

MODULO7 : A FULL STACK MUSIC INFORMATION RETRIEVAL AND STRUCTURED QUERYING ENGINE

First Author

Affiliation1

author@ismir.edu

Second Author

Retain these fake authors in

submission to preserve the formatting

Third Author

Affiliation3

author3@ismir.edu

ABSTRACT

This paper describes a novel Music Information Retrieval and Structured Querying framework named Modulo7. Modulo7 is a full stack implementation (both client and server side software) which facilitates indexing variegated sources of music (midi, mp3, music xml and digitized sheet music files). It implements a similarity search engine based on customized vector space representations of songs based music theory, an efficient indexing and persistent storage mechanism and an interface for querying relevant songs based on SQL(Structured Querying Language) like principles and a search engine functionality. This papers describes the design goals and implementation details of the framework, outlines the speed up and scale up results over input sources and other MIR frameworks and discusses potential improvements.

Keywords : MIDI, Music XML, MP3, Music Retrieval, SQL

1. INTRODUCTION

Given the explosive growth of Music Information Retrieval research, several approaches and software suites have been designed to tackle generic problems such as efficient song indexing, music recommendation, archival methods, structured and un-structured querying, feature extraction, audio and digital signal processing etc. A vast majority of the MIR frameworks in academia tend to approach very specific problems and does not support scalability as a significant end goal in itself [6]. Moreover, industry applications predominantly treat MIR approaches based on collaborative filtering techniques [11] and/or manual annotation schemes [8] instead of structural analysis of music sources. However these frameworks scale very well to large datasets.

Modulo7 is an attempt to capture the best features of both worlds. Modulo7 converts different music sources (midi, musicxml files, sheet music png/jpeg files and

mp3) into a single unified symbolic representation. It indexes songs on global properties (artist name, key signature, time signature etc) and provides a vector space implementation based on music theory and facilitates similarity searches based on the SIMILIE [10] platform and extended to include chords and polyphony and a novel querying approach. As a consequence, Modulo7 acts as an end to end software suite for Music Information Retrieval which analyses song's musical content directly instead of music consumer's opinions about music. The broad goals of Modulo7 can hence be summarized as

1. To efficiently store and index variegated data formats into a unified format while keeping symbolic information intact.
2. To create a framework for music retrieval from multiple music formats and in effect facilitate comparing songs of different formats.
3. To query/search relevant songs based on a quantified version of music theoretic principles.
4. To build a full stack open source platform for end to end MIR tasks. Modulo7 is hosted at [2]

2. RELEVANT WORK

Music Information Retrieval is a vast and interdisciplinary body of work. In order to facilitate research in MIR, several software suites and frameworks have been developed in the academic community in the past. Notable amongst them are software suites like jMIR [14] for automatic feature extraction from audio and midi sources to be used as input into machine learning algorithms for genre classification, marsyas [19], essentia [1] and librosa [13] for audio processing, humdrum [3] for automated musicological research, gamera [9] for optical music recognition and symbol identification and SIMILIE [10] for melodic similarity analysis.

On top of the effort done by academia, there has been significant effort undertaken by companies in the field of MIR with an overarching emphasis on music recommendation as an end goal. For instance Amazon implements a collaborative filtering approach by quantifying a consumer's shopping behavior [11]. Pandora on the other hand utilizes an approach in which trained musicologists manually classify similarity between artists and songs [8]



© First Author, Second Author, Third Author. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** First Author, Second Author, Third Author. "Modulo7 : A full stack Music Information Retrieval And Structured Querying Engine", 17th International Society for Music Information Retrieval Conference, 2016.

which is subsequently consumed by Pandora’s APIs for similarity searches. Shazam [7] implements a database of music fingerprints and supports music ID as a service (by comparing a generated fingerprint of an input song to their existing database).

3. SOFTWARE DESIGN

This section details the software design of the Modulo7 most notably its architecture and a typical consumer workflow description(i.e. how Modulo7 is initialized and consumed as a service).

3.1 Software Architecture and Design

The Modulo7 architecture can be visualized in figure 1 and is broadly divided into the following components.

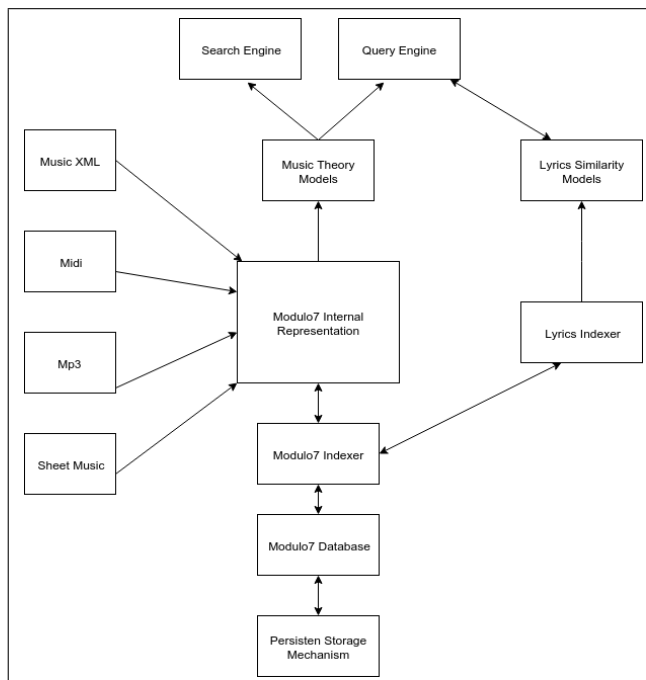


Figure 1. A block diagram of the Modulo7 Architecture.

1. **Music Parsers:** Components that individually parse different music sources into a unified symbolic representation. as visualized in figure 2. This object can be stored in memory or in a persistent serialized state via Apache Avro library. Depending on the source of the file, the parser utilizes a different approach.

- (a) For mp3 files, the input is parsed into chromagrams using echonest API and then converted into symbolic format via an in house implementation of the KKTonality profiles algorithm [12].
- (b) For Midi and Music XML files, a custom parser is built as a part of the Modulo7 code base which converts these formats.

- (c) For Sheet Music, the Audiveris library is used to convert digitized sheet music to Music XML files and subsequently consumed by the Modulo7 Music XML parser.

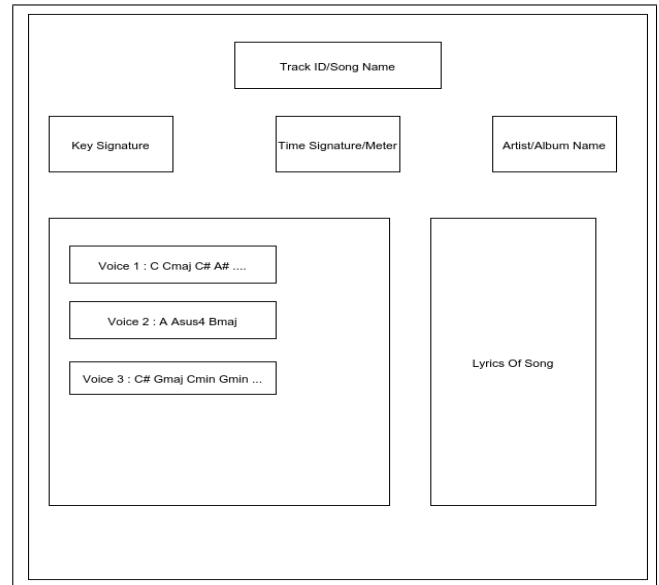


Figure 2. Modulo7 internal representation.

2. **Music theory Models:** Modulo7 implements vector space models derived from [10] with extensions to polyphonic music and chords using a template matching algorithm described in [18]. This forms a mathematical formulation of music sources. Its described in detail in section 4.
3. **Client side Querying:** Modulo7 implements client side abstractions that exposes two modes of information retrieval. Both these modes leverage the vector space models described in section 4.

- (a) **Structured Querying:** The structured querying mode allows a user to make SQL like queries based on certain criteria or statistic that is a quantified version of a music theoretic criteria. This is described in section 5.2
- (b) **Custom Search Engine:** Modulo7 exposes a standard search engine functionality(an input query song against a database) based on customized similarity metrics as described in section 5.1

4. **Lyrics Engine:** On top of providing support for music sources, Modulo7 also implements a standard text based retrieval model on lyrics of songs written in Apache Lucene. It is implemented as a standard text search engine.

3.2 Consumer Workflow

The consumer workflow can be defined as the steps/operations the user takes to initialize the Modulo7

server and client components and the output generated on search and structured queries.

1. Modulo7 server is initialized by specifying a source directory. Every music source file inside the directory is parsed and converted to either an in memory and/or persistent store version. It also indexes these songs on global properties (artist, key-signature, time-signature etc). It also guesstimates certain properties if not explicitly present (e.g key signature via the KK Tonality profiles key finding [12] algorithm)
2. It then exposes a prompt which asks the user for two modes (search engine/structured querying). (In either case the appropriate vector space model is built based on the query)
 - (a) If the user selects the search engine functionality, the engine then asks for an input query song and a similarity metric (from one of the choices described in section 5.1). Then the engine proceeds to rank the query song with each of the indexed songs based on the similarity metric and then returns an ordered list of "relevant" song(based on the similarity scores for each of the indexed songs against the query song).
 - (b) If the user selects the structured querying mode, Modulo7 asks for an input query - list of relevant song formats along with a set of music theoretic predicates songs must satisfy (details in section 5.2). Any song that satisfies this query is deemed relevant and returned to the user.
3. Modulo7 optionally caches results in memory so that it does not have to recompute queries.

4. VECTOR SPACE MODELS FOR MUSIC SOURCES

The vector space models form the core of the search and query functionality of Modulo7. Vector space models are derived from quantification of music theoretical representations of songs and extends functionality of the SIMILIE [10] melodic similarity library.

4.1 Building Blocks of unified internal representation

All vector space models are built from the unified internal representation of songs in Modulo7. This unified representation is based on certain building blocks described as follows

1. **Note:** A note is a pair of pitch (deterministic frequency played by an instrument) and onset (time instance at which a pitch is played). A **Chord** is a set of pitches played at the same onset.
2. **Interval:** The difference between two note pitches.

3. **Voice:** A set of notes or chords played in succession (usually on a single instrument).
4. **Song:** A song is a set of interleaved voices, along with some global properties (key signature, time-signature, artist, album etc)

4.2 Global Vector Space Models

Now once the above building blocks are ascertained for songs via Modulo7's parsers, certain vector space models can be defined over the entire song. Modulo7 implements 6 such vector space models. Some notable examples are

1. **Normalized Tonal Histogram Vector :** For each of the 12 intervals in western music, compute the frequencies(number of occurrences) of each interval and divide by the total number of intervals over the song. This forms a 12 dimensional vector irrespective of length of song. Mathematically

$$NTH(S) = \langle \Delta p_1^f, \Delta p_2^f, \dots, \Delta p_{12}^f \rangle \quad (1)$$

Here Δp_i^f is the fraction of the number of occurrences of the i^{th} interval.

2. **Normalized Tonal Duration Histogram Vector :** For each of the 12 intervals in western music, compute the cumulative durations each of these intervals are played and divide by total duration of the song. This forms a 12 dimensional vector irrespective of length of song. Mathematically

$$NTDH(S) = \langle \Delta t_1^f, \Delta t_2^f, \dots, \Delta t_{12}^f \rangle \quad (2)$$

Here Δt_1^f is the fraction of the duration that the i^{th} interval is played.

4.3 Voice specific vector space models

Considering each voice as a separate entity, a voice can be modeled as a vector independent of other voices in a song. A particular song can then be represented as a set of these vectors. Some notable examples of these are

1. **Pitch Interval Vector:** If a voice has k pitches played one after the other, there are k - 1 intervals (for each pitch its immediate subsequent pitch). The pitch interval vector is the chronological succession of these intervals. Mathematically

$$PI = \langle \Delta p_1, \Delta p_2, \dots, \Delta p_n \rangle \quad (3)$$

Here Δp_i is the i^{th} interval which is the difference between the i^{th} and $(i - 1)^{th}$ pitches.

2. **Rhythmically weighted Pitch Interval Vector:** Similar to the pitch interval vector, except it also takes into account the time t_i an interval Δp_i is played for. Mathematically

$$RPI = \langle rp_1, rp_2, \dots, rp_n \rangle \quad (4)$$

Here $rp_i = \Delta p_i \times t_i$

5. CONSUMER FACING ABSTRACTIONS

Modulo7 exposes two consumer facing abstractions used to facilitate acquiring relevant songs.

5.1 Search Engine Functionality

Based on the vector space models defined in sections 4.2 and 4.3, customized similarity metrics can be defined to form ranked search functionality for a given query song against an indexed data set. These similarity metrics are derived from the SIMILIE project [10] and extended ideas in [17]. Given two songs S_1 and S_2 , we can define similarity to be

$$Sim(S_1, S_2) \in \{0, 1\} \quad (5)$$

The definition of $Sim(S_1, S_2)$ can be implemented in various ways and Modulo7 allows the consumer to choose the similarity metric as an input parameter. So for a given query song S_{query} and a given similarity metric Sim , we can define the order list $Or(S_{query}, Sim)$ given the indexed set of songs S_{ind}

$$Or(S_{query}, Sim) = desc_sort_k \{S | Sim(S_{query}, S) = k, S \in S_{ind}\} \quad (6)$$

This order list is presented as an output to the consumer. Modulo7 implements 7 different similarity measures. Following are some notable examples

1. **Smith Waterman Song Similarity:** Given two voices V_1 and V_2 , Modulo7 uses the smith waterman genome alignment algorithm [17] to estimate similarity based on notes as units. In order to extend this as a song similarity metric, pair wise voice similarities are calculated and the max similarity is returned. Mathematically

$$SM(S_1, S_2) = arg_{max} k \{SMV(V_i, V_j) = k | V_i \in S_1, V_j \in S_2\} \quad (7)$$

The method $SMV(V_1, V_2)$ is identical to the smith waterman similarity algorithm as described in [17].

2. **N Grams Similarity :** Voices can be defined as N grams (strings of n length n of continuous notes/chords) as defined in [10]. These n grams can be aggregated over all voices and measures like **ukkonen**, **count distance**, **scm trigram** etc (defined in [15])

5.2 Structured Querying Functionality

Modulo7 implements a custom SQL like language with its own grammar definition and compiler implementation(facilitated via the Antlr library). Any query can be constructed as

$$\begin{aligned} &select\ input_types\ from \\ &DATABASE_NAME \\ &where(boolean_formula) \end{aligned} \quad (8)$$

Here the individual components stand for

1. **Input Types:** A list of the input types we are interested in. For example if we are only interested in midi files then this value is "midi". If we are interested in more than one type(e.g midi and musicxml) then Input Types = "midi, musicxml".
2. **DATABASE_NAME :** A placeholder for the of the data base (the name is created when Modulo7 server starts up)
3. **boolean_function :** A boolean function is a well formed propositional formula (chain of propositions or predicates connected by AND/ORs). Each predicate can either be defined as a criteria or statistic condition.

5.3 Criteria

A criteria is a particular music theoretic condition that a song either satisfies or not. Modulo7 implements 6 such criteria. Some notable examples are

1. Polyphony criteria : Whether a song is polyphonic or not
2. Key Signature criteria : Whether a song has a particular key signature
3. STAB Voice Criteria : Whether a song follows the Soprano, Alto, Tenor, Bass Voice structures as defined in [4]

5.4 Statistic

A statistic is a real number generated when applied to a particular song (effectively a quantified attribute of a song). A statistic condition is defined when this generated statistic satisfies a relational operator with a given threshold value. For example $PowerIndex(song) < 0.5$ is a valid statistic condition. Modulo7 implements 11 such statistics. Some notable examples are

1. Average Pitch Duration over the song duration.
2. Happiness Index : Fraction of major intervals over the total number of intervals
3. Power Index : Fraction of perfect intervals over the total number of intervals
4. Sadness Index : Fraction of minor intervals over the total number of intervals

Given the definition of statistic or criteria (cumulatively called conditions) a boolean formula can be defined as

$$boolean_formula = cond_1\ op_1\ cond_2\ ..\ op_{n-1}\ cond_n \quad (9)$$

Here op_i is a boolean operator (and/or).

6. EXPERIMENTAL EVALUATION

This section includes experiments conducted and improvements that have been attained over existing frameworks and input sources. For the purposes of this experiment, existing academic datasets such as the Saarland dataset [16], the Wikifonia Dataset [21] and the million song data set [5] and experiments were run on identical hardware for competing implementations.

6.1 Indexing and Persistent Store Improvements

Modulo7 achieves significant storage improvements over the input formats, while not losing any symbolic information. Figure 3 illustrates this for the Saarland dataset [16]

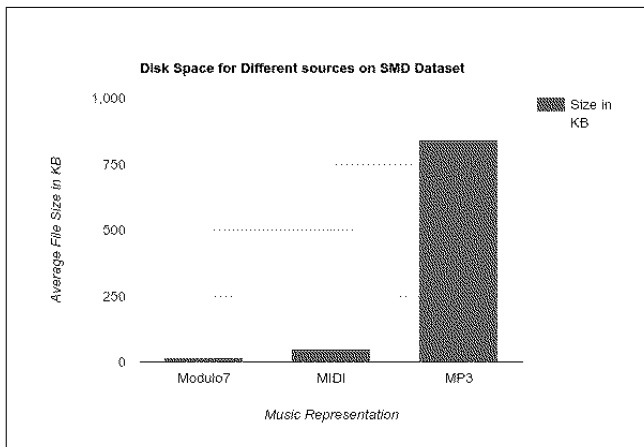


Figure 3. Space storage improvements over Saarland Dataset.

Similarly an experiment was carried out to ascertain compression as a function of data set size for the Wikifonia Dataset [21]

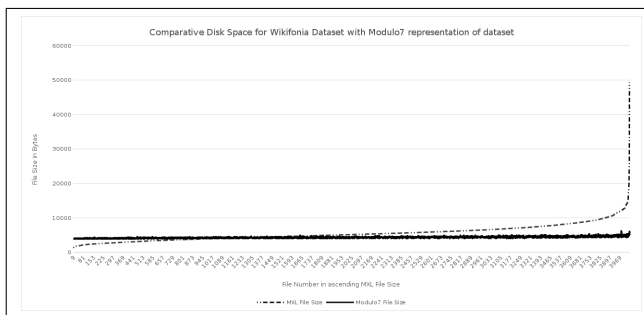


Figure 4. Modulo7 compression of Wikifonia Dataset over increasing subset sizes

Its evident from these graphs that Modulo7 achieves significant storage improvements and scales well over increased data set size.

6.2 Comparisons against JSymbolic

Some experiments were conducted against the jSymbolic component in the jMIR [14] software suite. Like the

jMIR suite, Modulo7 extracts features from songs, but with a different goal (relevance querying and similarity against jMIRs genre classification goals). Modulo7 implements **23** of jSymbolic’s **111** features. Moreover both frameworks are written in Java and hence language features do not influence the evaluation.

The experiment involved comparing memory footprints and time taken for extraction of 23 features for both frameworks over increasing subset sizes of the Saarland [16] dataset.

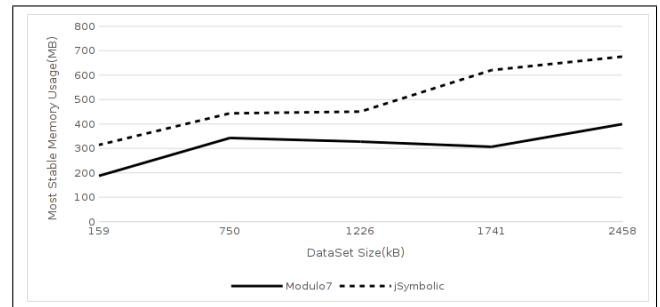


Figure 5. Average Memory taken for feature extraction, jSymbolic vs Modulo7

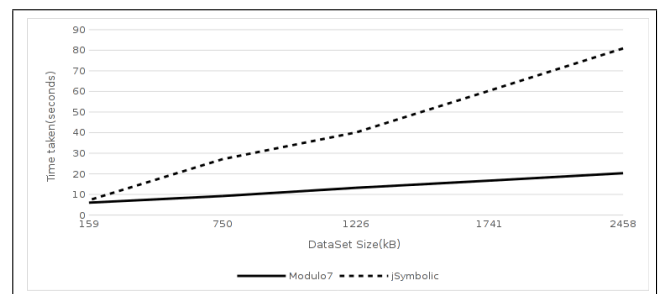


Figure 6. Total time taken for feature extraction, jSymbolic vs Modulo7

It can be inferred from these charts that Modulo7 outperforms jSymbolic in feature extraction speed and memory consumed and scales better with increasing data set sizes. Since jSymbolic is Open Source a quick review noted the following reasons for speedups

1. JSymbolic is single threaded whereas Modulo7 is multi threaded.
2. jSymbolic uses statically sized arrays, Modulo7 uses dynamically resizing maps for storing features.

6.3 Similarity Search Experiments

In order to estimate the efficacy of the similarity search algorithms, the million song dataset or MSD [5] was chosen to act as ground truth. MSD contains tags associated with songs (genre, comments about mood of the song and other such keywords) which was extracted by the researchers using the Last FM API. The ground truth of similarity was

based on the number of matching tags between two songs. Consider tag sets T_i and T_j for songs in the MSD S_i and S_j . The ground truth similarity is defined as

$$GTSim(S_1, S_2) = \sum_{t_i \in T(S_1)} \sum_{t_j \in T(S_2)} \begin{cases} 1 & t_i == t_j \\ 0 & otherwise \end{cases} \quad (10)$$

The measure was not normalized as MSD had a large number of "junk" tags (irrelevant comments about songs) which would have influenced the ground truth otherwise.

A direct download-able subset of the million song data set was used which contained 10000 songs. A ten fold cross validation was performed (1000 query songs, 9000 data hold out) and average precision and recall are calculated. A separate parser was developed for extracting symbolic information from MSD files. The following table illustrates average precision and recall calculated for different similarity measures.

| Similarity Measure | Average Recall | Average Precision |
|--------------------|----------------|-------------------|
| SCM Trigram | 0.308 | 0.299 |
| Ukkonen | 0.339 | 0.291 |
| Count Distance | 0.294 | 0.283 |
| Tonal Histogram | 0.341 | 0.362 |

While the performance of these measures are not very high, its comparable to the MIREX evaluation [20] on melodic similarity on smaller data sets.

Table 1. Average Precision and Recall for Similarity Measures

6.4 Exploratory Query Search Experiments

In order to ascertain the effectiveness of the query engine, certain custom experiments were conducted. Modulo7 queries were built which logically correlated with a music theoretic statement as ground truth. For example, rock songs have higher number of power chords and hence a higher number of perfect intervals - i.e. a higher power index. The following table shows Average Precision and Recall values for the million song data set [5] with a ten fold cross validation on the MSD subset [5] used in the similarity search experiments.

| Purpose | Query | Precision | Recall |
|---------------|-------|-----------|--------|
| Rock Song ID | Q1 | 0.13 | 0.98 |
| Sad Song ID | Q2 | 0.02 | 0.44 |
| Happy Song ID | Q3 | 0.018 | 0.4 |

Table 2. Results for the exploratory query analysis

We define Q1, Q2 and Q3 as follows

Q1 select mp3 from default_database where powerindex > 0.61; with a ground truth estimate as song tags being either "rock" or "pop_rock"

Q2 select mp3 from default_database where sadnessindex > 0.15 and scale = minor; with the ground truth estimate being song tags either "sad" or "sad_song"

Q3 select mp3 from default_database where happinessindex > 0.11 and scale = major; with ground truth being song tags as either "happy" or "happy_song"

It can be inferred that while the queries get the relevant information (i.e. high recall), it also acquires irrelevant songs (i.e. low precision). A logical solution would be to improve queries with more conditions to achieve higher precision.

7. CONCLUSION AND FUTURE WORK

This paper introduced the design and implementation of a novel Music Information Retrieval Framework called Modulo7. Unlike other MIR frameworks, Modulo7 allows indexing songs from heterogeneous sources of data and facilitates comparison of songs from these sources. It features a novel structured querying functionality for retrieving relevant songs on music theoretic criteria and facilitating customized search and querying based on these criteria. Modulo7's implementation indexes data very efficiently and provides significant speedups over existing frameworks like jMIR.

While Modulo7 is a novel full stack MIR framework, several modifications and additions can be done on top of the existing implementation to improve scalability, efficiency and functionality. Some notable extensions could be

1. Adding distributed computing support to allow parallel processing of query processing and intermediate steps in indexing songs (i.e. feature extraction). This can be facilitated by technologies like Hadoop and Spark.
2. Implementing additional query conditions and search schemes (vector space models and similarity metrics). This would allow for more specialized searching.
3. Adding a relevance feedback mechanism to dynamically modify queries.

8. REFERENCES

- [1] *ESSENTIA: an Audio Analysis Library for Music Information Retrieval*.
- [2] Modulo7 open source code base link: Omitted for peer review, February 2016.
- [3] Anon. The humdrum toolkit: Reference manual. menlo park, california: Center for computer assisted research in the humanities, 552 pages, isbn 0-936943-10-6.
- [4] Anon. Satb wikipedia link: <https://en.wikipedia.org/wiki/satb>.

- [5] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The million song dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.
- [6] Donald Byrd and Tim Crawford. Problems of music information retrieval in the real world. *Inf. Process. Manage.*, 38(2):249–272, March 2002.
- [7] Avery Li chun Wang and Th Floor Block F. An industrial-strength audio search algorithm. In *Proceedings of the 4th International Conference on Music Information Retrieval*, 2003.
- [8] W.T. Glaser, T.B. Westergren, J.P. Stearns, and J.M. Kraft. Consumer item matching method and system, February 21 2006. US Patent 7,003,515.
- [9] Ichiro Fuginaga Karl MacMillan, Micheal Droettbroom. Gamera: Optical music recognition in a new shell.
- [10] Daniel Mllensiefen Klaus Frieler. The simile algorithms for melodic similarity.
- [11] Greg Linden, Brent Smith, and Jeremy York. Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet Computing*, 7(1):76–80, January 2003.
- [12] Sren Tjagvad Madsen, Gerhard Widmer, and Johannes Kepler. Key-finding with interval profiles.
- [13] Brian McFee, Matt McVicar, Colin Raffel, Dawen Liang, Oriol Nieto, Eric Battenberg, Josh Moore, Dan Ellis, Ryuichi YAMAMOTO, Rachel Bittner, Douglas Repetto, Petr Viktorin, Joo Felipe Santos, and Adrian Holovaty. librosa: 0.4.1, October 2015.
- [14] Cory Mckay. Automatic music classification with jmir, 2010.
- [15] Daniel Mullensiefen and Klaus Frieler. *The Simile algorithms documentation 0.3*. 2006.
- [16] Meinard Müller, Verena Konz, Wolfgang Bogler, and Vlora Arifi-Müller. Saarland music data (SMD). In *Late-Breaking and Demo Session of the 12th International Conference on Music Information Retrieval (ISMIR)*, Miami, USA, 2011.
- [17] Pascal Ferraroa Pierre Hannaa and Matthias Robinea. On optimizing the editing algorithms for evaluating similarity between monophonic musical sequences. 2007.
- [18] Adam M. Stark and Mark D. Plumbley. Real-time chord recognition for live performance. 2009.
- [19] George Tzanetakis and Perry Cook. Marsyas: A framework for audio analysis. *Org. Sound*, 4(3):169–175, December 1999.
- [20] Julin Urbano. Mirex 2013 symbolic melodic similarity: A geometric model supported with hybrid sequence alignment.
- [21] Wikifonia. The wikifonia lead sheet collection <http://www.synthzone.com/files/wikifonia/wikifonia.zip/>, 2013.