
Introduction to Musical Corpus Studies

Release 0.0.1

Fabian C. Moss

Aug 25, 2020

CONTENT

1 Organization	3
2 Introduction to Musical Corpus Studies	7
2.1 What are Musical Corpus Studies?	8
2.2 Epistemological goals	8
2.3 Issues	8
2.4 MCS and traditional musicology	8
3 Folk Songs and the Melodic Arc	9
4 Solos in the <i>Weimar Jazz Database</i>	11
5 Harmony in Beethoven's String Quartets	13
5.1 Access the data	13
5.2 Harmonic Annotations	13
5.3 Chord Transitions	13
6 Harmonic Clusters in Pop Charts	15
7 Harmony and Form in Brazilian Choro	17
7.1 Data	17
8 A Data-Driven History of Tonality	19
8.1 Musical Pieces as Tonal Pitch-Class Distributions	19
8.2 Historical Development	22
9 Malian Percussion Music	25
10 Electronic Music 1950–1990	27
11 Conclusion	29
Bibliography	31

Warning: This material is still (heavily) under construction and might change throughout the course!

You can help improving the course and [let me know](#) about any errors and inconsistencies that you find or suggest other ways of improving the course.

Welcome!

These pages present the content of the course “Introduction to Musical Corpus Studies” at the [Institute of Musicology](#), given at [University of Cologne](#) in Fall 2020.

In the last two decades *Musical Corpus Studies* evolved from a niche discipline into a veritable research area. The growing availability of digital and digitized musical data as well as the application and development of modern methodologies from computer science, machine learning, and data science cast new light on old musicological questions and generate entirely novel approaches to empirical music research.

Moreover, the general methodological and epistemological approach of Musical Corpus Studies allows to transcend traditional intra-musicological boundaries between its sub-disciplines (historical/systematic/ethnological/...) without sacrificing the respective specific viewpoints and perspectives.

This course offers a fundamental and practical introduction into these topics. It demonstrates, explores, and critically reflects central thematic areas and methods by means of a number of case studies. In the engagement with these topics the course also introduces elementary methods from natural language and music processing, as well as statistics, data analysis and visualization.

The course is aimed at students at the undergraduate level who have little or no empirical background and are curious about quantitative approaches to musicology.

**CHAPTER
ONE**

ORGANIZATION

No.	Date	Time	Room	Topic	Corpus	Methods
1	Fr., 13.11.20	16:00-2017:20 Uhr	Neuer Seminarraum 1.315	Introduction / Back-ground		
2		17:40-19:00 Uhr		Folk Songs, Melodies, Pitches and Intervals	Essen Folk Song Collection	frequencies, mean, variance
3	Sa., 14.11.20	09:00-2010:20 Uhr	Neuer Seminarraum 1.315	Jazz Solos, Melodies	Weimar Jazz Database	Regular Expressions
4		10:40-12:00 Uhr		Beethoven's string quartets Harmony	Annotated Beethoven Corpus	<i>n</i> -grams, Markov models
		12:00-13:00 Uhr		Lunch Break		
5		13:00-14:20 Uhr		Pop Charts Billboard 100, harmony,	McGill Billboard Dataset	Hidden Markov Models
6		14:40-16:00 Uhr		Free group work		
7	Fr., 11.12.20	10:00-2011:20 Uhr	Alter Seminarraum 1.408	Brazilian Choro, harmony, form,	Choro Songbook Corpus	Context-Free Grammars
8		11:40-13:00 Uhr		19th century piano music, harmony	DCML Piano Corpus	Probabilistic CFGs
9	Sa., 12.12.20	09:00-2010:20 Uhr	Neuer Seminarraum 1.315	Malian Percussion Music, rhythm, meter	Interpersonal Entrainment in Music Performance: Malian Jembe	
10		10:40-12:00 Uhr		Electronic Music 1950-1990	Curated Corpus of Historical Electronic Music	
		12:00-13:00 Uhr		Lunch Break		
11		13:00-14:20 Uhr		Free group work		
12		14:40-16:00 Uhr		Recapitulation and conclusion		
4					Chapter 1. Organization	

Important: Ich gehe in der Seminarplanung von 12 Semesterwochen à 2 SWS aus, für das gesamte Blockseminar also 24 SWS. Das Seminar wird mit 3 CP bewertet, was 90 Stunden aktiver Arbeit entspricht. Davon entfallen 24 SWS an die Präsenzzeit im Seminar plus 48 SWS an Vor- und Nachbereitung der Seminarsitzungen. Die verbleibenden 18 SWS sind für die Lektüre der Fachliteratur vorgesehen.

CHAPTER
TWO

INTRODUCTION TO MUSICAL CORPUS STUDIES



2.1 What are Musical Corpus Studies?

tbc... (text from diss?)

2.2 Epistemological goals

tbc...

2.3 Issues

tbc [Coo06][Hon06][Hur13][Mar16][NR16][Pug15][Sch16][TV13]

2.4 MCS and traditional musicology

tbc

CHAPTER
THREE

FOLK SONGS AND THE MELODIC ARC

Huron... / MusThe Tutorial

SOLOS IN THE WEIMAR JAZZ DATABASE



Fig. 4.1: Photo by Janine Robinson on Unsplash

The first project we will have a look at is the [Jazzomat](#) project. Transcriptions of Jazz solos [PFAbesser+17]. The *Weimar Jazz Database* (WJD) consists of 456 transcriptions of Jazz solos from diverse substyles. As all the corpora that we deal with here, it is freely available on the internet.¹

The WJD contains a number of tables:

¹ <https://jazzomat.hfm-weimar.de/dbformat/dboverview.html>

Table 4.1: Tables in the *Weimar Jazz Database*

Table name	Description
<i>beats</i>	Table for beat annotation of WJD melodies, referenced by <i>melody(melid)</i>
<i>composition_info</i>	Infos regarding the underlying composition of a WJD solo, referenced by <i>melody(melid)</i>
<i>db_info</i>	Information regarding the distributed database file like version information, license, etc
<i>esac_info</i>	EsAC infos for EsAC melodies, referenced by <i>melody(melid)</i>
<i>melody</i>	Main table for all melody events
<i>melody_type</i>	Indicated type of melody: WJD solos or EsAC (Folk songs using Essen Associative Code), referenced by <i>melody(melid)</i>
<i>popsong_info</i>	Pop song infos, referenced by <i>melody(melid)</i>
<i>record_info</i>	Infos regarding the specific audio recording of a WJD solo was taken from, referenced by <i>melody(melid)</i>
<i>sections</i>	All sections (phrase, chorus, form, chords, etc.), referenced by <i>melody(melid)</i>
<i>solo_info</i>	Solo infos for WJD solos, referenced by <i>melody(melid)</i>
<i>track_info</i>	Information specific to a track on a record (or CD)
<i>transcription_info</i>	Transcription infos for WJD solos, referenced by <i>melody(melid)</i>

Here, we focus on the main table *melody*. First, we download the entire database from <https://jazzomat.hfm-weimar.de/download/download.html> (under “Weimar Jazz Database”) and save it as the file *wjazz.db*.

```
import sqlite3
import pandas as pd

# create connection to database
conn = sqlite3.connect("wjazzd.db")

# read all entries of the `melody` table into a pandas DataFrame
df = pd.read_sql("SELECT * FROM melody", con=conn)

df.head()
```

The part of the code `SELECT * FROM melody` reads “Select all entries from the table ‘melody’”.

HARMONY IN BEETHOVEN'S STRING QUARTETS

5.1 Access the data

The data lies on the GitHub repository [DCMLab/ABC](#). Either download the `.tsv` file directly and open it in pandas or load it from the URL as follows:

```
import pandas as pd  
  
df = pd.read_csv("https://github.com/DCMLab/ABC/corpus.tsv", sep="\t")
```

The corpus is now stored in the variable `df`.

5.2 Harmonic Annotations

- regular expressions

[NHMR18][MNHM19]

5.3 Chord Transitions

- n-grams

CHAPTER
SIX

HARMONIC CLUSTERS IN POP CHARTS



Clustering analysis in [SVJ+20].

HARMONY AND FORM IN BRAZILIAN CHORO



Fig. 7.1: A 2019 musical tribute to renowned choro composer Jacob do Bandolim, on the 50th anniversary of his death in Sao Paulo, Brazil. (c) Governo do Estado de Sao Paulo / CC BY 2.0

7.1 Data

The *Choro Songbook Corpus* contains transcriptions of the chord symbols of 295 Choro pieces [MSFR20]. It is available on [Zenodo](#) or [Github](#). See this link for a performance of the piece “Lamentos”

As before, we can import our data from a public resource.

```
df = pd.read_csv("https://raw.githubusercontent.com/DCMLab/choro/1.1/data/choro.tsv",  
sep='\t')
```

And display the first rows of the data frame:

```
df.head()
```

CHAPTER
EIGHT

A DATA-DRIVEN HISTORY OF TONALITY



8.1 Musical Pieces as Tonal Pitch-Class Distributions

Musical pieces are composed of notes which have a certain pitch, a certain duration, and are located in a specific position in a piece. Here, we disregard both the location and the duration of notes and consider only the pitch dimension of musical notes when counting them. It is moreover common to consider notes to be equivalent if their respective pitches are related by one or multiple octaves, and thus to speak of *pitch classes*. Pitch classes come in two varieties. The first, most commonly used representation in computational musicology and music information retrieval, distinguishes twelve different pitch classes. This representation system also assumes the enharmonic equivalence of certain notes, e.g. F♯ and G♭, C♯ and B, etc. The assumption of enharmonic equivalence is usually a consequence of music encoding formats that assume twelve-tone equal temperament, e.g. MIDI, in which enharmonically equivalent notes are indistinguishable. The assumptions of octave and enharmonic equivalence allow to represent pitch classes

as residuals in \mathbb{Z}_{12} and to arrange them on a circle. This arrangement of pitch classes is shown in Fig. 8.1, called the *circle of fifths*.

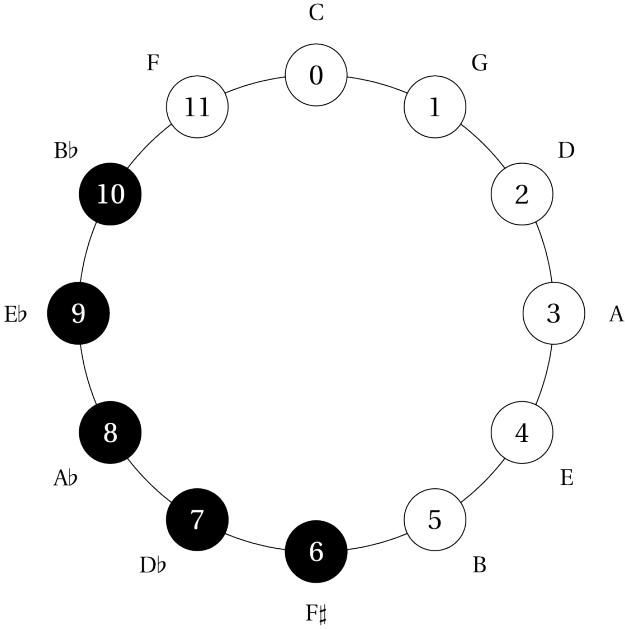


Fig. 8.1: Schematic depiction of the twelve neutral pitch-classes in \mathbb{Z}_{12} on the circle of fifths. One representative of each neutral pitch-class is shown as a tonal pitch-class label next to the node. The coloring of the nodes corresponds to the colors of the keys on the piano.

The numbers correspond to pitch classes in the order of fifths that can be transformed to chromatic ordering, the chromatic circle, by the mapping $t \mapsto 7t \bmod 12$. Pitch classes that correspond to white keys on the piano are shown in white and pitch classes that correspond to black keys on the piano are shown in black. The second variant of pitch classes does not assume enharmonic equivalence but only that octave-related notes are equivalent, and is hence more general. This representation allows to arrange pitch classes on the *line of fifths* [Tem00]. This linear ordering of tonal pitch-classes has been used by a number of music theorists, e.g. [Web51][Rie00][Han48][Bar56]. It is shown in Fig. 8.2.

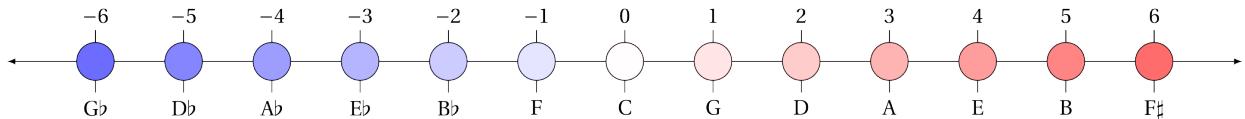


Fig. 8.2: Schematic depiction of the tonal pitch-classes on the line of fifths mapped to integers in \mathbb{Z} .

Tip: Think of other possible equivalence relations for pitch classes. E.g. which relation do we need to distinguish only two classes, high and low with respect to middle C?

Following [Tem00], we call the first representation *neutral pitch-classes* and the second one *tonal pitch-classes*. The linear structure of the line of fifths allows to associate each tonal pitch-class with an integer $k \in \mathbb{Z}$ such that this integer represents the number of perfect fifths that lie between this pitch-class and C (Gárdonyi and Nordhoff, 2002). In other words, the integer $k \in \mathbb{Z}$ corresponds to the number of flats (negative integers) or sharps (positive integers) that the major key has in which this tonal pitch-class is the root. For example, D is mapped to +2 because D major has a key signature with two sharps, Ab is mapped to -4 because the key signature for Ab major has four flats, and C is mapped to 0 because its key signature does not have any accidentals. Yet another benefit of this representation

can be seen in the way in which we associate each tonal pitch-class with a color. Positive integers ('sharpened' tonal pitch-classes) are associated with increasingly darker shades of red, negative integers are associated with increasingly darker shades of blue ('flattened' tonal pitch-classes), and C is associated with white as the neutral origin of the line of fifths. This color mapping is used throughout this part. The line of fifths does not only contain all tonal pitch-classes but also a number of central musical scales. For example, pentatonic scales are segments of length 4 (containing five pitch classes) on the line of fifths, e.g. from G \flat to B \flat , diatonic scales are defined by segments of length 6, e.g. from F to B, the early extensions of the natural diatonic scale by B \flat and F \sharp correspond to the segment spanning the eight fifths on the line between these two tonal pitch-classes, and the two whole-tone scales correspond to the odd and even numbers, respectively. Theoretically, the line of fifths extends to infinity in both directions but in actual compositions only a small segment of it is used. In the corpus that is used here, we consider only the segment from F $\flat\flat$ to B $\sharp\sharp$ because no piece in the corpus contains tonal pitch-classes outside this range. The vocabulary size of the corpus is thus $V = 35$, consisting of the seven natural pitch-classes F, C, G, D, A, E, B with two, one, or zero sharps or flats, respectively. While the transformation of tonal into neutral pitch-classes is achieved by deterministically mapping a tonal pitch-class $t \in \mathbb{Z}$ to a neutral pitch-class $t \mapsto t \bmod 12 \in \mathbb{Z}_{12}$ (in fifths ordering), the reverse direction involves some kind of inference and is called the problem of pitch spelling ([Tem01][Cam03][SRU04][CC05][Mer06]).

helix

Due to the encoding of the corpus, we adopt the representation of pieces as bags of notes (see Section~ref{sec:bagofnotes}) and represent each one as a distribution over the $V = 35$ tonal pitch-classes. That is, in the corpus used here with $D = 2012$ pieces, the tonal pitch-class distribution of a musical piece x_d is given by the relative frequencies of the tonal pitch class in that piece, for $d \in \{1, \dots, D\}$. Each piece is thus described by a V -dimensional vector, $x_d \in \mathbb{R}^V, \forall d \in \{1, \dots, D\}$, containing positive real numbers that sum to 1. In this chapter, we do not make any further assumptions about the process that generated this distribution and will postpone these considerations until Chapter~ref{chap:tonal_profiles}.

In this view, pieces simply correspond to points in the $V - 1$ -simplex

$$\Delta^{V-1} = \left\{ x_d \in \mathbb{R}^V \mid \sum_{i=1}^V x_{d,i} = 1; x_{d,i} \geq 0 \right\}.$$

In this space, those pieces with very different tonal pitch-class distributions will be very distant, whereas pieces that have similar tonal pitch-class distributions will be closer to one another and form clusters in the $V - 1$ -simplex. It is important to note that the bag-of-notes representation relies on the assumption that the V dimensions are independent, meaning that this model does not *a priori* assume any particular order between the tonal pitch-classes.

The average tonal pitch-class distribution of all pieces in the corpus is shown in Figure~ref{fig:pca_mean}. Figure~ref{fig:tpc_mean_dist_periods} in the appendix shows this distribution separately for each century from 1361 to 1943. Although we know that one can in principle order all tonal pitch-classes on the line of fifths (Fig. 8.2), we do not incorporate this assumption into the model but will show instead that it can be inferred from the data.

8.1.1 Distributions of TPCs

What are (discrete) distributions?

entropy

Hypothesis: similar pieces have similar distributions

operationalization: distances in vector space

Assumption (for coloring): tonal center is approximated by most frequent note

Taking tonal spaces into account: the line of fifths (dimensionality reduction, clustering, t-SNE, PCA)

8.1.2 Principle Component Analysis

Since the dimensionality of Δ^{V-1} is $|\{C, D, E, F, G, A, B\}| \cdot |\{\flat\flat, \flat, \natural, \sharp, \sharp\sharp\}| = 35$, it is impossible to visualize this space and pieces in it to see whether their arrangement contains any meaningful information. We address this problem by using a method for *dimensionality reduction* called *Principal Component Analysis* (PCA) [Bis06] that projects the data into a lower-dimensional space while at the same time maintaining as many characteristic properties of the original space as possible. PCA thus can aid to achieve a better understanding of the global structure of the space.¹

PCA considers the data to be represented as a matrix

$$X = \begin{bmatrix} x_1^\top \\ \vdots \\ x_D^\top \end{bmatrix} \in [0, 1]^{D \times V}, \quad x_d \in \Delta^{V-1},$$

where the rows are given by the D data points, the pieces in the corpus, and the columns are given by the V features, the number of distinct tonal pitch-classes in the vocabulary. All entries in X range between 0 and 1 and all rows of X sum to 1 because the pieces are represented as distributions. PCA determines the $M \leq V$ largest directions and magnitudes of the variance in the data in X by first calculating the *covariance matrix*

$$K_X = \text{cov}[X, X] = E[(X - \bar{x})(X - \bar{x})^\top] \in \mathbb{R}^{D \times D},$$

where E denotes the expected value and $\bar{x} \in [0, 1]^V$ is the mean of the columns of X . The main directions of the variance in the data and their magnitude is given by the eigenvectors w_i and eigenvalues λ_i of K_X which can be calculated by solving

$$K \cdot w_i = \lambda \cdot w_i.$$

The projection into the lower-dimensional space is then achieved by selecting the M largest eigenvalues and their corresponding eigenvectors and transform the data to X' , the dimensionality reduction of X by

$$X' = X \cdot [w_1, \dots, w_M] \in \mathbb{R}^{D \times M}$$

The sum of all eigenvalues λ_i is the total amount of variance in the data and the variance explained by each principal component is given by $\lambda_i / \sum_j \lambda_j$.

In the present context, X was transformed to have zero mean before applying PCA but that the variance was not standardized to 1. This was done because the features all are on the same scale and because the differences in the variance between the respective tonal pitch-classes is of particular interest here.

8.2 Historical Development

historical expansion on line of fifths,

¹ While PCA is one of the most commonly used methods for dimensionality reduction, there are many others (sometimes relying on particular assumptions about the distribution of the data). In a qualitative comparison, *Locally Linear Embedding* [RS00] achieved a similar picture, whereas *t-distributed Stochastic Neighbor Embedding* (t-SNE) [VDMH08] that does emphasize the local over the global structure of the data did not. We opted here for PCA because it preserves most of the global structure and the interpretation of the results is straight-forward.

8.2.1 LOWESS

**CHAPTER
NINE**

MALIAN PERCUSSION MUSIC

CHAPTER
TEN

ELECTRONIC MUSIC 1950–1990



Fig. 10.1: Photo by Antoine Julien on Unsplash.

**CHAPTER
ELEVEN**

CONCLUSION

Final thoughts, critical discussion...

[Some image]

BIBLIOGRAPHY

- [Bis06] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [Bar56] Layos Bárdos. Natürliche tonsysteme: methode ihrer messung. In B. Rajeczky and L. Vargyas, editors, *Studia Memoriae Belae Bartók Sacra*, pages 209—248. Akadémiai Kiadó, 1956.
- [Cam03] E. Cambouropoulos. Pitch spelling: a computational model. *Music Perception: An Interdisciplinary Journal*, 20(4):411–429, 2003. doi:[10.1525/mp.2003.20.4.411](https://doi.org/10.1525/mp.2003.20.4.411).
- [CC05] E. Chew and Y.-C. Chen. Real-time pitch spelling using the spiral array. *Computer Music Journal*, 2005. doi:[10.1162/0148926054094378](https://doi.org/10.1162/0148926054094378).
- [Coo06] Nicholas Cook. Border Crossings: A Commentary on Henkjan Honing’s “On the Growing Role of Observation, Formalization and Experimental Method in Musicology”. *Empirical Musicology Review*, 1(1):7–11, 2006.
- [Han48] Jacques Handschin. *Der Toncharakter: Eine Einführung in die Tonpsychologie*. Atlantis Verlag, 1948.
- [Hon06] Henkjan Honing. On the Growing Role of Observation, Formalization and Experimental Method in Musicology. *Empirical Musicology Review*, 1(1):2–6, 2006.
- [Hur13] David Huron. On the Virtuous and the Vexatious in an Age of Big Data. *Music Perception: An Interdisciplinary Journal*, 31(1):4–9, 2013.
- [Mar16] Alan Marsden. *Music Analysis by Computer: Ontology and Epistemology*. Springer, 2016.
- [Mer06] David Meredith. The ps13 pitch spelling algorithm. *Journal of New Music Research*, 35(2):121–159, 2006. doi:[10.1080/09298210600834961](https://doi.org/10.1080/09298210600834961).
- [MNHM19] Fabian C. Moss, Markus Neuwirth, Daniel Harasim, and Rohrmeier Martin. Statistical characteristics of tonal harmony: a corpus study of beethoven’s string quartets. *PLoS ONE*, 14(6):e0217242, 2019. doi:[10.1371/journal.pone.0217242](https://doi.org/10.1371/journal.pone.0217242).
- [MSFR20] Fabian C. Moss, Willian Souza Fernandes, and Martin Rohrmeier. Harmony and form in Brazilian Choro: A corpus-driven approach to musical style analysis. *Journal of New Music Research*, 2020. doi:[10.1080/09298215.2020.1797109](https://doi.org/10.1080/09298215.2020.1797109).
- [NHMR18] Markus Neuwirth, Daniel Harasim, Fabian C. Moss, and Martin Rohrmeier. The Annotated Beethoven Corpus (ABC): A Dataset of Harmonic Analyses of All Beethoven String Quartets. *Frontiers in Digital Humanities*, 5(July):1–5, 2018. doi:[10.3389/fdigh.2018.00016](https://doi.org/10.3389/fdigh.2018.00016).
- [NR16] Markus Neuwirth and Martin Rohrmeier. Wie wissenschaftlich muss Musiktheorie sein? Chancen und Herausforderungen musikalischer Korpusforschung. *Zeitschrift der Gesellschaft für Musiktheorie*, 13(2):171–193, 2016. doi:[10.31751/915](https://doi.org/10.31751/915).
- [PFAbesser+17] Martin Pfleiderer, Klaus Frieler, Jakob Abeßer, Wolf-Georg Zaddach, and Benjamin Burkhardt, editors. *Inside the Jazzomat - New Perspectives for Jazz Research*. Schott Campus, 2017.

- [Pug15] Laurent Pugin. The challenge of data in digital musicology. *Frontiers in Digital Humanities*, 2(4):1–3, 2015. doi:10.3389/fdigh.2015.00004.
- [Rie00] Hugo Riemann. *Musik-Lexikon*. Max Hesse Verlag, 5 edition, 1900.
- [RS00] S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000. doi:10.1126/science.290.5500.2323.
- [Sch16] Kris Schaffer. What is computational musicology? online, 2016. URL: <https://medium.com/@krishaffer/what-is-computational-musicology-f25ee0a65102>.
- [SVJ+20] Kris Schaffer, Esther Vasiete, Brandon Jacquez, Aaron Davis, Diego Escalante, Calvin Hicks, Joshua McCann, Camille Noufi, and Paul Salminen. A cluster analysis of harmony in the McGill Billboard dataset. *Empirical Musicology Review*, 14(3–4):146–162, 2020. doi:10.18061/emr.v14i3-4.5576.
- [SRU04] J. Stoddard, C. Raphael, and P. E. Utgoff. Well-tempered spelling: a key-invariant pitch spelling algorithm. In *Proceedings of the 5th International Conference on Music Information Retrieval (ISMIR 2004)*. 2004. doi:10.1.1.100.9663.
- [Tem00] David Temperley. The line of fifths. *Music Analysis*, 19(3):289–319, 2000. doi:10.2307/854457.
- [Tem01] David Temperley. *The Cognition of Basic Musical Structures*. MIT Press, 2001.
- [TV13] David Temperley and Leigh VanHandel. Introduction to the Special Issue on Corpus Methods. *Music Perception: An Interdisciplinary Journal*, 31(1):1–3, 2013.
- [VDMH08] L. Van Der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(November):2579–2605, 2008.
- [Web51] Gottfried Weber. *The Theory of Musical Composition, Treated with a View to a Naturally Consecutive Arrangement of Topics*. Messrs. Robert Cocks and Co., 1851.