On the Application of Generic Summarization Algorithms to Music

Francisco Raposo, Ricardo Ribeiro, David Martins de Matos, Member, IEEE

Abstract

Several generic summarization algorithms were developed in the past and successfully applied in fields such as text and speech summarization. In this paper, we review and apply these algorithms to music. To evaluate this summarization's performance, we adopt an extrinsic approach: we compare a Fado Genre Classifier's performance using truncated contiguous clips against the summaries extracted with those algorithms on 2 different datasets. We show that Maximal Marginal Relevance (MMR), LexRank and Latent Semantic Analysis (LSA) all improve classification performance in both datasets used for testing.

I. Introduction

Several algorithms to summarize music have been published [1]–[8], mainly for popular music songs whose structure is repetitive enough. However, those algorithms were devised with the goal of producing a thumbnail of a song as its summary, the same way an image's thumbnail is that image's summary. Therefore, the goal is to output a shorter version of the original song so that people can quickly get the gist of the whole piece without listening to all of it. These algorithms usually extract continuous segments because of their human consumption-oriented purpose.

Francisco Raposo is with Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais, 1049-001 Lisboa, Portugal Ricardo Ribeiro is with Instituto Universitário de Lisboa (ISCTE-IUL), Av. das Forças Armadas, 1649-026 Lisboa, Portugal David Martins de Matos is with Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais, 1049-001 Lisboa, Portugal

Ricardo Ribeiro and David Martins de Matos are with L2F - INESC ID Lisboa, Rua Alves Redol, 9, 1000-029 Lisboa, Portugal

This work was supported by national funds through FCT – Fundação para a Ciência e a Tecnologia, under project PEst-OE/EEI/LA0021/2013.

Generic summarization algorithms have also been developed for and are usually applied in text summarization. Their application, in music, to extract a thumbnail is not ideal, because a "good" thumbnail entails requirements such as coherence and clarity. These summaries are composed of small segments from different parts of the song which makes them unsuitable for human enjoyment and thus may help evade copyright issues. Nevertheless, most of these algorithms produce summaries that are both concise and diverse.

We review several summarization algorithms, in order to summarize music for automatic, instead of human, consumption. The idea is that a summary clip contains more relevant and less redundant information and, thus, may improve the performance of certain tasks that rely on processing just a portion of the whole audio signal. We evaluate the summarization's contribution by comparing the performance of a Portuguese music style Fado Genre Classifier [9] using the extracted summaries of the songs against using contiguous clips (truncated from the beginning, middle and end of the song). We summarize music using MMR, LexRank, LSA and also with a method for music summarization called Average Similarity for comparison purposes. We present results on 2 datasets showing that MMR, LexRank and LSA improve classification performance under certain parameter combinations.

Section II reviews related work on summarization. Specifically, the following algorithms are reviewed: Average Similarity in section II-A, MMR in section II-B, LexRank in section II-C and LSA in section II-D. Section III describes the details of the experiments we performed for each algorithm and introduces the Fado Classifier. Section IV reports and discusses our classification results and section V concludes this paper with some remarks and future work.

II. SUMMARIZATION

Several algorithms for both generic and music summarization have been proposed. However, music summarization algorithms were developed to extract an audible summary so that any person can listen to it coherently. Our focus is on automatic consumption, so coherence and clarity are not mandatory requirements for our summaries.

LexRank [10] and TextRank [11] are centrality-based methods that rely on the similarity between every sentence. These are based on Google's PageRank [12] algorithm for ranking web pages and are successfully applied in text summarization. GRASSHOPPER [13] is another method applied in text summarization, as well as social network analysis, focusing on improving diversity in ranking sentences. MMR [14], [15], applied in speech summarization, is a query-specific method that selects sentences according to their similarity to the query and to the sentences previously selected. LSA [16] is another

method used in text summarization based on the mathematical technique Singular Value Decomposition (SVD).

Music-specific summarization structurally segments songs and then selects which segments to include in the summary. This segmentation aims to extract meaningful segments (e.g. chorus, bridge). [1] presents two approaches for segmentation: using a Hidden Markov Model (HMM) to detect key changes between frames and Dynamic Time Warping (DTW) to detect repeating structure. In [2], segmentation is achieved by correlating a Gaussian-tempered "checkerboard" kernel along the main diagonal of the similarity matrix of the song, outputting segment boundaries. Then, a segment-indexed similarity matrix is built, containing the similarity between every detected segment. SVD is applied to that matrix to find its rank-K approximation. Segments are, then, clustered to output the song's structure. In [3], [4], songs are segmented in 3 stages. First, a similarity matrix is built and it is analyzed for fast changes, outputting segment boundaries. These segments are clustered to output the "middle states". Finally, an HMM is applied to these states, producing the final segmentation. These algorithms then follow some strategies to select the appropriate segments. [5] groups (based on the Kullback-Leibler (KL) divergence) and labels similar segments of the song and then the summary is generated by taking the longest sequence of segments belonging to the same cluster. In [6], [7], a method called Average Similarity is used to extract a thumbnail L seconds long that is most similar to the whole piece. Another method for this task is the Maximum Filtered Correlation [8] which starts by building a similarity matrix and then a filtered time-lag matrix, which has the similarity between extended segments embedded in it. Finding the maximum value in the latter is finding the starting position of the summary.

To apply generic summarization algorithms to music, first we need to segment the song into musical words/terms. This fixed segmentation differs a lot from the structural segmentation used in music-specific algorithms. Fixed segmentation does not take into account the human perception of musical structure. It simply allows us to look at the variability and repetition of the signal and use them to find the most important parts. Structural segmentation aims to find meaningful segments (to people) of the song so that we can later select those segments to include in the summary. This type of segmentation often leads to audible summaries which violate copyrights of the original songs. Fixed segmentation combined with generic summarization algorithms may help evade those issues.

In the following sections we review the algorithms we chose to evaluate: Average Similarity, MMR, LexRank, and LSA.

A. Average Similarity

This approach to summarization has the purpose of finding a fixed-length continuous music segment, of duration L, most similar to the entire song. This method was introduced in [6] and later used in other research efforts such as [7].

The method consists of building a similarity matrix for the song and calculating an aggregated measure of similarity between the whole song and every L seconds long segment.

In [6], 45 Mel Frequency Cepstral Coefficient (MFCC)s are computed but only the 15 with highest variance are kept. The cosine distance is used to calculate pairwise similarities.

In [7], the first 13 MFCCs and the spectral centre of gravity (sound "brightness") are used. The Tchebychev distance was selected for building the similarity matrix.

Once the similarity between every frame is calculated, we build a similarity matrix S and embed the similarity values between feature vectors v_i and v_j in it: $S(i,j) = s(v_i, v_j)$.

The average similarity measure can be calculated by summing up columns (or rows, since the similarity matrix is symmetric) of the similarity matrix, according to the desired summary length L, starting from different initial frames. The maximum score will correspond to the segment that is most similar to the whole song. To find the best summary of length L, we must compute the score $Q_L(i)$:

$$Q_{L}(i) = \bar{S}(i, i + L) = \frac{1}{NL} \sum_{m=i}^{i+L} \sum_{n=1}^{N} S(m, n)$$
(1)

N is the number of frames in the entire piece. The index $1 \le i \le (N - L)$ of the best summary starting frame is the one that maximizes $Q_L(i)$.

The evaluations of this method in the literature are subjective (human) evaluations that take into account whether the generated summaries include the most memorable part(s) of the song [6]. Other evaluations are averages of scores given by test subjects, regarding specific qualities of the summary such as Clarity, Conciseness and Coherence [7].

B. Maximal Marginal Relevance

MMR [17], selects sentences from the signal according to their relevance and to their diversity against the already selected sentences in order to output low-redundancy summaries. This approach has been used in speech summarization [14], [15]. It is a query-specific summarization method, though it is possible to produce generic summaries by taking the centroid vector of all the sentences (as in [15]) as the query.

MMR iteratively selects the sentence S_i that maximizes the following mathematical model:

$$\lambda \left(Sim_1 \left(S_i, Q \right) \right) - \left(1 - \lambda \right) \max_{S_i} Sim_2 \left(S_i, S_j \right) \tag{2}$$

 Sim_1 and Sim_2 are the, possibly different, similarity metrics; S_i are the unselected sentences and S_j are the previously selected ones; Q is the query and λ is a configurable parameter that allows the selection of the next sentence to be based on its relevance, its diversity or a linear combination of both. Usually sentences are represented as Term Frequency - Inverse Document Frequency (TF-IDF) scores vectors. The cosine similarity is frequently used as Sim_1 and Sim_2 .

C. LexRank

LexRank [10] is a centrality-based method that relies on the similarity for each sentence pair. This centrality-based method is based on Google's PageRank [12] algorithm for ranking web pages. The output is a list of ranked sentences from which we can extract the most central ones to produce a summary.

First, we compare all sentences, normally represented as TF-IDF scores vectors, to each other using a similarity measure. LexRank uses the cosine similarity. After this step, we build a graph where each sentence is a vertex and edges are created between every sentence according to their pairwise similarity. Usually, the similarity score must be higher than some threshold to create an edge. LexRank can be used with both weighted and unweighted edges. Then, we perform the following calculation iteratively for each vertex until convergence is achieved (when the error rate of two successive iterations is below a certain threshold for every vertex):

$$S(V_i) = \frac{(1-d)}{N} + S_1(V_i)$$
 (3)

$$S_1(V_i) = d \times \sum_{V_i \in adj[V_i]} \frac{Sim(V_i, V_j)}{\sum_{V_k \in adj[V_j]} Sim(V_j, V_k)} S(V_j)$$

$$\tag{4}$$

d is a damping factor to guarantee the convergence of the method, N is the total number of vertices and $S(V_i)$ is the score of vertex i. This is the case where edges are weighted. When using unweighted edges, the equation is simpler:

$$S(V_i) = \frac{(1-d)}{N} + d \times \sum_{V_j \in adj[V_i]} \frac{S(V_j)}{D(V_j)}$$

$$(5)$$

 $D(V_i)$ is the degree (i.e., number of edges) of vertex i. We can construct a summary by taking the highest ranked sentences until a certain summary length is reached.

This method is based on the fact that sentences recommend each other. A sentence very similar to many other sentences will get a high score. Sentence score is also determined by the score of the sentences recommending it.

D. Latent Semantic Analysis

LSA is based on the mathematical technique SVD that was first used for text summarization in [16]. SVD is used to reduce the dimensionality of an original matrix representation of the text. To perform LSA-based text summarization, we start by building a T terms by N sentences matrix A.

Each element of A, $a_{ij} = L_{ij}G_i$, has two weight components: a local weight and a global weight. The local weight is a function of the number of times a term occurs in a specific sentence and the global weight is a function of the number of sentences that contain a specific term.

Applying SVD to matrix A will result in a decomposition formed by three matrices: U, a $T \times N$ matrix of left singular vectors (its columns); Σ , a $N \times N$ diagonal matrix of singular values; and V^T , a $N \times N$ matrix of right singular vectors (its rows): $A = U \Sigma V^T$.

Singular values are sorted by descending order in matrix Σ and are used to determine topic relevance. Each latent dimension corresponds to a topic. We calculate the Rank K approximation by taking the first K columns of U, the $K \times K$ sub-matrix of Σ and the first K rows of V^T . We can extract the most relevant sentences by iteratively selecting sentences corresponding to the indices of the highest values for each (most relevant) right singular vector.

In [18], two limitations of this approach are discussed: the fact that K is equal to the number of sentences in the summary, which, as it increases, tends to include less significant sentences; and that sentences with high values in several dimensions (topics), but never the highest, will never be included in the summary. To compensate for these problems, a sentence score was introduced and K is chosen so that the K^{th} singular value does not fall under half of the highest singular value: $score(j) = \sqrt{\sum_{i=1}^k v_{ij}^2 \sigma_i^2}$.

III. EXPERIMENTS

To evaluate these algorithms on music, we tested their impact on a Fado classifier. This classifier simply classifies a song as Fado or non-Fado. Fado is a Portuguese music genre whose instrumentation usually consists solely of stringed instruments, such as the classical guitar and the Portuguese guitar. The classifier is a Support Vector Machine (SVM) [19].

The features used by the SVM consist of a 32-dimensional vector per song, which is a concatenation of 4 features: average vector of the first 13 MFCCs of the song; Root Mean Square (RMS) energy; high frequencies 9-dimensional rhythmic features; and low frequencies 9-dimensional rhythmic features.

These rhythmic features are computed based on the Fast Fourier Transform (FFT) coefficients on the 20 Hz to 100 Hz range (low frequencies) and on the 8000 Hz to 11025 Hz range (high frequencies). Assuming v is a matrix of FFT coefficients with frequency varying through columns and time through lines, each component of the 9-dimensional vector is: maxamp: max of the average v along time; minamp: min of the average v along time; number of v values above 80% of maxamp; number of v values above 15% of maxamp; number of v values above maxamp; number of v values below minamp; mean distance between peaks; standard deviation of distance between peaks; max distance between peaks.

These features capture rhythmic information in both low and high frequencies. Fado does not have much information in the low frequencies as it does not contain, for example, drum kicks. However, due to the string instruments used, Fado information content is higher in the high frequencies, making these features good for distinguishing it from other genres.

We used 2 datasets in our experiments which consist of 500 songs from which half of them are Fado songs and the other half are not. The 250 Fado songs are the same in both datasets. The datasets are encoded in mono, 16-bit, 22050 Hz Microsoft WAV files. We will make the post-summarization datasets available upon request.

We used 5-fold cross validation when calculating classification performance. The classification performance was calculated first for the beginning, middle and end sections (of 30s) of the songs to get a baseline and then we compared it with the classification using the summaries (also 30s) for each parameter combination and algorithm.

For feature extraction we used OpenSMILE's [20] implementation, namely, to extract MFCC feature vectors. We also used the Armadillo library [21] for matrix operations and the Marsyas library [22] for synthesizing the summaries.

For Average Similarity, we experimented with 3 different frame sizes (0.25, 0.5, and 1 s) with both 50% and no overlap. We also experimented with MFCC vector sizes of 12 and 24.

To use the generic summarization algorithms, however, we need additional processing steps. We adapted those algorithms to the music domain by mapping the audio signal frames (represented as MFCC vectors) to a discrete representation of words and sentences. For each piece being summarized, we cluster all of its frames using the mlpack's [23] K-Means algorithm implementation which calculates the vocabulary for that song (i.e., each frame is now a word from that vocabulary). Then, we segment the whole piece into fixed-size sentences (e.g., 5-word sentences). This allows us to represent each sentence as a vector of word occurrences/frequencies (depending on the type of weighting chosen) which lets us compare sentences with each other using the cosine distance.

In our implementation of MMR, we calculate the similarity between every sentence only once and then apply the algorithm until the desired summary length is reached. We experimented using 3 different values for λ (0.3, 0.5 and 0.7) and 4 different weighting types: raw (counting of the term), binary (presence of the term), TF-IDF and "dampened" TF-IDF (same as TF-IDF but takes logarithm of TF instead of TF itself).

The damping factor used in LexRank was 0.85 and the convergence threshold was set to 0.0001. We also calculated the similarity between every sentence only once, applying the iterative algorithm and picking sentences until the desired summary length is reached. We also tested LexRank using the same weighting types as for MMR.

We used Armadillo's [21] implementation of the SVD operation to implement LSA. After sentence/word segmentation, we apply SVD to the term by sentences matrix (column-wise concatenation of all sentence vectors). We then take the rank-K approximation of the decomposition where the Kth singular value is not smaller than half of the (K-1)th singular value. Then, we calculate the sentence score (as explained in section II-D) for each sentence and pick sentences according to that ranking until the desired summary length is reached. We tested LSA with both raw and binary weighting.

We tested MMR, LexRank, and LSA, with all combinations of the following parameter values: frame size of 0.5s with no overlap and with 50% (0.25s hops) overlap; vocabulary size of 25, 50, and 100 words; and sentence size of 5, 10, and 20 words. We used MFCC vectors (of size 12) as features for these experiments, they are widely used in many MIR tasks including music summarization in [2], [5]–[7].

IV. RESULTS

We present only the most interesting results, since we tried many different parameter combinations for each algorithm. The Frame/Hop Size columns indicate the frame/hop sizes in seconds, which can be interpreted as overlap (e.g., the pair 0.5, 0.25 stands for frames of 0.5s duration with a hop size of 0.25s, which corresponds to a 50% overlap between frames). The classification accuracy results for the 30s contiguous segments which constitute the baseline are 95.8%, 96.2%, and 94% for the beginning, middle, and end sections, respectively, on dataset 1 and 85.2%, 92%, and 90.4%, on dataset 2.

The Average Similarity algorithm was successful in improving classification performance on dataset 1 (98.8% as maximum accuracy obtained with frame size of 0.5 s, no overlap, 24 MFCCs), but not on dataset 2 (90.8% maximum accuracy with frame size of 0.25 s, no overlap, 12 MFCCs).

In table I, we can see that although not all parameter combinations for MMR yielded an increase in classification performance on both datasets, some combinations did do that. For example, the best

combination on the dataset 1 yielded 100% accuracy but on dataset 2 it yielded only 90.8% which is lower than the baseline (92%). However, all other parameter combination presented in those tables yield a better result than the baseline for both datasets. We also noticed that smaller values of λ would result in worse accuracy scores.

We can also see that the best parameter combination for LexRank on dataset 1 was also the best on dataset 2. Besides that, all other presented combinations are better when compared to the corresponding baseline, which suggests that these parameter combinations might also be good for other datasets.

Our experiments show that LSA works best with binary weighting when applied to music. This has to do with the fact that some musical sentences, namely, at the beginning of the songs, are strings with very few repeating terms, which increases term-frequency scores. Moreover, those terms might not even appear anywhere else in the song which will, in turn, decrease the document frequency of the term, thus increasing the inverse document frequency score. These issues are detected when LSA chooses those (unwanted) sentences because they will have a high score on a certain latent topic. The binary weighting alleviates these problems because we only check for the presence of a term (not its frequency) and the document frequency of that term is not taken into account. LSA also achieved results above the baseline (table I).

V. CONCLUSIONS AND FUTURE WORK

We evaluated summarization through classification for MMR, LexRank, and LSA in the music domain. More experimenting should be done to find a set of parameter combinations that will work for most music contexts. Future work includes testing other summarization algorithms, other similarity metrics, other types of features and other types of classifiers. The use of Gaussian Mixture Models may also help in finding more "natural" vocabularies and Beat Detection might be used to find better values for fixed segmentation.

TABLE I $\label{eq:mmr} \mbox{MMR, Lexrank and LSA (\#MFCC = 12)}$

Frame	Нор	Vocab.	Sentence	Weighting	λ	Accuracy
Size	Size	Size	Size			
MMR on Dataset 1						
0.5	0.5	50	5	dampTF	0.7	100%
0.5	0.25	100	5	Binary	0.7	99.2%
0.5	0.5	25	5	Binary	0.5	97.2%
0.5	0.5	25	10	dampTF	0.7	97.6%
MMR on Dataset 2						
0.5	0.5	50	5	dampTF	0.7	90.8%
0.5	0.25	100	5	Binary	0.7	93.4%
0.5	0.5	25	5	Binary	0.5	93.4%
0.5	0.5	25	10	dampTF	0.7	93.4%
LexRank on Dataset 1						
0.5	0.5	25	5	dampTF	-	99%
0.5	0.25	100	20	Binary	-	97.4%
0.5	0.5	25	10	dampTF	-	97.6%
0.5	0.5	25	10	Raw	-	97.6%
LexRank on Dataset 2						
0.5	0.5	25	5	dampTF	-	94%
0.5	0.25	100	20	Binary	-	93.8%
0.5	0.5	25	10	dampTF	-	93.8%
0.5	0.5	25	10	Raw	-	93.4%
LSA on Dataset 1						
0.5	0.5	100	20	Binary	-	99.6%
0.5	0.5	25	10	Binary	-	99.4%
0.5	0.5	50	10	Binary	-	96.6%
0.5	0.25	25	20	Binary	-	97%
LSA on Dataset 2						
0.5	0.5	100	20	Binary	-	91.2%
0.5	0.5	25	10	Binary	-	93.4%
0.5	0.5	50	10	Binary	-	93.4%
0.5	0.25	25	20	Binary	-	92.8%

REFERENCES

- [1] W. Chai, "Semantic Segmentation and Summarization of Music: Methods Based on Tonality and Recurrent Structure," Signal Processing Magazine, IEEE, vol. 23, no. 2, pp. 124–132, March 2006.
- [2] M. Cooper and J. Foote, "Summarizing Popular Music via Structural Similarity Analysis," in *Applications of Signal Processing to Audio and Acoustics*, 2003 IEEE Workshop on., Oct 2003, pp. 127–130.
- [3] G. Peeters, A. L. Burthe, and X. Rodet, "Toward Automatic Music Audio Summary Generation from Signal Analysis," in *Proc. Intl. Conf. on Music Information Retrieval*, 2002, pp. 94–100.
- [4] G. Peeters and X. Rodet, "Signal-based Music Structure Discovery for Music Audio Summary Generation," in *Proc. of the Intl. Computer Music Conf. (ICMC)*, 2003, pp. 15–22.
- [5] S. Chu and B. Logan, "Music Summary Using Key Phrases," Hewlett-Packard Cambridge Research Laboratory, Cambridge MA 02139, Tech. Rep. CRL 2000/1, April 2000.
- [6] M. Cooper and J. Foote, "Automatic Music Summarization via Similarity Analysis," in *Proc. Int. Conf. Music Information Retrieval*, 2002, M. Fingerhut, Ed., 2002, pp. 81–85.
- [7] J. Glaczynski and E. Lukasik, "Automatic Music Summarization: A "Thumbnail" Approach," *Archives of Acoustics*, vol. 36, no. 2, pp. 297–309, Jan. 2011.
- [8] M. A. Bartsch and G. H. Wakefield, "Audio Thumbnailing of Popular Music using Chroma-based Representations," *IEEE Transactions on Multimedia*, vol. 7, no. 1, pp. 96–104, 2005.
- [9] P. G. Antunes, D. M. de Matos, R. Ribeiro, and I. Trancoso, "Automatic Fado Music Classification," *CoRR*, 2014, arXiv:1406.4447.
- [10] G. Erkan and D. R. Radev, "LexRank: Graph-based Lexical Centrality as Salience in Text Summarization," *Journal of Artificial Intelligence Research*, vol. 22, no. 1, pp. 457–479, Dec. 2004.
- [11] R. Mihalcea and P. Tarau, "TextRank: Bringing Order into Texts," in *Proc. of EMNLP 2004*. Barcelona, Spain: Association for Computational Linguistics, July 2004, pp. 404–411.
- [12] S. Brin and L. Page, "The Anatomy of a Large-Scale Hypertextual Web Search Engine," in *Seventh Intl. World-Wide Web Conf. (WWW 1998)*, 1998.

- [13] X. Zhu, A. B. Goldberg, J. V. Gael, and D. Andrzejewski, "Improving Diversity in Ranking using Absorbing Random Walks," *Proc. of NAACL HLT*, pp. 97–104, 2007.
- [14] K. Zechner and A. Waibel, "Minimizing word error rate in textual summaries of spoken language," in *Proc. of the 1st North American chapter of the Association for Computational Linguistics conference*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2000, pp. 186–193.
- [15] G. Murray, S. Renals, and J. Carletta, "Extractive Summarization of Meeting Recordings," in *Proc. of the 9th European Conf. on Speech Communication and Technology*, 2005, pp. 593–596.
- [16] Y. Gong and X. Liu, "Generic Text Summarization Using Relevance Measure and Latent Semantic Analysis," in *Proc. of the 24th Annual Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval*. New York, NY, USA: ACM, 2001, pp. 19–25.
- [17] J. Carbonell and J. Goldstein, "The Use of MMR, Diversity-based Reranking for Reordering Documents and Producing Summaries," in *Proc. of the 21st Annual Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval*. New York, NY, USA: ACM, 1998, pp. 335–336.
- [18] J. Steinberger and K. Jezek, "Using Latent Semantic Analysis in Text Summarization and Summary Evaluation," in *Proc.* of ISIM 2004, 2004, pp. 93–100.
- [19] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, pp. 27:1–27:27, 2011.
- [20] F. Eyben, F. Weninger, F. Gross, and B. Schuller, "Recent Developments in openSMILE, the Munich Open-source Multimedia Feature Extractor," in *Proc. of the 21st ACM Intl. Conf. on Multimedia*, 2013, pp. 835–838.
- [21] C. Sanderson, "Armadillo: An open source c++ linear algebra library for fast prototyping and computationally intensive experiments," NICTA, Australia, Tech. Rep., October 2010.
- [22] G. Tzanetakis and P. Cook, "MARSYAS: A Framework for Audio Analysis," Organised Sound, vol. 4, no. 3, pp. 169–175, Dec. 1999.
- [23] R. R. Curtin, J. R. Cline, N. P. Slagle, W. B. March, P. Ram, N. A. Mehta, and A. G. Gray, "MLPACK: A Scalable C++ Machine Learning Library," *Journal of Machine Learning Research*, vol. 14, pp. 801–805, 2013.