

Multimodal music information processing and retrieval: survey and future challenges

Federico Simonetta, Stavros Ntalampiras, Federico Avanzini
 LIM – Music Informatics Laboratory
 Department of Computer Science
 University of Milan
 Email: {name.surname}@unimi.it

Abstract—Towards improving the performance in various music information processing tasks, recent studies exploit different modalities able to capture diverse aspects of music. Such modalities include audio recordings, symbolic music scores, mid-level representations, motion and gestural data, video recordings, editorial or cultural tags, lyrics and album cover arts. This paper critically reviews the various approaches adopted in Music Information Processing and Retrieval, and highlights how multimodal algorithms can help Music Computing applications. First, we categorize the related literature based on the application they address. Subsequently, we analyze existing information fusion approaches, and we conclude with the set of challenges that Music Information Retrieval and Sound and Music Computing research communities should focus in the next years.

Index Terms—Multimodal music processing, music information retrieval, music description systems, information fusion

I. INTRODUCTION

Beginning with the oldest evidence of music notation, music has been described in several different forms [1]. Such descriptions have been used by computational systems for facilitating music information computing tasks. Interestingly, when observing the history of music, one can see how the various descriptive forms have gradually emerged with a strict dependence both on technology advancements and changes in music practices.

Initially, no written description systems for music existed besides text. Between the 6th-7th cen., Isidore of Seville, Archbishop and theologian, wrote that no melody could be written. Indeed, the first systems to memorize music were based solely on lyrics and only later some signs over the words appeared. Such notation, called *neumatic*, evolved in more complex forms, which differed from region to region. Due to the need of more powerful tools to express music features, new notation systems, called *pitch specific*, took place, such as the *alphabetic* and the *staff*-based notations. In particular, the system introduced by Guido d'Arezzo (10th-11th cen.) was particularly successful and similar conventions spread all over Europe. Music notation was now able to represent text, pitches and durations at the same time. During the following centuries, other types of symbols were introduced addressing directly the performer towards peculiar colors, or sentiment(s). At the crossing of the 16th and 17th cen., Opera was born in Italy, after a long tradition of plays, including Greek drama, medieval entertainers and renaissance popular plays (both

liturgic and profane) [2]. The tremendous success of the Opera in Italy and then in the rest of Europe, determined a fundamental way to connect music and visual arts for the future centuries. A turning point in the history of music description systems was the invention of the *phonograph cylinder* by Thomas Edison in 1877 and the *disc phonograph* diffused by Emile Berliner ten years later [3]. In the same years, Edison and the Lumière brothers invented the first devices to record video [4]. Since then, a number of technologies were born paving the way for new music description systems. With the invention of computers and the beginning of the digital era, the elaboration of sound signals highlighted the need for more abstract information characterizing audio recordings. Thus, researchers started proposing *mid-level* representations [5], with reference to *symbolic* and *physical* levels [6]. Nowadays, the availability of vast, easily accessible quantities of data, along with appropriate modern computational technologies, encourages the collection of various types of *meta-data*, which can be either *cultural* or *editorial* [7].

From a cognitive point of view, the connecting, almost evolutionary, element between the above-mentioned representations is that each one relates to a different abstraction level. Psychology, indeed, is almost unanimous in identifying an abstraction process in our music cognition [8]: we can recognize music played on different instruments, with different timings, intensity changes, various metronome markings, tonalities, tunings, background noises and so on. The different descriptions of music developed in different era or contexts, can be seen as an answer to the necessity of representing new modalities – such as the visual one – or new unrevealed abstraction levels – such as the audio recordings and the mid-symbolic levels, or the pitch specific notation compared to the neumatic one.

Aside from these historical and cognitive considerations, it is a fact that in the last two decades researchers have obtained better results through multimodal approaches in respect to single-modalities approaches [9], [10]. As Minsky said [11]:

To solve really hard problems, we'll have to use several different representations.

We argue that music processing tasks can benefit profoundly from multimodal approaches, and that a greater focus is needed by the research community in creating such a syn-

ergistic framework. A fundamental step would be the study and design of suitable algorithms through which different modalities can collaborate. Then, a particular effort should be devoted in developing the needed technologies. In fact, given the course of history summarized above, we could expect that in the future, new disparate music representations will be born.

In this paper, we review the existing literature about Music Information Retrieval techniques which exploit multiple descriptions of music to the end of *multimodal fusion* [12]. The paper is organized as follows: in section II, we give some basic definition and discuss previous reviews on similar topics to explain the categorization and the taxonomy we used. Sections III to VII describe the different tasks faced with multimodal approaches, the various features extracted, the preprocessing steps and the fusion approaches adopted in literature; in section VIII we express our idea about how the multimodal paradigm can be enhanced.

II. DEFINITIONS, TAXONOMY AND PREVIOUS REVIEWS

We have found no univocal definition of modality. In the music computing literature, authors use the word *multimodal* in two main contexts:

- in computational psychology, where *modality* refers to a human sensory channel;
- in music information retrieval, where *modality* usually refers a source of music information;

Since we are focusing on music information retrieval methods, to the purpose of the present paper, with *modality* we mean a specific way to digitize music information. Different modalities are obtained through different transducers, in different places or times, and/or belong to different media. Examples of modalities that may be associated to a single piece of music include audio, lyrics, symbolic scores, album covers, and so on.

Having defined what we mean by modality, we define *multimodal music information processing* as an MIR [13] approach which takes as input multiple modalities of the same piece of music. All the papers which we are going to discuss show methods which take as input various music representations. Conversely, we are not considering those approaches which exploit features derived through different methods from the same modality: an example is pitch, rhythmic and timbral features, when they are all derived from the audio [14]. Similarly we are not considering approaches which process multiple representations of the same modality: an example is spectrograms (treated as 2D images) and traditional time-domain acoustic features [15], which are both derived from the audio. Moreover, we do not focus on general multimodal sound processing: the idea which moves our effort is that music is characterized by the *organization* of sounds in time; thus, we are interested in exploiting this organization, which is not available in general sound processing.

One previous review on multimodal music processing was written in 2012 [12]. However, that work was more focused on a few case studies rather than on an extensive survey. The

authors recognized a distinction between “the effort of characterizing the *relationships* between the different modalities”, which they name *cross-modal processing*, and “the problem of efficiently combining the information conveyed by the different modalities”, named *multimodal fusion*. To our analysis, this distinction is useful if with *cross-modal processing* we mean the end-user systems which offer an augmented listening experience by providing the user with additional information. If this is the case, we are primarily interested in *multimodal fusion*; nevertheless, some synchronization algorithms, which are classified as *cross-modal processing* by the previous authors [12], are used as pre-processing steps in other works. Because of this ambiguous distinction, we base our classification on the performed task rather than on the processing stage – see section III.

Almost all authors dealing with multimodal information fusion talk about two approaches: *early fusion* and *late fusion*. Figure 1 shows the main difference between the two approaches: in *early fusion*, data is used “as is” in one single processing algorithm which fuse the data representation, while in *late fusion* data from each modality is first processed with specific algorithms and then all the output are merged, so that it is the output to be fused and not the data. Because of this, *early fusion* is also called *feature-level fusion*, and *late fusion* is also called *decision-level fusion*, even if features extraction and decision algorithms are not the only approaches for multimodal processing. Some reviews [9] also talk about *hybrid fusion* for multimedia analysis, but we have found no example in the music domain.

Finally, we have found useful to introduce a new diagram to represent the data flow in retrieval systems (see fig. 2). Indeed, in most of these systems, one modality is used to query a database for retrieving another modality; in such cases, no fusion exists, but just a data conversion and a similarity computation.

An exhaustive and continuously updated table, which summarizes all the works reviewed in this paper, is available online.¹

III. MULTIMODAL MUSIC PROCESSING TASKS

To date, several tasks have been experimented in multimodal approaches. We found two possible categorizations for the application level:

- **less vs more** studied tasks: some tasks have been extensively studied with a multimodal approach, such as *audio-to-score alignment*, *score-informed source separation*, *music segmentation*, *emotion* or *mood* recognition; other tasks, instead, have been little explored and are worth of more attention.
- **macro-task** based categorization: we identified 4 different macro-tasks, that are a partial re-elaboration of a previous effort [13]: *classification* of music, *synchronization* of different representations, *similarity* computation

¹Link: <https://frama.link/multimodal-MIR>

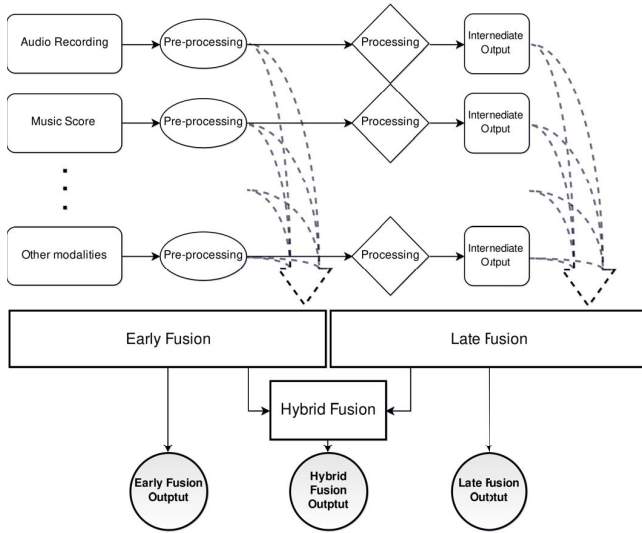


Figure 1. Diagram showing the flow of information in *early-fusion* and *late-fusion*. Early fusion process takes as input the output of the pre-processing of the various modalities, while the late fusion takes as input the output of specific processing for each modality. *Hybrid fusion*, instead, uses the output of both *early* and *late* fusion.

between two or more modalities, and *time-dependent representation*.

Figure 3 outlines all the tasks that we found in the literature. Here, instead, we are going to briefly describe each task and how it has been fulfilled by exploiting a multimodal approach.

A. Synchronization

Synchronization algorithms aim at aligning in time or space different modalities of music, i.e. creating associations between points in different modalities. They can be performed both in real-time and offline. In the real-time case, the challenge is to predict if a new event discovered in a real-time modality – e.g. an onset in the audio – corresponds to an already known event in another off-line modality – e.g. a new note in the score. Off-line synchronization, instead, is usually referred to as *alignment* and involves the fusion of multiple modalities by definition. Well-studied alignment algorithms include *audio-to-score* alignment [17], *audio-to-audio* alignment [17] and *lyrics-to-audio* alignment [18]. An interesting task is to align the audio recording to the images, without using any symbolic data [19]. Very often, alignment algorithms are a fundamental pre-processing step for other algorithms – see section IV.

B. Similarity

With *similarity*, we mean the task of computing the amount of similarity between the information content of different modalities. Often, this task has the purpose of retrieving documents from a collection through a query, which can be explicitly expressed by the user or implicitly deduced by the system. The multimodal approach, here, can exist either in the different modalities between the query and the retrieved

documents or in the query itself. A very common example of explicit queries for retrieving another modality is *query-by-humming* or *query-by-example*, in which the query is represented by an audio recording and the system retrieves the correct song; this task is usually performed with two main approaches: by using a collection of recordings in a single-modality fashion, or by exploiting multimodality with a collection of symbolic data [20], [21]. An example of *implicit* query systems, instead, are recommender systems and playlist generators, where the user is usually not aware of which specific parameters are used for the recommendations; most of the recent research in this field tries to exploit multimodal approaches – also called *hybrid* – involving *metadata*, user *context*, *audio* features [22], [23]. An emerging field in the retrieval context is the so-called *multimodal queries*, where the user can explicitly create a query by using different parameters for different modalities [16], [24]. Following this line of thought, some researchers devised and studied novel tasks in the context of multimodal music retrieval. Some example are: a system for retrieving music score images through audio queries [19]; an algorithm to retrieve the cover of a given song [25]; systems to retrieve audio recordings through symbolic queries [26], [27]; an approach to query a music video database with audio queries [28].

C. Classification

The *classification* process consists in taking as input a music document and returning one or more labels. A popular multimodal classification task is the *mood* or *emotion* recognition [29], while an emerging one is *genre* classification [30]–[37]. Both these two tasks can take advantage of audio recordings, lyrics, cover arts and meta-tags. Additionally, emotion recognition can exploit EEG data, while for genre classification one can use music video and generic text such as critic reviews. Usually, just one modality is considered in addition to audio recordings, but an interesting work [37] tries to exploit more than two modalities. Other multimodal classification tasks found in the literature are:

- *artist* identification, through lyrics and audio fusion [38];
- *derivative works* classification of youtube video through audio, video, titles and authors [39];
- *instrument* classification by exploiting audio recordings and performance video [40], [41];
- *tonic* identification, that is: given an audio recording and the note level, find the tonic [42];
- *expressive musical description*, which consists in associating a musical annotation to an audio recording by extracting features with the help of symbolic level [43].

D. Time-dependent representation

With *time-dependent representation*, we mean the creation of a time-dependent description of the music data, created by merging and processing multiple modalities. Possibly the most studied task within this family is *score-informed source separation* [17], in which symbolic music data and audio recordings of a musical ensemble are used to create different

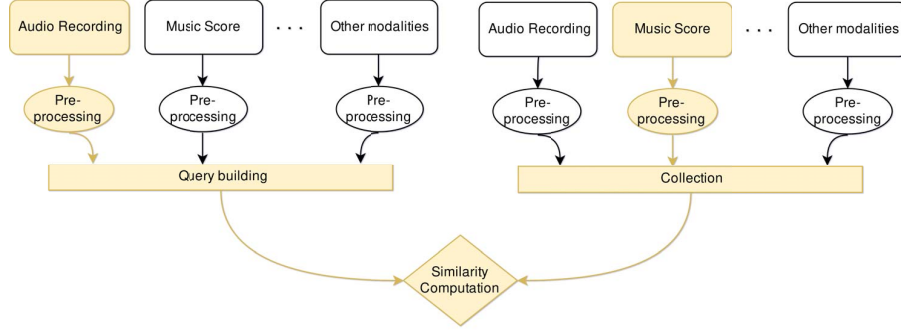


Figure 2. Multimodal retrieval: usually, the query and the collection contain different modalities, so that the diagram should be collapsed to the highlighted elements; however a more general case is possible [16], in which both the query and the collection contain multiple modalities.

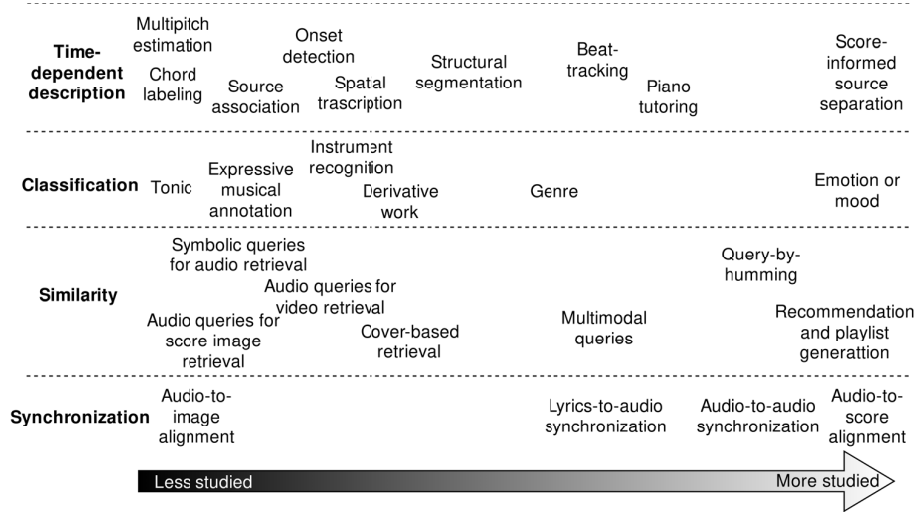


Figure 3. The tasks identified in literature, divided in 4 macro-tasks and plotted along a *less - more* studied axis. Tasks for which only one paper has been found appear at the left-side (*less studied*); at the rightmost side are tasks for which extensive surveys are already available; the other tasks are placed in the remaining space proportionally to the number of corresponding papers found in literature. All references to these tasks can be found in the discussion and in the online spreadsheet – see footnote 1. Note that labels refer to the multimodal approach at hand and not to generic MIR tasks – e.g. *genre* classification task is intended to be performed with a multimodal approach and thus it has been less studied than *emotion or mood* classification in the context of multimodal approaches.

audio recordings for each different instrument. A number of researchers have also tried to use audio and video recordings of a music performance or of a dancer to extract *beat tracking* information [44]–[48]. An emerging task is *piano tutoring*, which consists in the tracking of errors in a piano performance: to this end, the audio recording, the knowledge about the instrument timbre and the symbolic score can be exploited [49]–[55]. Less studied tasks are:

- *music segmentation*, in which audio and video, lyrics or note level can be exploited to identify the music piece structure [56]–[58];
- *spatial transcription*, that is the inference, starting from audio and video, of the note level of songs for fretted instruments, so that the resulting score includes the annotation of fingering [59], [60];
- *onset detection* through audio and performer video [61] or rhythmic structure knowledge;

- *chords* labeling, by comparing multiple audio recordings of the same work [62];
- *source association*, that is the detection of which player is active time by time by exploiting audio, video and music scores [63], [64];
- *multi-pitch* estimation, that is the transcription of parts being played simultaneously, with the help of performance video to detect play-nonplay activity of the various instruments [65].

IV. DATA PRE-PROCESSING

Data pre-processing is the elaboration of data to the end of transforming their representation to a more suitable format for the subsequent steps. We have identified a number of possible non-exclusive types of pre-processing :

- *Synchronization*: the synchronization process described in section III-A is sometime used as pre-processing step

to align multiple modalities; thus, the pre-processing itself can be multimodal. For example, in *piano tutoring* and *score-informed source separation*, an *audio-to-score* alignment is performed; *audio-to-audio* synchronization is a fundamental pre-processing step in tasks requiring comparison of multiple recordings of the same piece [62]; *audio-to-score* alignment is also used in several previously cited works [26], [27], [43], [58];

- *Feature extraction*: usually, music representations are not used as they are, but a number of features are extracted – see section V.
- Other pre-processing steps include:
 - *conversion* from one modality to the other, such as in *query-by-humming* – which includes a conversion from audio to the symbolic level – or in *audio-to-score* alignment where symbolic scores can be converted to audio through a synthesis process.
 - *feature selection* through *Linear Discriminant Analysis* (LDA) [28] or *ReliefF* [43]
 - *normalization* of the extracted features [48]
 - *source-separation* in lyrics-to-audio alignment and source association [63], [64]
 - chord labeling on audio only [62]
 - multi-pitch estimation on audio only [65]
 - video-based hand tracking [59]
 - *tf-idf*-based statistics – see section V-C – adapted for audio [38]

Finally, we think that a step worthy of a particular attention is the *conversion to a common space* of the extracted features, to make them comparable. We will talk about this step in section VI. The accompanying online table (see footnote 1) contains a short description of the pre-processing pipeline adopted in each cited paper.

V. FEATURE EXTRACTION IN MULTIMODAL APPROACHES

Various types of features can be extracted from each modality. In this section, we provide a general description for audio, video, textual and symbolic score features.

A. Audio features

This section is mainly written with reference to a previous review [66]. Audio features can be broadly subdivided in *physical* and *perceptual*.

1) *Physical features*: these can be computed in various domains, such as time, frequency or wavelet. Time-domain features can be computed directly on the digitally recorded audio signal and include *zero-crossing rate*, *amplitude*, *rhythm* and *power-based* features, such as the *volume*, the *MPEG-7 temporal centroid* or the *beat histogram*. Frequency-domain features are the richest category; they are usually computed through a Short-Time Fourier Transform (STFT) or an autoregression analysis and can be subdivided in: *autoregression-based*, *STFT-based* and *brightness*, *tonality*, *chroma* or *spectrum shape* related. Features in the Wavelet-domain are computed after a Wavelet transform, which has the advantage of being able to represent discontinuous, finite, non-periodic or

non-stationary functions. Image-domain features are computed through a graphic elaboration of the spectrogram, that is a matrix that can be represented as a one-channel image computed with the STFT; often, spectrogram is used as input for a convolutional neural network (CNN), which is trained to compute *ad-hoc* features, which lack straightforward interpretation.

2) *Perceptual features*: these try to integrate human sound perception processing in the feature extraction stage or in the elaboration of physical audio features. Most of them aim at mapping certain measurements to a perceptual-based scale and/or metrics. For example, *Mel Frequency Cepstral Coefficients* (MFCC) are derived by mapping the Fourier transform to a Mel-scale, thus improving the coherence with human perception. *Perceptual wavelet packets* [67] employ a perceptually motivated critical-band based analysis to characterize each component of the spectrum using wavelet packets. *Loudness* is computed from the Fourier transform with the aim of providing a psychophysically motivated measure of the intensity of a sound.

B. Video and image features

This section is mainly written with reference to a previous work [48]. Video features used in the music domain are very similar to visual features used in general purpose video analysis. Image features can be based on the *color space* (RGB or HSV), on *edges* detection, on the *texture* – such as the LBP –, or on the *moment* of a region. In video, motion detection is also possible and can be performed with *background detection* and *subtraction*, *frame difference* and *optical flow*. *Object tracking* has been also used to detect hand movements, for example in *piano-tutoring* applications. *Object tracking* can happen by exploiting the difference between frames of the detected object contours, by using deviations frame-to-frame of whole regions or generic features. In video, one can also detect *shots*, for example by analyzing the variation of the color histograms in the video frames, using the Kullback-Leibler distance [68] or other metrics.

In genre and mood related analysis, other features can also be exploited [69]. The use of *tempo* is essential to express emotions in video clips, and can be analyzed through features related to motion and length of video shots. Another relevant factor is lighting, that can be measured through brightness-based features. Colors have an affective meaning too, and color features are consequently useful for genre or emotion recognition.

Finally, images can also be used as they are as input of CNNs.

C. Text features

This section is written with reference to a previous review [70]. The most common text representations are based on *tf-idf*. In this context, $tf(d, t)$ is the *term frequency* and is computed as the number of occurrences of a term t in a document d . Instead, $idf(d, t)$ is a short for *inverse document frequency* and is needed to integrate the discrimination power

of the term t for the document d , considering the whole collection; it is related to the inverse ratio between the number of documents containing t at least once and the total number of documents in the considered collection:

$$\text{idf} = \frac{\text{docs in collection}}{\text{docs containing } t} \quad (1)$$

Usually, *tf-idf* takes the following form:

$$\text{tf-idf}(d, t) = \text{tf}(d, t) \times \log[\text{idf}(d, t)] \quad (2)$$

Features based on *tf-idf* are often used in Bag-of-Words (BoW) models, where each document is represented as a list of words, without taking care of the cardinality and order of words. In order to make BoW and *tf-idf* models effective, a few preliminary steps are usually performed, such as *punctuation* and *stop-words* removal and *stemming*. More sophisticated methods are also available, allowing for topic- or semantics-based analysis, such as *Latent Dirichlet Allocation* (LDA), *Latent Semantic Analysis* (LSA), *Explicit Semantic Analysis* (ESA) [71] and CNN feature extraction.

For lyrics analysis, other types of features can be extracted, like rhymes or positional features. Finally, when the available text is limited, one can extend it with a semantic approach consisting in *knowledge boosting* [37].

D. Symbolic score features

Symbolic music scores have been rarely used in feature extraction approaches. Most of the papers which deal with symbolic scores use MIDI-derived representations, such as the pianoroll [17] or inter-onset intervals (IOI) [58]. To the end of audio-symbolic comparison, one can compute chromograms, that are also computable from the audio modality alone. However a number of representation exist and have been tested in Music Information Retrieval applications, such as *pitch histograms*, *Generalized Pitch Interval Representation* (GPIR), *Spiral Array*, *Rizo-Iñesta trees*, *Pinto graphs*, *Orio-Rodà graphs* and others. A brief review of the music symbolic level representations is provided in a previous work [72].

VI. CONVERSION TO COMMON SPACE

The conversion of the extracted features to a *common space* is often a mandatory step in *early fusion* approaches. Nevertheless, almost no authors emphasize this aspect. Thus, we think that greater attention should be posed on this step of the pre-processing pipeline.

The conversion to a common space consists in the mapping of the features coming from different modalities to a new space where they are comparable. This can be needed in single-modality approaches too, when the features refer to very different characteristics of the signal. Indeed, many papers describe techniques which include a mapping of the features to a common space, both in the pre-processing and in the processing stages, but no particular attention to the conversion itself. Common methods include:

- *normalization*, that is the most basic approach;
- *conversion from one modality to another*, so that features can be computed in the same units;

- *machine learning algorithms* such as CNNs or SVMs: SVMs compute the best parameters for a kernel function that is used to transform the data into a space where they are more easily separable; CNNs, instead, can be trained to represent each input modality in a space such that the last network layers can use as input the concatenation of these representations;
- *dimensionality reduction* algorithms, which usually search for a new space where data samples are representable with a fewer number of dimensions without losing the ability to separate them; examples are *Principal Component Analysis* (PCA) and *Linear Discriminant Analysis* (LDA).

It must be said that some types of features are suitable for multimodal fusion without any conversion step. For example, *chroma features* can be computed from both the audio recordings and the symbolic scores and thus can be compared with no additional processing.

A possible categorization of the conversion to common space methods is between *coordinated* and *joint*: in the former type, the mapping function takes as input a unimodal representation, while in the latter type it takes as input a multimodal representation [10]. In other words, *coordinated* conversion learns to map each modality to a new space trying to minimize the distance between the various descriptions of the same object, while *joint* conversion learns the best mapping function which uses all the modalities and optimizes the subsequent steps – e.g. SVM.

VII. INFORMATION FUSION APPROACHES

Two major information fusion approaches exist: *early fusion* and *late fusion* – see fig. 1. Some authors also report a *hybrid* approach [9], which consists in fusing information both in a *early* and *late* fashion and in adding a further step to fuse the output of the two approaches. Nevertheless, we did not find any existing application to the music domain. Before discussing in detail the two approaches, we recall that no fusion is usually needed in *similarity* tasks, but just a comparison of the various modalities and, thus, a conversion to a common space. The accompanying online table (see footnote 1) contains a short description of the fusion approach used in all the cited papers. To our understanding the main difference between *early* and *late* fusion is about their efficiency and ease of development; however authors disagree about which one is the more effective.

A. Early fusion

Early fusion consists in the fusion of the features of all the modalities, using them as input in one single processing algorithm. Although the development of such techniques is more straightforward, they need a more careful treatment because the features extracted from various modalities are not always directly comparable.

To the end of *synchronization*, the most used approach exploits Dynamic Time Warping algorithms (DTW) [73]. DTW is a well-known technique based on a similarity matrix

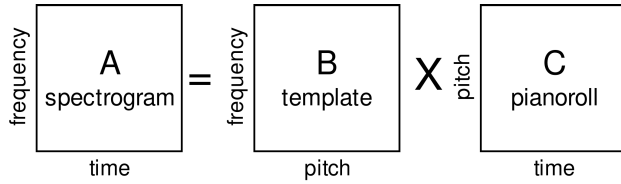


Figure 4. Exemplification of Non-negative Matrix Factorization for music transcription.

between two sorted sets of points, for example two time-sequences. By using a dynamic programming algorithm, one can exploit the similarity matrix to find the best path which connects the first point in one modality to the last point in the same modality and which satisfies certain conditions. This path will indicate the corresponding points between the two modalities. Other common methods for synchronization purposes are Hidden Markov Models (HMMs) [18], [58] where hidden states represent points in one modality and observations represent points in a second modality; this is particularly effective for real-time alignment or generic sequence fusion such as in *time-dependent descriptions*.

Aside HMMs, many additional machine learning [74] approaches are used to perform early fusion: Support Vector Machines (SVMs), Gaussian Mixture Models (GMMs), Convolutional Neural Networks (CNNs) and Particle Filters are the most used techniques.

Another interesting method is Non-negative Matrix Factorization (NMF), through which audio and symbolic scores can be exploited to the end of precise performance transcription, as in *score-informed source separation* and *piano tutoring* applications [17]. In NMF, a matrix A is decomposed in two components C and B , so that $A = B \times C$. If A is a spectrogram and B is a *template matrix* dependent on the instrumentation, then we can think to C as a pianoroll matrix – see fig. 4. Consequently, one can use an optimization algorithm to minimize the function $f(B, C) = A - B \times C$, by initializing C with a symbolic score; at the end of the optimization, C will be a precise transcription of the performance contained in A .

Finally, feature fusion can also happen at the feature selection stage [38], [43].

B. Late fusion

Unlike *early fusion*, *late fusion* is the fusion of the output of various *ad-hoc* algorithms, one for each modality. It is also called *decision-level* fusion, even if a decision process is not mandatory. The main advantage of *late fusion* is that it allows for a more adjustable processing of each modality. However, it is usually more demanding in terms of development costs.

In *classification* and *time-dependent description* tasks, the most used types of late fusion are *rule-based*. Rules can include voting procedures [62], [64], linear and weighted combinations [41], [75], maximum and minimum operations [41], [75]. Many authors have developed sophisticated algorithms to execute this step, such as in *beat tracking*, *piano tutoring*

and *structural segmentation* [57], multi-pitch estimation [65] and tonic identification [42].

In *synchronization* tasks, instead, no *late-fusion* approach is possible, since the task consists in creating associations between points in different modalities and, thus, the process must take as input all the modalities, eventually in some common representation.

VIII. FUTURE DIRECTIONS

In this paper, we have analyzed the literature on multimodal music information processing and retrieval. Based on our study, we propose the following concluding remarks.

First of all, we note the unavailability of datasets of suitable size. This issue is usually addressed with various methods such as *co-learning* approaches [10], which has the side-effect of impoverishing the goodness of the developed algorithms. Although a few datasets have been recently created [37], [76]–[78], a great effort should still be carried out in this direction. Indeed, existing multimodal music datasets are usually characterized by limited sizes and only rarely include a wide range of modalities. However an exhaustive list of the available datasets is out of the scope of this paper. We argue that this limit is due to two main reasons: first, the precise alignment of various modalities is a hard computational task and should be controlled by human supervision; second, no largely adopted standard exists for multimodal music representation. About the first point, more effort should be devoted to the development of algorithms for the alignment of various sequences. The representation of the intrinsic music multimodality, instead, is faced by the *IEEE 1599*² standard and the *Music Encoding Initiative*³; moreover, the *W3C* group is currently working on a new standard with the purpose of enriching *MusicXML* with multimodal information⁴. The course of history described in section I and the rapid technology advancements of our times suggest that new representation modalities could be needed in the future and that multimodal representation standards should also focus on this challenge.

Another challenge that multimodal music researchers should face in the next years is the exploration of various techniques already used in multimodal processing of multimedia data, that have not been tested in the musical domain. According to previous surveys [9], [10], multimodal methods never applied to the music domain include: the hybrid approach, the Dempster-Shafer theory, Kalman filters, the maximum entropy model, Multiple Kernel Learning and Graphical Models. Moreover, we have found only one paper in which the information fusion happens during the feature extraction itself [58] and not afterwards. This approach should be explored more deeply.

Finally, we suggest that the conversion to common space should be more rigorously addressed. To this end, transfer learning technologies could be explored towards forming a

²IEEE 1599 website: <http://ieee1599.lim.di.unimi.it/>

³MEI website: <https://music-encoding.org/>

⁴W3C music notation group website: <https://www.w3.org/community/music-notation/>

synergistic feature space able to meaningfully represent multiple modalities [79], [80]. Such a direction may include the use of an existing feature space characterizing a specific modality, or the creation of a new one where multiple modalities are represented. Such a space could satisfy several desired properties, such as sparseness, reduced dimensionality, and so on.

REFERENCES

- [1] I. D. Bent, D. W. Hughes, R. C. Provine, R. Rastall, A. Kilmer, D. Hiley, J. Szendrei, T. B. Payne, M. Bent, and G. Chew, "Notation," in *Grove Music Online*. Oxford University Press, 2001.
- [2] H. M. Brown, E. Rosand, R. Strohm, M. Noiray, R. Parker, A. Whittall, R. Savage, and B. Millington, "Opera (i)," in *Grove Music Online*. Oxford University Press, 2001.
- [3] G. Mumma, H. Rye, B. Kernfeld, and C. Sheridan, "Recording," in *Grove Music Online*. Oxford University Press, 2003.
- [4] D. A. Cook and R. Sklar, "History of the motion picture," in *Britannica Academic*. Encyclopædia Britannica, 2018.
- [5] T. Kitahara, "Mid-level representations of musical audio signals for music information retrieval," in *Advances in Music Information Retrieval*. Springer Berlin Heidelberg, 2010, pp. 65–91.
- [6] H. Vinet, "The Representation Levels of Music Information," in *Computer Music Modeling and Retrieval*, U. K. Wil, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, pp. 193–209.
- [7] F. Pachet, "Musical metadata and knowledge management," in *Encyclopedia of Knowledge Management, Second Edition*, D. G. Schwartz and D. Te'eni, Eds. IGI Global, 2005, pp. 1192–1199.
- [8] D. Deutsch, *The Psychology of Music (third Edition)*, third edition ed., D. Deutsch, Ed. Academic Press, 2013.
- [9] P. K. Atrey, M. A. Hossain, A. E. Saddik, and M. S. Kankanhalli, "Multimodal fusion for multimedia analysis: A survey," *Multimedia Systems*, vol. 16, no. 6, pp. 345–379, Apr. 2010.
- [10] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, "Multimodal machine learning: A survey and taxonomy," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2018.
- [11] M. Minsky, "Logical versus analogical or symbolic versus connectionist or neat versus scruffy," *AI Magazine*, vol. 12, no. 2, pp. 34–51, 1991.
- [12] S. Essid and G. Richard, "Fusion of multimodal information in music content analysis," in *Multimodal Music Processing*. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik GmbH, Wadern/Saarbruecken, Germany, 2012.
- [13] M. Schedl, E. Gómez, and J. Urbano, "Music information retrieval: Recent developments and applications," *Foundations and Trends in Information Retrieval*, vol. 8, no. 2-3, pp. 127–261, 2014.
- [14] G. Tzanetakis and P. Cook, "Musical genre classification of audio signals," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 5, pp. 293–302, July 2002.
- [15] L. Nanni, Y. M. G. Costa, A. Lumini, M. Y. Kim, and S. Baek, "Combining visual and acoustic features for music genre classification," *Expert Syst. Appl.*, vol. 45, pp. 108–117, 2016.
- [16] L. Zhonghua, "Multimodal music information retrieval: From content analysis to multimodal fusion," Ph.D. dissertation, 2013.
- [17] M. Müller, *Fundamentals of Music Processing: Audio, Analysis, Algorithms, Applications*, 1st ed. Springer Publishing Company, Incorporated, 2015.
- [18] H. Fujihara and M. Goto, "Lyrics-to-Audio Alignment and its Application," in *Multimodal Music Processing*, ser. Dagstuhl Follow-Ups, M. Müller, M. Goto, and M. Schedl, Eds. Dagstuhl, Germany: Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2012, vol. 3, pp. 23–36.
- [19] M. Dorfer, A. Arzt, and G. Widmer, "Learning audio-sheet music correspondences for score identification and offline alignment," in *Proceedings of the 18th International Society for Music Information Retrieval Conference, ISMIR 2017, Suzhou, China*, S. J. Cunningham, Z. Duan, X. Hu, and D. Turnbull, Eds., 2017, pp. 115–122.
- [20] A. Kotsifakos, P. Papapetrou, J. Hollmén, D. Gunopulos, and V. Athitsos, "A survey of query-by-humming similarity methods," in *Proceedings of the 5th International Conference on Pervasive Technologies Related to Assistive Environments*, ser. PETRA '12. New York, NY, USA: ACM, 2012.
- [21] MIREX Community. (2016) 2016:query by singing/humming. [Online]. Available: https://www.music-ir.org/mirex/wiki/2016:Query_by_Singing/Humming
- [22] G. Bonnin and D. Jannach, "Automated generation of music playlists: Survey and experiments," *ACM Comput. Surv.*, vol. 47, no. 2, pp. 26:1–26:35, Nov. 2014.
- [23] P. Deshmukh and G. Kale, "A Survey of Music Recommendation System," *International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRCSEIT)*, vol. 3, no. 3, pp. 1721–1729, Mar. 2018.
- [24] M. Müller, A. Arzt, S. Balke, M. Dorfer, and G. Widmer, "Cross-modal music retrieval and applications: An overview of key methodologies," *IEEE Signal Processing Magazine*, vol. 36, no. 1, pp. 52–62, Jan. 2019.
- [25] A. Correya, R. Hennequin, and M. Arcos, "Large-Scale Cover Song Detection in Digital Music Libraries Using Metadata, Lyrics and Audio Features," *Arxiv E-prints*, Aug. 2018.
- [26] I. S. H. Suyoto, A. L. Uittenboger, and F. Scholer, "Searching musical audio using symbolic queries," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, no. 2, pp. 372–381, Feb. 2008.
- [27] S. Balke, V. Arifi-Müller, L. Lamprecht, and M. Müller, "Retrieving audio recordings using musical themes," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Mar. 2016, pp. 281–285.
- [28] O. Gillet, S. Essid, and G. Richard, "On the correlation of automatic audio and visual segmentations of music videos," *Ieee Transactions on Circuits and Systems for Video Technology*, vol. 17, no. 3, pp. 347–355, 2007-03.
- [29] Y. E. Kim, E. M. Schmidt, R. Migneco, B. G. Morton, P. Richardson, J. J. Scott, J. A. Speck, and D. Turnbull, "State of the art review: Music emotion recognition: A state of the art review," in *Proceedings of the 11th International Society for Music Information Retrieval Conference, ISMIR 2010, Utrecht, Netherlands, August 9-13, 2010*, J. S. Downie and R. C. Veltkamp, Eds. International Society for Music Information Retrieval, 2010, pp. 255–266.
- [30] T. Li and M. Ogihara, "Music artist style identification by semi-supervised learning from both lyrics and content," in *Proceedings of the 12th Annual ACM International Conference on Multimedia*, ser. MULTIMEDIA '04. New York, NY, USA: ACM, 2004, pp. 364–367.
- [31] R. Mayer, R. Neumayer, and A. Rauber, "Combination of audio and lyrics features for genre classification in digital audio collections," in *Proceedings of the 16th ACM International Conference on Multimedia*, ser. MM '08. New York, NY, USA: ACM, 2008, pp. 159–168.
- [32] —, "Rhyme and style features for musical genre classification by song lyrics," in *ISMIR 2008, 9th International Conference on Music Information Retrieval*, Drexel University, Philadelphia, PA, USA, September 14-18, 2008, 2008, pp. 337–342.
- [33] R. Mayer and A. Rauber, *Multimodal Aspects of Music Retrieval: Audio, Song Lyrics – and Beyond?* Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 333–363.
- [34] C. Zhen and J. Xu, "Multi-modal music genre classification approach," in *Proc. 3rd Int. Conf. Computer Science and Information Technology*, vol. 8, Jul. 2010, pp. 398–402.
- [35] R. Mayer and A. Rauber, "Musical genre classification by ensembles of audio and lyrics features," in *Proceedings of the 12th International Society for Music Information Retrieval Conference (ISMIR 2011)*. University of Miami, 2011, pp. 675–680, vortrag: 12th International Society for Music Information Retrieval Conference (ISMIR 2011), Miami; 2011-10-24 – 2011-10-28.
- [36] A. Schindler and A. Rauber, "Harnessing music-related visual stereotypes for music information retrieval," *ACM Transactions on Intelligent Systems and Technology*, vol. 8, no. 2, pp. 1–21, Oct. 2016.
- [37] S. Oramas, F. Barbieri, O. Nieto, and X. Serra, "Multimodal deep learning for music genre classification," *Transactions of the International Society for Music Information Retrieval*, vol. 1, no. 1, pp. 4–21, 2018.
- [38] K. Aryafar and A. Shokoufandeh, "Multimodal music and lyrics fusion classifier for artist identification," in *2014 13th International Conference on Machine Learning and Applications*. IEEE, Dec. 2014.
- [39] J. B. L. Smith, M. Hamasaki, and M. Goto, "Classifying derivative works with search, text, audio and video features," in *Proc. IEEE Int. Conf. Multimedia and Expo (ICME)*, Jul. 2017, pp. 1422–1427.
- [40] O. Slizovskaia, E. Gómez, and G. Haro, "Musical instrument recognition in user-generated videos using a multimodal convolutional neural network architecture," in *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval - ICMR '17*. ACM Press, 2017.
- [41] A. Lim, K. Nakamura, K. Nakadaï, T. Ogata, and H. G. Okuno, "Audio-visual musical instrument recognition," 2011.

- [42] S. Sentürk, S. Gulati, and X. Serra, "Score informed tonic identification for makam music of turkey," in *Proceedings of the 14th International Society for Music Information Retrieval Conference, ISMIR 2013, Curitiba, Brazil, November 4-8, 2013*, A. de Souza Britto Jr., F. Gouyon, and S. Dixon, Eds., 2013, pp. 175–180.
- [43] P.-C. Li, L. Su, Y.-H. Yang, and A. W. Y. Su, "Analysis of expressive musical terms in violin using score-informed and expression-based audio features," in *Proceedings of the 16th International Society for Music Information Retrieval Conference, ISMIR 2015, Málaga, Spain, October 26-30, 2015*, M. Müller and F. Wiering, Eds., 2015, pp. 809–815.
- [44] G. Weinberg, A. Raman, and T. Mallikarjuna, "Interactive jamming with shimon: A social robotic musician," in *Proceedings of the 4th ACM/IEEE International Conference on Human Robot Interaction*, ser. HRI '09. New York, NY, USA: ACM, 2009, pp. 233–234.
- [45] A. Lim, T. Mizumoto, L. Cahier, T. Otsuka, K. Takahashi, K. Komatani, T. Ogata, and H. G. Okuno, "Robot musical accompaniment: integrating audio and visual cues for real-time synchronization with a human flutist," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems*, Oct. 2010, pp. 1964–1969.
- [46] T. Itoharu, T. Otsuka, T. Mizumoto, T. Ogata, and H. G. Okuno, "Particle-filter based audio-visual beat-tracking for music robot ensemble with human guitarist," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems*, Sep. 2011, pp. 118–124.
- [47] D. R. Berman, "AVISARME: Audio Visual Synchronization Algorithm for a Robotic Musician Ensemble," Ph.D. dissertation, University of Maryland, 2012.
- [48] M. Ohkita, Y. Bando, Y. Ikemiya, K. Itoyama, and K. Yoshii, "Audio-visual beat tracking based on a state-space model for a music robot dancing with humans," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems (IROS)*, Sep. 2015, pp. 5555–5560.
- [49] S. Wang, S. Ewert, and S. Dixon, "Identifying missing and extra notes in piano recordings using score-informed dictionary learning," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 10, pp. 1877–1889, 2017.
- [50] T. Fukuda, Y. Ikemiya, K. Itoyama, and K. Yoshii, "A score-informed piano tutoring system with mistake detection and score simplification," in *Proc of the Sound and Music Computing Conference (SMC)*. Zenodo, 2015.
- [51] E. Benetos, A. Klapuri, and S. Dixon, "Score-informed transcription for automatic piano tutoring," in *European Signal Processing Conference*, 2012, pp. 2153–2157.
- [52] O. Mayor, J. Bonada, and A. Lascos, "Performance analysis and scoring of the singing voice," in *AES 35th International Conference: Audio for Games*, 2009.
- [53] W. Tsai and H. Lee, "Automatic evaluation of karaoke singing based on pitch, volume, and rhythm features," and *Language Processing IEEE Transactions on Audio, Speech*, vol. 20, no. 4, pp. 1233–1243, May 2012.
- [54] J. Abeßer, J. Hasselhorn, C. Dittmar, A. Lehmann, and S. Grollmisch, "Automatic quality assessment of vocal and instrumental performances of ninth-grade and tenth-grade pupils," in *International Symposium on Computer Music Multidisciplinary Research*, 2013.
- [55] J. Devaney, M. I. Mandel, and I. Fujinaga, "A study of intonation in three-part singing using the automatic music performance analysis and comparison toolkit (AMPACT)," in *Proceedings of the 13th International Society for Music Information Retrieval Conference, ISMIR 2012, Mosteiro S.Bento Da Vitória, Porto, Portugal, October 8-12, 2012*, F. Gouyon, P. Herrera, L. G. Martins, and M. Müller, Eds. FEUP Edições, 2012, pp. 511–516.
- [56] Y. Zhu, K. Chen, and Q. Sun, "Multimodal content-based structure analysis of karaoke music," in *Proceedings of the 13th Annual ACM International Conference on Multimedia*, ser. MULTIMEDIA '05. New York, NY, USA: ACM, 2005, pp. 638–647.
- [57] 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 *2009 IEEE International Symposium on Circuits and Systems*. IEEE, May 2009.
- [58] J. Gregorio and Y. Kim, "Phrase-level audio segmentation of jazz improvisations informed by symbolic data," in *Proceedings of the 17th International Society for Music Information Retrieval Conference, ISMIR 2016, New York City, United States*, 2016, pp. 482–487.
- [59] M. Paleari, B. Huet, A. Schutz, and D. Slock, "A multimodal approach to music transcription," in *Proc. 15th IEEE Int. Conf. Image Processing*, Oct. 2008, pp. 93–96.
- [60] A. Hrybyk, "Combined audio and video analysis for guitar chord identification," in *11th International Society for Music Information Retrieval Conference (ISMIR 2010)*, 2010.
- [61] B. Marengo, M. Fuentes, F. Lanzaro, M. Rocamora, and A. Gómez, "A multimodal approach for percussion music transcription from audio and video," *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, Jan. 2015.
- [62] V. Konz and M. Müller, "A cross-version approach for harmonic analysis of music recordings," in *Multimodal Music Processing*, ser. Dagstuhl Follow-Ups, M. Müller, M. Goto, and M. Schedl, Eds. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, Germany, 2012, vol. 3, pp. 53–72.
- [63] B. Li, C. Xu, and Z. Duan, "Audiovisual source association for string ensembles through multi-modal vibrato analysis," *Proc. Sound and Music Computing (smc)*, 2017.
- [64] B. Li, K. Dinesh, Z. Duan, and G. Sharma, "See and listen: Score-informed association of sound tracks to players in chamber music performance videos," in *Proc. Speech and Signal Processing (ICASSP) 2017 IEEE Int. Conf. Acoustics*, Mar. 2017, pp. 2906–2910.
- [65] K. Dinesh, B. Li, X. Liu, Z. Duan, and G. Sharma, "Visually informed multi-pitch analysis of string ensembles," in *Proc. Speech and Signal Processing (ICASSP) 2017 IEEE Int. Conf. Acoustics*, Mar. 2017, pp. 3021–3025.
- [66] F. Alías, J. Socoró, and X. Sevillano, "A review of physical and perceptual feature extraction techniques for speech, music and environmental sounds," *Applied Sciences*, vol. 6, no. 5, p. 143, May 2016.
- [67] S. Ntalampiras, I. Potamitis, and N. Fakotakis, "Exploiting temporal feature integration for generalized sound recognition," *EURASIP Journal on Advances in Signal Processing*, vol. 2009, no. 1, p. 807162, Dec 2009.
- [68] D. Brezeale and D. J. Cook, "Automatic video classification: A survey of the literature," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 38, no. 3, pp. 416–430, May 2008.
- [69] S. Wang and Q. Ji, "Video affective content analysis: A survey of state-of-the-art methods," *IEEE Transactions on Affective Computing*, vol. 6, no. 4, pp. 410–430, Oct. 2015.
- [70] W. B. Croft, D. Metzler, and T. Strohman, *Search engines: Information retrieval in practice*. Pearson Education, Inc., 2015, vol. 283.
- [71] E. Gabrilovich and S. Markovitch, "Computing semantic relatedness using wikipedia-based explicit semantic analysis," in *IJCAI*, vol. 7, 2007, pp. 1606–1611.
- [72] F. Simonetta, "Graph based representation of the music symbolic level. A music information retrieval application," Master's thesis, Università di Padova, Apr. 2018.
- [73] M. Müller, "Dynamic Time Warping," in *Information Retrieval for Music and Motion*. Springer Berlin Heidelberg, 2007, pp. 69–84.
- [74] C. M. Bishop, *Pattern recognition and machine learning*, 5th Edition, ser. Information science and statistics. Springer, 2007.
- [75] N. Degara, A. Pena, M. E. P. Davies, and M. D. Plumbley, "Note onset detection using rhythmic structure," in *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, Mar. 2010, pp. 5526–5529.
- [76] G. Meseguer-Brocal, A. Cohen-Hadria, and P. Geoffroy, "Dali: A large dataset of synchronized audio, lyrics and notes, automatically created using teacher-student machine learning paradigm," in *19th International Society for Music Information Retrieval Conference, ISMIR*, Ed., Sep. 2018.
- [77] E. Maestre, P. Papiotis, M. Marchini, Q. Llimona, O. Mayor, A. Pérez, and M. M. Wanderley, "Enriched multimodal representations of music performances: Online access and visualization," *Ieee Multimedia*, vol. 24, no. 1, pp. 24–34, Jan. 2017.
- [78] B. Li, X. Liu, K. Dinesh, Z. Duan, and G. Sharma, "Creating a multi-track classical music performance dataset for multi-modal music analysis: Challenges, insights, and applications," *IEEE Transactions on Multimedia*, p. 1, 2018.
- [79] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, oct 2010.
- [80] S. Ntalampiras, "A transfer learning framework for predicting the emotional content of generalized sound events," *The Journal of the Acoustical Society of America*, vol. 141, no. 3, pp. 1694–1701, mar 2017.