

The MeloSol Corpus

David John Baker¹

¹ Louisiana State University

Author Note

David John Baker now works at Flatiron School in London, England.

Correspondence concerning this article should be addressed to David John Baker, .

E-mail: davidjohnbaker1@gmail.com

Abstract

9 This paper introduces the *MeloSol* corpus, a collection of 783 Western, tonal monophonic
10 melodies. We first begin by describing the overall structure of the corpus, then proceed to
11 detail its contents as they would be helpful for researchers working in the field of
12 computational musicology or music psychology. In order to contextualize the MeloSol
13 corpus, compare descriptive statistics generated using the FANTASTIC feature extraction
14 toolkit with that of the Essen Folk Song Collection as well as The Densmore Collection of
15 Native American Songs. We suggest possible uses of this corpus including extending
16 research which investigates Western tonality, perceptual experiments needing novel
17 ecological stimuli, or work involving the musical generation of monophonic melodies in the
18 style of Western tonal.

19 *Keywords:* corpus studies, FAIR data, kern

20 Word count: X

The MeloSol Corpus

Introduction

This data report introduces the *MeloSol* corpus, a collection of 783 monophonic melodies taken from *A New Approach to Sight Singing: Fifth Edition* (Berkowitz, Fontrier, Kraft, Goldstein, & Smaldone, 2011). The title *MeloSol* derives from a combination of the corpus' content—*Mel*odic data—and the first name of the original author of the collection, *Sol* Berkowitz.

The corpus is divided into two major sections: a collection of sight singing melodies composed specifically for pedagogical purposes ($n = XXX$) and examples from the Western Classical Music canon ($n = XXX$).

- Point of Edit

Within each of the two larger sections exists FIVE further subdivisions. These five subdivisions tend to be mapped in conjunction with aural skills classroom. For example, the first section of both the sight singing melodies and the first section of the Literature align with melodies that a first semester undergraduate student would be expected to learn in their first semester of college in an aural skills classroom. Each section is meant to increase in difficulty. The fifth and final section of both the sight singing melodies and examples from the literature contains melodies either meant to be atonal or have some sort of unstable tonality (bi-tonality/modality). A visual depiction of the breakdown of melodies from the two larger sections in terms of count data is presented IN FIGURE HERE.

- FIGURE HERE

In terms of analyzable data, the 783 melodies are all encoded in **kern** format with each file containing metadata listing the excerpt's LIST HERE. Overall, the corpus consists

of XXXXX digital tokens, a subset of which are XXX note heads. All melodies in the corpus were encoded by hand by the author using MUSE SCORE 3, initially saved as XML, then converted to kern using the HUMDRUM EXTRAS `xml2hum` with the current meta data added using the `name-of-script.R` file. Further addition to the metadata can be added with modifications to `name-of-script.R`.

From a more meaningful point of view, the descriptive statistics of the corpus are displayed in FIGURE TWO and FIGURE THREE.

- FIGURE TWO (subset out Section Five)
- FIGURE THREE (Section Five)

Comparison

Further descriptive statistics of the corpus generated from MULLENSIEFEN'S FANTASTIC TOOLBOX can help contextualize the *MELOSOL* corpus in context with other corpora commonly used in the literature. One of the most cited corpora in the field of computational musicology is ESSEN. ESSEN contains XYZ and is often taken as proxy for representing implicit understanding via statistical learning (HURON) and PEARCE and OTHERS. Essen also has Chinese songs. The *MeloSol* corpus also falls under umbrella of Western Music and as discussed below, may be helpful novel corpus for continuing to investigate claims. Since publication of ESSEN there has also been DENSMORE CITE. DENSMORE is collection of melodies encoded by Shanahan and Shanahan. From musicological point of view, both DENSMORE and CHINESE are expected to be different for reasons of both location as well as style. Here we compare the two to get high level reduction of idea.

First in FIGURE FOUR we compare high level descriptive statistics between WESTERN ESSEN, CHINA, DENSMORE, and MELOSOL. Figure contains comparative

overlap of LIST OF FEATURES HERE. Note there is a very large difference in the size of MELOSOL (and others) compared to ESSEN.

- FIGURE FOUR

Second in FIGURE FIVE we compare MORE ABSTRACT FEATURES

- FIGURE FIVE

Here is some small comparison on the differences in features.

Useful

As the *MeloSol* corpus is made of Western music, can be used to continue research that has made claims about certain features of Western music if need proxy. For example there are a lot of claims made by HURON about contour class that have been initially explored by BAKER. There also have been lots of modeling of expectation using IDyOM by Pearce that have used ESSEN. If buy the idea of sample population as generation, this could be taken forward in that area.

Also note that now have dataset was initially generated using pedagogical materials and might be helpful in that domain. For example, extending work of MY DISSERTATION could look at proxies of difficulty using FANTASTIC. Could also see if enough data here can be used for generative data analyses using LSTM.

Data analysis

We used R (Version 3.6.2; R Core Team, 2019) and the R-package *papaja* (Version 0.1.0.9942; Aust & Barth, 2020) for all our analyses.

References

89

90 Aust, F., & Barth, M. (2020). *papaja: Create APA manuscripts with R Markdown*.

91 Retrieved from <https://github.com/crsh/papaja>

92 Berkowitz, S., Fontrier, G., Kraft, L., Goldstein, P., & Smaldone, E. (2011). *A new*

93 *approach to sight singing* (5th ed). New York: W.W. Norton.

94 R Core Team. (2019). *R: A language and environment for statistical computing*. Vienna,

95 Austria: R Foundation for Statistical Computing. Retrieved from

96 <https://www.R-project.org/>