using Jamendo," in *MediaEval 2019 Workshop*, 2019.
- [10] 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.
- [11] J. A. Russell and A. Mehrabian, "Evidence for a three-factor theory of emotions," *Journal of Research in Personality*, vol. 11, no. 3, pp. 273–294, 1977.
- [12] D. M. Greenberg, M. Kosinski, D. J. Stillwell, B. L. Monteiro, D. J. Levitin, and P. J. Rentfrow, "The song is you: Preferences for musical attribute dimensions reflect personality," *Social Psychological and Personality Science*, vol. 7, no. 6, pp. 597–605, 2016.
- [13] 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.
- [14] J. A. Russell, "A circumplex model of affect," *Journal of Personality and Social Psychology*, vol. 39, no. 6, pp. 1161–1178, 1980.
- [15] T. Eerola, O. Lartillot, and P. Toivainen, "Prediction of multidimensional emotional ratings in music from audio using multivariate regression models," in *International Society for Music Information Retrieval Conference (ISMIR 2009)*, 2009.
- [16] M. Soleymani, A. Aljanaki, and Y. Yang, "DEAM: MediaEval database for emotional analysis in music," 2016. [Online]. Available: <https://cvml.unige.ch/databases/DEAM/manual.pdf>
- [17] M. Soleymani, M. N. Caro, E. M. Schmidt, C.-Y. Sha, and Y.-H. Yang, "1000 songs for emotional analysis of music," in *ACM International Workshop on Crowdsourcing for Multimedia (CrowdMM 2013)*, 2013.
- [18] J. Grekow, "Music emotion maps in arousal-valence space," in *IFIP International Conference on Computer Information Systems and Industrial Management (CISIM 2016)*, 2016.
- [19] J. Bai, J. Peng, J. Shi, D. Tang, Y. Wu, J. Li, and K. Luo, "Dimensional music emotion recognition by valence-arousal regression," in *IEEE International Conference on Cognitive Informatics & Cognitive Computing (ICCI\*CC 2016)*, 2016.
- [20] C. Akiki and M. Burghardt, "MuSe: The musical sentiment dataset," *Journal of Open Humanities Data*, vol. 7, no. 10, 2021.
- [21] K. W. Cheuk, Y.-J. Luo, B. Balamurali, G. Roig, and D. Herremans, "Regression-based music emotion prediction using triplet neural networks," in *International joint conference on neural networks (IJCNN 2020)*, 2020.
- [22] A. Livshin, "The importance of cross database evaluation in musical instrument sound classification," in *International Conference on Music Information Retrieval (ISMIR 2004)*, 2003.
- [23] D. Bogdanov, A. Porter, P. Herrera, and X. Serra, "Cross-collection evaluation for music classification tasks," in *International Society for Music Information Retrieval Conference (ISMIR 2016)*, 2016.

- [24] P. Alonso-Jiménez, D. Bogdanov, J. Pons, and X. Serra, “TensorFlow audio models in Essentia,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2020)*, 2020.
- [25] Yi-Hsuan Yang and H. H. Chen, “Ranking-based emotion recognition for music organization and retrieval,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 762–774, 2011.
- [26] J. Fan, K. Tatar, M. Thorogood, and P. Pasquier, “Ranking-based emotion recognition for experimental music,” in *International Society for Music Information Retrieval Conference (ISMIR 2017)*, 2017.
- [27] A. Metallinou and S. Narayanan, “Annotation and processing of continuous emotional attributes: Challenges and opportunities,” in *IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG 2013)*, 2013.
- [28] G. N. Yannakakis and H. P. Martinez, “Grounding truth via ordinal annotation,” in *IEEE International Conference on Affective Computing and Intelligent Interaction (ACII 2015)*, 2015.
- [29] R. T. Dean, F. Bailes, and E. Schubert, “Acoustic intensity causes perceived changes in arousal levels in music: An experimental investigation,” *PLOS One*, 2011.
- [30] P. Kuppens, F. Tuerlinckx, J. A. Russell, and L. F. Barrett, “The relation between valence and arousal in subjective experience,” *Psychological Bulletin*, vol. 139, no. 4, pp. 917–940, 2013.
- [31] P. Kuppens, F. Tuerlinckx, M. Yik, P. Koval, J. Coosemans, K. J. Zeng, and J. A. Russell, “The relation between valence and arousal in subjective experience varies with personality and culture,” *Journal of Personality*, vol. 85, no. 4, pp. 530–542, 2017.
- [32] J. S. Gómez-Cañón, E. Cano, T. Eerola, P. Herrera, X. Hu, Y.-H. Yang, and E. Gómez, “Music emotion recognition: Toward new, robust standards in personalized and context-sensitive applications,” *IEEE Signal Processing Magazine*, vol. 38, no. 6, pp. 106–114, 2021.
- [33] J. Gómez-Cañón, E. Cano, Y. Yang, P. Herrera, and E. Gómez, “Let’s agree to disagree: Consensus entropy active learning for personalized music emotion recognition,” in *International Society for Music Information Retrieval Conference (ISMIR 2021)*, 2021.
- [34] T. Li and M. Ogihara, “Detecting emotion in music,” in *International Conference on Music Information Retrieval (ISMIR 2003)*, 2003.
- [35] 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.
- [36] F. Weninger, F. Eyben, and B. Schuller, “On-line continuous-time music mood regression with deep recurrent neural networks,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2014)*, 2014.
- [37] R. Panda, R. Malheiro, and R. P. Paiva, “Novel audio features for music emotion recognition,” *IEEE Transactions on Affective Computing*, vol. 11, no. 4, pp. 614–626.
- [38] B. Bhattarai and J. Lee, “Automatic music mood detection using transfer learning and multilayer perceptron,” *International Journal of Fuzzy Logic and Intelligent Systems*, vol. 19, no. 2, pp. 88–96, 2019.
- [39] M. Soleymani, M. N. Caro, E. M. Schmidt, and Y.-H. Yang, “The MediaEval 2013 brave new task: Emotion in music,” in *MediaEval 2013 Workshop*, 2013.
- [40] A. Aljanaki, Y.-H. Yang, and M. Soleymani, “Emotion in Music Task at MediaEval 2014,” in *MediaEval 2014 Workshop*, 2014.
- [41] —, “Emotion in Music Task at MediaEval 2015,” in *MediaEval 2015 Workshop*, 2015.
- [42] —, “Developing a benchmark for emotional analysis of music,” *PLOS One*, vol. 12, no. 3, p. e0173392, 2017.
- [43] S. Zhao, Y. Li, X. Yao, W. Nie, P. Xu, J. Yang, and K. Keutzer, “Emotion-based end-to-end matching between image and music in valence-arousal space,” in *ACM International Conference on Multimedia (MM 2020)*, 2020.
- [44] H. Liu, Y. Fang, and Q. Huang, “Music emotion recognition using a variant of recurrent neural network,” in *International Conference on Mathematics, Modeling, Simulation and Statistics Application (MMSSA 2018)*, 2019.
- [45] A. B. Warriner, V. Kuperman, and M. Brysbaert, “Norms of valence, arousal, and dominance for 13,915 English lemmas,” *Behavior Research Methods*, vol. 45, no. 4, pp. 1191–1207, 2013.
- [46] J. Fan, Y.-H. Yang, K. Dong, and P. Pasquier, “A comparative study of western and chinese classical music based on soundscape models,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2020)*, 2020.
- [47] J. Fan, M. Thorogood, and P. Pasquier, “Emo-soundscapes: A dataset for soundscape emotion recognition,” in *IEEE International Conference on Affective Computing and Intelligent Interaction (ACII 2017)*, 2017.
- [48] D. Yang and W.-S. Lee, “Disambiguating music emotion using software agents,” in *International Conference on Music Information Retrieval (ISMIR 2004)*, 2004.



- [49] European Broadcasting Union (EBU), “EBU Tech 3341. Loudness metering: ‘EBU mode’ metering to supplement loudness normalisation in accordance with EBU R 128.”
- [50] D. Bogdanov, N. Wack, E. Gómez, S. Gulati, P. Herrera, O. Mayor, G. Roma, J. Salamon, J. Zapata, and X. Serra, “ESSENTIA: An audio analysis library for music information retrieval,” in *International Society for Music Information Retrieval Conference (ISMIR 2013)*, 2013.
- [51] N. F. Gutiérrez Páez, J. S. Gómez-Cañón, L. Porcaro, P. Santos, D. Hernández-Leo, and E. Gómez, “Emotion annotation of music: A citizen science approach,” in *International Conference on Collaboration Technologies and Social Computing (CollabTech 2021)*, 2021.
- [52] J. Pons and X. Serra, “musicnn: Pre-trained convolutional neural networks for music audio tagging,” in *International Society for Music Information Retrieval Conference (ISMIR 2019) Late Breaking Demo*, 2019.
- [53] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold *et al.*, “CNN architectures for large-scale audio classification,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2017)*, 2017.
- [54] P. Alonso-Jiménez, X. Serra, and D. Bogdanov, “Music representation learning based on editorial metadata from Discogs,” in *International Society for Music Information Retrieval (ISMIR 2022)*, 2022.
- [55] C. J. Willmott, “On the validation of models,” *Physical geography*, vol. 2, no. 2, pp. 184–194, 1981.
- [56] D. Chicco, M. J. Warrens, and G. Jurman, “The coefficient of determination r-squared is more informative than smape, mae, mape, mse and rmse in regression analysis evaluation,” *PeerJ Computer Science*, vol. 7, p. e623, 2021.