# Valerio Velardo,\* Mauro Vallati,† and Steven Jan\*

\*School of Music, Humanities and Media

<sup>†</sup>School of Computing and Engineering

\*†University of Huddersfield

Queensgate, Huddersfield, HD1 3DH, UK velardovalerio@gmail.com {m.vallati, s.b.jan}@hud.ac.uk

# Symbolic Melodic Similarity: State of the Art and Future Challenges

**Abstract:** Fostered by the introduction of the Music Information Retrieval Evaluation Exchange (MIREX) competition, the number of systems that calculate symbolic melodic similarity has recently increased considerably. To understand the state of the art, we provide a comparative analysis of existing algorithms. The analysis is based on eight criteria that help to characterize the systems, highlighting strengths and weaknesses. We also propose a taxonomy that classifies algorithms based on their approach. Both taxonomy and criteria are fruitfully exploited to provide input for new, forthcoming research in the area.

The advent of the Internet has made a large quantity of audio and symbolic musical data freely available. The analysis of these data can provide useful insights into several aspects of music. By comparing many musical pieces, it is possible to abstract relevant rules and processes that characterize a particular style. Also, the analysis of large databases can improve our understanding of the generative process, shedding light on the evolutionary path undergone by music over time. To capitalize upon the significant body of knowledge currently stored within online music data sets, a number of reliable and efficient automatic tools have been developed over the last decades. Melodic similarity-detection algorithms are an instance of such tools. When used on online musical data sets, they can provide valuable information on intra- and interwork melodic relationships and on the underlying melodic structures of the pieces analyzed.

Given two or more sequences of notes, symbolic melodic similarity (SMS) aims to evaluate their degree of likeness, as human listeners are able to do. This task has relevance both within the academy and in industry. For instance, beyond the purely academic benefits of identifying the degree of likeness between musical pieces and composer practices afforded by melodic similarity systems, the detection of plagiarism constitutes an example of a practical application of this task with clear legal and commercial implications. Many algorithms for judging melodic similarity have been introduced

Computer Music Journal, 40:2, pp. 70–83, Summer 2016 doi:10.1162/COMJ\_a\_00359 © 2016 Massachusetts Institute of Technology.

over the years. Even though such tools perform essentially the same task, they may be based on theories and methods that belong to radically different disciplines. For example, there are some algorithms based on principles from music theory, others based on cognitive constraints, and others that implement notions from pure mathematics.

Thanks to the "Symbolic Melodic Similarity" track of the Music Information Retrieval Evaluation Exchange (MIREX) competition (Downie 2008), introduced in 2005, the number of tools in this field has increased dramatically. The last published surveys on SMS, however, only consider algorithms developed up to 2004 (Müllensiefen and Frieler 2006; Hofmann-Engl 2010). This lack of review of the state of the art is the first motivation for this article. Accessibility is a second motivation. The literature on SMS is distributed across many different sources, which cover numerous topics from computer science to music theory. This survey brings together recent studies on SMS, describing them in a concise way so that researchers can form an initial overview of the approaches used by other scholars. Our article proposes a highly modular taxonomy that allows the effective categorization of techniques according to the approach they exploit. We identify eight criteria, which summarize the most relevant aspects of each melodic similarity algorithm. By analyzing the state of the art, we are also able to provide guidelines and recommendations for a further development of the field. Moreover, the aforementioned taxonomy and the criteria facilitate the classification and comparison of future systems.

The article is structured as follows. First, we outline relevant background information and related

works. Second, we outline the identified criteria. Next, the taxonomy is presented and the considered systems are briefly described and categorized. Then, algorithms are compared with regard to the eight criteria. Finally, by synthesizing the knowledge gained with the analysis, we propose new directions for research and offer conclusions.

# **Background**

Stephen Downie defined music information retrieval (MIR) as a "multidisciplinary research endeavor that strives to develop innovative content-based searching schemes, novel interfaces, and evolving networked delivery mechanisms in an effort to make the world's vast store of music accessible to all" (Downie 2004, p. 12). The interdisciplinary environment of MIR encompasses many fields, such as computer science, psychology, musicology, music cognition, and signal processing. Music information retrieval combines these disciplines to create realworld applications that are capable of extracting relevant information from music. Its techniques have been applied to solve a large number of tasks, such as music recommendation, automatic music transcription, and track separation. There are two ways MIR systems can represent music: audio and symbolic. Systems that adopt audio representation directly encode musical information through digital audio formats such as WAV and MP3. Applications that use symbolic representation are usually based on MIDI and MusicXML formats. Symbolic encoding allows the system to manipulate musical items (such as notes that each have a specified pitch and duration) without recourse to signal processing, while having a clear representation of the composition.

Symbolic melodic similarity is a central issue of MIR. Many applications exploit musical similarity to retrieve pieces from a database, to perform musical analysis, and to categorize music. Essentially, all systems that are based on melodic similarity try to find musical utterances that match the information needed by the user, expressed in a query. Applications that evaluate melodic similarity improve the accessibility of musical databases, allowing

users to efficiently retrieve the musical information they need. Furthermore, such systems can enhance the understanding of the structure of music itself. Indeed, musicologists can exploit applications to track stylistic traits of musical pieces and to trace the occurrences of musical patterns both within and between musical works. This can deepen our understanding of musical style, while actively promoting a new quantitative, empirical analytical approach among musicologists and music theorists.

A major issue with melodic similarity is that assessing the likeness between two musical phrases is an extremely difficult process, which reflects the cognitive complexity of the task. Most of the complexity associated with melodic similarity detection arises from the multidisciplinary nature of this process. Melodic similarity spans several elements of music theory, ethnomusicology, cognitive science, and computer science, all of which have to be considered simultaneously. Music theory suggests how to identify syntactically relevant musical structures. Ethnomusicology accounts for the variety and the cultural dependency of melodies from distinct geographical regions. Music cognition describes the basic cognitive processes that humans deploy to recognize melodies as similar. Finally, computer science affords a means to create "intelligent" computational systems able to embed the insights provided by the other fields. Although melodic similarity has been extensively investigated in all of these disciplines, there is no agreed, clear-cut definition of the field yet. Even scholars from the same background disagree on the ambit and methodologies of melodic similarity.

This survey focuses on MIR systems performing melodic similarity analysis that have been presented since 2004. Systems introduced before have already been reviewed (Müllensiefen and Frieler 2006; Hofmann-Engl 2010). For the sake of completeness, we briefly report the main strategies developed before 2004, which strongly affected the later research environment. In 1996, McNab and coworkers introduced an algorithm based on the edit distance between motives (McNab et al. 1996). Emilios Cambouropoulos (1998) devised a system based on melodic contrast that dynamically creates motivic categories containing similar melodic structures. In

the same year, Donncha O'Maidín (1998) developed a system that exploits the distance in pitch between two melodies, weighted through correlation and difference coefficients. Later, Downie (1999) presented a system that assesses melodic similarity based on *n*-grams. Finally, Meek and Birmingham (2002) developed an algorithm that exploits hidden Markov chains.

MIREX was established in 2005 (Downie 2008). This organization runs an important annual competition that aims to compare state-of-the-art algorithms and systems relevant for MIR. One of the several tracks of the contest is SMS, and in this respect MIREX helps to facilitate cross-fertilization among researchers, while constantly fostering the improvement of the techniques and strategies adopted in melodic similarity research. MIREX has become the main forum for researchers and practitioners interested in evaluating and comparing algorithms that span several tasks of MIR. As a consequence, many systems considered in this article have been tested or trained on MIREX benchmarks. Indeed, this competition has been providing the MIR community with a large set of useful benchmarks for ten years. [Editor's note: This article was written prior to the 2015 MIREX competition.]

#### **Criteria**

We have formulated eight criteria that are useful for analyzing melodic similarity systems. They have been designed to investigate the functionality of systems from a number of different perspectives: flexibility, similarity evaluation, training, and validation. It should be noted that the proposed criteria are not meant to be used for evaluating the considered systems, but for characterizing them.

- 1. *Polyphony* indicates the ability of a system to deal with musical pieces that include one or more voices. Monophonic systems can evaluate melodic similarity only between pieces containing a single melodic line.
- 2. *Scope* denotes the musical genres and styles a system is able to investigate. Some methods

- have a general scope, being able to analyze a wide range of musical styles. Others evaluate melodic similarity only in a specific type of music (e.g., folk, classical, or pop).
- 3. Similarity function indicates the methodology used to calculate melodic similarity. To calculate melodic similarity, algorithms rely on functions that provide a quantitative measure. Usually, such functions are based on well-known geometrical, mathematical, cognitive, or musical notions.
- 4. Musical parameters list the parameters that have been taken into account by a system. Because melodic similarity is a multidimensional problem, tools can consider several parameters (such as pitch, duration, etc.) to evaluate the similarity between melodies.
- 5. Musical representation covers the encoding used to represent music, which can vary considerably between systems. Each representation shows strengths and weaknesses that affect the operation of the algorithm. Pieces have been variously encoded as strings, numbers, trees, or graphs.
- 6. Experiment-based denotes whether or not the similarity function of a system has been designed based upon the results of some empirical, cognitive research. Melodic similarity is deeply rooted in perception and cognition. To develop efficient algorithms, it is sometimes necessary to draw upon experiments in human perception.
- 7. Training indicates whether or not an algorithm has been trained—using some machine-learning techniques—on some specific data set. Trained systems are expected to have good performance on musical pieces that are comparable to those used for training.
- 8. Empirical validation investigates whether or not the similarity function has been validated. Algorithms exploit a similarity function for obtaining a quantitative value of likeness between two musical pieces.

# **Taxonomy**

In this section, we briefly describe the 15 algorithms for melodic similarity detection that we considered. The systems are organized into four categories, according to the strategies they exploit, namely: cognition, music theory, mathematics, and hybrid. Because the melodic similarity task is intrinsically interdisciplinary, researchers have addressed it by exploiting techniques that derive from very different areas. This taxonomy reflects the four main areas of investigation adopted in the field of music information retrieval.

It is worth noting that some of these disciplines overlap, thus it can sometimes be difficult to unequivocally classify techniques based on approaches at the edge of two areas. For example, several theories of music such as pitch-class set theory (Forte 1973) and the generative theory of tonal music (GTTM, Lerdahl and Jackendoff 1985) are strongly founded on mathematics. Therefore, any algorithm built upon one of these theories could be safely classified both as a member of the music theory and the mathematics categories. In such cases, and for the sake of clarity, we have classified systems according to the category that fits them better qualitatively.

This taxonomy aims to provide a first step in the direction of a comprehensive ontology for melodic similarity techniques. The proposed framework can be extended straightforwardly by adding new categories as well as by identifying meaningful subcategories. Indeed, additional categories might be needed when new approaches are developed.

#### **Systems Based on Cognition**

Cognitive constraints have only recently been used for evaluating melodic similarity. Algorithms that rely on this approach are usually based on one or both of two methods: (1) the linear combination of cognition-based metrics and (2) human-tailored pattern recognition. The work of Roig and colleagues (2013) and that of Vempala and Russo (2015) exploit the linear combination of different metrics, based on the evaluation of differences between pairs of musical features extracted from the melodies, such

as pitch distance, pitch direction, and rhythmic salience. On the other hand, de Carvalho and Batista (2012) propose a system—one based on a form of Prediction by Partial Matching (Cleary and Witten 1984)—that simulates the operation of the human auditory cortex in assessing the similarity between given MIDI files.

#### **Systems Based on Music Theory**

Music theory has long provided tools for understanding the structure of melodies. Therefore, several melodic similarity tools are built upon theories such as the GTTM (Lerdahl and Jackendoff 1985), the implication/realization (I/R) model of Narmour (1992), and Schenkerian analysis (Forte and Gilbert 1982).

Grachten et al. (2004) propose a melodic similarity measure based on the I/R model. Specifically, the algorithm tries to implement most of the innate processes presented by the I/R model. Melodies are annotated by performing I/R analysis. Annotations are then used for comparing different melodies. For overcoming one major problem of the I/R model—the impossibility of unambiguously differentiating the intervallic direction of a sequence of notes—Yazawa et al. (2013) designed an algorithm based on an extended I/R model that includes a number of new symbols for improving expressivity.

In a Schenkerian vein, Orio and Rodà (2009) introduced an approach, based on the representation of musical pieces as hierarchical graphs, to identify the relative structural importance of notes. The most relevant notes in a melody are then compared in order to find similarities between melodic segments.

## **Systems Based on Mathematics**

Systems in this category use a number of mathematical approaches that serve as a basis for evaluating the degree of likeness between melodies. Many of these algorithms represent musical data as functions within an abstract space and exploit notions from geometry as a means of comparison. Other common mathematical strategies are

based on statistical analysis or information-retrieval techniques.

Aloupis and colleagues (2006) present two algorithms that exploit a geometrical approach. Melodies are represented as polygonal chains within a pitch-time abstract plane, and their similarity is calculated as the minimum area between polygonal chains. In 2010, the S2 and W2 algorithms were introduced (Laitinen and Lemström 2010; Lemström 2010). Both are geometric algorithms that work with point-set representations of music and are designed to process polyphonic music. The similarity between the query pattern and the target melody is defined by the number of elements that match between them, after the application of some invariants. Finally, Julián Urbano (2013) proposes three different systems: ShapeH, ShapeTime, and Time. All three rely on a geometric model that encodes melodies as curves in a pitch/time space, that are then compared.

Bohak and Marolt (2009) investigate how melodybased features relate to folk-song variants. Specifically, the authors extract 94 melody-based features from each melody that are then used for performing comparisons.

Wolkowicz and Kešelj (2011) propose six different algorithms (WK1-6) that exploit text-based information-retrieval approaches. They extract features from an input data set of MIDI files. Stringbased methods can be directly applied to such input files, since none of the notes is concurrent or overlapping. Then they build *n*-grams (substrings of *n* consecutive tokens [i.e., notes]), which are used to calculate similarity between melodies. In the framework proposed by Frieler (2006), melodies are represented by series of arbitrary length in an abstract space of events. *n*-grams are then used for measuring the similarity of two melodies. *n*grams have an identity that allows one to assess their degree of similarity, thus providing further information.

### **Hybrid Systems**

To maximize the efficiency of algorithms, hybrid systems combine several techniques that usually

belong to two or more of the aforementioned categories. The SIMILE toolbox (Müllensiefen and Frieler 2004; Frieler and Müllensiefen 2005) is based on a linear combination of about 50 algorithms covering different aspects of the SMS task. Authors gathered rating data by human experts, which were considered as the ground truth for evaluating and optimizing the linear combination of the considered techniques' similarity values.

Fanimae (Suyoto and Uitdenbogerd 2010) combines two similarity measurements: the NGR5 approach, which exploits 5-grams for matching melodies by considering only the pitch; and the newly introduced PIOI technique, which evaluates similarity by considering either pitch or duration.

Rizo and Inesta (2010) introduced the UA systems—four methods designed with the objective of obtaining a good trade-off between accuracy and processing time. Two of these methods are based on tree representation, one is based on the quantized point-pattern representation, and one is the combination of the three other methods.

#### **Comparison of Systems**

By analyzing according to the criteria previously introduced in the corresponding section, it is possible to infer general trends and exceptions in how researchers conceived, designed, and implemented their systems. A detailed mapping of each system against our criteria is provided in Table 1. Each criterion gives an insight into one specific aspect of a melodic-similarity algorithm. Rather than listing decisions taken by authors with regard to each criterion, we outline trends and relevant deviations from trends that have been highlighted by the criteria.

The majority of the systems calculate melodic similarity between monophonic sequences only. The system developed by Laitinen and Lemström (Laitinen and Lemström 2010; Lemström 2010), as well as the algorithm built by Suyoto and Uitdenbogerd (2010), can additionally evaluate melodic similarity between polyphonic pieces. It should also be noted, however, that some of the systems' authors do not provide information about their systems' polyphonic analysis capabilities.

**Table 1. Description of Systems** 

System	Category	Polyphony	Scope	Similarity Function	Musical Parameters	Musical Representation	Based on Experiments	Trained	Empirical Validation
Frieler and Müllensiefen (2004)	Hybrid	No	General	Linear combination of metrics (edit distance, N-grams, geometric distance, correlation coefficient)	Pitch, duration, contour, tonality, accent structure	Not specified	Yes	Yes	Musical incipits (RISM)
Grachten et al. (2004)	Music theory	-	General	Edit distance	Pitch, contour	String of symbols (I/R symbols, note sequences)	No	No	Jazz songs
Aloupis et al. (2006)	Mathematics	No	General	Minimum area between polygonal chains	Pitch, duration	Geometric (polygonal chains)	No	No	Not specified
Frieler (2006)	Mathematics	_	General	Linear combination of N-gram similarity measures	Pitch, duration	Sequence of symbols	No	No	No
Bohak and Marolt (2009)	Mathematics	No	Folk songs	Statistical differences	Melodic complexity, meter, entropy, pitch, duration	96 statistical melodic features	No	yes	Folk songs
Orio and Rodà (2009)	Music theory	-	Music based on harmony	Shortest path between two nodes in a graph	Harmony, metre, pitch	Graph	No	No	Musical incipits (RISM)
Laitinen and Lemström (2010)	Mathematics	Yes	General	Number of matching elements between melodies	Pitch, duration	Point-set representation	No	No	Musical incipits
Rizo and Iñesta (2010)	Hybrid	No	General	Linear combination of metrics	Pitch, duration	Tree	No	No	Not specified
Suyoto and Uitdenbogerd (2010)	Hybrid	Yes	General	Linear combination of metrics (N-grams, dynamic programming)	Pitch, duration	String of symbols	No	No	Musical incipits
Wolkowicz and Kešelj (2011)	Mathematics	No	General	Text retrieval methods	Pitch	String of symbols	No	Yes	Musical incipits
de Carvalho and Batista (2012)	Cognition	No	General	Match frequent sequences of pitch intervals, duration ratios between melodies	Pitch intervals, duration ratio	String of symbols	No	No	Musical incipits

Table 1. Continued.

System	Category	Polyphony	Scope	Similarity Function	Musical Parameters	Musical Representation	Based on Experiments	Trained	Empirical Validation
Roig et al. (2013)	Cognition	No	General	Linear combination of metrics	Downbeat onset, passing note onset, pitch direction, pitch intervals, transposition	Not specified	No	No	Not specified
Urbano (2013)	Mathematics	No	General	Change in shape of curves	Pitch, duration	Geometric (curve in pitch/time plane)	No	No	Musical incipits
Vempala and Russo (2015)	Cognition	No	Tonal melodies	Linear combination of metrics	Pitch distance, pitch direction, duration, contour, tonal stability	Not specified	Yes	Yes	Tonal melodies
Yazawa et al. (2013)	Music theory	No	General	Matched N-gram sequences	Pitch direction, duration	String of symbols (extended I/R symbols)	Yes	No	Folk songs (Essen collection)

Most systems have a general scope, since they can be used to calculate the likeness of any melody belonging to any style. On the other hand, the algorithm developed by Vempala and Russo (2015) focuses on melodies that are rooted in tonality, the system by Orio and Rodà (2009) analyzes music in which harmony plays a functional role, and the system built by Bohak and Marolt (2009) is designed for folk songs only.

There is no dominant trend in similarity functions exploited by algorithms. Combining several similarity measures together, by means of a linear combination, is a popular approach (Müllensiefen and Frieler 2004; Frieler 2006; Suyoto and Uitdenbogerd 2010; Roig et al. 2013; Vempala and Russo 2015). In such cases, however, the specific measures that contribute to form the linear combination change from system to system. Other systems (Aloupis et al. 2006; Urbano 2013) calculate the differences of shape and area between curves that represent melodies in abstract mathematical spaces. There are also approaches based on statistics (Bohak and Marolt 2009), the shortest path between two nodes in a graph (Orio and Rodà 2009), and text-retrieval methods (Wolkowicz and Kešelj 2011). Also, some traditional methods, such as edit distance (Grachten et al. 2004) and n-grams (Yazawa et al. 2013), are exploited.

Musical parameters considered for calculating similarity are almost always pitch and duration. Information about time and frequency seem to be expressive enough to allow a precise judgement of likeness between musical excerpts. Some systems also rely, however, on parameters such as harmony (Orio and Rodà 2009), contour (Grachten et al. 2004; Müllensiefen and Frieler 2004; Vempala and Russo 2015), and pitch direction (Yazawa et al. 2013; Vempala and Russo 2015). It should be noted that only the algorithm developed by Bohak and Marolt (2009) includes a number of unusual and interesting parameters, such as melodic complexity, metric accent, and entropy.

The diversity of approaches used to calculate melodic similarity is matched by (and perhaps even springs from) the heterogeneity of musical representations. A common method of encoding musical information is through sequences of symbols (Frieler 2006; Rizo and Inesta 2010; Wolkowicz and Kešelj 2011; de Carvalho and Batista 2012). Other strategies explored to represent music are trees (Rizo and Inesta 2010), graphs (Orio and Rodà 2009), statistical features (Bohak and Marolt 2009), and I/R symbols (Grachten et al. 2004; Yazawa et al. 2013). Finally, two algorithms that rely on mathematics (Aloupis et al. 2006; Urbano 2013) represent melodies as curves in an abstract mathematical space.

Unfortunately, the large majority of systems reviewed in this survey are not based on first-hand experiments in music perception. Indeed, only three algorithms (Müllensiefen and Frieler 2004; Yazawa et al. 2013; Vempala and Russo 2015) are guided by direct experimental enquiry. Other algorithms either implement theoretical hypotheses, or follow well-established notions in music perception.

Moreover, only 4 algorithms out of the 15 considered have been trained on a set of melodies (Müllensiefen and Frieler 2004; Bohak and Marolt 2009; Wolkowicz and Kešelj 2011; Vempala and Russo 2015). The other systems do not benefit from the use of such machine-learning methods. As empirically observed in many areas of computer science and artificial intelligence, machine learning is a promising technique that could enhance the performance of algorithms considerably.

We observe that some of the reviewed approaches lack empirical validation. Although most methods have been tested on groups of melodies, the lack of empirical validation for some of the algorithms (Aloupis et al. 2006; Frieler 2006; Rizo and Inesta 2010; Roig et al. 2013) makes their assessment difficult, if not impossible. In the case of the method proposed by Frieler (2006), however, the lack of validation process is due to the theoretical nature of the article, which introduces innovative mathematical notions that the author acknowledges need future evaluation. Among the systems that do provide empirical validation, most have been tested on musical incipits (Müllensiefen and Frieler 2004; Orio and Rodà 2009; Lemström 2010; Suyoto and Uitdenbogerd 2010; Wolkowicz and Kešelj 2011; de Carvalho and Batista 2012; Urbano 2013). Other

**Table 2. System Rankings** 

	MIREX Edition (No. of Participants)						
System	2010 (13)	2011 (10)	2012 (6)	2013 (5)	2014 (4)		
Laitinen and Lemström (2010)	4	_	_	_	_		
Rizo and Iñesta (2010)	7	_	_	_	_		
Suyoto and Uitdenbogerd (2010)	8	_	_	_	_		
Wolkowicz and Kešelj (2011)	_	4	_	_	_		
de Carvalho and Batista (2012)	_	_	6	_	_		
Roig et al. (2013)	_	_	_	3	_		
Urbano (2013)	1	1	1	1	1		
Yazawa et al. (2013)	_	_	_	4	_		

Rankings in the Melodic Similarity track of different editions of MIREX.

algorithms have been tested with folk songs (Bohak and Marolt 2009; Yazawa et al. 2013), jazz melodies (Grachten et al. 2004), and tonal melodies (Vempala and Russo 2015).

Although not conclusive, an analysis of the results of the MIREX competitions in the SMS track is informative about the strengths and weaknesses of algorithms. It is worth noting that the best algorithm according to one particular MIREX competition is not necessarily the best system in all situations. From Table 2, it is possible to see that the systems proposed by Urbano (2013) have been the most successful, winning five editions of the MIREX competition. To compare the results achieved by the different algorithms proposed we consider their  $F_1$ -score, which is the harmonic mean of precision and recall ranging from 0 to 1. The performance of all the algorithms that competed in the 2010 edition of MIREX was close. None of the systems had a  $F_1$ -score greater than 0.30: Urbano (2013)  $F_1 = 0.30$ ; Laitinen and Lemström (2010)  $F_1 = 0.27$ ; Suyoto and Uitdenbogerd (2010)  $F_1 = 0.26$ ; and Rizo and Inesta (2010)  $F_1 = 0.26$ . In the 2011 campaign there was a significant increment in the performance of some of the algorithms, which doubled their accuracy. Both the systems proposed by Urbano (2013) and Wolkowicz and Kešelj (2011) achieved the same result, with  $F_1 = 0.64$ . After 2011, improvements have slowed down. So far, the best performance has been achieved by Urbano (2013) in the 2014

MIREX competition, with  $F_1 = 0.77$ . It is difficult to compare these systems with those that did not enter the MIREX competitions, since they have not been assessed in accordance with the same standardized evaluation procedure.

Although there has been a significant improvement in the methods introduced in the last decade, not all new systems are better than those developed before 2004. This is the case with the algorithm developed by Yazawa et al. (2013), which perhaps pays the price for experimenting with a new strategy. Nonetheless, it is safe to assert that pre-2004 algorithms have been outperformed in all aspects by newer methods.

Regardless of the rankings obtained by the systems in MIREX, it is possible to identify situations in which algorithms based on different methods perform best. In the case of melodies that differ by only a few pitches, intervals, and rhythms, an approach based on the number of matching elements in the two melodies (e.g., Laitinen and Lemström 2010; de Carvalho and Batista 2012) is likely to be the most effective. This situation is common in both folk and classical styles, where the composer introduces variety by altering a few musical elements of a melody (see Figure 1b). In the case of melodies that are substantially different but that occasionally share similar musical fragments, algorithms based on matched n-gram sequences (e.g., Frieler 2006; Yazawa et al. 2013) have an edge, because they are

Figure 1. Melodies with different categories of similarity: the model (a); a melody obtained by changing few pitches and durations of the model (b);

a melody that conserves only the head and the tail of the model (c); and an ornamented version of the model (d).



able to detect similarities even if only tangential. This scenario is often found in contrapuntal music (e.g., fugues and motets), where it is common that only the head and the tail of a theme are maintained throughout different repetitions (see Figure 1c). Systems based on a linear combination of different metrics and statistical musical differences (e.g., Suyoto and Uitdenbogerd 2010; Roig et al. 2013; Vempala and Russo 2015) are most effective when used on melodies that are loosely similar. In this category are those melodies that share only a similar contour, or in which ornamental notes are intercalated between the fundamental notes (see Figure 1d). This compositional technique is used in classical music and it is frequently encountered in pieces built around a single melodic idea. Finally, when the relationship between two melodies results from a mixture of these three situations, the most effective approach is that of Urbano (2013), which calculates differences in the shape of geometric representations of a melody.

## **Guidelines and Recommendations**

The work done by researchers in the last decade has significantly increased the number of algorithms that calculate melodic similarity, and it has improved the variety and quality of approaches. Systems are now capable of finding melodies similar to a target sequence of notes, and searching through large databases of melodies in a more reliable, efficient and precise way than before. There are still some limitations, however, and several steps can be taken to improve the performance and functionality of these systems.

As already mentioned in the Background section, there is no clear-cut definition of melodic similarity. As Alan Marsden (2012) points out, the notion of melodic similarity is highly dependent on context. Different models are required to account for similarity judgments carried out by people. For instance, to reduce the time needed by the similarity process carried out by human judges, in MIREX each

person considers different sets of melodies. These melodies can have very different backgrounds, however, thus providing different evaluations. This strongly affects both the way in which systems are evaluated and the development of new algorithms. The lack of a widely approved definition of melodic similarity is also reflected in the lack of a common ground truth, shared by all researchers interested in this field, against which the performance of algorithms can be tested. The evaluation framework provided by MIREX is a first attempt to solve this issue. As Urbano (2013) points out, however, there are some problems with this musical data set. The MIREX evaluation framework failed to give a consistent measure of performance for the same system over a number of years. As a consequence, this framework cannot be adopted to benchmark performances of different systems over time. The solution for a reliable ground truth for melodic similarity is to have an extremely large database of melodies, in which each melody has a series of ratings that indicate the degree of similarity between it and all other melodies in the collection. To develop this evaluation framework, a number of experiments with human listeners are needed. Because it is necessary to design a large database, however, these experiments could be too onerous to conduct in a traditional laboratory-based setting. This issue can be overcome by conducting experiments on the Internet, which allows exploitation of the very large number of people available online.

A second issue affecting current systems is that most of them focus on monophonic music only. We believe that next-generation algorithms should be able to analyze polyphonic music and thus find the degree of similarity between two polyphonic excerpts. Polyphonic music in the Western tradition considerably outweighs monophonic music. So the impact of polyphonic similarity systems on musicology and music theory would be significantly increased. Developing polyphonic similarity analysis tools would also encourage researchers to discover the main differences between evaluating likeness between monophonic and polyphonic musics. This research should be informed by experiments in music perception and, in turn, would

provide new insights to musicologists and music theorists.

There is a close relationship between music cognition and melodic similarity algorithms. To develop more efficient systems, it is of paramount importance to conduct focused experiments in music perception. The results of these studies could eventually suggest innovative notions and techniques still overlooked in the design of computational systems.

Another relevant point to consider for enhancing the performance of algorithms is their scope. Most of the methods reviewed have a general scope. In other words, these systems can be used to analyze melodies belonging to any musical style. It is well known, however, that every style has specific rules that help create its unique sound. This is true for melody generation as well. Therefore, if we want to obtain greater efficiency, we need to concentrate on specific styles, rather than developing algorithms able to parse any kind of music. While building these style-specific tools, we might discover some of the deepest rules that contribute to define a style. Understanding these rules would not only improve the efficiency of systems, but would also help musicologists and music theorists understand the complex phenomenon of musical style. Stylespecific rules can be extracted, for instance, by using techniques of machine learning. Only a few systems analyzed in this survey are based on training strategies. This implicitly indicates that most of the systems developed so far have a general scope or have been manually configured following developer intuitions. Clearly, training on specific styles will make systems less general. To avoid this pitfall, compound tools, made up of a series of style-specific subsystems, can be designed. In this way, systems would operate distinct subsystems, depending on the style of melodies they are analyzing.

A weakness of using machine learning is that there is a trade-off between performance and interpretability (James et al. 2013). The more complex the machine-learning technique used, the looser the relationship between the predictors (i.e., musical features) and melodic similarity. This phenomenon is because of the risk for overfitting in highly flexible methods. If, however, such systems aim only

at providing a score of likeness between melodies, and are not meant to be used as tools for studying the general relations between musical features comparison and similarity, the interpretability of the predictive model is not of interest. In this case, the trade-off between performance and interpretability is not a significant concern.

Melodic similarity algorithms based on mathematics have been dominating MIREX competitions for the last five years. In our opinion, this does not mean that algorithms based on mathematics are intrinsically superior to those based on cognition or on music theory. Rather, this should encourage the improvement of tools based upon different theories, in order to close the performance gap with mathematical systems. Indeed, the task of melodic similarity is highly interdisciplinary and should be tackled from diverse perspectives. For this reason, it would be advisable to develop tools that merge different approaches. By creating hybrid systems, researchers can ameliorate the weaknesses of specific techniques, while augmenting the overall strength of the system.

#### **Conclusions**

Melodic similarity systems have a wide range of applications. For instance, they can be used to perform information retrieval on large music data sets, and they can help identify music plagiarism. Recently, because of the introduction of the MIREX competition, a large number of algorithms have been developed. In this article, we have briefly described existing approaches, and we have presented a new modular taxonomy and eight criteria that are instrumental in classifying and comparing melodic similarity algorithms. The taxonomy classifies algorithms based on their approach, i.e., whether they are based on cognition, music theory, mathematics, or some hybrid of these. This article fills the gap since the previous surveys in the area (Müllensiefen and Frieler 2006; Hofmann-Engl 2010), and provides a clear overview of existing techniques. The analysis of the 15 systems considered has allowed us to identify the strengths and weaknesses of the algorithms as well as wider trends in the field.

Starting from this point, we have been able to recommend avenues of future research that one hopes will lead to further improvement in the area. In this regard, we highlight the lack of a widely approved definition of melodic similarity, which strongly affects both the development of algorithms and the outcomes of competitions. We recommend the development of a common ground truth that can be used as a standard corpus for comparing systems. Also, we observed a worryingly small number of systems that are able to analyze polyphonic pieces. The capacity to analyze such music is of primary importance in fostering the exploitation of melodic similarity tools in real-world applications. Finally, because musical works belonging to distinct styles are often very significantly different, we have found that it is extremely difficult to develop systems that perform well across a range of styles. Therefore, we believe that future algorithms should have a relatively limited scope, and that they should be configured on a specific style of music through machine-learning techniques.

### References

- Aloupis, G., et al. 2006. "Algorithms for Computing Geometric Measures of Melodic Similarity." *Computer Music Journal* 30(3):67–76.
- Bohak, C., and M. Marolt. 2009. "Calculating Similarity of Folk Song Variants with Melody-based Features." In *Proceedings of the International Conference on Music Information Retrieval*, pp. 597–602.
- Cambouropoulos, E. 1998. "Towards a General Computational Theory of Musical Structure." PhD dissertation, University of Edinburgh.
- Cleary, J. G., and I. Witten. 1984. "Data Compression Using Adaptive Coding and Partial String Matching." *IEEE Transactions on Communications* 32(4):396–402.
- de Carvalho, A. D., Jr., and L. V. Batista. 2012. "SMS Identification Using PPM, Psychophysiological Concepts, and Melodic and Rhythmic Elements." In *Proceedings of the Annual Music Information Retrieval Evaluation Exchange*. Available online at music-ir.org/mirex/abstracts/2012/DB1.pdf. Accessed February 2015.
- Downie, J. S. 1999. "Evaluating a Simple Approach to Music Information Retrieval: Conceiving Melodic N-Grams as Text." PhD dissertation, The University of Western Ontario, London, Ontario.

- Downie, J. S. 2004. "The Scientific Evaluation of Music Information Retrieval Systems: Foundations and Future." *Computer Music Journal* 28(2):12–23.
- Downie, J. S. 2008. "The Music Information Retrieval Evaluation Exchange (2005–2007): A Window into Music Information Retrieval Research." Acoustical Science and Technology 29(4):247–255.
- Forte, A. 1973. *The Structure of Atonal Music.* New Haven, Connecticut: Yale University Press.
- Forte, A., and S. E. Gilbert. 1982. *Introduction to Schenkerian Analysis: Form and Content in Tonal Music.* New York: Norton.
- Frieler, K. 2006. "Generalized N-Gram Measures for Melodic Similarity." In Bagatelj, V., et al., eds. *Data Science and Classification*. Berlin: Springer, pp. 289–298.
- Frieler, K., and D. Müllensiefen. 2005. "The Simile Algorithm for Melodic Similarity." In *Proceedings of the Annual Music Information Retrieval Evaluation Exchange*. Available online at music-ir.org/mirex /abstracts/2005/frieler.pdf. Accessed February 2015.
- Grachten, M., et al. 2004. "Melodic Similarity: Looking for a Good Abstraction Level." In *Proceedings of the International Conference on Music Information Retrieval*. Available online at ismir2004.ismir.net /proceedings/p040-page-210-paper166.pdf. Accessed February 2015.
- Hofmann-Engl, L. 2010. "An Evaluation of Melodic Similarity Models." Available online at www.chameleongroup.org.uk/research/evaluation.pdf. Accessed January 2016.
- James, G., et al. 2013. An Introduction to Statistical Learning. Berlin: Springer.
- Laitinen, M., and K. Lemström. 2010. "Geometric Algorithms for Melodic Similarity." In *Proceedings of the Annual Music Information Retrieval Evaluation Exchange*. Available online at music-ir.org/mirex /abstracts/2010/LL1.pdf. Accessed February 2015.
- Lemström, K. 2010. "Towards More Robust Geometric Content-Based Music Retrieval." In *Proceedings of the Conference of the International Society for Music Information Retrieval*, pp. 577–582.
- Lerdahl, F., and R. Jackendoff. 1985. A Generative Theory of Tonal Music. Cambridge, Massachusetts: MIT Press.
- Marsden, A. 2012. "Interrogating Melodic Similarity: A Definitive Phenomenon or the Product of Interpretation?" *Journal of New Music Research* 41(4):323–335.
- McNab, R. J., et al. 1996. "Towards the Digital Music Library: Tune Retrieval from Acoustic Input." In Proceedings of the ACM International Conference on Digital Libraries, pp. 11–18.

- Meek, C., and W. P. Birmingham. 2002. "Johnny Can't Sing: A Comprehensive Error Model for Sung Music Queries." In *Proceedings of the International Conference on Music Information Retrieval*. Available online at ismir2002.ismir.net/proceedings/02-FP04-4.pdf. Accessed February 2015.
- Müllensiefen, D., and K. Frieler. 2004. "Optimizing Measures of Melodic Similarity for the Exploration of a Large Folk Song Database." In *Proceedings of the International Conference on Music Information Retrieval*. Available online at ismir2004.ismir.net /proceedings/p052-page-274-paper178.pdf. Accessed February 2015.
- Müllensiefen, D., and K. Frieler. 2006. "Evaluating Different Approaches to Measuring the Similarity of Melodies." In V. Batagelj, et al., (eds.) *Data Science and Classification*. Berlin: Springer, pp. 299–306.
- Narmour, E. 1992. The Analysis and Cognition of Melodic Complexity: The Implication-Realization Model. Chicago: University of Chicago Press.
- O'Maidín, D. 1998. "A Geometrical Algorithm for Melodic Difference." Computing in Musicology: A Directory of Research 11:65–72.
- Orio, N., and A. Rodà. 2009. "A Measure of Melodic Similarity Based on a Graph Representation of the Music Structure." In *Proceedings of the International* Conference for Music Information Retrieval, pp. 543– 548.
- Rizo, D., and J. M. Inesta. 2010. "Trees and Combined Methods for Monophonic Music Similarity Evaluation." In Proceedings of the Annual Music Information Retrieval Evaluation Exchange. Available online at music-ir.org/mirex/abstracts/2010/RI1.pdf. Accessed February 2015.
- Roig, C., et al. 2013. "Submission to MIREX 2013 Symbolic Melodic Similarity." In *Proceedings of* the Annual Music Information Retrieval Evaluation Exchange. Available online at www.music-ir.org/mirex /abstracts/2013/RTBB1.pdf. Accessed February 2013.
- Suyoto, I. S. H., and A. L. Uitdenbogerd. 2010. "Simple Orthogonal Pitch with IOI Symbolic Music Matching." In *Proceedings of the Annual Music Information Retrieval Evaluation Exchange*. Available online at music-ir.org/mirex/abstracts/2010/SU1.pdf. Accessed February 2015.
- Urbano, J. 2013. "A Geometric Model Supported with Hybrid Sequence Alignment." In *Proceedings of the Annual Music Information Retrieval Evaluation Exchange*. Available online at music-ir.org/mirex /abstracts/2013/JU1.pdf. Accessed February 2013.

- Vempala, N. N., and F. A. Russo. 2015. "An Empirically Derived Measure of Melodic Similarity." *Journal of New Music Research* 44(4):391–404.
- Wolkowicz, J., and V. Kešelj. 2011. "Text Information Retrieval Approach to Music Information Retrieval." In *Proceedings of the Annual Music Information* Retrieval Evaluation Exchange. Available online at
- music-ir.org/mirex/abstracts/2011/WK1.pdf. Accessed February 2013.
- Yazawa, S., et al. 2013. "Melodic Similarity Based on Extension Implication-Realization Model." *MIREX Symbolic Melodic Similarity Results* Available online at music-ir.org/mirex/abstracts/2013/YHKH1.pdf. Accessed February 2013.