

3. METHODS

Next, we describe the key components of our model, starting with the dataset collection, the network architecture of the form analyzer, as well as a peak-picking scheme needed as input for the phrase analyzer.

3.1 Data Collection

We constructed a dataset — the **SMFSA Database** — consisting of 200 pieces of classical music in MIDI format² (for ease of score conversion, error correction, and signal processing). For each piece, we have an accompanying text document containing the form classification and part/phrase labels with their respective timestamps (obtained from the sheet music analysis performed by a human annotator),³ written in a format similar to the SALAMI dataset. To represent the audio signals, the Mel Spectrogram was first generated by resampling the audio with a sample rate of 44.1 kHz, then computing the Short-Time Fourier Transform (STFT) over the entire signal with a hop length of 6144 (0.139 seconds per frame, see [28]) and 8192 samples per frame (a window size of 0.209 seconds per frame multiplied by the sampling rate). Further, to extract meaningful representations from the raw data, we obtain the 2D Self-Similarity Matrix (SSM) [25] of the Mel Spectrogram, constructed from various attributes and similarity metrics as well as the duration of the piece to create the set of features. The final tabulated dataset contains the mean and variance arrays for the following features: Mel Spectrogram SSM, Mel-Frequency Cepstral Coefficient (MFCC) Spectrogram SSLMs (Euclidean and Cosine distances), Chromagram SSLMs (Euclidean and Cosine distances) [24], as well as 15 other features. Since only 200 data samples may not be enough for many training tasks, we extend our dataset to 1200 by augmenting the dataset using speed shifting, pitch shifting, time shifting, and noise injection [12, pp. 41-46]. The features were extracted similarly for the augmented dataset as well.

3.2 Form Analyzer Architecture

Our goal is to identify a piece of classical music as one of the following 12 forms: Arch, Bar, Binary, Minuet & Trio, Ritornello, Rondo, Sonata (-allegro), Ternary, Theme & Variation, Through Composed, Unary/Strophic, and Unique. Note that not all of the forms are equally distributed in the dataset (reflecting real-world musical pieces) and this leads to a class-imbalanced multi-class classification problem (see Supplement Figure 3 for the relative proportion of each form in the dataset). To address this problem, we adopt a framework that implements an ensemble of decision trees using a neural network, which has been shown to work well for class-imbalanced multi-class datasets [29, 30]. We note that other alternatives such as

² While VSTs may differ across sequencing software, the Form Analyzer and Peak-Picking algorithm are intended to be instrument-neutral methods. As such, the timbral difference is negligible — including running the algorithm on different rendered sequences of the same pieces; we only seek to discover how the pitches are distributed structurally and temporally.

³ We simply follow the human annotations from the original analyzed sheet music and copy these exactly at the time of their occurrence (in seconds, i.e., Part *B* with phrase *d* occurring at 15.27 seconds would be written as "15.270 B, d"), of course omitting the implicit labels as well.

focal loss [31] or imposing fairness in terms of model performance across each pair of classes can also be used [32].

In our implementation, we use the TreeGrad approach by [33], which implements Gradient Boosted Decision Trees as Neural Networks and has been shown to have a small computational footprint. TreeGrad is an extension of Deep Neural Decision Forests (DNDF), which treats the node split structure of a decision tree as a neural network architecture search problem. Similar to DNDF, they have three layers: a *decision node* layer, a *routing* layer, and a *prediction* layer [34]. The decision layer consists of decision stumps, each of which computes the probability of routing the data node x to (left/positive or right/negative) child nodes for a node n and is formulated as (shown here only for the positive route) $d_n^+(x; \theta_+) = \sigma(\theta_+^T x + b)$, where θ_+ are the learnable parameters and σ is a temperature-controlled softmax function. These routing probabilities (both positive and negative) of all nodes are then concatenated to yield $D(x; \tilde{\theta})$, which is the output of the first layer. The parameter vector in $\tilde{\theta}$ forms the linear decision boundary that results in how each node is routed. This is followed by the routing layer, which computes the probability that each node uses to route a data point to a particular leaf. For this step, we use a binary preconstructed matrix Q of size $l \times 2n$, where l is the number of leaves, and n is the number of internal nodes. This allows for enumerating the relationship between nodes and leaves. This layer computes μ_l which indicates the probability of reaching a leaf l and can be written as

$$\begin{aligned} \mu_l(x|\theta) &= \prod_{j=1}^n (D(x; \tilde{\theta}_j) \odot Q_l + (1 - Q_l)) \quad (1) \\ &= \exp(Q_l^T \log(D(x; \tilde{\theta}))) \end{aligned}$$

where Q_l is the l -th row of Q and \odot is the Hadamard Product [33]. We then pass it through a softmax activation function to produce the output of the second layer. The final output layer is the leaf layer which is a fully connected layer. The activation function used here is $\exp(x)$. See Figure 2 for an illustration of the network architecture.

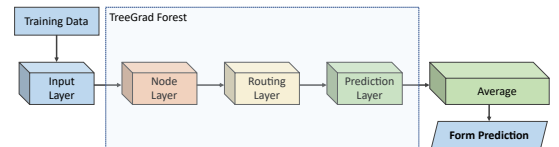


Figure 2. Architecture diagram of Form Analyzer

3.3 Phrase Analyzer Architecture

Recall that phrase analysis is a significantly more challenging task than form. To be able to learn phrase analysis similar to a musician, the model must be able to differentiate between musical events that can be labeled as having only the part label (*CODA*, *Episode 1*, *Middle Entry*, etc.), or only the phrase label (*a*, *at*, *b*, *transition*, *melodic (variation)*, etc.) or both part and phrase labels. This presents several problems — the first deals with making the subtle distinction of when to assign both labels or just one, as well as choosing a suitable metric that captures the similarities correctly. Moreover, labels for varied phrases are also extremely difficult to grade, as a phrase (or part) may be repeated with a different harmonic goal (phrase *a* to *a'*

or Parts A and A'). Variations of a can be either the same variation a' or a new variation a'' or a''' (these varied repetitions sometimes continue beyond 3, such as a^4 and so on). Since scoring such variations can be different based on the annotator, for experimentation on our dataset, we simply reduced the label set to remove prime marks and retain the phrase/part letter(s). Likewise, large and/or hybrid forms (such as Sonata-allegro form) in our dataset which have their own unique labels are simplified (during experimentation) to the letter label set with parts A, B, A' .

The task of Phrase Analysis requires labels to be provided with the form of the piece as well as the knowledge of where the musical phrase starts — a problem potentially solved using Onset or Peak Detection [25]. To do this, we implement a simple algorithm for peak-picking based on existing literature which is described next. The timestamps of the peaks, along with the output of the Form Analyzer (which is added as a feature), are provided as input to the Phrase Analyzer, where labels are classified sequentially (shown in Figure 4). We describe its components next.

3.3.1 Onset Detection – Peak-Picking for Phrase Events

To find the timestamps (in seconds) of the musical phrases, we use a reduced version of the peak-picking algorithm described in [35]. This algorithm computes the Mel Spectrogram and Self-Similarity Lag Matrix (SSLM) Chromagram (a group of pitch class distributions over time) to perform the onset detection. The Chromagram SSLM is computed using the Mel Spectrogram, which is clustered by pitch-class using k-Nearest Neighbors. The audio frames captured by the Short-Time Fourier Transform (STFT) are represented as a computed vector of peaks, which is returned as an array of timestamps. The choice of this scheme was based on preference for ideas that are common in the community, and we verified that it was comparable to other CNN architectures [36] and also greatly reduced system design time, given that training was not needed (since the approach was unsupervised).

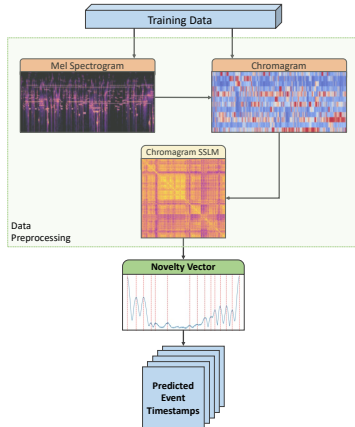


Figure 3. Schematic diagram of Peak-Picking algorithm

3.3.2 Bidirectional LSTMs for Phrase Classification

Using the onset detection method, the musical piece is essentially segmented into smaller phrases (i.e., the musical content between timestamps t and $t + 1$ is denoted as feature vector x_t). Therefore, we think of a piece X as a collection of phrases $\{x_1, \dots, x_T\}$, which occur in sequence.

The task is to tag each phrase with a label for the phrase and/or the part. In all, we use 9 phrases, 9 parts, and 5 variations (phrase labels used for Theme & Variation pieces) for the labeling; specific names of these are provided in the [Supplement](#). Clearly, the specific labeling associated with a phrase is dependent on what came before it, as well as what comes after. Therefore, we propose to use Bidirectional LSTMs to capture the dependence in such sequence data and replicate segment-segment similarity analysis.

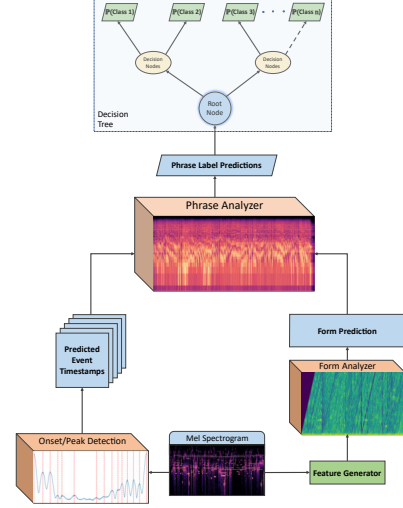


Figure 4. Architecture diagram of Phrase Analyzer

We now discuss the network architecture. X is passed as input to a Bidirectional Long Short-Term Memory (LSTM) [37] network. LSTM networks can capture long-term dependencies in temporal data and have been successfully used for a number of time series classification problems. LSTM (similar to Recurrent Neural Networks (RNNs) [38]) contains loops in its architecture that allow it to memorize previous states such that the network can effectively process temporal data. A typical LSTM layer consists of a set of recurrently connected blocks, known as memory blocks. Each block contains one or more recurrently connected memory cell (c^t) and three multiplicative units - the input (i^t), output (o^t), and forget gates (f^t) which regulate the extent to which data is propagated through the LSTM unit. The operations inside an LSTM block can be formulated using:

$$\begin{aligned} f_t &= \rho(W_f x_t + U_f h_{t-1} + b_f) \\ i_t &= \rho(W_i x_t + U_i h_{t-1} + b_i) \\ o_t &= \rho(W_o x_t + U_o h_{t-1} + b_o) \\ \hat{c}_t &= \phi(W_c x_t + U_c h_{t-1} + b_f) \\ c_t &= f_t \circ c_{t-1} + i_t \circ \hat{c}_t \\ h_t &= o_t \circ \phi(c_t) \end{aligned} \quad (2)$$

Here ρ and ϕ are activation functions, \circ denotes the element-wise product operation, x_t is the input vector, W and U are weights, and h_t is the hidden state vector — also known as the output vector of the LSTM unit. Because of the dependence on both past and future phrases, we use a Bidirectional LSTM (Bi-LSTM), which includes both a forward and backward layer of LSTMs. Both the forward and backward layer outputs are calculated by using the standard LSTM update equations (2). Then, Bi-LSTM

connects the two hidden layers to the same output layer. More details about Bi-LSTMs can be found in [39, 40].

To perform the final classification, the Bi-LSTM is connected to a Decision Tree to provide the label(s) for each timestamp. To accomplish this, we used the feature vector output from the (Time Distributed) Dense layer as the training set for the tree, which is fit along with the original set of labels. Once the tree is combined with the Bi-LSTM, it can be used to provide the final prediction for new timestamps. Using this approach also helped greatly decrease the training time of the Phrase Analyzer model and reduced the overfitting that the LSTM would suffer from by itself.

4. EXPERIMENTS

We evaluate each component of our system individually as well as a whole for the Phrase Analyzer, which utilizes all the components. For all the experiments, we utilize the entire augmented dataset, where 85% of the data was used for training, and the rest is used for testing.

4.1 Form Analyzer Evaluation

Setup: The Form Analyzer model takes in the features as described in Section 3.1 (the duration and Mel Spectrogram SSM) and outputs the predicted classification, which is compared against the “ground truth” label in the dataset. To turn the Mel Spectrogram SSM into a usable feature vector, we calculate the mean of the SSM to obtain a 1D array, which the model receives with the duration appended. Note that we can use the Form Analyzer architecture as a standalone model to identify the form of the musical piece. This can be useful to search a musical database by forms or for a musicologist to understand when and why a composer may choose a particular form over another. On the other hand, it can also be used to provide input to the Phrase Analyzer. Therefore, we evaluate its efficacy independently to evaluate its performance compared to other methods. The augmented dataset was split into 85% training and 15% testing for cross-validation,⁴ and the analysis model was evaluated using the classification accuracy in addition to other metrics such as Precision, Recall, and F1 scores.

Parameters: The TreeGrad model was tuned to 31 leaves per tree, an unbound max depth, 100 estimators (trees in the forest) with refit splits enabled, and the learning parameters include $\alpha = 0.1$ and a batch size of 32.

Results: The final Form Analyzer model achieved a classification accuracy of 83.9% — a surprisingly good performance given the subjective nature of the form classification, as many of the “inaccurate” classifications may be subjectively true based on personal bias. We compared our model with a number of other classification methods — 2D CNN [41], 2D CNN Ensemble [42], 1D CNN [43], Autokeras [44], Neural-SVM [45], XBNNet [46], DJINN [47], and an implementation of Neural Decision Trees [48]. In the plot in Figure 5 (left), our model outperforms other methods significantly. The next closest approach is the decision tree method, illustrating that the decision tree-based approaches indeed works better for this data. Note that the

model can be further fine-tuned by using a different boosting method. Furthermore, our method reported a Precision value of 85%, whereas both Recall and F-Score were 84%.

In addition to numerical accuracy, we wanted to study the mistakes in our method and how others performed in identifying the forms, and whether some specific forms were commonly misclassified as others. For this, we compared the confusion matrices from each approach. To compare the confusion matrices numerically, we compute the confusion entropy [49] for each — this is shown in Figure 5 (center and right). Confusion entropy is computed by exploiting the class distribution information of misclassifications of all classes as other classes (off-diagonal elements of the confusion matrices) and is an appropriate measure for this purpose. The results show overall, our method has the lowest entropy. The confusion entropy for each class (Figure 5 (right)) shows which classes were harder to classify. Here we found that our method has an entropy of .24 even for the worst-performing class. In addition, by quantitative analysis, a common theme we found is that most models tended to confuse forms that are typically similar in length (binary/unary, theme & variation/sonata, etc.) or hybrid/compound forms and their derived form (e.g., sonata/binary) even though for our model, such misclassifications are minimal.

4.2 Peak-Picking Algorithm Evaluation

Results: While the onset detection algorithm was not evaluated by a formal metric, it was tested against the training data, and the output timestamps were found to be nearly identical or had a low enough difference to be subjectively true (similar to the bias of a human analyst). The output of the algorithm was compared to numerous hand-labeled pieces from the dataset, and the difference was found to be negligible for most pieces (around $\approx \pm 5.269$ seconds on average), with the greatest absolute time difference being $\approx \pm 14.823$ seconds for pieces much greater than 6 minutes in duration. This metric was based on the approach found in [25]. The algorithm was also found to be comparable to other machine learning approaches such as SOMs [23] and CNNs [24]. Our result was based on a comparison between the “ground truth” and predicted timestamps in the validation dataset by index pairs rather than by the pair of closest timestamps since the algorithm may report more or fewer events than our own analysis.

4.3 Phrase Analyzer Evaluation

This model is much more difficult to evaluate using a classification evaluation metric due to numerous factors that vary from conventional machine learning classification models. First, the model gets the timestamps from the Peak-Picking algorithm, which is difficult to compare to human-annotated scores (our “ground truth”) due to integrative disagreement (i.e., a group of analysts would need to collectively agree on a ground truth as their analyses will likely vary). The labels are also often highly subjective and vary by the judgment of the annotator. In addition, some of the labels are implicit (for example, Part A continues until timestamp n but is only labeled at the first occurrence).

⁴ The dataset is sorted alphabetically by form per augmentation; other permutations typically performed the same or worse.

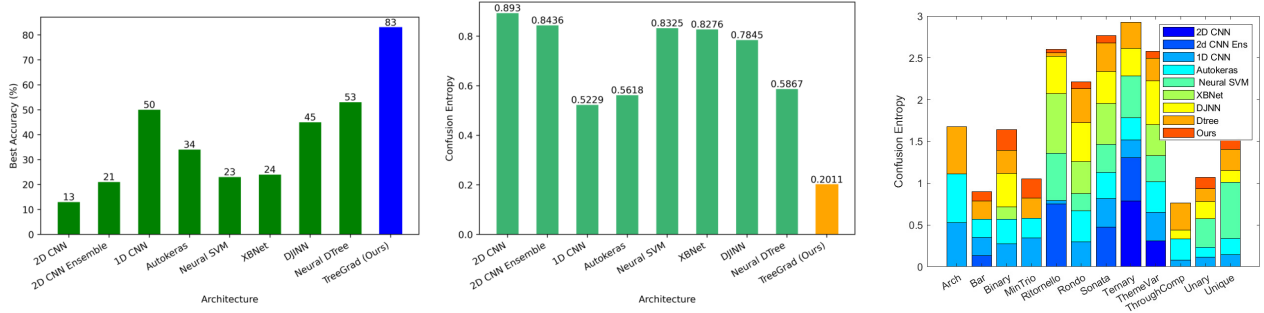


Figure 5. Form Analyzer Architecture Comparison with other methods (left), Confusion Entropy calculated for each method (center), and the class-wise Confusion Entropy for each method (right)

Therefore, we decided not to use a 0/1 accuracy metric but a more realistic one based on how a music theory expert would grade the labeling of other musicians (such as in a music theory class) based on a rubric. Such rubric-based grading is commonly used in case of challenging problems such as automated essay and clinical note grading [50, 51].

Setup: We design a rubric as follows: we score each labeling in $[0, 1]$ giving a score of 1 if the part and/or phrase match perfectly, but also have a variety of other values in $[0, 1)$ for partially correct labeling. These partial score scenarios cover a range of possibilities, such as if either the part or phrase but not both has been guessed correctly, the degree of similarity of the predicted phrase to the actual phrase, and accounting for the fact that time stamping may be off by a few seconds. These partial scores are ordered from high to low depending on the extent and multiplicity of the issues mentioned above. The rubric is presented in the [Supplement](#), and the code to do this grading on the output of the system will be provided with the dataset [13].

Parameters: The LSTM-Tree model has 4 LSTM units with a dropout rate of 0.2, a batch size of 10, sigmoid activation, and the learning parameters use the binary cross-entropy loss function, Adam optimizer, and 5 epochs.

Results: We evaluated the phrase analyzer using other state-of-the-art methods such as 1D CNN, RNN, and Feed-forward Neural Network and used the same rubric to come up with a weighted assessment score of the predicted phrases. We found that our LSTM-Tree architecture outperformed other NN architectures — our model reported a weighted score of 79.16%, whereas the best of the comparable methods reported a score of 77.75% (see [Supplement Figure 1](#) for a plot related to this). A qualitative evaluation of the results reveals our LSTM-Decision Tree method produced labels that are more consistent with a human analyst (who is prone to some errors) even on one attempt, whereas for some of the other methods, the labelings are more random and less interpretable. These results show that the LSTM-Decision Trees not only substantially increased accuracy quantitatively for the Phrase Analyzer but also gave reasonably realistic phrase labels.

5. DISCUSSION

Here we briefly discuss some possible extensions and improvements to the dataset as well as methods presented in

the paper. On the dataset front, we hope to expand the number of pieces as well as add additional annotators. The current dataset also features class imbalance; anthologies of classical music classified by form are lacking,⁵ though our system could be employed to assist in the compilation of such a database. A more sophisticated system may also greatly benefit from restructuring the analysis labels in the database to the standards specified in [54]. On the method front, an obvious area of improvement is the Peak-Picking algorithm, which can be replaced by a deep learning approach such as CNN or LSTM so that the whole network can be considered as one large deep learning system and parameter weights learned jointly. The Phrase Analyzer model can also benefit from the inclusion of Curriculum Learning or a Human-in-the-Loop type of approach, much like that of a traditional Form and Analysis class — though a Sequence-to-Sequence or Autoencoder model may also be useful in developing a faster/more accurate system.

6. CONCLUSION AND FUTURE WORK

In this paper, we introduce a new dataset, the **SMFSA Database**, accompanied by deep learning benchmark methods to analyze classical music forms and phrases. To the best of our knowledge, this is the most comprehensive study of a computational approach used toward classical music structure analysis. While the current system is specific to classical music, it could be extended to classify numerous additional forms (e.g., through transfer learning or an extended architecture), such as those found in popular music, ethnomusicology, and more complex hybrid forms. Another difficult task lacking substantial research is Optical Music Recognition (OMR) — our methods could potentially be extended to perform both visual music analysis and the classification/segmentation on the score directly (i.e., in the style of a human analyst). The system may also be extendable for use in Forensic Musicology and Copyright Detection systems, using the output analysis from the system to compare multiple pieces of music for potentially similar or exact replications of musical phrases.

⁵ Green’s book on Form Analysis [52] was found to be the most useful and accurate resource resembling such an anthology, whereas other resources were presented as anthologies for analysis rather than classified works. As such, our dataset was built primarily on the pieces (and musical forms) selected in this book (including some from [53] as well) to serve as a common ground for analysis.

7. REFERENCES

- [1] C. Burges, D. Plastina, J. Platt, E. Renshaw, and H. Malvar, "Using audio fingerprinting for duplicate detection and thumbnail generation," vol. 3, Apr 2005, pp. iii/9 – iii/12.
- [2] J. Valinsky, "Facebook now lets users share music, listen to 30-second song snippets," *Digiday*, Nov 2015. [Online]. Available: <https://digiday.com/?p=145138>
- [3] M. Müller, *Fundamentals of Music Processing: Using Python and Jupyter Notebooks*. Springer Nature, Apr 2021.
- [4] R. Liao, "How china's acrccloud detects copyrighted music in short videos," *TechCrunch*, Aug 2020. [Online]. Available: <https://techcrunch.com/2020/08/12/acrccloud-profile/>
- [5] T. Li, M. Ogiwara, and G. Tzanetakis, *Music Data Mining*. CRC Press, Jul 2011.
- [6] —, "Special section on music data mining," *Multi-media, IEEE Transactions on*, vol. 16, pp. 1185–1187, Aug 2014.
- [7] M. Gotham, "Moments musicaux," in *6th International Conference on Digital Libraries for Musicology*, ser. DLFM '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 70–78. [Online]. Available: <https://doi.org/10.1145/3358664.3358676>
- [8] R. Team, "Teaching the elements of music - form," Feb 2016. [Online]. Available: <https://www.connollymusic.com/stringovation/teaching-the-elements-of-music-form>
- [9] M. L. Maher and D. Fisher, "Using ai to evaluate creative designs," in *Proceedings of the 2nd International Conference on Design Creativity (ICDC)*, vol. 1, Sep 2012, p. 45–54.
- [10] H.-T. Cheng, Y.-H. Yang, Y.-C. Lin, and H. H. Chen, "Multimodal structure segmentation and analysis of music using audio and textual information," in *IEEE International Symposium on Circuits and Systems*, 2009, pp. 1677–1680.
- [11] J. B. L. Smith, J. A. Burgoyne, I. Fujinaga, D. De Roure, and J. S. Downie, "Design and creation of a large-scale database of structural annotations," in *ISMIR*, vol. 11. Miami, FL, 2011, pp. 555–560.
- [12] D. Szelogowski, "Deep learning for musical form: Recognition and analysis," Master's thesis, Apr 2022. [Online]. Available: <https://doi.org/10.13140/RG.2.2.33554.12481>
- [13] —, "Form-nn," Apr 2022. [Online]. Available: <https://github.com/danielathome19/Form-NN>
- [14] C. M. Resource, "Welcome to the classical midi and mp3 resource." [Online]. Available: <http://www.classicalmidiresource.com/>
- [15] B. Krueger, "Classical piano midi page - main page," 2018. [Online]. Available: <http://www.piano-midi.de/midicoll.htm>
- [16] K. der Fuge Project, "kunstderfuge.com - the largest classical music midi collection." [Online]. Available: <https://www.kunstderfuge.com/>
- [17] S. Ritchie, 2022. [Online]. Available: <https://www.classicalmidi.co.uk/page7.htm>
- [18] T. C. A. Team, 2022. [Online]. Available: <https://www.classicalarchives.com/midi.html>
- [19] E. Peiszer, T. Lidy, and A. Rauber, "Automatic audio segmentation: Segment boundary and structure detection in popular music," Tech. Rep., 2007.
- [20] J. Yang, "Music genre classification with neural networks: An examination of several impactful variables," 2018.
- [21] Y. Guan, J. Zhao, Y. Qiu, Z. Zhang, and G. Xia, "Melodic phrase segmentation by deep neural networks," 2018. [Online]. Available: <https://arxiv.org/abs/1811.05688>
- [22] D. Hörnel and W. Menzel, "Learning musical structure and style with neural networks," *Computer Music Journal*, vol. 22, no. 4, pp. 44–62, 1998. [Online]. Available: <http://www.jstor.org/stable/3680893>
- [23] P. J. Ponce de León and J. M. Iñesta, "Pattern recognition approach for music style identification using shallow statistical descriptors," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 37, no. 2, pp. 248–257, 2007.
- [24] K. Ullrich, J. Schlüter, and T. Grill, "Boundary detection in music structure analysis using convolutional neural networks," in *ISMIR*, 2014.
- [25] T. Grill and J. Schlüter, "Structural segmentation with convolutional neural networks mirex submission," 2015.
- [26] T. O'Brien, "Musical structure segmentation with convolutional neural networks," 2016.
- [27] J. de Berardinis, M. Vamvakaris, A. Cangelosi, and E. Coutinho, "Unveiling the hierarchical structure of music by multi-resolution community detection," *Transactions of the International Society for Music Information Retrieval*, vol. 3, no. 1, 2020.
- [28] carlosholivan, "Selfsimilaritymatrix-serra.ipynb," Sep 2020. [Online]. Available: https://github.com/carlosholivan/SelfSimilarityMatrices/blob/master/SelfSimilarityMatrix_Serra.ipynb
- [29] T. R. Hoens, Q. Qian, N. V. Chawla, and Z.-H. Zhou, "Building decision trees for the multi-class imbalance problem," in *Advances in Knowledge Discovery and Data Mining*, P.-N. Tan, S. Chawla, C. K. Ho, and J. Bailey, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 122–134.

- [30] X. Yuan, L. Xie, and M. Abouelenien, "A regularized ensemble framework of deep learning for cancer detection from multi-class, imbalanced training data," *Pattern Recognition*, vol. 77, pp. 160–172, 2018.
- [31] K. Pasupa, S. Vathavanavaro, and S. Tungjitnob, "Convolutional neural networks based focal loss for class imbalance problem: A case study of canine red blood cells morphology classification," *Journal of Ambient Intelligence and Humanized Computing*, pp. 1–17, 2020.
- [32] V. S. Lokhande, A. K. Akash, S. N. Ravi, and V. Singh, "Fairalm: Augmented lagrangian method for training fair models with little regret," in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Cham: Springer International Publishing, 2020, pp. 365–381.
- [33] C. Siu, "Transferring tree ensembles to neural networks," in *Neural Information Processing*, T. Gedeon, K. W. Wong, and M. Lee, Eds. Cham: Springer International Publishing, 2019, pp. 471–480.
- [34] P. Kontschieder, M. Fiterau, A. Criminisi, and S. Rota Bulò, "Deep neural decision forests," Dec 2015, pp. 1467–1475.
- [35] J. Serrà, M. Müller, P. Grosche, and J. L. Arcos, "Unsupervised music structure annotation by time series structure features and segment similarity," *IEEE Transactions on Multimedia*, vol. 16, no. 5, pp. 1229–1240, 2014.
- [36] C. Hernandez-Olivan, J. R. Beltran, and D. Diaz-Guerra, "Music boundary detection using convolutional neural networks: A comparative analysis of combined input features," Aug 2020. [Online]. Available: <https://arxiv.org/abs/2008.07527>
- [37] Y. Yu, X. Si, C. Hu, and J. Zhang, "A review of recurrent neural networks: Lstm cells and network architectures," *Neural computation*, vol. 31, no. 7, pp. 1235–1270, 2019.
- [38] A. Sherstinsky, "Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm) network," *Physica D: Nonlinear Phenomena*, vol. 404, p. 132306, 2020.
- [39] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional lstm networks," in *Proceedings of the IEEE International Joint Conference on Neural Networks.*, vol. 4, 2005, pp. 2047–2052.
- [40] Z. Cui, R. Ke, Z. Pu, and Y. Wang, "Stacked bidirectional and unidirectional lstm recurrent neural network for forecasting network-wide traffic state with missing values," *Transportation Research Part C: Emerging Technologies*, vol. 118, p. 102674, 2020.
- [41] P. Nandi, "Cnns for audio classification," *Towards Data Science*, Dec 2021. [Online]. Available: <https://towardsdatascience.com/cnns-for-audio-classification-6244954665ab>
- [42] L. Nanni, Y. M. G. Costa, R. L. Aguiar, R. B. Mangolin, S. Brahnam, and J. Silla, Carlos N., "Ensemble of convolutional neural networks to improve animal audio classification," *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2020, no. 1, May 2020.
- [43] S. Kiranyaz, T. Ince, O. Abdeljaber, O. Avci, and M. Gabbouj, "1-d convolutional neural networks for signal processing applications," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 8360–8364.
- [44] J. Brownlee, "How to use autokeras for classification and regression," Sep 2020. [Online]. Available: <https://machinelearningmastery.com/autokeras-for-classification-and-regression/>
- [45] M. A. Wiering, M. H. v. d. Ree, M. Embrechts, and L. Schomaker, "The neural support vector machine," Nov 2013. [Online]. Available: https://www.researchgate.net/publication/258435394_The_Neural_Support_Vector_Machine
- [46] T. Sarkar, "Xbnet - an extremely boosted neural network," Jun 2021. [Online]. Available: <https://arxiv.org/abs/2106.05239>
- [47] K. D. Humbird, L. Peterson, and R. McClarren, "Deep jointly-informed neural networks," Jul 2017. [Online]. Available: https://www.researchgate.net/publication/318205627_Deep_Jointly-Informed_Neural_Networks
- [48] Y. Yang, I. G. Morillo, and T. M. Hospedales, "Deep neural decision trees," Jun 2018. [Online]. Available: <https://arxiv.org/abs/1806.06988>
- [49] J.-M. Wei, X.-J. Yuan, Q.-H. Hu, and S.-Q. Wang, "A novel measure for evaluating classifiers," *Expert Systems with Applications*, vol. 37, no. 5, pp. 3799–3809, 2010.
- [50] H. T. T. Nguyen, "Neural networks for automated essay grading," 2016.
- [51] W.-w. Yim, A. Mills, H. Chun, T. Hashiguchi, J. Yew, and B. Lu, "Automatic rubric-based content grading for clinical notes," in *Proceedings of the Tenth International Workshop on Health Text Mining and Information Analysis (LOUHI 2019)*. Association for Computational Linguistics, 2019. [Online]. Available: <http://dx.doi.org/10.18653/v1/d19-6216>
- [52] D. M. Green, *Form in Tonal Music: An Introduction to Analysis*. Cengage Learning, 1979.
- [53] W. E. Caplin, *Classical Form: A Theory of Formal Functions for the Instrumental Music of Haydn, Mozart, and Beethoven*. Oxford University Press, Dec 2000.
- [54] M. R. H. Gotham and M. Ireland, "Taking form: A representation standard, conversion code, and example corpora for recording, visualizing, and studying analyses of musical form," in *ISMIR*, 2019.

BAF: AN AUDIO FINGERPRINTING DATASET FOR BROADCAST MONITORING

Guillem Cortès [♫]

Alex Ciurana [♫]

Emilio Molina [♫]

Marius Miron [♫]

Owen Meyers [#]

Joren Six [♭]

Xavier Serra [♫]

[♫] BMAT Licensing S.L., Barcelona

[#] Epidemic Sound, Stockholm

[♫] MTG, Universitat Pompeu Fabra, Barcelona

[♭] IPEM, Ghent University, Ghent

ABSTRACT

Audio Fingerprinting (AFP) is a well-studied problem in music information retrieval for various use-cases e.g. content-based copy detection, DJ-set monitoring, and music excerpt identification. However, AFP for continuous broadcast monitoring (e.g. for TV & Radio), where music is often in the background, has not received much attention despite its importance to the music industry. In this paper (1) we present BAF, the first public dataset for music monitoring in broadcast. It contains 74 hours of production music from Epidemic Sound and 57 hours of TV audio recordings. Furthermore, BAF provides cross-annotations with exact matching timestamps between Epidemic tracks and TV recordings. Approximately, 80% of the total annotated time is background music. (2) We benchmark BAF with public state-of-the-art AFP systems, together with our proposed baseline *PeakFP*: a simple, non-scalable AFP algorithm based on spectral peak matching. In this benchmark, none of the algorithms obtain a F1-score above 47%, pointing out that further research is needed to reach the AFP performance levels in other studied use cases. The dataset, baseline, and benchmark framework are open and available for research.

1. INTRODUCTION

Audio Fingerprinting (AFP) is the information retrieval task of identifying audio recordings in a given database of reference songs. The task is based on extracting content-based signatures that summarize an audio recording (*extraction*) [1], storing them in a database or in hash tables (*indexing*), and efficiently linking short snippets of unlabeled audio to the same content in the database (*matching*).

AFP has been successfully applied to different tasks such as query by example [2], advertisement tracking [3],

integrity verification [4], and data deduplication [5]. A relevant AFP application is broadcast monitoring for its crucial role in royalties distribution. In 2021, US\$ 2.9 billion were distributed among rights holders [6], representing 11.5% of the global recorded music industry revenues.

A great AFP system for broadcast monitoring must provide exact start and end timestamps within long audio recordings. It also must be robust against common distortions in broadcasting like music being in the background and with low SNR. To the best of our knowledge, the literature has not addressed the AFP use case of broadcasting monitoring with a heavy presence of background music. In this scenario, music may have a very low Signal-to-Noise ratio (SNR) that makes it difficult to identify [7]. Moreover, music may be masked by a large variety of non-musical sounds, such as speech, applause, laughter, urban and nature sounds, etc. [8].

To promote research in AFP for broadcast monitoring we propose BAF: a Broadcast Audio Fingerprinting dataset with TV recordings in which production music is played. The references are part of Epidemic Sound’s private catalog [9], a collection of music for content creators (production music). BAF reflects a challenging (but realistic) scenario [7] for broadcast AFP systems: low sample rate monaural audio, background music of variable SNR, a large variety of contexts, long queries with true negative sections, and matches of multiple durations. More details about the dataset can be found in Section §3.

In order to show the challenges that BAF presents, in Section §4 we present a benchmark with 4 of the available AFP algorithms: Audfprint [10], Panako [11, 12], Olaf [13], and NeuralFP [14]. These are open-source implementations that represent different approaches to AFP. In addition, we propose an open implementation of a simple baseline based on spectral peak matching and the evaluation framework used in the benchmark.

2. RELATED DATASETS

In this section, we present public datasets that have been used in the AFP literature. We also gather information about the contents of private datasets and depict them in Table 1. We include information about the queries, references, and the goal of the dataset or publication. We have found in the literature 21 datasets that have been used for

Corresponding author: Guillem Cortès (cortes.sebastia@gmail.com)



© G. Cortès, A. Ciurana, E. Molina, M. Miron, O. Meyers, J. Six and X. Serra. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** G. Cortès, A. Ciurana, E. Molina, M. Miron, O. Meyers, J. Six and X. Serra, “BAF: an Audio Fingerprinting dataset for Broadcast Monitoring”, in *Proc. of the 23rd Int. Society for Music Information Retrieval Conf.*, Bengaluru, India, 2022.

Work	(Synthetic) Queries	References	Goal	Used in
2002 Haitsma, J. [15]	(✓) 4 excerpts	20,000	Scalable, robust to noise	
2003 Wang, A. [2]	250 excerpts	10,000	Scalable, robust to noise	
2005 Bartsch, M.A. [16]	93	93	Structural redundancy	
2008 Baluja, S. [17]	(✓) 1,000	10,000	Robust to pitch and time	
2008 Bellettini, C. [18, 19]	(✓) 15,000	15,000	Robust to pitch	
2011 Fenet, S. [20]	7d radio broad.	7,309 (60s)	Scalable, robust to pitch	
2011 Fenet, S. [20]	5d radio broad.	30,000	Scalable, robust to pitch	
2014 Malekesmaeili, M. [21]	(✓) 200	200	Robust to noise, pitch and time	
2015 Zhang, X. [22]	(✓) 10s excerpts	2,075	Robust to pitch and time	
2016 Sonnleitner, R. [23]	(✓) 300	20,000	Robust to noise, pitch and time	
2016 Sonnleitner, R. [24]	8 DJ-mixes	296 \approx 7h	Robust to pitch and time	
2016 Walter, T. [25]	10s excerpts	10M	Efficient and scalable AFP	
2017 Gfeller, B. [26]	12,000 excerpts	450h	Low-power music recognizer	
2020 Son H.-S. [27]	(✓) 100	100	Robust to pitch	
2020 Yu, Z. [28]	(✓) 5,000 (10s)	345,000	Robust to any degradation	
2009 MagnaTagATune [29]		25,863 \approx 208h	Music Tagging	[30]
2006 TRECVID [31, 32]	(✓) 201	11,200 \approx 400h	Content-based Copy Detection	[33–35]
2014 Panako [11]	(✓) 600 excerpts	30,000 \approx 277h	Robust to pitch and time	[11]
2016 Mixotic [24]	10 DJ-mixes	723 \approx 11h	Robust to pitch and time	[24]
2016 QuadFP [23]	(✓) 450,000	100,011 \approx 6,899h	Robust to pitch and time	[23]
2021 NeuralFP [14]	(✓) short excerpts	100k \approx 8,000h	High-specific audio retrieval	[14]

Table 1. Private (top \uparrow) and public (bottom \downarrow) datasets that have been used for AFP. For each dataset information about the queries, references and the goal of the original work is given. Synthetic (✓) queries have been created by applying transformations to reference audios. We also indicate if the queries are excerpts and the number of original audio pieces that have been transformed, not meaning the total number of queries after the transformation.

AFP. Even though the nature of the data varies from each dataset, most of them rely on the same principle: a private collection of tracks that constitute a reference set and a query set formed by applying transformations to some of the references. This is a good way to test the limits of the robustness to degradations like pitch-shifting or time-scaling. Still, it does not target the characteristics of realistic broadcast monitoring, where the music is in the background, masked by speech and a wide variety of sounds and noises (See Section §3.3.1 for a detailed analysis).

As Table 1 reflects, only 6 out of the 21 datasets are public, mainly due to the difficulty of legally publishing copyrighted music. In other cases, data remains private to protect intellectual property, as with private companies. To that extent, private datasets hinder reproducibility and slow-down scientific progress in AFP research. Of all private datasets, only the ones built by Fenet et al. [20] reflect the use case of broadcast monitoring: they use real radio broadcasted emissions as queries and a reference set of 459 songs. Moreover, only Fenet’s [20], Sonnleitner et al., [24] and Walter and Gould [25] used unknown audios as queries. All other works used a version of the reference songs often modified with some degradation: echo, pitch and/or tempo alteration, reverb, etc. Besides these private datasets, there are some other public datasets that have been used in AFP works, even though none of them fits the broadcast monitoring task.

MagnaTagATune [29] is a public dataset created for Music Tagging. Each audio clip has associated a vector

of binary annotations of 188 tags that describe the music piece. The dataset was used for AFP by Ramona and Peeters [30] in which they followed the experimental protocol of Haitsma and Kalker [15] applying a series of distortions like Amplitude dynamic compression, MP3 encoding, time-shifting, or equalization to 500 music clips from MagnaTagATune.

NIST-TRECVID [31, 32]. One of the tasks the TREC Video Retrieval Evaluation (TRECVID) proposed until 2011 was Content-Based Copy Detection (CCD) [36]. Various works [33–35] used the dataset provided with the task to evaluate AFP algorithms. The length of queries varies from 3 to 180 seconds and comprises multiple transformed fragments from 201 unique audio recordings.

Mixotic [24] was created by Sonnleitner et al. to test the robustness of QuadFP, Panako, and Audfprint in real DJ mixes. It was generated from free, CC-licensed DJ mixes that were published on mixotic netlabel.

Panako [11], **QuadFP** [23, 37] and **NeuralFP** [14] present public datasets that were curated to test their algorithms. Since these datasets are reproducible we list them as public AFP datasets, but in the case of Panako and QuadFP they share a script that downloads tracks from Jamendo instead of the audio files. It can happen that some audio works are not available anymore thus impeding the reconstruction of the dataset. NeuralFP shares the audio files since they come from the FMA dataset [38]. It was trained on 10k FMA songs and tested on a larger set of 100k songs (\approx 8,000 hours).