

## Searching Online for Early Music

It could be observed that music history consists of many strands of information that by turns ebb into a pool and flow apart. Writing from the perspective of someone who has been engaged in methods of searching written scores of music (as one would search text files for words or phrases), my aim here is to share some observations on the searchability of music from the fifteenth through the eighteenth centuries. There are now quite a few search engines for encoded music.<sup>1</sup> Each one has its own datasets, but the overall level of cooperation among those discussed here is considerable.

As dry and dehumanized as computer-assisted research is sometimes perceived to be, no account of a subject such as this one can be entirely value-neutral. Despite a legacy of more than a half century of effort, the rapid evolution of technology has forced this field, which is of little interest to industry, to fend for itself through countless changes of computer operation and endless rounds of data migration. These comments are derived from involvement in our work (from 1984) at the Center for Computer Assisted Research in the Humanities (CCARH); our graduate courses in computer-assisted music topics at Stanford University (from 1997); and our various collaborative involvements with multiple spokes of the RISM project (principally from about 2000); and our in-house projects MuseData (from 1984) and Themefinder (from 1996). The last-named now have progeny of their own, which is mentioned as appropriate below. A list of web links for sites mentioned is provided at the end.

Score searching must be distinguished from audio searching, which is usually intended by the term “music search” in common parlance. Audio search is heavily supported by commercial interests with an inclination towards seeking market advantage as opposed to musical validity.<sup>2</sup> Having been developed almost entirely in academic and non-profit quarters, the progress of score searching (generally called “symbolic search”) has inevitably been slower. Its needs are in many

---

<sup>1</sup> I am immensely grateful to Craig Stuart Sapp for preparing the musical examples and for help in illuminating some of the features of the Josquin Research Project.

<sup>2</sup> Audio search can produce results that are more refined for the details of timing *in recordings* but far less competent in efforts to identify melodies, themes, harmonic patterns, or procedures. Much of the audio-search software in current use relies at least partly on meta-data (text fields such as title, author, performer, and total length of “track”), sometimes with audio samples of dynamic level and timbral composition.

respects more complex, however. Apart from many well documented problems of a technical nature, the effectiveness of a user-initiated search for music depends both on the accuracy of the user's recall but also on the user's general appreciation of materials available and search strategies appropriate for the music sought. There has been a growing appreciation of a collateral need for studies in music cognition, which are regularly reported by professional societies in Europe, Asia, and North America but which rarely have interchanges with the early music community.

The musical examples discussed here are intended more to entice new users to explore on their own than to answer specific research questions. We begin with a few useful principles. Most music that may interest performers may not be searchable. Digitized scores are of enormous value but encoded data is more specific and provides a "higher-resolution" impression than a photographic image. This kind of information will be of value to those seeking to explore questions of musical similarity, which are endemic in music history and of considerable interest to early musicians because so much of early music provides ready examples.

A central area of complexity, however, is that two pieces of music can be similar in countless different ways. There is no ranking system for degrees of similarity in music. Certain "tune families" are well explored in folk music research, and studies in chant centonization are also at hand. A vast middle ground stands apart these instances—music derived from pre-existing songs, from subjects derived from names, and from chants associated with specific feasts; repertories for lute and keyboard based on songs of many kinds and, later, on *bassi ostenati*; cyclically organized masses and keyboard suites; arias in free circulation that varied from singer to singer.

Among these conceivable contexts an important question to pose is where one should be as specific as possible and where as general as possible. Search for music can be literal and exacting. Many levels of detail are possible. "Fuzzy" search that allows for variability (by key, by elimination of the articulation of offbeats, and sometimes by meter) can be valuable in appropriate circumstances. In the space at hand it is possible to consider in detail only one variable, that of pitch.

Pitch is the most useful feature by which encoded music is searched and the one that in most cases best identifies characteristic passages of a particular work. In written music pitch has three parameters—a note name (A, B, C, D, E, F, G), an inflection (#b-), and an octave. In pitch searches the level of detail desired by the user must be supported both by the underlying encoding of the music and by the search software itself. Musicians can easily be deceived about the capability of the tool in use, because they are often conversant with a written or aurally remembered example

of the music sought. We call search capabilities that are limited to pitch names a diatonic or base-7 pitch representation. A full chromatic roster of written notes within an octave (C, C#, Db, D etc.) would be a base-21 system.<sup>3</sup> Octaves are numbered, although numbering systems vary from one to the next.<sup>4</sup> Diatonic systems are almost always adequate for chant repertoires and medieval monophony.

In schemes of pitch representation the ubiquity of MIDI instruments, files, and software have imposed another representation on the tones of the octave. These are “key numbers” (invented to address the physical keys of an electronic instrument). They work well for those with an aural orientation whose interests lies in pitch classes, since in an equal-tempered scale there is no sounding difference between C# and Df. This means that software can not distinguish between the two; it assigns them the same number. This can be a problem for hexachordal repertoires, since a long-standard practice in the MIDI trade has been to render all black notes as A#s. A soft hexachord is therefore impossible to reproduce correctly.<sup>5</sup>

Since most underlying score data available today comes from MIDI files, it suffers from a paucity of information to define pitch enharmonically. This would require a base-21 representation of differently named tones within an octave. MIDI has also been used as a default interchange format, and because of its limitations it can mean that a file acquired from another user or program is likely to have passed through a filter that eliminates pitch refinements and much else. MIDI is also the default basis for transposition in commercial software. This produces reasonable results in

---

<sup>3</sup> Extended systems accommodate double sharps and double flats, which can be useful for enharmonic notations of Baroque keyboard music. Walter B. Hewlett’s base-40 system goes one step further by inserting five null tokens (empty slots) to preserve correct enharmonic spellings in automatic transposition and harmonic assessment. Hewlett, Walter B. (1992). “A Base-40 Number-line Representation of Musical Pitch,” *Musikometrika* 4, 1-14 (reproduced at <http://www.ccarh.org/publications/reprints/base40/>) and “Method for Encoding Music Printing Information in a MIDI Message,” U. S. Patent 5,675,100 (October 7, 1997). The MIDI implementation later became known as MIDIPlus and is described in E. Selfridge-Field, *Beyond MIDI: The Handbook of Musical Codes* (MIT Press, 1997); reproduced with permission at <http://beyondmidi.ccarh.org/beyondmidi-600dpi.pdf>.

<sup>4</sup> Variability in octave labeling can be an issue for MIDI implementations in commercial systems. For example Peachnote (<http://peachnote.com>), which has numerous valid uses, currently produces playback that is an octave too high when the virtual keyboard is used, but can produce the correct pitch when the alphanumeric search box is reset to Middle C = 60.

<sup>5</sup> Over time commercial notation programs such as Finale and Sibelius (among many others) have developed algorithm solutions that remedy a high percentage of the errors inherent in the keyboard capture. The algorithms are usually optimized for works in major keys. Music in minor keys and unusual modes will still incur a noticeable number of errors.

conventional works in major keys not too distant from C Major, but degrades as one moves further away from convention.

### *Search Engines for Music*

The following search engines discussed here come are related to two large clusters of projects—those originating or maintained at the Center for Computer Assisted Research in the Humanities (CCARH) at Stanford University and those related to the Répertoire International des Sources Musicales (RISM). These project constellations have entirely different purposes but support search technology for encoded musical data. *MuseData*, focused on the development of means to encode, print, and archive full scores to a scholarly standard, operates on in-house software developed by Hewlett. Since 1984 a team of data specialists has encoded more roughly 1200 works (1680-1850) including significant quantities of orchestral and chamber repertory plus selected operas and oratorios. The data is freely downloadable and a few hundred of the scores are also downloadable as PDFs. Most of the data has been translated into several other formats for music data.

*Themefinder*, which stores musical incipits of diverse repertories totaling slightly more than 100,000 items, was initiated by David Huron, a CCARH visitor in 1996 and subsequently. His aim was to monitor how users go about searching for music. Since that time the search capability has been expanded as the repertory has grown. The sample repertories it includes come from classical works (1650-1900), Latin motets, and folksongs from Europe and Asia. Five levels of search are supported: (1) pitch (A..G plus differentiated sharps and flats); (2) intervals (including direction); (3) scale degree (1..7); (4) gross contour (Up, Down, Repeat); and (5) refined contour (up / down by step or leap, or repeat).<sup>6</sup> *Themefinder* also has filters for meter and mode.

The Josquin Research Project (JRP), under the direction of Jesse Rodin with implementation by Craig Sapp, is a project of recent vintage (2010) currently holds about 500 works. Its search engine is adapted from *Themefinder* to better suit a repertory in mensural notation. These include such fields as pitch, interval, and rhythm. Filters for genre and mensuration are among those

---

<sup>6</sup> Most of the software in current use has been designed and implemented by Craig Sapp. Andreas Kornstädt designed the original interface. Roughly 20 former Stanford students have worked on various aspects of the project. Datasets have also been contributed by various colleagues.

currently supported. Holdings are not limited to Josquin (whose works are subdivided into those with secure attributions and those without), Ockeghem, Obrecht, and Dufay<sup>7</sup>, and several others.

The RISM project, which can be traced to early projects in the 1950s, supports and coordinates music bibliography of many kinds. The most visible project is the music manuscript inventory ("A II" in the original nomenclature), which has been searchable online since 2011. Its aim is to inventory all manuscripts throughout the world containing music from the seventeenth and eighteenth centuries. It originated (as did *MuseData*) long before the internet, but RISM was from a start a computer-based project. However a lot of data it contains was first transcribed manually and later encoded. The manuscript inventory is rich in text fields, but its one field for encoding musical incipits is our focus here.<sup>8</sup> It was an early hope that the music-transcription field would facilitate the ready identification of anonymous works.<sup>9</sup>

The manuscript inventory is based on the transcription code Plaine & Easie, which was developed in 1966 by Barry Brook and Murray Gould.<sup>10</sup> Although the format has been modified and the data translated many times, the encoded incipits remain reasonably transparent. At this writing the central collection serves almost 900,000 listings. The completion of the composite (international) collection managed by the RISM editorial office in Frankfurt, with online access facilitated by the Bavarian State Library in Munich, is still some years away. Full synchronization is the end goal, but since cataloguing and software development are still in progress, some cooperating RISM national repositories provide their own search engines. These include those of Ireland, Switzerland, the United Kingdom, and the United States (herein cited as RISM EI, RISM CH, RISM UK, and RISM US).

The small differences between them are quite instructive to those interested in further exploring the myriad possibilities and pitfalls of music search. In 2004 CCARH set up, with the

---

<sup>7</sup> The Dufay holdings have been contributed by Alejandro Planchart and are currently in process of data-translation to the JRP music format.

<sup>8</sup> Other RISM projects focus on printed music by single composers, anthologies of printed music, music theory, and so forth.

<sup>9</sup> The last time its efficacy was assessed, the results had greater implications for dis-attribution than for attribution. See Joachim Schlichte, "Der automatische Vergleich vom 83,243 Musikincipits aus der RISM Datenbank: Ergebnisse—Nutzen—Perspektiven," *Fontes artis musicae* 37/1 (1990), 35-46, and John Howard, "Strategies for Sorting Melodic Incipits," *Melodic Comparison: Concepts, Procedures, and Applications, Computing in Musicology* 11 (1998), 119-128.

<sup>10</sup> The most recent description is by Howard, "Plaine and Easie Code: A code for music bibliography" in *Beyond MIDI*, pp. 326-72.

cooperation of the central editorial offices and the US RISM office, an experimental website for searching US music manuscripts. Based on the Themefinder search engine, it preserved Themefinder's multi-tiered search functions but also supports text search of some principal fields. This example inspired a virtual-keyboard adjunct to Laurent Pugin's Swiss RISM search site, parts of which were later incorporated in the UK RISM website. Pugin's music search page replaces Themefinder's search boxes with sliders for pitch and duration (available at the advanced search page <http://www.rism-ch.org/manuscripts/search?strategy=index> after clicking "incipits"). It also has filters for meter. Since many medieval manuscripts have been included at the Swiss RISM website, this ordinary fielded search supports searches by liturgical feast and other fields appropriate for such repertories.

### *Sample Searches*

A test set of incipits was assembled for the purpose of exploring the capabilities of these various search engines. Direct comparisons are not possible. The fact that there are few overlaps in data holdings means that some kinds of queries are destined to work only in one of situations. The six themes chosen are grouped into three categories: (1) Renaissance dances, (2) authored polyphony (Ockeghem, Josquin) from the Renaissance, and (3) late Baroque music (Handel, Bach). Some Renaissance dances had a long afterlife that carried through the eighteenth century. In a study that was nothing short of heroic, Luigi Ferdinando Tagliavini assembled a stunning collection of instances of the "Ballo di Mantova".<sup>11</sup>



Simple as the melodic outline is and quick though hearers are to find in it an antecedent of Smetana's "Má Vlast" (1879) and the Israeli national anthem (adopted in 1948), Tagliavini's manual cull of loose matches, which now ranges far beyond 100 examples, is concentrated in the 16<sup>th</sup>-through 18<sup>th</sup> centuries. It is not difficult to find among them pieces that deviate by key, meter, mode, or rhythmic detail in ways that some would consider to invalidate the match. Yet the study is enormously valuable for demonstrating that a determined human being can outperform a

---

<sup>11</sup> Tagliavini, L. F. "Il ballo di Mantova, ovvero, *Fuggi, fuggi da questo cielo*, ovvero, *Cecilia*, ovvero..." , *Max Lütolf zum 60. Geburtstag: Festschrift*, (Basel: Wiese, 1994). Tagliavini's collection is now being maintained and expanded by Liuwe Tamminga.

computer search by a very large distance. In most of them a musical ear will hear a melody similar to Smetana's piece.

A melodic (note-by-note) search for the first 9 notes of the "Ballo di Mantova" in the RISM DE (international) database located only 11 instances and, as with most melodic searches, produced some hits of questionable value. Yet a melodic search should be desirable in cases such as this one where titles are almost as numerous as instances of the melody. In a scale-degree search with a minor mode filter and two wildcards in 9-note profile *Themefinder* found 30 matches, some of which fail for rhythmic reasons. An incidental match occurred for what is labeled "Schöne Minka", otherwise identified as Beethoven's "Schöne Minka, ich muss scheiden", said to be on a Cossack air.

with held to be a Russian folksong ("Schöne Minka" in the Essen database, an encoded folksong repository). A US RISM search located 10 hits from manuscripts in the US, including a transcription for viol from the collection of the Harvard Music Club. Musicians familiar with other tune families of the sixteenth and seventeenth centuries will appreciate that the Ballo di Mantova phenomenon (and its attendant search problems) finds numerous parallels in such popular items as La Girometta, chaconne, Folia, and others.

### ***RENAISSANCE POLYPHONY***

Searching for Renaissance polyphony raises different issues. It is widely held view and the foundation of generations of research that much of the period's more elaborate music, especially sacred vocal music, is based on earlier secular songs and includes melodies once associated with particular texts, often in tightly knit schemes of rhyme and repetition. The degree of matching can revolve in the polyphonic context on other kinds of detail, among them mensuration, number of voices, rhythmic properties (which can vary greatly from one part to another). Some of this kind of variance can be judged from a comparison of Ockeghem's three-voice song "Ma Maistresse" and his four-voice mass.

Discantus

Tenor

Contra

When we turn to the start of the mass setting we find a significant redeployment of voices such that they are treated in pairs and in a more clearly organized imitative manner. However the new discantus interleaves elements of the old discantus with the old contra voice.

Discantus

Contratenor

Contra

Bassus

This comparison barely scratches the surface of an important implication for automatic melodic searches: not only may the line sought wander from part to part but it may also decay and be recomposed. Searches for melodic matches may confront a new battery of obstacles.

However, other music of the period present much more straightforward possibilities. An almost contrary example is Josquin's well-known *soggetto cavato* "Hercules Dux Ferrariae", which he used as the basis for a mass. The eight tones of the subject were derived by vowel substitution. A search through the JRP located 11 instances of its use in the like-named mass. They heavily concentrated in the Sanctus. In searches elsewhere RISM DE provides one quasi-match by Cristóbal de Morales in a Latin motet for Annunciation. The melody seems not to have been one that natural powers of invention could have conjured up. The "L'homme arme" melody, in contrast, elicits 175 matches in 71 works through its intervallic search (with the expression 1 4 1 -2 -2 -2).

### ***THE BAROQUE***

To represent the baroque era we selected contrasted search examples—the aria "Lascia ch'io piango" from Act Two of Handel's *Rinaldo* and the Bach fugal subject B-A-C-H. The aria is a rhythmically simple diatonic one, while the Bach *soggetto cavato* is not simply chromatic but may require enharmonic definition in some search contexts. The difficulty one may encounter in searching for the aria is that because of its enduring popularity all the liberties of baroque interpretation have accrued to its performing history over three centuries. Although historians

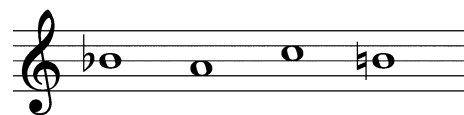


generally see music as fixed to its time of origin and therefore an authoritative version to produce legitimate matches, the nineteenth century produced far more musical manuscripts than earlier centuries and modifications were rampant. A text search on the aria in RISM DE cited among other sources two early ones (GB-Ob Mus. Sch. C. 41, GB-Lam MS 90), one from later in the eighteenth century (US-NYp Mus. Res. JOG 72-138), and one from the nineteenth (I-BRc 19<sup>th</sup> cent, Fondo Pasini 1). The subtle differences are easy to identify in comparison.



The surprising truth is that in a melodic search “Lascia ch’io piango” yielded the lowest hit rate across all the search engines mentioned here. (RISM DE identified 25 instances in title search. With a wild-card search Themefinder offered turned up a loose relationship to two movements in Handel harpsichord site, and with a meter filter it located a copy of the aria in the New York Public Library.) Because the singing of baroque arias still varies as much today as it did in Handel’s time, pitch searching that is too faithful to any one iteration is likely to miss others that are legitimate. Although the addition of rhythmic search features is a widely held goal of music search engines, in this case rhythmic features carry things further afield because of this repertory’s well known susceptibility to variations in dotting patterns.

The Prelude and Fugue on the name B-A-C-H (the authenticity of which is opened to question, is in Bb Major) is a useful test case of a different kind. It is short (too short, it seems, to RISM UK’s expectations of 5-7 note searches) and distinctive.



The German music spelling B [Bb] – A – C – H [B<sup>n</sup>] proved to be as difficult to locate with virtual keyboards as without them. Virtual keyboards usually rely on MIDI, because when a user plays a black note the software is unable to distinguish A# from Bb. (Trained musicians tends not to realize this lapse, because they have a clear conception of the configuration.) The DE RISM search engine offers a pop-up window for refine pitch spelling for each note of the virtual keyboard. The user can select three possible interpretations for of the twleve physical keys. (The black note between A and B can be defined as xA / bB / bbC; the B natural as B / xxA / bC.) It locates the four-note theme in 149 works, but the number of hits drops rapidly as the theme is extended (a phenomenon common to all melodic searches). Among its more likely matches is a theme in 34 from Act III of Kurt Stiegler’s parody *Der Thomaskantor* (1928) on a text by F. A. Geissler. US RISM, utilizing the *Themefinder* search engine, finds five examples, among them the same J. L. Krebs “Fuga alla breve” (Yale School of Music Library) as RISM reports in the Benedictine abbey of St. Boniface (Munich), the Prussian State Library (Berlin).

A postscript to the B-A-C-H search comes from an experimental search through all the MuseData holdings for J. S. Bach in 1996 by Hewlett.<sup>12</sup> The aim was to demonstrated the viability of direct data searches in encoded data, in this case given complete encodings at a high level of pitch specification. While some of the results were entirely coincidental, others challenge us to consider whether or not they are so. For example, in the duet “Du wahrer Gott und Davids Sohn” of the cantata BWV 23) the local context is chromatic, while the specific match coincides with the words “erbarm’ dich” (“Have mercy”). Clearly the prospective of deeper penetration into repertories polyphonic repertories will perplex with further questions of interpretation.

The image shows a musical score for a duet. The Soprano part (Sop.) and Alto part (Alto) are written in treble clef. The Soprano part has lyrics: "barm' dich mein, er - barm' dich mein!". The Alto part has lyrics: "dich, er - barm' dich mein!". Above the Soprano part, the letters "B A C H" are written above the notes. Below the Alto part, there is a bass line with figured bass notation. The figures are: 7 6 6b 7 6 7b 5, 6 4 5 6 6, # 6 4+ 6 4, 7b 6 3 6 4 7, and 5 4 3. The score is in 2/4 time and has a key signature of one flat (Bb).

<sup>12</sup> Walter B. Hewlett, “A Derivative Database Format for High-speed Searches,” *Computing in Musicology* 10 (1995-96), 131-42.

## ***Conclusions and Caveats***

Since melodic searching is still in its infancy and since no dataset contains more than a modest portion of all the repertoires that figure in the wider purview of early music, users can expect to experience many partial successes in music search, irrespective of the search engine and its interface. Through a better understanding of the diverse methods of pitch searching they can optimize their changes by a making accommodations in query construction. Other frontiers loom on the horizon—methods of metrical and rhythmic search, coordinated pitch and lyrics searches, and support for seeking musical features pertinent to particular repertoires. For the time being they are invited to explore the search sites listed below.

### **Music-search websites mentioned (open-access)**

CCARH-maintained websites:

CCARH: <http://www.ccarh.org>

Josquin Research Project: <http://jrp.stanford.edu>

KernScores: <http://kern.humdrum.org>

MuseData: <http://musedata.stanford.edu/>

Themefinder: <http://themefinder.stanford.edu/>

RISM-related websites:

RISM CH: <http://rism-ch.ch>

RISM DE: <https://opac.rism.info/metaopac/start.do?View=rism>

RISM IE: <http://www.rism-ie.org/> (no incipit search yet)

RISM UK: <http://www.rism.org.uk/>

RISM US: <http://rism.themefinder.org>

Teaching websites on musical search and analysis at Stanford University:

[http://wiki.ccarh.org/wiki/Music\\_253](http://wiki.ccarh.org/wiki/Music_253)

[http://wiki.ccarh.org/wiki/Music\\_253/CS\\_275a\\_Syllabus](http://wiki.ccarh.org/wiki/Music_253/CS_275a_Syllabus)

[http://wiki.ccarh.org/wiki/Music\\_254](http://wiki.ccarh.org/wiki/Music_254)

[http://wiki.ccarh.org/wiki/Music\\_254/CS\\_275b\\_Syllabus](http://wiki.ccarh.org/wiki/Music_254/CS_275b_Syllabus)

<http://kern.humdrum.org>