



Norwegian University of  
Science and Technology

# Combining Audio Fingerprints

Vegard Andreas Larsen

Master of Science in Computer Science

Submission date: June 2008

Supervisor: Trond Aalberg, IDI

Norwegian University of Science and Technology  
Department of Computer and Information Science



# Problem Description

The thesis will investigate the possibility of combining two or more existing acoustic fingerprinting solutions into a common solution, to generate a more universal fingerprint. By combining the probabilities from each fingerprinting solution in various ways into a combined score, the solution will be more adept at finding equivalent recordings than a single solution will be. The thesis will test the fingerprinting systems accuracy separately, and then combine them in different ways to see if a better set of matches can be found.

Assignment given: 15. January 2008  
Supervisor: Trond Aalberg, IDI



# Abstract

Large music collections are now more common than ever before. Yet, search technology for music is still in its infancy. Audio fingerprinting is one method that allows searching for music.

In this thesis several audio fingerprinting solutions are combined into a single solution to determine if such a combination can yield better results than any of the solutions can separately. The solution is used to find duplicate music files in a personal collection.

The results show that applying the weighted root-mean square (WRMS) to the problem most effectively ranked the results in a satisfying manner. It was notably better than the other approaches tried. The WRMS produced 61% more correct matches than the original FDMF solution, and 49% more correct matches than libFoolD.



# Contents

<b>Contents</b>	<b>iii</b>
<b>List of Figures</b>	<b>v</b>
<b>List of Tables</b>	<b>vii</b>
<b>Acknowledgments</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Objectives . . . . .	1
1.3 Approach . . . . .	2
1.4 Results . . . . .	2
1.5 Structure . . . . .	2
<b>2 Pre-study</b>	<b>5</b>
2.1 Sound theory . . . . .	5
2.2 Music Information Retrieval . . . . .	10
2.3 Audio fingerprinting . . . . .	12
2.4 State of the art . . . . .	16
2.5 Software . . . . .	17
<b>3 Method</b>	<b>21</b>
3.1 Music collection . . . . .	21
3.2 Descriptor architecture . . . . .	26
3.3 Software . . . . .	26
3.4 System architecture . . . . .	28
3.5 Combining results . . . . .	29
<b>4 Results</b>	<b>39</b>
4.1 Evaluating results . . . . .	39
4.2 Individual descriptors . . . . .	39
4.3 Combined solutions . . . . .	50
4.4 Examples . . . . .	56
4.5 Improving search speed . . . . .	56

<b>5 Conclusion</b>	<b>59</b>
5.1 Results . . . . .	59
5.2 Evaluation . . . . .	59
5.3 Further work . . . . .	60
<b>Bibliography</b>	<b>61</b>
<b>Index</b>	<b>65</b>
<b>A Software used</b>	<b>69</b>
<b>B Source code</b>	<b>71</b>
B.1 fdmf . . . . .	71
B.2 libFooID . . . . .	81
B.3 Gunderson's descriptors . . . . .	88
B.4 Combining descriptors . . . . .	93
B.5 Licenses . . . . .	109
<b>C Music collection</b>	<b>123</b>



# List of Figures

2.1	Features of waves illustrated. . . . .	6
2.2	Combination of three sine waves results in complex wave. . . . .	7
2.3	Example of audio sampling . . . . .	9
2.4	Example of audio wave . . . . .	14
2.5	Example frequency spectrum . . . . .	15
3.1	FRBR terminology . . . . .	22
3.2	Illustration of clustered duplicate pairs . . . . .	24
3.3	Euclidean distance for the mfcc_avg descriptor . . . . .	28
3.4	Database architecture . . . . .	29
3.5	The flow of data through the system. . . . .	30
3.6	Weights by count, score, average and performance . . . . .	36
3.7	Three sets of manual weights. . . . .	36
4.1	Results from descriptors . . . . .	40
4.2	False positive rate for individual descriptors . . . . .	42
4.3	False positive rate from simple individual descriptors . . . . .	43
4.4	Simple overview of false positive rate in individual descriptors . . . . .	44
4.5	False positive rate in combined solutions . . . . .	51
4.6	False positive rate in combined solutions (close-up) . . . . .	52
4.7	Simple overview of false positive rate in combined solutions . . . . .	53
4.8	Potential index performance . . . . .	58



# List of Tables

2.1	The chromatic scale . . . . .	8
2.2	MD5 hashes of three similar text strings. . . . .	12
2.3	Failure modes. . . . .	13
3.1	Thresholds for various descriptors by Gunderson. . . . .	27
3.2	Weights used . . . . .	37
4.1	First three erroneous duplicate pairs from fdmf_0. . . . .	41
4.2	First three erroneous duplicate pairs from fdmf_1. . . . .	45
4.3	First three erroneous duplicate pairs from fdmf_2. . . . .	45
4.4	First three erroneous duplicate pairs from libFooID. . . . .	46
4.5	First three erroneous duplicate pairs from the centroid descriptor. . . . .	46
4.6	First three erroneous duplicate pairs from using track length as a descriptor. . . . .	47
4.7	First three erroneous duplicate pairs from the mean/square ratio descriptor. . . . .	47
4.8	First three erroneous duplicate pairs from the steepness descriptor. . . . .	48
4.9	First three erroneous duplicate pairs from the mfcc_avg descriptor. . . . .	48
4.10	First three erroneous duplicate pairs from the mfcc_delta_avg descriptor. . . . .	49
4.11	First three erroneous duplicate pairs from the flcc_avg descriptor. . . . .	49
4.12	First three erroneous duplicate pairs from the flcc_delta_avg descriptor. . . . .	49
4.13	False positive rate of combinations . . . . .	54
4.14	Weights by performance for various descriptors. . . . .	56
4.15	Scores for example 1 . . . . .	56
4.16	Scores for example 2 . . . . .	57
4.17	Scores for example 3 . . . . .	57
4.18	Performance characteristics of scalar descriptors as index . . . . .	58



# Acknowledgments

Thanks to Trond Aalberg for supervising me throughout the semester with this thesis.

Thanks to Gian-Carlo Pascutto for providing me with helpful hints regarding libFooID, even though he no longer actively maintains the project. Thanks to Kurt Rosenfeld for giving me insight into the inner workings of fdmf. Thanks to Steinar Gunderson for providing helpful pointers on how to modify his system.

Thanks to (in alphabetical order) Trond Aalberg, Anne Cecilie Burhol, Hilde Halland, Francis Rath and Frode Sandholtbråten for proof-reading drafts of this thesis.

Vegard Andreas Larsen

Trondheim, June 3rd 2008



# Chapter 1

## Introduction

### 1.1 Motivation

Music is everywhere. In the 150 years since its invention recorded music has come into everyday use. Music players are now so compact that mobile phone producers include them in almost every phone they sell.

It is now normal to have large digital music collections that were unfeasible even ten years ago. When computers became a commodity, it was mostly used for textual documents. It took many years before it was also used for music. Music is different in fundamental ways from text in that it is not easily searchable without associated metadata. In the extreme case of a music collection without any associated metadata it can be very difficult to locate a song in that collection.

With a large music collection it is not uncommon to have multiple copies of songs or even entire albums. Without even knowing it, you may have copies of a certain album in two different formats, meaning that the size of the files, or even the quality, may be different.

Audio fingerprinting systems are in use today for a variety of different uses, such as automatically finding metadata for a given recording, intellectual property rights management or clustering of similar music to find music suggestions. When you listen to a song on your car stereo, a modern fingerprinting system could tell you the performing artist and the song title if it has the fingerprint of the song in its database. Broadcast monitoring could be made easier if there were fingerprinting systems in place to automatically transcribe a list of music being played on the radio or on TV. This thesis confines the area of study to finding such duplicates in a music collection.

### 1.2 Objectives

This thesis focuses on finding audibly similar music in a collection. The main objective is better performing fingerprint solutions that correctly identifies audibly similar music. Statistical and probabilistic combination methods are applied to audio fingerprints from different systems to increase recognition rates.

This thesis seeks to answer the following questions:

- Can combining multiple descriptors from separate authors produce a better fingerprint?
- Which methods of combining descriptors produce consistently better results?
- How do the descriptors tested rank in performance?

### 1.3 Approach

This thesis uses a music collection with 9536 music files. The music collection is a real-world example of what might be in someone's personal music collection. The music collection consists mainly of MP3 and WMA files in a ratio of approximately 8:1.

This thesis uses the `fdmf` and `libFooID` software, and in addition includes fingerprinting software from a masters thesis. `fdmf` and `libFooID` (in the form of `Foosic`) are both projects that can quickly be tested by downloading them from the Internet. An individual audio fingerprint, known as a descriptor, for a file can be compared to the same descriptor for another file. The result is a percentage that indicates the similarity between the two files. Outputs from multiple systems are combined using several approaches, e.g. the average and a naïve Bayes classifier, to give a combined score for how similar the two files being combined are.

The first 1000 matches are compared to a reference list of verified equivalent files, and the failure rate is used as a measure to compare performance characteristics of a given combination.

### 1.4 Results

This thesis has proven that combining descriptors from multiple sources is a viable way to create a better fingerprint. It finds that a weighted root-mean square (WRMS) of a given set of 8 descriptors gives the best results with the given music collection. When picking descriptors carefully, other combination methods can be used with results nearly as good as the WRMS.

Among the tested descriptors it is found that `libFooID` is the most reliable, followed closely by `fdmf_0`, two of the MFCC descriptors, `F1CC average`, `fdmf_2` and `fdmf_1`. This thesis also found that several of the descriptors cannot be used to reliably identify recordings, amongst them are song length, centroid, steepness, mean/square ratio and the rate of zero crossings.

### 1.5 Structure

The report is divided into the following chapters.

- Chapter 2 presents basic sound theory, examines previous work in the field, and presents the software solutions to be tested.
- Chapter 3 presents the methods used to test and combine the various solutions into a working system. It shows how fingerprints are compared, how the results are combined, and how the accuracy of a system is measured.



- Chapter 4 presents the results of this thesis, and discusses any discrepancies seen.
- Chapter 5 presents the conclusion to this thesis, and outlines improvements that could be made to the system.



# Chapter 2

## Pre-study

This chapter examines the basis for audio fingerprints, and looks at previous work in the area of music information retrieval (MIR). It then looks at the individual descriptors that will be combined into a single solution.

### 2.1 Sound theory

Sound is physical waves moving through a medium such as air with a frequency between approximately 20 Hz and 20 kHz. The physical waves can be described as sequential compression and decompression of adjacent molecules in the medium. Sound waves are mostly started by vibration in a physical object, setting the medium into motion.

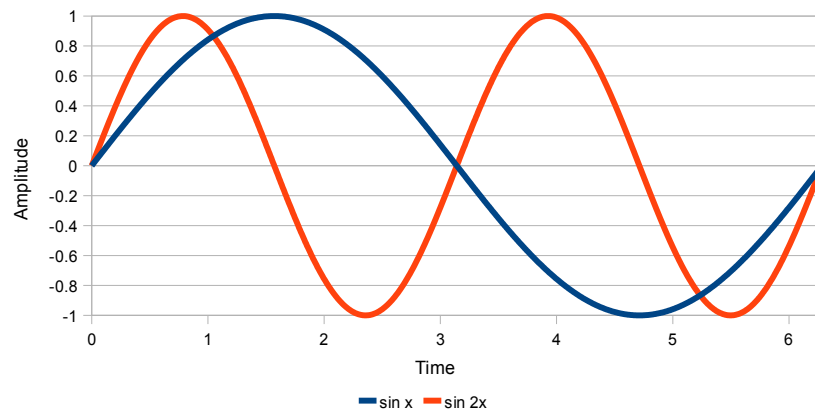
Interactions between waves grows complex very quickly. As can be seen in figure 2.2, the combination of three simple waves result in a wave that is very unpredictable. In the real world several hundred single waves combine to make the sound of a violin.

#### 2.1.1 Waves and interactions

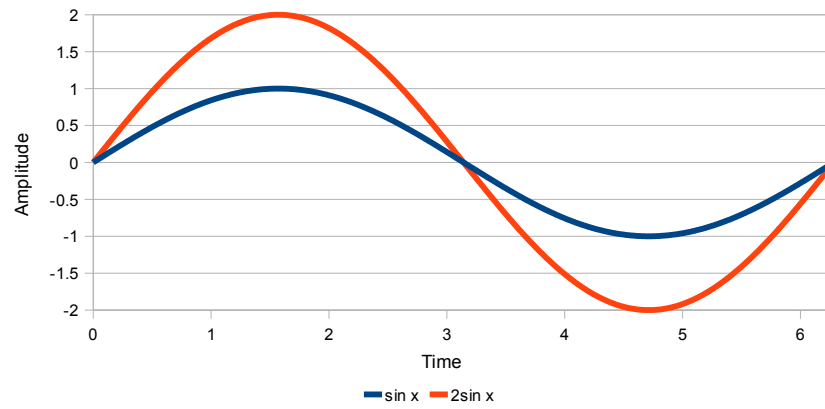
Pure sound waves are characterized by two main properties; the *frequency* and the *amplitude*. The frequency measures how often the wave repeats itself, and is measured in Hertz (Hz). 1 Hz indicates an event occurring once every second. Figure 2.1(a) shows two sine waves, one with a frequency of  $\frac{1}{2\pi}$ , the other with a frequency of  $\frac{1}{\pi}$ . Waves with a higher frequency has a shorter distance between two local maxima, a distance known as the *cycle* or *period* of a wave.

The amplitude of the wave is a measurement of how much of the medium the wave can displace during one period. When talking about amplitude, one usually means the peak amplitude, which is the distance from neutral to the maxima. The waves presented in figure 2.1(b) has peak amplitudes of 1 ( $\sin x$ ) and 2 ( $2\sin x$ ). A single sine wave will sound like a clear, steady tone.

Waves interact in very complex ways, which produce different sounds to the human ear. Waves can collide, and provide positive and negative feedback that either neutralizes the sound wave, or increases its amplitude. In extreme cases where two waves with the same frequency collide and maxima line up, the amplitude is added together for the two waves. This is known as the waves being



(a) Sine waves of different frequency



(b) Sine waves of different amplitude.

Figure 2.1: Features of waves illustrated.

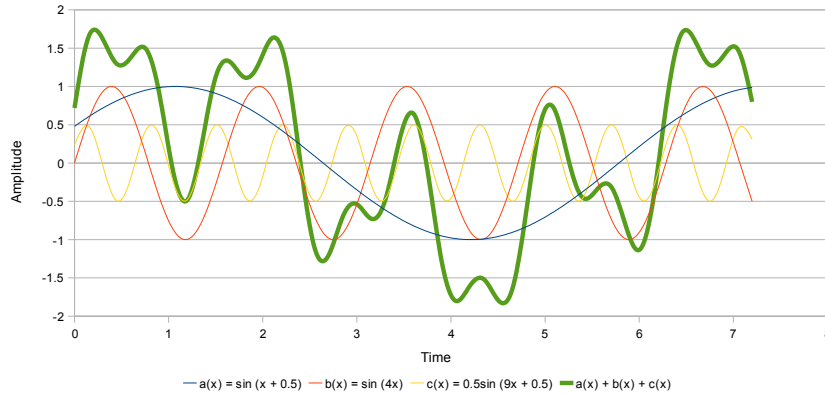


Figure 2.2: Combination of three sine waves results in complex wave.

*in phase*. When the waves are out of phase they cancel each other out. The combined wave is always the sum of the displacement of the individual waves at any given point.

Waves are defined by the wave equation, a second-order linear partial differential equation. It gives the propagation of waves with a given speed  $v$ , with  $\nabla^2$  being the Laplacian. For a detailed explanation of the wave equation, the reader is referred to an article in MathWorld [38].

$$\nabla^2 = \frac{1}{v^2} \frac{\partial^2 \psi}{\partial t^2} \quad (2.1)$$

### 2.1.2 Humans and sound

Sound is perceived by humans through the ears, where the sound waves hit the tympanic membrane<sup>1</sup>. The tympanic membrane transfers the sound waves' kinetic energy to the ossicles in the middle ear, which again transfers the energy to the cochlea. Inside the cochlea fluids are set in motion, and tiny hairs in the Organ of Corti register the movements and signals the brain. The brain then interprets these signals as sound. For more details, see [29].

The human auditory system does not perceive sound linearly, and the sound pressure level, the volume, is often quoted in decibel (dB). The sound pressure level is logarithmic, and measures the root-mean square change from the ambient pressure caused by a sound wave (equation 2.2). The reference sound pressure  $p_{ref}$  is  $20\mu Pa$ . Human perception of sound loudness roughly follows the sound pressure level on the decibel scale, meaning that a doubling of the sound pressure will be perceived as a constant increase, no matter what the previous sound pressure.

$$L_p = 10 \log_{10} \left( \frac{p_{rms}^2}{p_{ref}^2} \right) = 20 \log_{10} \left( \frac{p_{rms}}{p_{ref}} \right) \quad (2.2)$$

<sup>1</sup>Colloquially known as the *eardrum*.

Note	Frequency	n
G	392Hz	-2
G $\sharp$ / A4 $\flat$	415Hz	-1
A4	440Hz	0
A4 $\sharp$ / B $\flat$	466Hz	1
B	494Hz	2
C	523Hz	3
C $\sharp$ / D $\flat$	554Hz	4
D	587Hz	5
D $\sharp$ / E $\flat$	622Hz	6
E	659Hz	7
F	698Hz	8
F $\sharp$ / G $\flat$	740Hz	9
G	784Hz	10
G $\sharp$ / A5 $\flat$	831Hz	11
A5	880Hz	12

Table 2.1: The chromatic scale around A4, with frequencies and n-values (see equation 2.3).

Pitch is another facet of sound that the human auditory system perceives roughly logarithmically [10]. It is the perceived frequency of sound, and the perception can be changed by playing other frequencies simultaneously. Pitch is close to the fundamental frequency of a sound, yet the perceived pitch can change subtly when harmonics are introduced.

Harmonics and pitch are closely related. Harmonics are the effects produced when multiple sound waves, all with frequencies that are multiples of a common fundamental frequency, occur concurrently. The canonical example is the musical note A4, which in current Western music is defined to have a frequency of 440 Hz. The 440 Hz A4 note is considered the first harmonic for itself. The 880 Hz overtone is called the second harmonic, and corresponds to the musical note A5. 1320 Hz is the third harmonic (A6). Those familiar with music notation will notice that moving from A4 to A5 is the same as moving up an octave.

For Western music the range between two harmonics is usually divided into 12 half-steps, of which seven are assigned one-letter names [20]. When only using these seven denotations (A, B, C, D, E, F, G), you are using the diatonic scale. When also looking at the half-steps between some of these notes, you are using the chromatic scale. It is worth noting that notes in the diatonic scale is not evenly spaced. The frequency of a note is calculated as in equation 2.3, where  $n$  is the number of half-steps away from A4. The equation has been used to generate table 2.1.

$$f = 2^{n/12} \times 440 \quad (2.3)$$

### 2.1.3 Digital sound

For sound to be stored in a computer it must be digitized. Since computers operate using discrete numbers — contrasted to the continuous spectrum of sound — digitizing sound carries an inherent degradation. Sound is *sampled* at

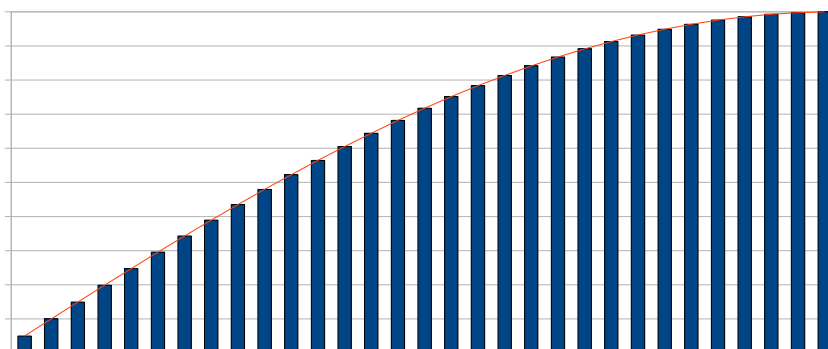


Figure 2.3: The line shows a continuous line, representing the continuous audio signal in the medium, and the vertical bars representing the sampled values.

regular intervals, and the computer stores the sound pressure. The sampling frequency is also measured in Hz, or samples per second, and specifies how often the samples are taken. Due to the Nyquist-Shannon sampling theorem<sup>2</sup>, the sampling frequency has to be at least twice the highest frequency of the signal you are trying to capture. Specifically;

If a function  $f(t)$  contains no frequencies higher than  $W$  cps<sup>3</sup>, it is completely determined by giving its ordinates at a series of points spaced  $\frac{1}{2W}$  seconds apart.[23]

Since human hearing is limited upwards at around 20kHz, most sound today is sampled at 44.1kHz or 48kHz using 16 bits to represent the sound pressure in each sample. Because of the high number of samples required, digitally stored sound has high storage requirements. For a single CD, which can contain up to 74 minutes of sound, a staggering 650 megabytes is required. This results in a bit rate of approximately 1411kbps; or around  $172 \frac{kB}{sec}$ .

There is an alternative way to store music using a discrete computer encoding. The best known is MIDI (Musical Instrument Digital Interface), where the individual notes and pauses are stored using a compact digital encoding. It can be thought of as a computer equivalent to sheet music, with instructions as to which instrument should be played at what time, using a given volume and pitch. A three minute song may therefore be encoded in as little as a few kilobytes of data, but cannot contain lyrics and can only use the instruments specified in the MIDI standard. Music encoded as MIDI is therefore perceptually very different to a human listener. Headway is being made in this kind of compact representation, using three-dimensional models of instruments and implementing physical laws to synthesize the sound [4]. This thesis is only concerned with the representation of sound as sampled audio.

<sup>2</sup>sometimes referred to as the Shannon sampling theorem, or simply the sampling theorem

<sup>3</sup>cycles per second

### 2.1.4 Compressing digital sound

Because of the high storage requirements for raw digital sound, several methods for compressing the sound has emerged. Traditional compression techniques, used for e.g. text documents, can also be used for compressing sound. The compression ratio gained from traditional compression techniques as applied to sound is inadequate for enabling storage of large amounts of digital sound.

Other techniques evolved as a response, with the two main branches being lossless and lossy compression techniques. With lossless compression, it is guaranteed that a file that has been compressed and then decompressed is bit for bit identical to the original. With lossy encoding, the compression algorithm can throw away non-significant data to decrease size. If the algorithm is programmed to discard data that is hard for humans to perceive, the loss of data can be acceptable to human listeners. A lossy algorithm works very badly with text, as letters, words or sentences may disappear completely from the text. Traditional text compression techniques are therefore lossless.

The most famous audio codec is MP3 (MPEG-1 Audio Layer 3), a lossy compression that was designed for the MPEG video encoding format to store audio [5]. Before it was chosen as an audio layer in the MPEG-1 standard, it was known as ASPEC (Adaptive Spectral Perceptual Entropy Coding) [6]. The most common MP3 file format has a bit rate of 128kbps, which is approximately 10 percent of the original file, assuming the original was from a CD. The MP3 format discards data or reduces precision in data that cannot be easily perceived by humans, based on psycho-acoustic models. Several other lossy formats, such as Ogg Vorbis, Windows Media Audio and AAC exist.

Lossless formats only look at efficient ways of storing the sound data, such as finding recurring patterns, and do not reduce precision in or discard data. These formats typically reduce the files size to 40 to 50 percent of the original file's size. Formats here include FLAC, Windows Media Audio Lossless and ATRAC.

## 2.2 Music Information Retrieval

The field of music information retrieval (MIR) is very broad, and encompasses a wide range of possible uses. The field is mainly concerned with extracting information from music contents, without relying on human-supplied metadata. Various research areas include:

**Speech recognition** Systems that automatically extract lyrics from a song could help index large amounts of audio recordings for search using standard information retrieval techniques. Problems with this approach include large amounts of background noise<sup>4</sup>, variations in speech patterns and the large amount of languages used in music.

**Automated transcription** Producing sheet music from any piece of music, suitable for an orchestra or soloist to play. The current models for understanding sound waves — the models that are used for extracting the individual notes for an individual instrument — are not adequate to understand the complex interactions between sound waves produced by multiple

---

<sup>4</sup>Noise in the context of extracting the spoken word.



instruments playing in unison. It can be hard for untrained humans to differentiate between different types of instruments in a sound recording.

**Query by humming (QBH)** Allows searching through a database of music by humming or whistling a song to the computer. The computer then identifies the fundamental frequencies of the humming, and looks them up in the indexed music. The music indexed by current QBH systems has to be in structured form and not as sampled audio.

**Audio fingerprinting** Automatically finds the identifying characteristics in a piece of music. Some audio fingerprint systems allow nearest neighbor searches for music, others concentrate solely on identifying a specific recording of a song. Audio fingerprint systems can e.g. be used to identify the currently playing song on the radio or to find duplicates in a personal music collection.

**Feature extraction** Used for identifying the features of a musical track, such as e.g. beats per minute. The Echo Nest has a free API that can be used to extract various features from audio, and has an example that allows a user to put together a seamless mix of several songs automatically [7].

**Genre classification** A wide variety of music genres is represented in present-day music. A lot of work has gone into producing a system that can automatically classify music into genres, and that can compare similarities between different genres of music.

### 2.2.1 Properties of music

Describing music is a very hard task, and attempts to reduce the complexity of this task is usually done by splitting the music into different aspects. Downie presents seven facets that can be used to describe music, listed here for convenience [9]:

**Pitch facet** The perceived frequency of a sound. Hard to determine from a recording.

**Temporal facet** Duration of a musical event, i.e. the length of holding a given note.

**Harmonic facet** Several instruments playing at once with pitch frequencies that are multiples of a given fundamental frequency for that instrument.

**Timbral facet** The property that differentiates the sound from two different instruments, i.e. how a clarinet and a trumpet is different.

**Editorial facet** Instructions given to the performers that may change between recordings.

**Textual facet** Lyrics of a song.

**Bibliographic facet** Any metadata, such as title, artist and composer. Usually textual.

Most MIR software only concern themselves with a subset of these facets. The techniques used in this thesis do not at all rely on the editorial, textual or bibliographic facets.

String	MD5
A MD5 hash value changes quickly.	eb0583fed3abca7103419cfce517046e
A MD5 hash changes quickly.	0e04f19f5c76bce72cf68e0c007346a0
A MD5 hash value changes drastically.	6cabee6499ab39254ae765dde785c7a9

Table 2.2: MD5 hashes of three similar text strings.

## 2.3 Audio fingerprinting

Audio fingerprinting can be loosely defined as a process that produces a small data sequence that uniquely can identify a specific piece of sound. Contrasted to a human fingerprint, an audio fingerprint can be thought of as a severely one-way compressed copy of the sound. Some features are desirable in fingerprints, and these features are highly interconnected. Some features are:

**Compactness** The fingerprints should be small, so they can be easily stored and searched through.

**Spatial placement** Fingerprints of songs that are related should be less different than fingerprints of songs that are not related using a given distance measure, such as Euclidean distance. This allows for the fingerprints to be used to find similar music, not only identical music, or clustering music into genres.

**Robustness** A small or medium amount of noise should not affect the generated fingerprint significantly. This also applies to transformations that might occur when using various different compression algorithms. If the fingerprint is not robust, this results in false negatives.

**Reliability** Songs should preferably never be mislabeled as other songs, an event known as a false positive.

**Granularity** How much audio is needed to construct a fingerprint. For applications where you have a recording of an entire song, the entire song can be used. In other cases you might only have a few seconds of audio.

**Destructiveness** An audio fingerprint does not have to be reconstructable into the original, or perceptually similar, audio.

Audio fingerprints are sometimes mistakenly compared to hashes, because hashes by definition produce a small fingerprint for a large amount of data. For example, a very common hash function is MD5, which can be used to verify the integrity of a downloaded file. Hash functions meant for use in cryptography — such as MD5 — produce completely different outputs if even a single bit is changed in the input value. As can be seen in table 2.2, a small change in the text string, changes the MD5 hash completely. This is a property that is undesirable in audio fingerprinting, as it has to be resistant to small changes in the output. Ideally, a small change in the audio will only result in a correspondingly small change in the audio fingerprint.

		Actual condition	
		Same song	Different song
Fingerprint indicates	Same song	True positive	False positive
	Different song	False negative	True negative

Table 2.3: Failure modes.

### 2.3.1 Use cases

Audio fingerprints can be used in a wide variety of situations, both by consumers and corporations. Consumers may want to figure out if they have a song in their music library, or if they have to purchase it. They may want to search through their music library for songs that are similar to the song they are listening to, or eliminate duplicates in their library. Consumers will also be interested in identifying music that is played on the radio, and there exists several solutions for cell phones. New SonyEricsson cell phones are delivered with a software called TrackID, which allows you to record a few seconds of a song using your cell phone, and send the sample fingerprint to an online service. SonyEricsson has partnered with Gracenote for the fingerprint identification. According to their website, Gracenote has a database of over 80 million tracks [12].

Corporations have different uses. The music industry will be interested in verifying that music shared on file-sharing networks are not copyrighted, and can use audio fingerprints to test against their entire music library. Google wants to choose advertisements based on the music it can identify from background noise in the room you are sitting in [1].

### 2.3.2 Misclassifications

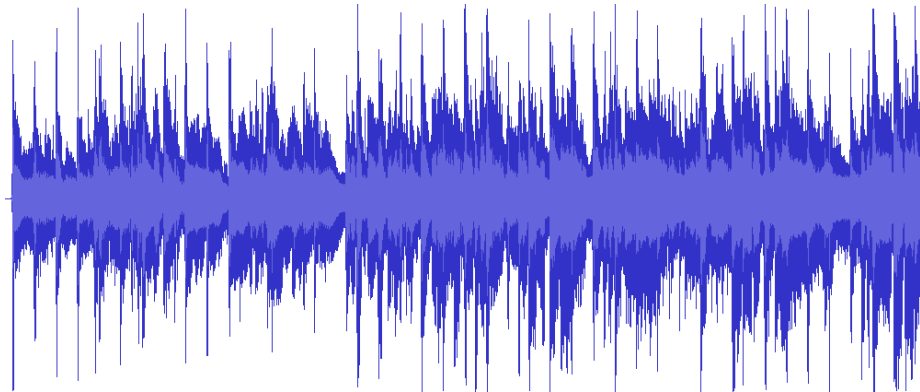
When comparing two fingerprints from two tracks, there are four possible outcomes. Note that equivalence is context-sensitive.

1. The two fingerprints are from the same track, and the comparison indicates that the tracks are equivalent. This is known as a *true positive*.
2. The two fingerprints are from different tracks, and the comparison indicates that the tracks are not equivalent. This is known as a *true negative*.
3. The two fingerprints are from different tracks, but the comparison indicates that the tracks are equivalent. This is known as a *false positive*<sup>5</sup>.
4. The two fingerprints are from the tracks song, but the comparison indicates that the tracks are not equivalent. This is known as a *false negative*<sup>6</sup>.

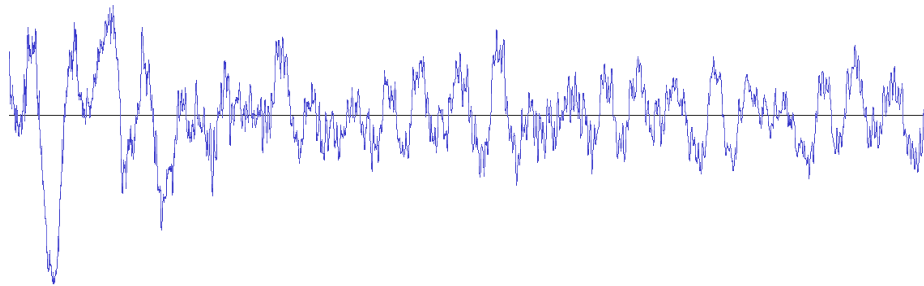
Table 2.3 shows these failure modes in relation to each other. False negatives and false positives are the failure modes of a classification. False positives and false negatives can be said to be errors of equal magnitude when looking at audio fingerprints. In the domain of spam filtering a false positive (classifying a non-spam mail as spam) is much worse than a false negative (classifying a spam mail as non-spam) [13].

<sup>5</sup>False positives are called *type I* or  $\alpha$  errors by statisticians.

<sup>6</sup>False negatives are called *type II* or  $\beta$  errors by statisticians.



(a) First 30 seconds



(b) 120ms extract

Figure 2.4: Time-based representation of the first 30 seconds and a 120ms extract of “deLillos - Fullstendig oppslukt av frykt”. Y-axis is power scale, X-axis is time.

The error analysis in this thesis concerns itself mainly with false positives for two main reasons:

1. Since the output is ranked according to probability, false positives will be very easy to spot in the top results.
2. Because this thesis operates on a real life music collection, the false negatives will be very hard to find in such a large collection.

### 2.3.3 Analyzing audio

There are multiple ways in which an audio fingerprint can be constructed. Three methods are outlined here: time-based representation, frequency spectrum-based representation and wavelets.

When looking at audio using a normal editor, one usually looks at the *time-based representation* in the shape of a waveform. In such a representation the power is represented on the Y-axis, while the time runs along the X-axis. This representation allows you to see features such as rhythm and beats per minute, BPM, by counting peaks. An example can be seen in figure 2.4.

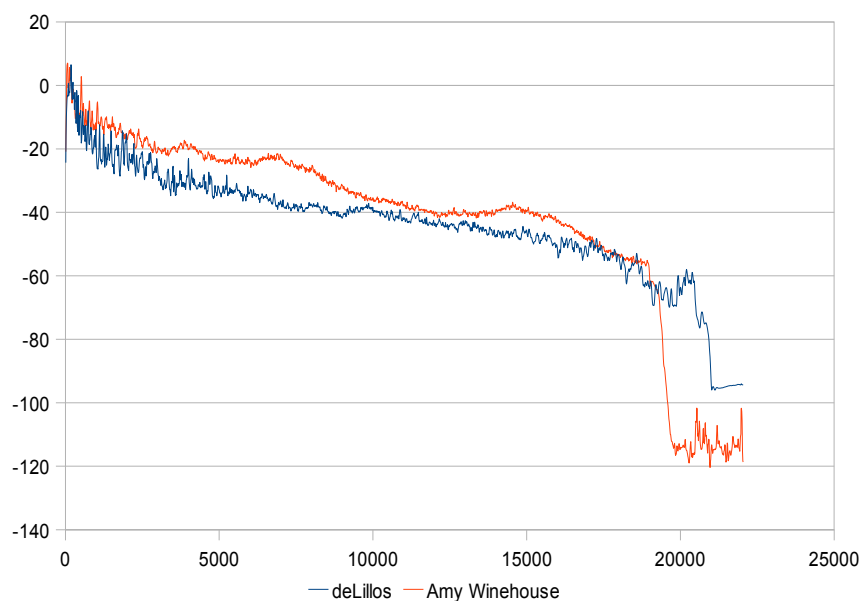


Figure 2.5: The frequency spectrum of the first 23.8 seconds of “deLillos - Fullstendig oppslukt av frykt” vs. “Amy Winehouse - Rehab” using a Hann window function. X-axis is frequency in Hz, Y-axis is power in dB.

The *frequency spectrum-based representation* is usually found using a short-time Fourier transform (STFT), and shows which frequencies are used at what power levels. The STFT is applied to short sequences of sound, depending on the size of the features one attempts to extract. The frequency spectrum can then be reduced using various algorithms to be used as a descriptor. One common algorithm is to establish frequency ranges called bins, and sum the power levels for frequencies in that bin.

The *wavelet representation* uses a wavelet transform to create a range of wavelets that allows examination of frequency components on a suitable scale [14]. Wavelets are also easier to compute than the Fourier transform, with an asymptotic run time of  $O(n)$ , compared to the fast Fourier transform’s  $O(n \log n)$ .

How? Construct a function, shift it by some amount, and change its scale. Apply that structure in approximating a signal. Now repeat the procedure. Take that basic structure, shift it, and scale it again. Apply it to the same signal to get a new approximation. And so on. It turns out that this sort of scale analysis is less sensitive to noise because it measures the average fluctuations of the signal at different scales.[14]

### 2.3.4 Intellectual property

The ownership of music is considered under the broad term *intellectual property* (IP). Buying a CD does not mean that you own the music, but that you own a license to play that music. IP law covers the legal aspects and varies from country to country. Most countries' copyright laws include a term known as *fair use*, which specifies under which conditions one can use a copyrighted work without paying royalty fees.

Audio fingerprints are usually destructive, and is therefore not covered by copyright law, as the result can not be used to infringe on the original work. Non-destructive fingerprints may be covered, and could probably not be shared freely. However, this thesis focuses solely on destructive audio fingerprints.

Several of the audio codecs mentioned and used in this thesis are patented. Norway does not recognize software patents, and such issues are therefore not considered [24].

## 2.4 State of the art

The field of music information retrieval progresses quickly, and there exists several solutions that are use different descriptors than the ones used in this thesis.

Several solutions have been based directly on the work of Ke, Hoiem and Sukthankar [18], which concentrates on applying computer vision techniques on the spectrogram of audio files. By using a STFT, the spectrogram contains the power of logarithmically spaced frequency bands, much like the process used in Gunderson [15]. Using machine learning they identify a set of filters that perform well. The implementation is freely available under the GPL, and this work has formed the basis for at least two systems. libFingerprint from last.fm is an adaptation of the work presented by Ke *et. al.*, with improvements mainly in the speed of lookups and featuring a cleaner API [17].

Google released a set of papers that details Waveprint, which do not rely on machine learning but rather on wavelets [8, 2]. Wavelets are introduced easily in [14], and allows examination of both the small and large features of a signal, and is considered an alternative to Fourier analysis. None of the software solutions examined in this thesis use wavelets.

Another type of audio fingerprinting models the audio signal by using a sinusoid generated by parameters such as amplitude, phase and frequency, and takes into account the residual noise. The sinusoidal models are extracted using Fourier analysis. Some sinusoidal peaks are selected, mainly those that have high amplitudes and conforms to the sinusoidal model. Betser *et. al.* details the selection process, and finds an increased recall compared to Haitsma [16]. The system can recognize segments as short as 1 second [28].

Some papers state the goal of finding audio recordings that descend from the same musical work, such as Miotto and Larsen [22, 19]. Larsen demonstrates that systems such as fdmf and libFooID cannot be used to identify a more abstract recording reliably. Miotto builds a statistical model of a performance that predicts possible alternative performances. The original performance is segmented and various audio features are extracted from the segments. Hidden Markov models are applied to enable identification.

MIRtoolbox — a package built for Matlab — can be used as an introduction to MIR [35]. It allows easy processing of many common MIR methods, such as finding the mean/square ratio, spectral measures such as MFCC and cepstrum, skewness and centroid. It also has an extensive filter bank, and can quickly be used to test the effects of filters. It can also be used to analyze the pitch of a song, resulting in a range of key candidates, or to find the rhythm in a song.

## 2.5 Software

For testing our hypothesis that several descriptors can be used to produce a more accurate fingerprint, a set of existing software is used. Throughout this thesis these are referred to as the individual descriptors. The choice of descriptors were made based on the availability of the software.

### 2.5.1 fdmf

The fdmf (find duplicate music files) package is a small software package designed explicitly to find duplicate music files in a collection. It is written mostly in GNU C, with parts written in Perl. fdmf consists of a program that fingerprints all the audio files, and stores it in an internal database. This database is then used by a second program that matches all fingerprints against all other fingerprints in the database, and prints out the results.

#### 2.5.1.1 Fingerprints

fdmf generates a combined fingerprint of 768 bits, that is in fact three separate descriptors of 256 bits. The first 256 fingerprint-bits are a summary of the energy spectrum of the audio file. The following 256 bits are a summary of the ratio spectrum, a mathematical equivalent to the power spectrum [32]. The final 256 bits are a summary of the twist spectrum.

Rosenfeld details the generation of the descriptors [31, 30]. The descriptors are generated from 250 ms segments of a mono-channel version of the original audio. After applying the STFT, four band energies are calculated for each segment. Equations 2.4 - 2.6 show how these band energies are converted into chunk metrics [31]:

$$cm_1 = be_1 + be_2 + be_3 + be_4 \quad (2.4)$$

$$cm_2 = (be_3 + be_4) / (be_1 + be_2) \quad (2.5)$$

$$cm_3 = (be_1 + be_3) / (be_2 + be_4) \quad (2.6)$$

A second STFT is applied to the chunk metrics, to limit the time misalignment that might occur in different recordings of files. 256 frequencies are chosen uniformly across the resulting spectra. These values are used to quantize the three chunk metrics into a 768 bit fingerprint.

#### 2.5.1.2 Comparing fingerprints

When comparing two audio fingerprints, fdmf looks at the three descriptors separately. For each of the three descriptors it calculates the Hamming distance

between the parts. If the Hamming distance is less than the threshold for that descriptor, the score for those two audio files matching is increased. If all three descriptors has a difference smaller than their threshold, the files are considered to be equivalent. The thresholds are specified at startup as a vector, and if not given, a default threshold vector is used.

### 2.5.2 libFooID

libFooID has not been documented in any research papers, but its website details the fingerprinting process [27]. The audio is normalized, and combined to a mono recording. Starting silence is skipped, and the first 100 seconds of the audio is processed. The audio is resampled to 8000Hz, and only 90 seconds of the audio is used.

A Hann-windowed discrete Fourier transform is applied to 8192 sample blocks, resulting in 87 frequency spectra frames. The frequency spectrum is partitioned into the Bark scale using a modified version of Traunmüller's formula [34]:

$$z = \left\lceil \frac{26.81f}{1960 + f} \right\rceil - 0.53 \quad (2.7)$$

The first Bark band is ignored, and the final band is enlarged, which leaves 16 Bark bands.

$$z_2 = \begin{cases} z + 0.15 \times (2 - z) & z < 2 \\ z & 2 \leq z \leq 20.1 \\ z + 0.22 \times (z - 20.1) & z > 20.1 \end{cases} \quad (2.8)$$

A least square regression line fit is used on the spectral data, and the correlation coefficient and the dominant spectral line per frame is stored using 2 and 6 bits, respectively. The fingerprint is stored using 424 bytes, and includes some metadata, such as the length of the song.

#### 2.5.2.1 Comparing fingerprints

Comparing fingerprints is done by using a simple Hamming distance for the dominant spectral lines, and uses a quadratically weighted Hamming distance for the correlation coefficient.

### 2.5.3 Gunderson's descriptors

A set of descriptors as presented in Gunderson's master thesis was also used in this thesis [15]. The thesis implements a variety of basic descriptors, such as the spectral centroid, play time, rate of zero crossings (with and without Schmitt triggering), steepness, and the mean/square ratio. Gunderson also implements mel frequency cepstral coefficients (MFCC) and floor-1 cepstral coefficients (F1CC), that are included in this thesis.

The descriptors implemented by Gunderson are briefly explained below, and are explained in depth in this thesis itself. Common to all of the descriptors implemented, with the exception of song length, is that they only look at the first 30 seconds of a song.



### 2.5.3.1 Centroid

A spectral measure that finds the “center of mass” for the audio signal. The audio signal is transformed to the frequency domain by a discrete Fourier transform.

### 2.5.3.2 Song length

The song length was implemented as the length of the track in seconds, but could also be the number of samples if the sample rate is equal in all songs.

When using lossy compression methods with variable bit rates, the entire track may have to be decoded to be able to determine the song length, if this is not given in any of the metadata. The very common MP3 file format, when encoded using a variable bit rate, does not store the song length automatically as metadata, and as such the song length can only be estimated until the entire file has been decoded.

### 2.5.3.3 Mean/square ratio

The mean/square ratio measures the average distance from zero, calculated after the audio signal has been normalized. Equation 2.9 is from [15].

$$r = \frac{\text{avg}(|x|)}{RMS} = \frac{\frac{\sum |x_i|}{N}}{\sqrt{\frac{\sum x_i^2}{N}}} = \frac{\sum |x_i|}{\sqrt{N \sum x_i^2}} \quad (2.9)$$

### 2.5.3.4 Steepness

Steepness is a frequency measure that is less affected by noise than the zero crossing rate. It is computed by taking the mean of the numerical derivative of the audio signal, as in equation 2.10, again from [15].

$$s = \frac{1}{N-1} \sum_{i=1}^{N-1} |x_i - x_{i-1}| \quad (2.10)$$

### 2.5.3.5 Rate of zero crossings (with and without Schmitt triggering)

The rate of zero crossings counts the number of times the signal changes sign divided by the number of samples. Schmitt triggering defines a guard band around the x-axis that needs to be crossed for a zero crossing to be counted. Schmitt triggering increases the resilience to noise.

### 2.5.3.6 Mel frequency cepstral coefficients

Mel frequency cepstral coefficients (MFCC) is based on the mel scale that attempts to mimic the human auditory system’s perception of varying pitches at constant loudness. The frequency is transformed to the mel scale using the following equation:

$$B(x) = 1127 \ln \left( 1 + \frac{x}{700} \right) \quad (2.11)$$

Related frequencies are linearly grouped into spectral bands on the mel scale. MFCC outputs several sets of individual descriptors, of which the important

ones are `mfcc_avg` and `mfcc_delta_avg`. `mfcc_avg` is titled the *first zero moment*, and is calculated separately for the 32 coefficients:

$$\hat{\mu}_m = \frac{1}{N} \sum_{i=0}^{N-1} c_{m,i} \quad (2.12)$$

`mfcc_delta_avg` is calculated as the average (as in equation 2.12) of the first derivative of the coefficients.

$$\Delta c_{m,i} = |c_{m,i+1} - c_{m,i}| \quad (2.13)$$

### 2.5.3.7 Floor-1 cepstral coefficients

Floor-1 cepstral coefficients (F1CC) is the descriptor tested by Gunderson to determine if audio masking techniques used in lossy compression could be reliably used when creating fingerprints. Floor-1 is based on the Ogg Vorbis psychoacoustical model, which models how certain weak tones are hidden by louder nearby tones [11, 36].

The calculation of fingerprints is otherwise equivalent to MFCC, and the descriptors named `f1cc_avg` and `f1cc_delta_avg` are considered by Gunderson to be the most accurate.

### 2.5.3.8 Comparing fingerprints

Gunderson's descriptors outputs a vector of floating point values. To compare two fingerprints to each other Gunderson experiments with using both the Euclidean distance and Mahalanobis distance, the latter of which is scale-invariant. For this thesis, only the Euclidean distance has been used. For details about how the distance measure is converted to a probability of match, see 3.3.3.1.

# Chapter 3

## Method

This chapter describes the work that was done to combine the three solutions chosen for this thesis. It outlines the structure of the value-added software that was written for this thesis.

### 3.1 Music collection

The music collection used consists of two merged personal collections, with contents in WMA, MP3 and Ogg Vorbis formats, using high bit rates. The collection contains in total 9536 music files, of which an estimated 2200 combinations of two files (see 3.1.3) are considered duplicates. In total, the files consume 57GB of disk space. It would take 26 days and nights to play at normal speed. The collection is strongly biased to contain MP3 files over WMA files, with 8429 MP3 files, 1027 WMA files and 80 Ogg Vorbis files.

The sample collection is meant to mimic an average personal collection, and no attempts were made to remove duplicate songs or albums when combining the two separate music collections. A listing of the albums in the collection can be found in appendix C. The collection is stored in a folder hierarchy with the naming scheme <artist>/<album>/<song>, but this naming scheme is not strictly enforced throughout the collection.

An alternative music collection could be created by using perfectly preserved originals, encoded in a lossless format, and re-encoding these into several popular formats. This approach was not used for two main reasons:

1. No access to a large lossless music collection.
2. In the real world, a collection does not usually contain several copies of the same song in different formats. The results would be skewed towards how resistant an algorithm is at ignoring encoding artifacts.

#### 3.1.1 Equivalent music

The question of what defines two songs as being equivalent has never been settled, and is usually considered to be a matter of personal preference, context and usage scenario. To examine this more closely, a terminology is needed to

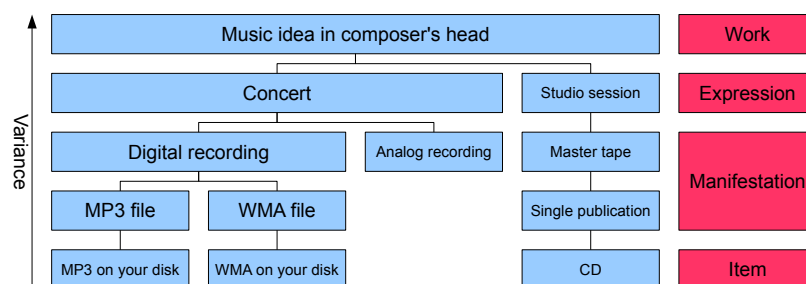


Figure 3.1: Simplistic illustration of the FRBR terminology as applied to the lifespan of a fictional song.

describe music. IFLA presents the FRBR<sup>1</sup> model, which can be used as such a terminology [25, 19]. It defines four entities that can be used to describe music. Originally it was intended to be used for more general bibliographic records to describe relationships between any items that can be stored in a library.

Behind every piece of music is the intellectual idea of the music, as imagined by the composer. This is the essence of that song, and is referred to as the *work*. The work can be expressed in a variety of ways, i.e. the composer might write sheet music for the song, or performing the song using an instrument. These are both considered *expressions* of that *work*. If someone uses a microphone to record the composer performing the song, the recording constitutes a *manifestation* of the expression.

The entity defined as *manifestation* encompasses a wide range of materials, including manuscripts, books, periodicals, maps, posters, sound recordings, films, video recordings, CD-ROMs, multimedia kits, etc. As an entity, *manifestation* represents all the physical objects that bear the same characteristics, in respect to both intellectual content and physical form. [25]

This means that our composer might write down the sheet music for the song, but if the sheet music is printed into many copies, those copies as a collection is the manifestation. A single copy is known as an *item*. The model also allows for works that are derivatives of each other, or when the changes are smaller, as different expressions of the same *work*.

In music a remix of a song can be said to be a derived *expression* of the original expression. In much the same way, a live performance of a song is a different *expression* than the original studio recording of a song. For the purposes of this thesis, live performances and remixes are considered as different enough to be nonequivalent with the original recording.

In several genres of music it is common to issue remastered versions of a song. When remastering, a sound technician has gone back to the original recording and removed noise or otherwise enhanced the audio. A remastering of a song is

<sup>1</sup>Functional Requirements for Bibliographic Records

considered a different *manifestation* of a song, and the audible difference from the original is sometimes hard for a human to detect. In this thesis such a remastered version is considered equivalent with the original. When digitally storing music, different encodings can be used. These encodings usually introduce various artifacts into the audio. Different encodings are also considered different *manifestations* of a song, and are considered to be equivalent in this thesis.

In figure 3.1 the FRBR hierarchy is illustrated with an example. The song is thought of by the composer, who plays it at a concert. At the concert it is recorded in both digital and analog format. For some reason, the analog recording gets ignored, while the digital recording is encoded into two formats (MP3 and WMA) and is distributed. The composer also records the song in a recording studio, resulting in a master tape, which subsequently is produced into a CD, which you can purchase copies of in your local record store.

When determining equivalence in this thesis, where the songs are not related in any of the ways mentioned above, the deciding factor for determining equivalence has been whether or not a human would say that the songs sound the same after having listened to the song once.

### 3.1.2 What are duplicate pairs?

For this thesis we define a *duplicate pair* as two tracks that cannot be easily distinguished from each other by a human listener. This definition might vary from study to study, and our definition is dependent on that this thesis chose a real-life music collection. In FRBR terminology, this will usually mean that the tracks are the same or derived manifestations. Yet in some cases, they might not be related manifestations, but still the same expression.

Because some of the descriptors used do not analyze more than 30 seconds of a track, a decision was made that two tracks will also be considered equivalent if they are audibly equivalent to a human for the first 30 seconds, as long as their total length is approximately equal.

This definition of a duplicate pair does not consider the fact that three or four tracks may be indistinguishable from each other, and that such clusters will in fact significantly impact the numbers of duplicate pairs. Assuming a cluster of four tracks that are related to each other, a total of 6 duplicate pairs can be found, as illustrated in figure 3.2. The number of duplicate pairs found for a cluster of  $n$  equivalent tracks is given by  $\binom{n}{2}$ , a number that increases quickly for large clusters.

### 3.1.3 Estimating duplicate pair count

A very simple process was used to estimate the number of duplicate pairs in the collection. The files in the collection have names that roughly corresponds to the title of the song. The file name is split into individual words, ignoring everything but letters and punctuation that naturally occurs within a word. Words shorter than three letters are ignored.

The tokenized version of two file names are then compared to each other, and the number of tokens found in both file names is counted. This value is then divided by the number of tokens in the longest file name, and if the result is

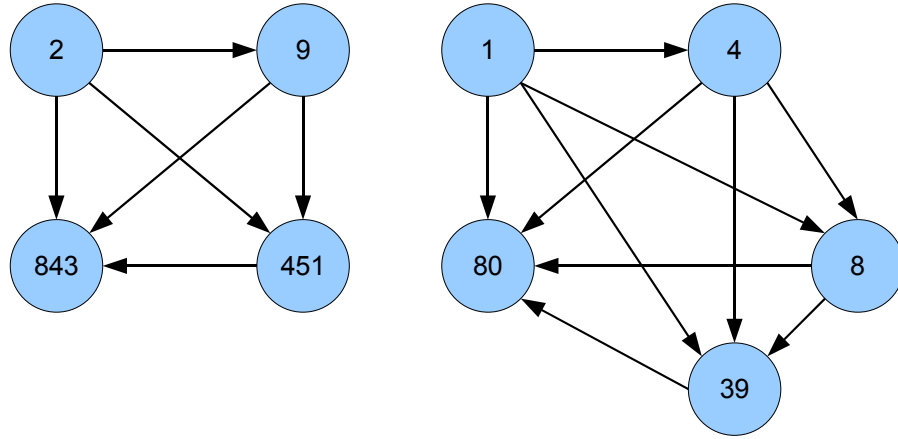


Figure 3.2: The edges represent a single duplicate pair, while the circles represent a track. The numbers indicate the identifier of the song, and note that the arrows always direct to a larger identifier. Left: the 6 combinations found when 4 identical songs are found. Right: the 10 combinations found when 5 identical songs are found.

---

**Algorithm 1** Pseudo-code for checking if two file names are close enough to be considered a potential duplicate pair.

---

```

1  if artistname1 <> artistname2:
2      return false;
3  tokens1 = tokenize(filename1)
4  tokens2 = tokenize(filename2)
5  all_tokens = merge(tokens1, tokens2)
6  common_tokens_count = 0
7  for token in all_tokens:
8      if token in tokens1 and token in tokens2:
9          increase(common_tokens_count)
10 tokens_count = max(count(tokens1), count(tokens2))
11 if tokens_count < 1:
12     return false;
13 ratio = common_tokens_count / tokens_count
14 return ratio > 0.6

```

---

greater than 0.6, and the files are from the same artist (indicated by the folder they are located in), the two file names are considered a match.

The number gained from this method is likely to be grossly over-dimensional because songs might have the same file name even though they are not audibly similar. For example, a significant number of music files in the collection contains live versions of songs which are not counted as duplicates of a studio recording. This error is compounding, as it will increase the count greatly when there are clusters of matching songs. This method of estimation resulted in 6991 combinations of two files listed as duplicate pairs.

It is also very likely that the comparison of just file names will produce a number of false positives. To estimate the amount of false positive, 500000 random combinations of file names were sampled. Using the method above produced 129 combinations that it believed to be the same song. When verifying these 129 combinations by hand, only 47 were likely to be equivalent based on the file name. We can therefore conclude that the actual number of hits is less than 35 percent of the number of hits, ergo around 2400 combinations.

When sampling all the combinations were the file names match and checking for the word “live” in the two file paths in the combination, around 10 percent of the combinations turned out to contain one song from a live recording and one from a studio recording. After this very simple check, we are left with an estimate of 2200 duplicate pairs.

### 3.1.4 Verifying duplicate pairs

Since the test is based around a non-controlled, real life collection of music, knowing in advance how many and which combinations should be considered duplicate pairs is not trivial. Verifying the identified duplicate pairs therefore becomes a highly labor-intensive manual process. However, since we know that we have approximately 2200 duplicate pairs (see 3.1.3), the amount of work is somewhat reduced.

When the various averages had been computed, an administration interface was set up to verify each duplicate pair. To aid in the process, the same code that was used in 3.1.3 to find duplicate pairs was used to highlight combinations that based on their file names could be actual duplicates. The first<sup>2</sup> few thousand potential duplicate pairs were then read through, and a human indicated to the system which were actual duplicate pairs.

This process might be error prone, and surely some of the duplicate pairs verified by a human will be wrong. The administration interface showed potential duplicate pairs with high rankings from both separate descriptors and the combined solutions, so the duplicate pairs that were verified is concentrated heavily in the top ranked results for each combined solution.

An assumption was made that perceived equality when talking about music is transitive. That is, given that song A is a duplicate of song B, and song B is a duplicate of song C, we can say that A is a duplicate of song C. This assumption might not hold in all cases, where subtle changes may make song C more different from A than B was. For a more concrete, yet fictional example, take three different recordings of the Billie Holiday song “Georgia On My Mind”, titled “Take 1”, “Take 2” and “Take 3”. The “Take 1” and “Take 2” recordings

---

<sup>2</sup>as ranked by the individual descriptors

might be similar enough that they can be listed as duplicate pairs, and likewise for “Take 2” and “Take 3”, while the “Take 1” and “Take 3” might be different enough for them not to be considered a duplicate pair. Please note that the example is fictional, and that no examples of non-transitive duplicates are known to exist in our collection.

In total 2057 duplicate pairs were verified.

## 3.2 Descriptor architecture

The combined solution should be able to quickly find equivalent audio files in a large collection. Several aspects have to be taken into consideration when designing the solution.

**Decoding** The most time-consuming process when creating the fingerprints is decoding the audio, and as such, the solution is designed to only decode the audio once and keep the audio data in memory. Then all the fingerprinting solutions can fetch the data from main memory. Once the fingerprints have been generated, the decoded audio can be discarded.

**Output** The systems should all be able to compare two fingerprints and produce a probability that the files are equivalent. This makes it possible to average the probabilities from multiple engines.

## 3.3 Software

Three different software packages were combined to form a single software system, that in one pass generates fingerprints for a set of audio files. The combined solution needs every system to output a probability of two songs being equivalent, and some of the systems had to be modified to output such a probability.

### 3.3.1 fdmf

By examination of the source code, fdmf in its original version returns only a list of the audio files that has a Hamming distance less than the threshold vector. The Hamming distance counts the number of positions in which two signals differ. We used the Hamming distance  $d_i$  and the threshold value  $t_i$  for each fingerprint part  $i$  to generate a probability  $P$  according to equation 3.1.

$$P_i(d_i, t_i) = 1 - \frac{d_i}{t_i} \quad (3.1)$$

The probability measure thus has a resolution of  $t_i$ , and ensures that the entire range will be used. The threshold values are used since a threshold of 255 would mean that the fingerprints are complete inverses of each other. fdmf’s default threshold values were used.



Descriptor	$t_i$
Centroid	50
Song length	5
Mean/square ratio	0.005
Steepness	25.0
Rate of zero crossings	2000
Rate of zero crossings (w/ Schmitt triggering)	1000
MFCC average	40
MFCC delta average	3.0
F1CC average	20
F1CC delta average	5

Table 3.1: Thresholds for various descriptors by Gunderson.

### 3.3.2 libFooID

Software was written around libFooID to create fingerprints from PCM audio data, since such software was not available for Linux with source code. Source code for our software can be found in appendix B. The software stores the fingerprints in the database table tblFingerprint, as described in section 3.4.1.

The software supplied with libFooID to compare fingerprints was adapted to automatically compare all the fingerprints in the database against each other, in much the same way as with fdmf.

### 3.3.3 Gunderson’s descriptors

The fingerprinting software was used without modifications to generate fingerprints for all of the audio files. A separate parser was then written to read the files produced into our fingerprint table (see 3.4.1). A script compares all of these fingerprints to each other, generating a probability that the two files match. The probability is calculated as in equation 3.2, where  $d(f_1, f_2)$  is the Euclidean distance between the fingerprints  $f_1$  and  $f_2$ , and  $t_i$  is the minimum Euclidean distance required for that descriptor according to table 3.1.

$$P_i(f_1, f_2) = 1 - \frac{d(f_1, f_2)}{t_i} \quad (3.2)$$

#### 3.3.3.1 Determining thresholds

The thresholds presented in table 3.1 was determined by pulling a random sample of 10 000 Euclidean distances for a given descriptor, and then sort the distances ascendingly. Due to the enormous amount of combinations we are only interested in Euclidean distances that are significantly smaller than the rest. With 9536 songs in our database, there are  $\binom{9536}{2} = \frac{9536 \times 9535}{2} = 45462880$  possible ways to combine two files, and we are only interested in identifying the duplicate pairs in our collection, which we have estimated to be around 2200. We are therefore only interested in 0.005 percent of the combinations presented.

Figure 3.3 shows a plot of the samples taken from the mfcc\_avg descriptor. By looking at similar plots for the different descriptors, we determined appropri-

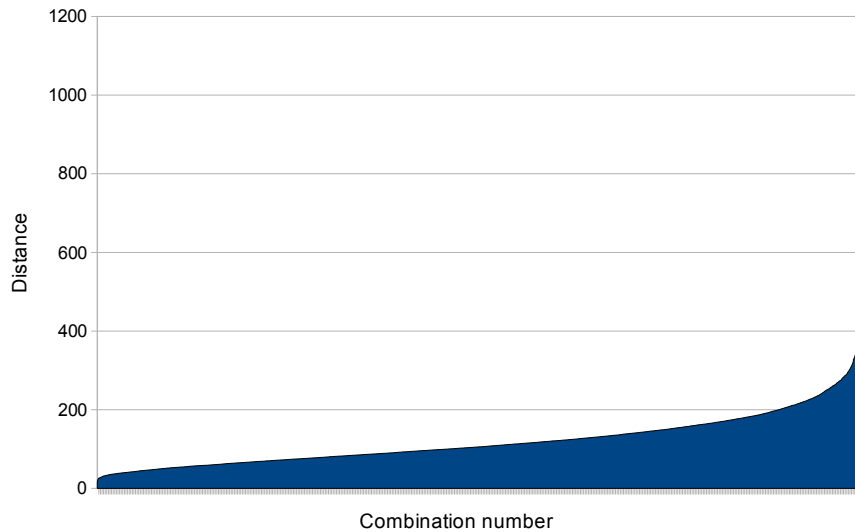


Figure 3.3: 10000 random Euclidean distances sampled from the `mfcc_avg` descriptor.

ate thresholds that at would result in around 50000 matches for each descriptor. 50000 was chosen to allow for a large safety margin.

## 3.4 System architecture

The system uses a MySQL database to store intermediate data, and several processing steps to combine the data in meaningful ways.

### 3.4.1 Database architecture

The fingerprints and scores from the various descriptors are fed into a common database to allow for easy querying. The database structure is presented in Figure 3.4.

It consists of four tables. Every file that is fingerprinted is present in `tblFile`, where the path to the file is stored, along with the size of the file. A unique identifier is assigned to each file. For quick lookups the MD5 sum of the file's path is used as the index. Since all of the paths have a common prefix and a MD5 causes a shorter prefix lookup in the index.

`tblType` contains a list of the different descriptors used. Since `fdmf` produces three separate fingerprints, it has three entries in this table. `tblFingerprint` contains the fingerprints for files for those descriptors that do not have their own internal fingerprint database. `tblScore` contains the likelihood that two files are equivalent when comparing with a given descriptor. If a descriptor finds no similarity between two files, it will not insert a row into `tblScore` to save calculation time and space. The subsequent calculation will account for missing 0-score rows. `tblSummedScore` is used for storing the final likelihood

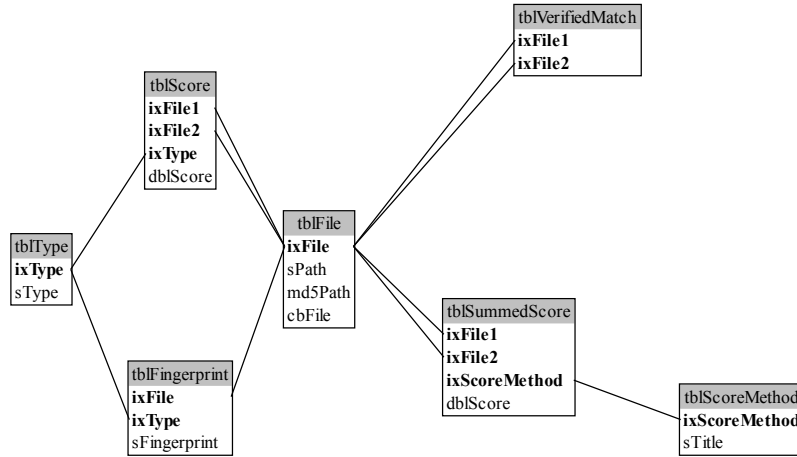


Figure 3.4: Database architecture of the combined solution. Primary keys are marked with a bold font.

that two files are equivalent, using different combination methods, as outlined in section 3.5.

It is worth noting that for all the tables that contain the two identifiers `ixFile1` and `ixFile2`, it is always the case that `ixFile1` is smaller than `ixFile2`. This assertion holds as long as all descriptor matching is guaranteed to be symmetric. All the descriptor matching is done using the Euclidean distance. The Euclidean distance is a metric, and metrics are symmetric by definition [37].

### 3.4.2 Data flow

The combined system modifies the software packages used to input data in several processing steps. There are three separate fingerprint processes that are run on the audio data; `fdmf`, `libFooID` and Gunderson’s descriptors. `fdmf` uses its own internal database to store the fingerprint data, while both `libFooID` and Gunderson’s descriptors relies on our MySQL database.

There are three comparison processes, that compares fingerprints for all files against each other using their descriptors, and stores the results in `tblScore`. A single process combines these results (as presented in 3.5), and enters the results into `tblSummedScore`. The entire process is outlined in figure 3.5.

## 3.5 Combining results

Several methods exists for finding a central tendency among a range of values, and the goal of this thesis is to highlight the merits of each of these methods. In the following, the likelihood for two songs matching as given by descriptor  $i$  is designated  $p_i$ .

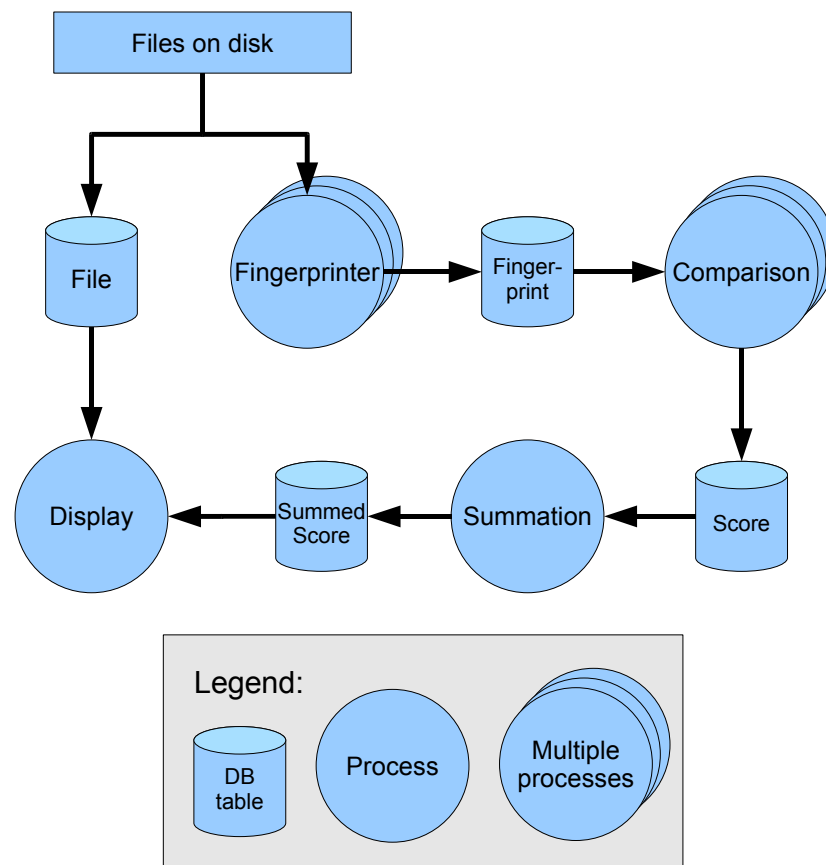


Figure 3.5: The flow of data through the system.

In addition to a likelihood measure given by various averages, the standard deviation is calculated when we are dealing with statistical measures. The standard deviation using average  $x$  is given in equation 3.3, where  $P_x$  is the probability given by e.g. the arithmetic mean or the truncated mean. It is worth noting that calculating the standard deviation is very similar to taking the root-mean square, except that you sum the squares of deviation. Also worth noting is that the standard deviation can not be compared directly unless the same averaging method is used.

$$\sigma_x = \sqrt{\frac{1}{I} \sum_{i=0}^I (p_i - P_x)^2} \quad (3.3)$$

Results are then sorted by descending values of probability, and then by ascending values of standard deviation. This ensures that high-probability duplicate pairs are sorted first, and when two duplicate pairs have the same probability, the standard deviation is the tie-breaker. If both the probability and standard deviation match, the files are sorted according to which duplicate pair's first file was first entered into the system.

Some of these methods are highly affected by outliers, and some descriptors produces lots of outliers. Therefore, in addition to running on the full set of data, some where also run on a set of the top performing descriptors, while ignoring the other descriptors. These were the arithmetic mean, the root-mean square, the truncated arithmetic mean and the Bayesian classifier. The descriptors that were chosen for this set was `fdmf_0`, `fdmf_1`, `fdmf_2`, `libfooid`, `mfcc_avg`, `mfcc_delta_avg`, `mfcc_f1_avg` and `mfcc_f1_delta_avg`. The median was used only for the top performing descriptors.

### 3.5.1 Arithmetic mean

The arithmetic mean (AM), commonly referred to as the *mean*, is the simplest and most common averaging method. It is one of several methods to indicate the central tendency of a range of numbers, and is calculated as in equation 3.4

$$P_{am} = \frac{1}{I} \sum_{i=0}^I p_i \quad (3.4)$$

The arithmetic mean is a good measure for distributions that have a small number of outliers. Outliers can significantly impact the result. A high standard deviation would usually indicate that the arithmetic mean is unsuitable.

### 3.5.2 Root-mean square

The root-mean square (RMS) of a set of values is the square root of the mean of values squared. It is most useful when values shift from positive to negative, and is well known in electrical engineering and with audio enthusiasts under the abbreviation RMS.

$$P_{rms} = \sqrt{\frac{1}{I} \sum_{i=0}^I p_i^2} \quad (3.5)$$

### 3.5.3 Weighted arithmetic mean

The weighted arithmetic mean (WAM) is an averaging method used when different values are assigned different weights. If we know that one fingerprinting method has a higher likelihood of producing false positives, we can reduce its weight, and thus reduce its final impact on the score. For the weighted arithmetic mean, each likelihood  $p_i$  has an associated weight  $w_i$ .

$$P_{wam} = \frac{\sum_{i=0}^I w_i p_i}{\sum_{i=0}^I w_i} \quad (3.6)$$

The arithmetic mean is a special case of the weighted mean, where all the weights are equal. The weights were adjusted manually for some runs, based on the reputation of individual descriptors.

### 3.5.4 Weighted root-mean square

A weighted version of the root-mean square measure was derived from the ordinary RMS measure after seeing the performance of the weighted arithmetic mean. Using the weights from weighing by performance and the second set of manual weights, equation 3.7 was used to compute a weighted root-mean square (WRMS).

$$P_{rms} = \sqrt{\frac{\sum_{i=0}^I w_i p_i^2}{\sum_{i=0}^I w_i}} \quad (3.7)$$

The formula was adapted from the root-mean square for this thesis, but has been used in a similar form in e.g. [3].

### 3.5.5 Truncated mean

A truncated mean discards some data values to get rid of outliers. The truncated arithmetic mean was used to get rid of the worst outliers, by removing an equal number of values on each end of the input data, and calculating the arithmetic mean of the remaining data. We only truncated two values, one from each side of the data set.

### 3.5.6 Median

The median is a simplistic measure that does not work well at all when all the descriptors are present. With a large number of descriptors returning perfect matches for non-related songs, the median will be more or less random in many situations.

When working with the set of known good descriptors, it could very well indicate a tendency among the values. The median is not affected as considerably by outliers as other averages.

### 3.5.7 Naïve Bayes

Bayesian analysis lends itself elegantly to the classification problem and is used with a high degree of success in spam filters. By starting with Bayes theorem,

a formula for calculating the probability that two songs match given a set of clues will be derived. This approach is called naïve Bayes (NB) [21].

A prerequisite for using Bayesian analysis is that the scores returned by the varying engines is the actual probability of the two songs matching. It is a well known fact that this assumption does not hold for many of the descriptors, and as such, the Bayesian analysis has been limited to only encompass the 8 best-ranked descriptors from weighing by performance (see 3.5.9).

In the following,  $P(M)$  is the probability of a correct match.  $P(C_i)$  is the probability for the two songs matching given by descriptor  $i$ . The formulas assume only two clues, while the final formula is easily expandable.

$$P(M|C_1 \cap C_2) = \frac{P(C_1 \cap C_2|M) P(M)}{P(C_1 \cap C_2)} \quad (3.8)$$

$$= \frac{P(C_1 \cap C_2|M) P(M)}{P(M) P(C_1 \cap C_2|M) + P(\neg M) P(C_1 \cap C_2|\neg M)} \quad (3.9)$$

Bayes' theorem was used to get to 3.8, and use the usual expanded form of Bayes' to get 3.9. To continue the assumption that our clues are independent is needed, which turns the algorithm into a naïve Bayes algorithm. We can now substitute  $P(C_1 \cap C_2|M) = P(C_1|M) P(C_2|M)$ :

$$= \frac{P(C_1|M) P(C_2|M) P(M)}{P(M) P(C_1|M) P(C_2|M) + P(\neg M) P(C_1|\neg M) P(C_2|\neg M)} \quad (3.10)$$

Applying Bayes' theorem three times gives, after some calculation:

$$= \frac{\frac{P(M|C_1)P(M|C_2)}{P(M)}}{\frac{P(M|C_1)P(M|C_2)}{P(M)} + \frac{P(\neg M|C_1)P(\neg M|C_2)}{P(\neg M)}} \quad (3.11)$$

Substituting  $a = P(M|C_1)$ ,  $b = P(M|C_2)$  and  $S = P(M)$  leaves a simpler formula:

$$L(F_1, F_2) = \frac{\frac{ab}{S}}{\frac{ab}{S} + \frac{(1-a)(1-b)}{1-S}} \quad (3.12)$$

This formula expands as you might expect:

$$L(F_1, F_2) = \frac{\frac{abc}{S}}{\frac{abc}{S} + \frac{(1-a)(1-b)(1-c)}{1-S}} \quad (3.13)$$

Depending on the number of clues, a generalized formula is presented below:

$$L(F_1, F_2) = \frac{\prod C_i}{\prod C_i + \frac{\prod (1-C_i)}{1-S}} \quad (3.14)$$

In our database, a total of 31 million potential duplicate pairs will be evaluated using this algorithm. An estimated 2200 correct duplicate pairs exist. This means  $P(M) = S \approx \frac{2200}{31000000} = 8 \times 10^{-5}$ , and that the probability for picking a correct duplicate pair from the potential duplicate pairs is 0.008 percent. However, when using this value, the formula proved to be to extremely sensitive

to outliers. After repeated experimentation the value was therefore manually adjusted to 0.8. The value had to be adjusted because otherwise a tremendous amount of evidence would be required for the naïve Bayes algorithm to classify two songs as a match. A value of 0.8 results in the values being more spread out across the available spectrum.

### 3.5.7.1 Problems with naïve Bayes

There are two main problems with Bayesian analysis as applied to our system. The first is the reason for excluding descriptors. As an example, assume that we have three descriptor scores, length with a probability of 1.0, fdmf\_0 with a score of 0.2 and mfcc\_avg with a score of 0.1. Since the probability of the two songs being similar is 1 as given by the length descriptor, the combined probability is also 1. This means that any descriptor producing a perfect score will completely distort the result. Inversely, if any descriptor returns zero, the combined score will also be zero.

Two solutions can be used to prevent this problem. First, the scores are capped at both ends, and is always between 0.01 and 0.99. This ensures that a single descriptor cannot override the rest of the results. Second, one can pick scores from accurate descriptors liberally, which we are doing when we are selecting the top performing descriptors.

The second problem lies with the assumption that the scores from all the descriptors are statistically independent. This assumption is likely not to hold, as the algorithms used by the various descriptors are similar. This problem is frequently ignored when writing Bayesian classifiers and is why this is considered a *naïve* Bayes algorithm.

### 3.5.8 Discarded approaches

Several averaging methods were not selected for this approach, mostly due to the the large quantity of zero values in the data set. For example, the harmonic mean, defined as in equation 3.15 will result in division by zero, and distort the result. The harmonic mean is therefore not applicable in our tests.

$$P_{hm} = \frac{I}{\sum_{i=0}^I \frac{1}{p_i}} \quad (3.15)$$

The geometric mean is another average where the  $I$ -th root of the product of values is taken. Equation 3.16 shows the calculation, but due to the multiplication of probabilities, we will always get zero probability if only one of the results is zero.

$$P_{gm} = \sqrt[I]{\prod_{i=0}^I p_i} \quad (3.16)$$

### 3.5.9 Determining weights

Some of the descriptors tested produced a larger number of potential duplicate pairs than others. Some even two orders of magnitude apart. It is likely that a potential duplicate pair from a descriptor that produces few potential duplicate



pairs, will be more likely to be correct than a potential duplicate pair from a descriptor that produces one hundred times as many.

By this assumption, four different ways of calculating the weights were devised. In one run, the weight  $w_i$  was calculated using the number of non-zero scores  $n_i$ , compared to the maximum number of non-zero scores produced by a single descriptor  $n_{max}$ . It is shown as equation 3.17. The multiplication by 2, is to ensure that the lowest assigned weight is 0.5, instead of 0.

$$w_i = 1 - \frac{n_i}{2 \times n_{max}} \quad (3.17)$$

This method, referred to as weighing by count, only looks at the number of potential duplicate pairs found to determine which descriptors are trustworthy. So, if a descriptor produces a lot of very low-ranked potential duplicate pairs (e.g. with a score of 0.01 for 99% of the comparisons), it will be ranked as unreliable, even if it is not. Another weighing scheme looks at the sum of the of scores instead, so that a descriptor that produces very many low-ranking results will not be penalized as harshly, while a descriptor that produces many high-ranking results compared to the amount of results will be penalized harder.

More specifically, the weight for a given descriptor  $i$  was calculated using the sum of scores  $s_i$ , as presented in equation 3.18, where  $s_{max}$  is the highest  $s_i$ .

$$w_i = 1 - \frac{s_i}{2 \times s_{max}} \quad (3.18)$$

The third method uses the arithmetic mean of the scores given by one descriptor. A high average is assumed to indicate that the descriptor is less trustworthy. Equation 3.19 shows the calculation, where  $s_{j,i}$  is the  $j$ -th score for descriptor  $i$ , and  $J_i$  is the number of scores for descriptor  $i$ .

$$w_i = 1 - \frac{\sum_{j=0}^{J_i} s_{j,i}}{J_i} \quad (3.19)$$

The most accurate way to get a measure of the usefulness of a descriptor, is to look at its actual performance. After having finished the database of verified duplicate pairs, the performance was measured using the false positive rate in the highest-ranked one thousand matches. This method was titled weighing by performance.

Using these equations on the fingerprint database produced weights as can be seen in figure 3.6. It is apparent from weighing by performance that a set of the descriptors are nearly useless, an indication that weighing by average also seems to support. Only weighing by performance recognizes `mfcc_delta_avg` as a good indicator, with the exception of weighing by average. Almost all of these methods seem to agree that `fdmf_1` is a bad estimator as well.

In addition to the automatically computed weights, several sets of weights were created manually, based on the performance of the different weighing schemes. These weights are shown in figure 3.7. The first set of manual weights was chosen before the relative merits of the different descriptors were known, and was based on their reputation. Manual weights method 2 and 3 are based on weighing by performance, where weights less than 0.05 and 0.1 respectively are ignored. This makes the manual weights comparable to the other means when only a chosen set of descriptors is used, as it is the same set of descriptors that is used.

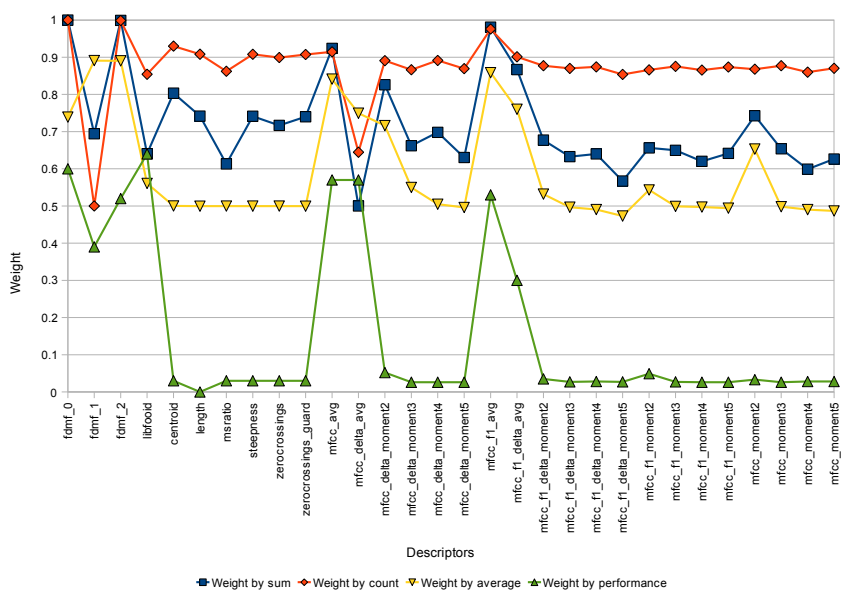


Figure 3.6: Weights as generated by equation 3.17, 3.18, 3.19 and weighing by performance. Note that the weights are relative to weights within its series, and as such, can not be compared directly to weights given by another series.

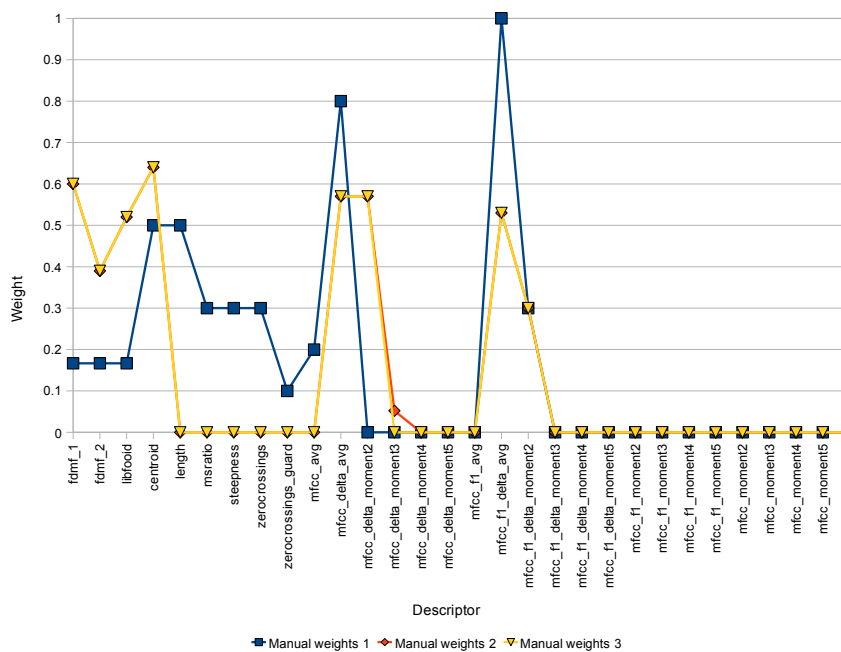


Figure 3.7: Three sets of manual weights.

Descriptor	Weights						
	Sum	Count	Avg.	Perf.	M. 1	M. 2	M. 3
fdmf_0	1	1	0.74	0.6	0.17	0.6	0.6
fdmf_1	0.69	0.5	0.89	0.39	0.17	0.39	0.39
fdmf_2	1	1	0.89	0.52	0.17	0.52	0.52
libfooid	0.64	0.85	0.56	0.64	0.5	0.64	0.64
centroid	0.8	0.93	0.5	0.03	0.5	0	0
length	0.74	0.91	0.5	0	0.3	0	0
msratio	0.61	0.86	0.5	0.03	0.3	0	0
steepness	0.74	0.91	0.5	0.03	0.3	0	0
zerocrossings	0.72	0.9	0.5	0.03	0.1	0	0
zerocrossings_guard	0.74	0.91	0.5	0.03	0.2	0	0
mfcc_avg	0.92	0.91	0.84	0.57	0.8	0.57	0.57
mfcc_delta_avg	0.5	0.64	0.75	0.57	0	0.57	0.57
mfcc_delta_moment2	0.83	0.89	0.72	0.05	0	0.05	0
mfcc_delta_moment3	0.66	0.87	0.55	0.03	0	0	0
mfcc_delta_moment4	0.7	0.89	0.5	0.03	0	0	0
mfcc_delta_moment5	0.63	0.87	0.5	0.03	0	0	0
mfcc_f1_avg	0.98	0.98	0.86	0.53	1	0.53	0.53
mfcc_f1_delta_avg	0.87	0.9	0.76	0.3	0.3	0.3	0.3
mfcc_f1_delta_moment2	0.68	0.88	0.53	0.04	0	0	0
mfcc_f1_delta_moment3	0.63	0.87	0.5	0.03	0	0	0
mfcc_f1_delta_moment4	0.64	0.87	0.49	0.03	0	0	0
mfcc_f1_delta_moment5	0.57	0.85	0.47	0.03	0	0	0
mfcc_f1_moment2	0.66	0.87	0.54	0.05	0	0	0
mfcc_f1_moment3	0.65	0.88	0.5	0.03	0	0	0
mfcc_f1_moment4	0.62	0.87	0.5	0.03	0	0	0
mfcc_f1_moment5	0.64	0.87	0.49	0.03	0	0	0
mfcc_moment2	0.74	0.87	0.65	0.03	0	0	0
mfcc_moment3	0.65	0.88	0.5	0.03	0	0	0
mfcc_moment4	0.6	0.86	0.49	0.03	0	0	0
mfcc_moment5	0.63	0.87	0.49	0.03	0	0	0

Table 3.2: The set of weights used by the weighted arithmetic mean. The weighted root-mean square only use the sets labeled manual 2 and performance.



# Chapter 4

## Results

This chapter presents the results from the experiments done in chapter 3.

### 4.1 Evaluating results

Several measurements are common when evaluating search results. A very common approach is to use a precision versus recall curve to outline how the system performs in retrieving relevant results. Since our system does not have a complete list of which results are relevant, the precision is impossible to compute. The F-measure, defined as  $\frac{2pr}{p+r}$ , where p is precision and r is recall, is excluded for the same reason.

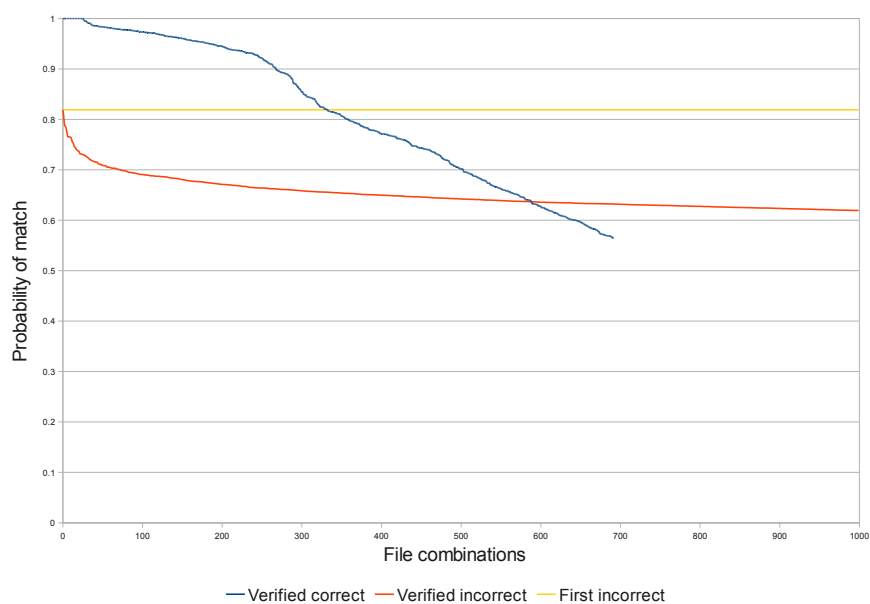
Instead an approach that measures the number of false positives in the first thousand results were chosen. For each descriptor, we extract the thousand highest-ranked results, and compare these results to our list of verified duplicate pairs. By grouping the results into groups of one hundred, we can see roughly how many trustworthy results a descriptor produces. The descriptors are then ranked by the percentage of false positives in the first thousand positive results. After the first thousand results none of the individual descriptors or combined solutions produced less than 50% error rates when divided into partitions of one hundred.

This method of evaluation takes into account the fact that results have limited use if not sorted correctly. Google became popular because they realized that the important thing about search was putting relevant results first [33, 26]. This also applies to finding duplicates in a music collection, the most relevant results should appear first.

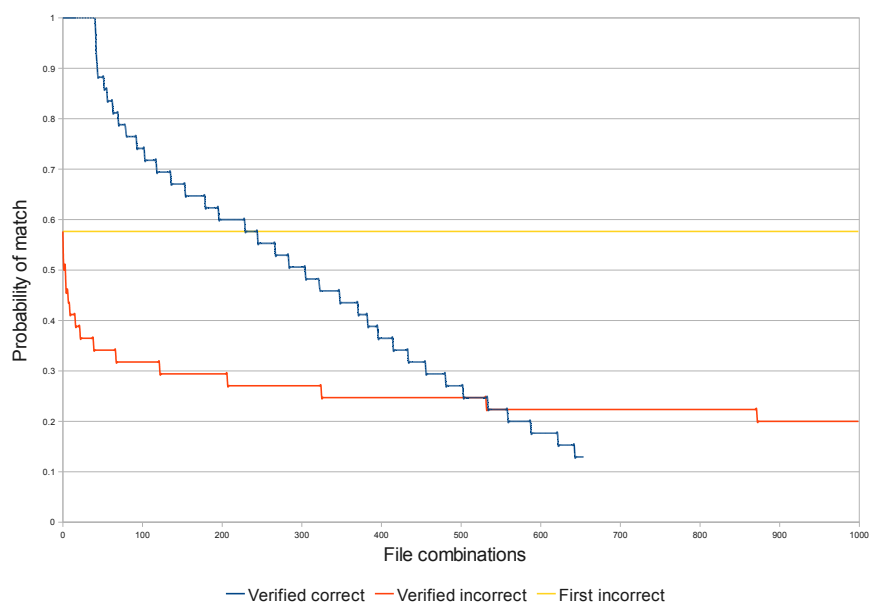
### 4.2 Individual descriptors

The results from the individual descriptors were not impressive. To examine the results closely, the first thousand results were ranked and then partitioned into 10 partitions of 100 results. Figures 4.2 and 4.3 show the results from the individual descriptors, with the x-axis being the partitions, and the y-axis being the number of false positives.

As can be seen in the graph, several of the descriptors do reasonably well. `mfcc_avg` and `mfcc_f1_avg` do not produce a single false positive in the first 300



(a) mfcc avg



(b) fdmf 2

Figure 4.1: Results from two individual descriptors graphed. Data generated from list of potential duplicate pairs ordered by the likelihood given by the individual descriptors. The top line shows verified duplicate pairs, while the lower line shows the incorrect matches. The middle, straight line highlights the first error.

ranked results. `libfooid`, `fdmf_0`, `fdmf_2`, `mfcc_delta_avg` and `mfcc_f1_delta_avg` also show promising results, and produced few false positives in the first 400 results.

Several of the descriptors, such as song length, perform badly. The first 3 599 results returned by that descriptor all indicate perfect matches, and ranking results is therefore reduced to which files were entered first into the database. It is therefore obvious that song length as a descriptor is completely inadequate for its intended purpose, and it can at best be used as an optimization detail.

Figure 4.4 shows the total false positive rate for the first thousand for each individual descriptor. These numbers were used to generate the weights by performance as detailed in section 3.5.9.

The results from a subset of the descriptors were analyzed to determine how they fail, and if there are patterns in the results that can be exploited. The top 100 erroneous results from each descriptor was, and the three highest ranking erroneous matches was excerpted for the reader's reference.

### 4.2.1 `fdmf_0`

This descriptor represents the energy spectrum of the audio signal, and is one of the better-performing descriptors. Considering the first few erroneous matches, only two has a probability higher than 0.5.

File 1	File 2	Score
Billie Holiday - Spoken Introduction	Sound of electric drill	0.86667
Billie Holiday - A Sunbonnet Blue	Sound of electric drill	0.57333
Billie Holiday - Just One of Those Things	Taube - Sa länge skutan kan gå	0.44

Table 4.1: First three erroneous duplicate pairs from `fdmf_0`.

None of these audio recordings are in any way perceptually similar, which indicates that the energy spectrum does not look at the audio in the same way as a human would. The electric drill recording is a recording that was inadvertently included in the collection, but illustrates that music collections very often has short recordings that might not be music.

Worth noting with this descriptor, is that it has not given a full score to two files of different file formats. This could be because the music collection contains more MP3 files than WMA files, but the absence of WMA files in the first 50 results is noteworthy, as they appear regularly from that point onwards, always with a score lower than 0.9.

### 4.2.2 `fdmf_1`

Looking at the ratio spectrum fingerprint generated by `fdmf`, it is the poorest performing of the `fdmf` descriptors, however it does better than some of the other descriptors.

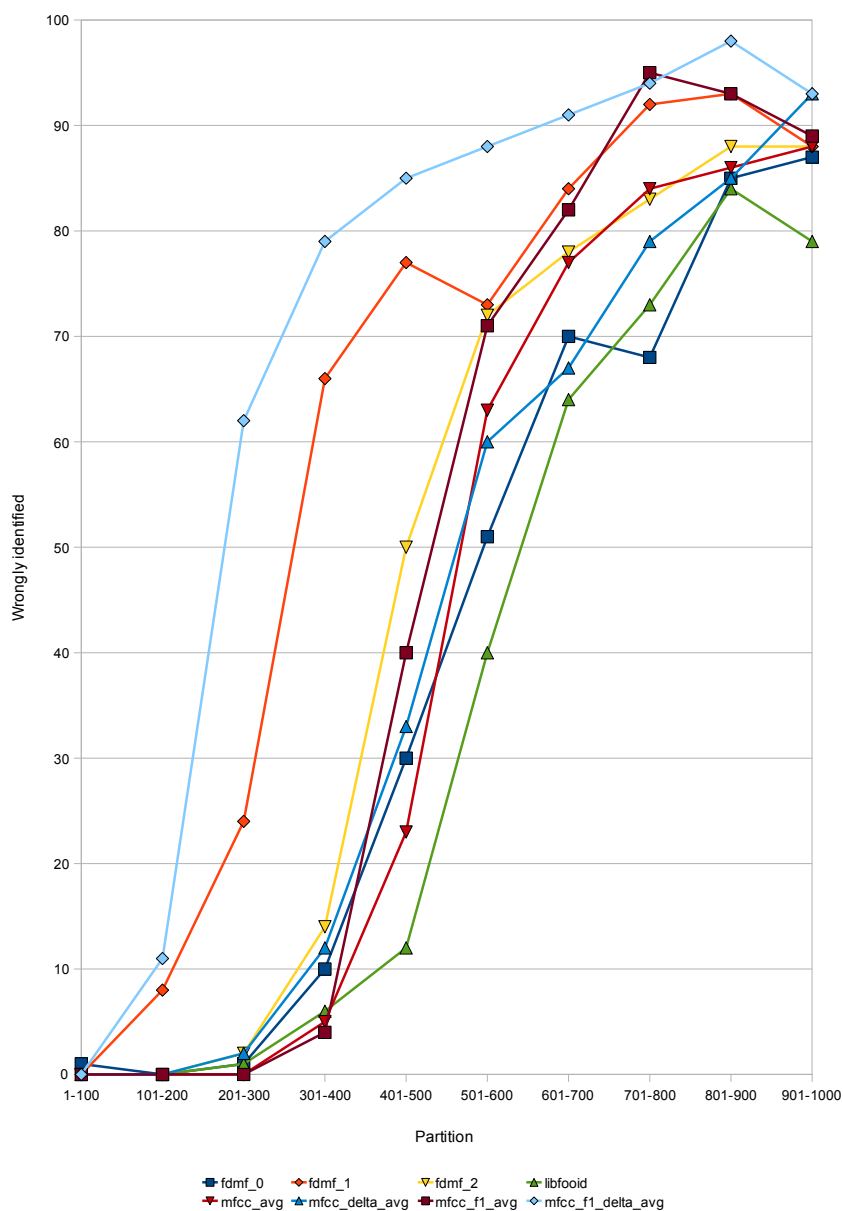


Figure 4.2: Number of false positives in the first 1000 results for the best performing descriptors, divided into partitions of 100 results. x-axis shows partition number, y-axis shows the number of false positives.



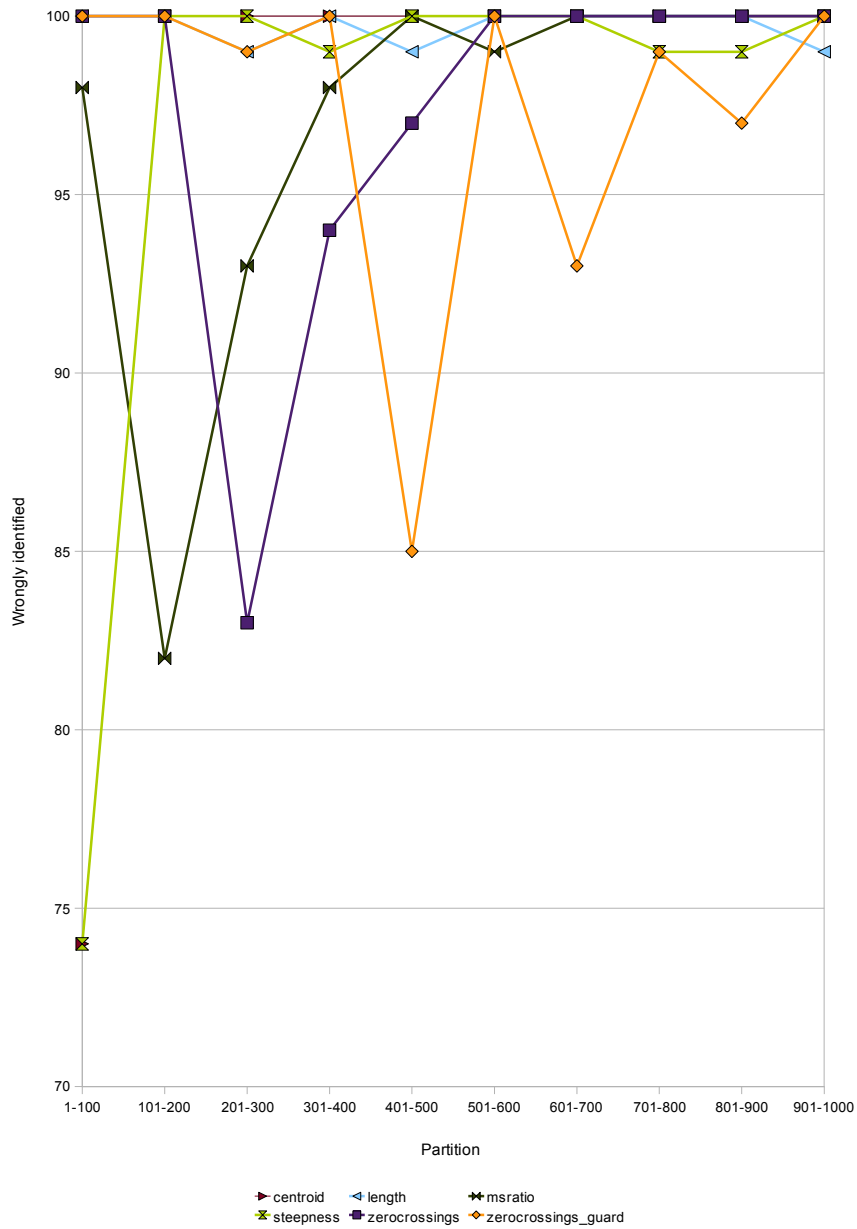


Figure 4.3: Number of false positives in the first 1000 results for the simple descriptors, divided into partitions of 100 results. x-axis shows partition number, y-axis shows the number of false positives. Note that the scale varies from figure 4.2.

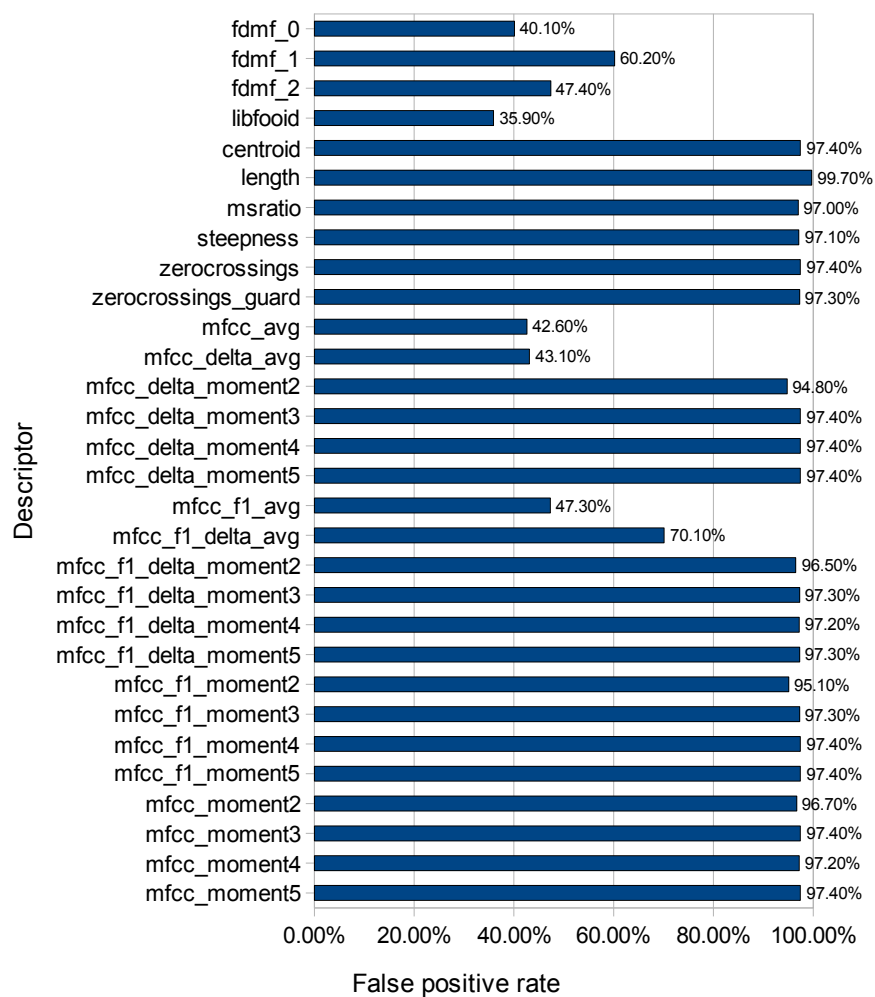


Figure 4.4: False positive rate from the individual descriptors for the first thousand results.

File 1	File 2	Score
The Doors - Hello, I Love You	Taube - Stockholmsmelodi	0.791304
The Doors - Hello, I Love You	Vamp - Ingeborg	0.773913
Hugh Cornwell - Snapper	Taube - Stockholmsmelodi	0.773913

Table 4.2: First three erroneous duplicate pairs from `fdmf_1`.

Unlike `fdmf_0`, `fdmf_1` produces a large number of high-probability matches that are erroneous. Among the wrong results, some songs are repeatedly matched against a variety of songs, while the songs themselves sound very different to a human.

The same tendency as for `fdmf_0` in detecting duplicate pairs with different file types is present in `fdmf_1`. The first cross-match occurs at position 50, with a score of less than 0.9.

### 4.2.3 `fdmf_2`

The twist spectrum fingerprint from `fdmf` produces few high-probability erroneous matches. Only the first 16 erroneous matches have a score higher than 0.4.

File 1	File 2	Score
Billie Holiday - Spoken Introduction	Sound of electric drill	0.576471
Billie Holiday - Yesterdays	Billie Holiday - Our Love Is Different	0.505882
Billie Holiday - Jeepers Creepers	Billie Holiday - Our Love Is Different	0.505882

Table 4.3: First three erroneous duplicate pairs from `fdmf_2`.

Again we see the same audio matching in number 1 for `fdmf_0`, and a predominance for Billie Holiday recordings. Common for those recordings are the poor quality and the noise level, which may confuse the descriptor.

Again, `fdmf_2` seems to match worse on files with different file formats. This might indicate that `fdmf`'s spectrum analysis is sensitive to encoding artifacts in various file formats.

### 4.2.4 `libFooID`

`libFooID` produces a range of high-probability erroneous duplicate pairs, while having the lowest false-positive rate of the tested descriptors. However, the first three erroneous duplicate pairs are all moody, string-heavy music with no drums, which could indicate that `libFooID` picks up the lack of percussion as a similarity.

This descriptor deserves a special mention, as it is the only complete fingerprinting solution that is used without decomposing it into its internal parts. It is also the best performing single descriptor tested.

File 1	File 2	Score
Keith Jarret - Hymn For Good Friday	Michael Andrews - Carpathian Ridge	0.852641
Brian Eno - Harmonic Studies	Michael Andrews - Cellar Door	0.851268
Michael Andrews - Carpathian Ridge	Michael Andrews - The Tangent Universe	0.505882

Table 4.4: First three erroneous duplicate pairs from libFooID.

Also libFooID seems to rank differently encoded files somewhat lower than files in the same file format.

### 4.2.5 Centroid

The centroid does not perform well, and a large amount of audio files produce very high probability matches. The erroneous duplicate pairs identified by the centroid do not appear to have anything in common when listened to by a human.

File 1	File 2	Score
Talking Heads - Listening Wind	The Triffids - Jerdacuttup Man	1
Neil Young - Trans Am	Rush - Fly By Night	1
Apoptygma Berzerk - Black Pawn	The Waterboys - Fisherman's Blues	0.999999

Table 4.5: First three erroneous duplicate pairs from the centroid descriptor.

### 4.2.6 Song length

The song length is such an obviously flawed descriptor to use for audio identification of a single audio file. It therefore comes as no surprise that it is the worst performing descriptor.

The usefulness of the track length as a descriptor is highly dependent on the use case in question: In discerning different tracks from CD, it can be an excellent measure, in genre classification, it is nearly useless, and finally, in some cases (such as when identifying what music is being played on the radio at a given instant), it might not be available at all. [15]

It is hard to imagine the song length as a reliable indicator of the contents of the music, as a lot of the songs produced today is between 3 and 4 minutes long. It can be useful in detecting outliers, such as quickly skipping a comparison of two songs if the songs are 10 minutes different in length.

	Song name	Length	File size
File 1	Björk - Hunter	4:15	6 160 832 bytes
File 2	Hugh Cornwell - Always The Sun (live)	4:15	8 390 786 bytes
File 1	Björk - Venus As A Boy	4:41	6 796 152 bytes
File 2	Enya - Water Shows The Hidden Heart	4:41	6 796 956 bytes
File 1	Björk - Army of Me	3:55	5 695 546 bytes
File 2	The Rolling Stones - It Won't Take Long	3:55	8 330 019 bytes

Table 4.6: First three erroneous duplicate pairs from using track length as a descriptor.

If the track length is available, using it as an optimization to quickly remove tracks that vary greatly in length could work, and is done by e.g. libFooID. It cannot be trusted as an individual descriptor due to the very high likelihood that other songs fall within the same range.

#### 4.2.6.1 Song length as an optimization

The song length is a descriptor that is very easy to understand, and so are the reasons for why it will not work as a descriptor for identifying single audio files. The average song length in our collection is 230 seconds (3 minutes and 50 seconds). In our collection, a total of 2826 songs are within 30 seconds of the average song length. 1465 songs are within 15 seconds from the average track length.

If using the song length as an optimization and automatically failing comparisons of two songs with a difference in length exceeding 30 seconds; you can skip approximately 85% of the actual fingerprint comparisons on average. In our system, this excluded 231 verified duplicate pairs, as our criteria specifies only looking at the first minute of audio.

#### 4.2.7 Mean/square ratio

The mean/square ratio performs on an equal level with the centroid, but has a large number of erroneous duplicate pairs with a perfect score. The erroneous duplicate pairs found do not carry a perceptual similarity when listened to by a human.

File 1	File 2	Score
Björk - Joga	Led Zeppelin - Moby Dick	1
Billy Idol - Crank Call	David Bowie - Strangers When We Meet	1
Billie Holiday - Romance in the Dark	Tom Waits - My Baby Left Me on A Trash Day	1

Table 4.7: First three erroneous duplicate pairs from the mean/square ratio descriptor.

### 4.2.8 Steepness

Steepness is another badly performing descriptor, and reports many high-probability erroneous duplicate pairs. The erroneous duplicate pairs again show no discernible pattern in what triggers a match using that descriptor.

File 1	File 2	Score
Apoptygma Berzerk - More Serotonin Please	Billie Holiday - Interview 1a	0.999998
Annie Lennox - Loneliness	deLillos - Fange i ditt eget bur	0.999997
Kevin Bloody Wilson - Flowers	XTC - Knights In Shining Karma	0.999996

Table 4.8: First three erroneous duplicate pairs from the steepness descriptor.

### 4.2.9 Rate of zero crossings

Both with and without a guard band, the ratio of zero crossings produce a large amount of erroneous duplicate pairs with perfect scores. This means the results are as random as the song length, and no discernible pattern can be found in the duplicate pairs it finds.

### 4.2.10 Mel frequency cepstral coefficients

mfcc\_avg is one of the best performing descriptors for our task, and has an apparent tendency in its erroneous duplicate pairs to link two songs by the same artist. In fact, of the top ten erroneous duplicate pairs produced, only two were different artists. In the eight cases, all the songs matched were found on the same album.

File 1	File 2	Score
Tom Waits - I'll Shoot The Moon	Tom Waits - Pony	0.818893
Dane Cook - Someone Shit On The Coats	Dane Cook - Driveway Intruder	0.806178
Tom Waits - Eyeball Kid	Tom Waits - Pony	0.79033

Table 4.9: First three erroneous duplicate pairs from the mfcc\_avg descriptor.

This clearly indicates that mfcc\_avg identifies songs using the same instruments and having the same remastering as more similar.

When looking at mfcc\_delta\_avg, it is apparent that it follows some of the same pattern. It performs on the same level as mfcc\_avg, but produces different erroneous duplicate pairs.

File 1	File 2	Score
Johnny Cash - Where We'll Never Grow Old	Johnny Cash - I'm Bound For The Promised Land	0.857711
Kevin Bloody Wilson - The Browneye Medley	Kevin Bloody Wilson - Australian Anthems	0.851686
Perssons Pack - I all vår tid	deLillos - Nå lever den av seg selv	0.843354

Table 4.10: First three erroneous duplicate pairs from the `mfcc_delta_avg` descriptor.

The various `mfcc_moment` and `mfcc_delta_moment` descriptors are as inaccurate as the centroid and the steepness. Examples of how they fail have therefore not been included.

#### 4.2.11 Floor-1 cepstral coefficients

`f1cc_avg` performs slightly worse for our tasks than `mfcc_avg`. Considering the improvements made to F1CC over MFCC, it was intended to increase the resilience to encoding artifacts. This might have changed the algorithm sufficiently to make it perform worse on other tasks.

File 1	File 2	Score
Two To Tango - Vision	deLillos - Den feite mannen	0.725673
Dane Cook - Struck By A Vehicle	Dane Cook - Abducted	0.71712
Kevin Bloody Wilson - Take It Like A Man	Tom Waits - Get Behind The Mule	0.716681

Table 4.11: First three erroneous duplicate pairs from the `f1cc_avg` descriptor.

The tendency to identify songs from the same artist as duplicate pairs is decreased in `f1cc_avg` compared to `mfcc_avg`, but still present. `f1cc_avg` produces a significantly different set of erroneous duplicate pairs than does `mfcc_avg`.

File 1	File 2	Score
Astor Piazzolla - Nuevo Tango	Nick Cave & The Bad Seeds - Sleeping Annaleah	0.88834
Billie Holiday - Stars Fell on Alabama	Billie Holiday - 'Deed I Do	0.882058
Johnny Cash - Where We'll Never Grow Old	Johnny Cash - If We Never Meet Again This Side Of Heaven	0.881461

Table 4.12: First three erroneous duplicate pairs from the `f1cc_delta_avg` descriptor.

Considering `f1cc_delta_avg`, it too produces significantly different erroneous duplicate pairs from `mfcc_delta_avg`, yet persists in the tendency of identifying duplicate pairs where both songs are from the same artist.

### 4.3 Combined solutions

When looking at figure 4.5 it is immediately evident that the applied methods works as a method of producing more stable results.

The best performing combination was the weighted root-mean square, which only had a failure rate of 4.6% among the first thousand results. This means it correctly identified 11 more duplicate pairs than the runner-up.

#### 4.3.1 Arithmetic mean

The arithmetic mean is highly affected by the large number of always high-probability descriptors. It therefore has a higher failure rate than many of the individual descriptors, and can therefore not be used without an algorithm for choosing reliable descriptors.

When using the arithmetic mean on the chosen descriptors only, the results markedly improve. The total error rate is now only 8%, and it has no erroneous results in the first five hundred. This suggests that it works well as a method of producing a stable fingerprint.

#### 4.3.2 Root-mean square

As with the arithmetic mean, the root-mean square produces useless results when not limiting the set of descriptors. However, when choosing the descriptor set carefully, the root-mean square performs slightly better than the arithmetic mean on the same descriptors, with a failure rate of 6.6%.

#### 4.3.3 Weighted arithmetic mean

The weighted arithmetic mean is among the most interesting combination methods in our set, and produces some of the lowest failure rates when certain sets of weights are used. First, weighing by the number of results produced by a descriptor did not provide a significant change in the results, and neither did weighing by the sum of probabilities. Weighing by the average score did not produce a large variation in the number of correct results found.

Weighing by performance did produce very interesting results, and is the sixth best performing method. Our manual weight adjustments, that removed descriptors with weights less than 0.05 and 0.10 produced a decrease in failure rates by 1 and 0.9 percentage points respectively.

It is noteworthy that the WAM with manual weights set 2 and 3 did produce its first erroneous duplicate pair later than any of the other methods used.

#### 4.3.4 Weighted root-mean square

The weighted root-mean square (WRMS) measure was introduced after seeing the success rates of the weighted arithmetic mean.

WRMS in this thesis used the same weights as weighing by performance for the weighted arithmetic mean, and the second set of manual weights because it ranked better in the weighted arithmetic mean. When weighing by performance, WRMS achieved a failure rate of 6.4%, which is 0.5 percentage points lower than



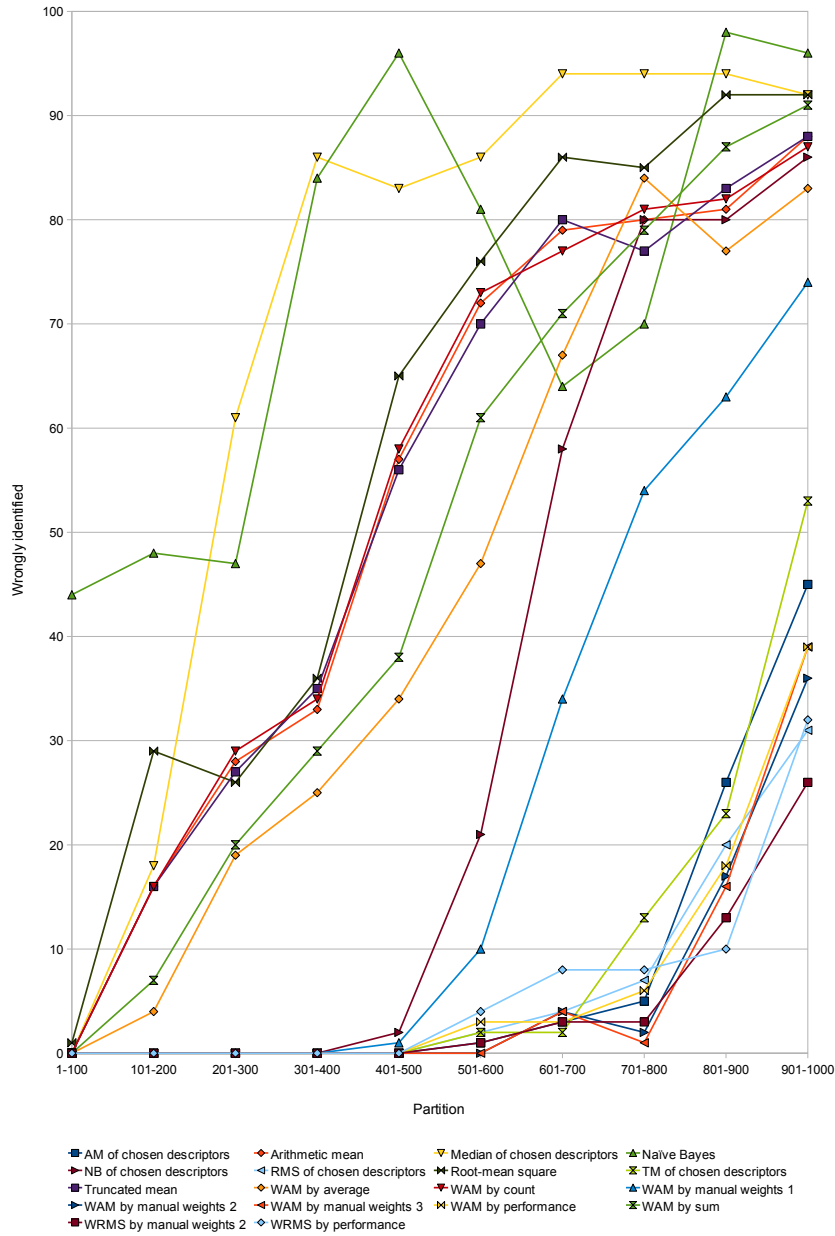


Figure 4.5: Number of false positives in the first 1000 results from the combined solutions, divided into partitions of 100 results. x-axis shows partitions, y-axis shows the number of false positives.

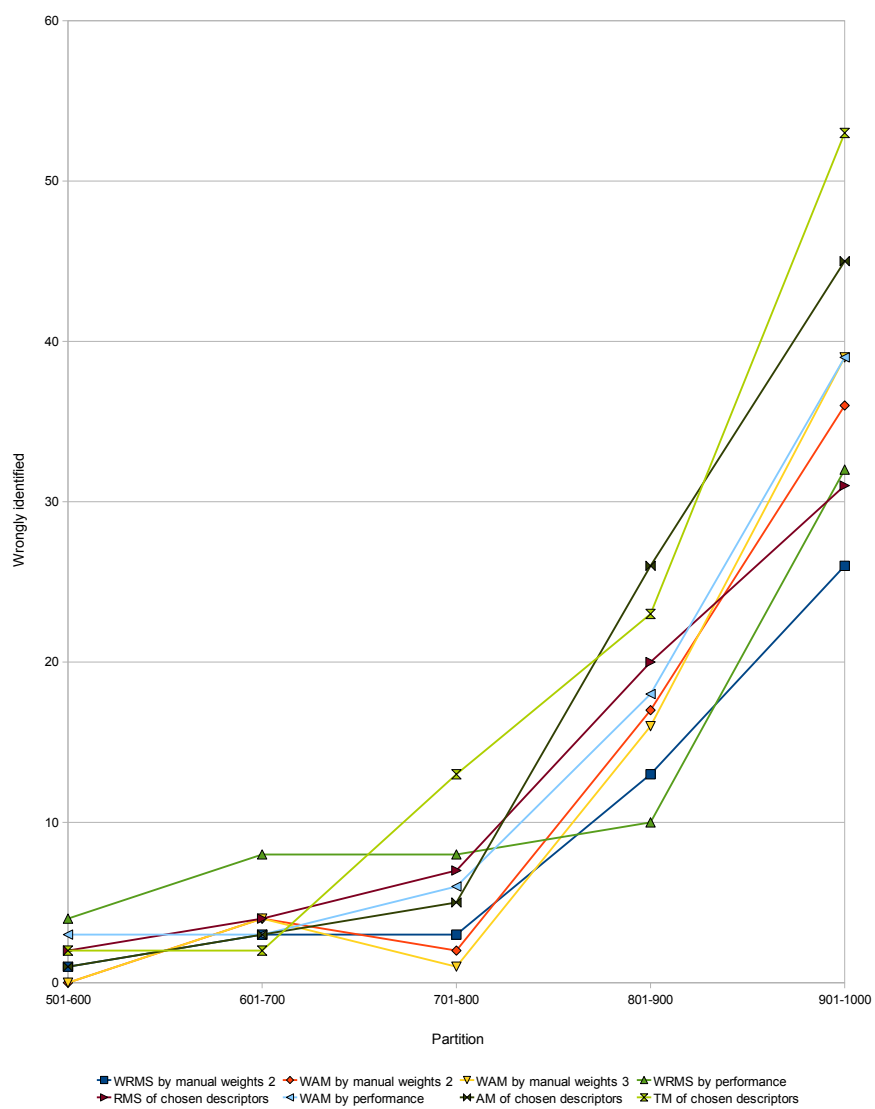


Figure 4.6: Zoomed in view of the lower right part of figure 4.5 for a subset of the combination methods.

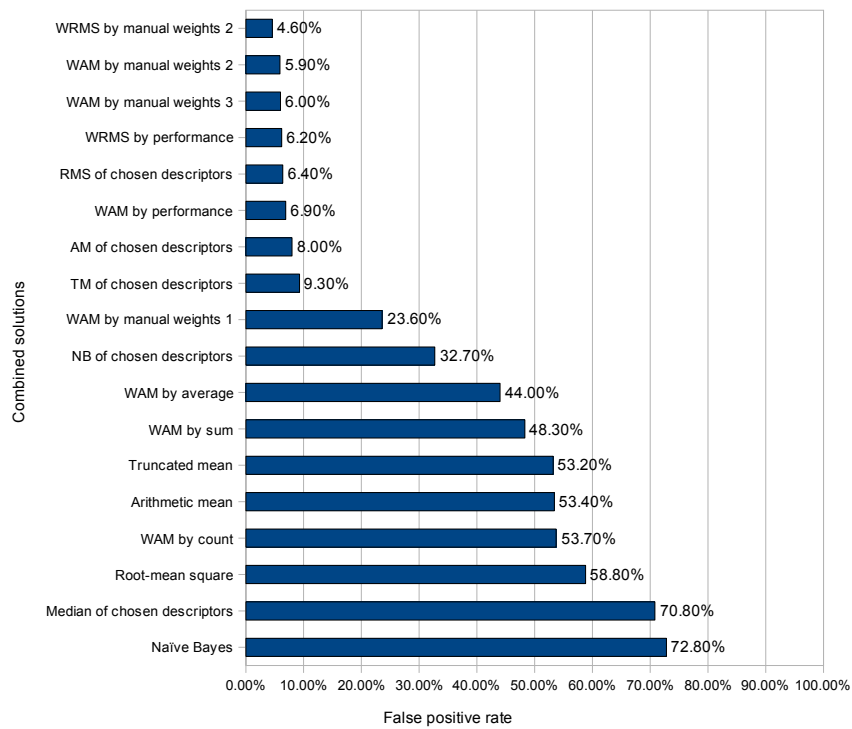


Figure 4.7: False positive rate from the combined solutions for the first thousand results.

Name	Partitions										False positive rate	First error at result #
	1	2	3	4	5	6	7	8	9	10		
WRMS using manual weights 2	0	0	0	0	0	1	3	3	13	26	4.60%	596
WAM using manual weights 2	0	0	0	0	0	0	4	2	17	36	5.90%	603
WAM using manual weights 3	0	0	0	0	0	0	4	1	16	39	6.00%	624
WRMS by performance	0	0	0	0	0	4	8	8	10	32	6.20%	513
RMS of chosen descriptors	0	0	0	0	0	2	4	7	20	31	6.40%	560
WAM by performance	0	0	0	0	0	3	3	6	18	39	6.90%	525
AM of chosen descriptors	0	0	0	0	0	1	3	5	26	45	8.00%	583
TM of chosen descriptors	0	0	0	0	0	2	2	13	23	53	9.30%	547
WAM using manual weights 1	0	0	0	0	1	10	34	54	63	74	23.60%	448
NB of chosen descriptors	0	0	0	0	2	21	58	80	80	86	32.70%	489
WAM by average	0	4	19	25	34	47	67	84	87	83	44.00%	168
WAM by sum	0	7	20	29	38	61	71	79	87	91	48.30%	161
TM	0	16	27	35	56	70	80	77	83	88	53.20%	136
AM	0	16	28	33	57	72	79	80	81	88	53.40%	128
WAM by count	0	16	29	34	58	73	77	81	82	87	53.70%	124
RMS	1	29	26	36	65	76	86	85	92	92	58.80%	96
Median of chosen descriptors	0	18	61	86	83	86	94	94	94	92	70.80%	113
NB	44	48	47	84	96	81	64	70	98	96	72.80%	1

Table 4.13: The false positive rate of combinations, sorted from lowest to highest. The position of the first erroneous duplicate pair was included for reference.

WAM. When weighing using the second manual set of weights, it only had a 4.8% error rate, which is clearly the best achieved in this thesis.

### 4.3.5 Truncated mean

The truncated mean is a small improvement in results over the arithmetic mean when run on all the descriptors, with a difference of only two more correctly identified results. This is not a statistically significant difference from the arithmetic mean.

When run on the chosen set of descriptors, it performs worse than the arithmetic mean, with a failure rate of 9.3%. Discarding values therefore seems to be a bad idea when working with a carefully chosen set of descriptors.

### 4.3.6 Median

Of the combination methods working on the reduced set of descriptors, the median performs worst of all. It is simply not a good indicator of central tendency for this use, with a failure rate of 70.8%.

### 4.3.7 Naïve Bayes

Results-wise, the naïve Bayes approach does not perform as well as its competitors. Even when considering only the top descriptors, it had a failure rate of 32.8%, which is worse than most other methods. When used on all descriptors, it has a failure rate of 72.8%.

It does have one redeeming feature, that can come in handy in certain cases. By changing the value of  $P(M)$  manually, you can adjust how much a negative or positive probability affect the final score, and this allows you to produce a smaller number of high-probability matches that are almost guaranteed to be correct. This is a trade-off that will limit the number of false positives, but it will dramatically increase the number of false negatives as well.

Because of the results seen, it is likely that the assumption that the score given by a descriptor can be interpreted as a measure of probability is false.

### 4.3.8 Original FDMF

The results from the original FDMF package was easy to reproduce using the new architecture. FDMF does not provide any ranking of its results, and simply lists them in the order the files are encountered by the fingerprinter. The same conditions were applied here.

FDMF found 627 duplicate pairs, of which 593 were verified to be correct. This means that compared to the best combined solutions presented, FDMF results in 38% fewer correct results.

Because of the lack of a sorting algorithm in FDMF, a simple arithmetic mean was applied to the scores given by the three descriptors to give it a ranked score. When using this sorting, FDMF had only 1 erroneous match in the first 500 scores.

Descriptor	Weight
fdmf_0	0.6
fdmf_1	0.39
fdmf_2	0.52
libfooid	0.64
mfcc_avg	0.57
mfcc_delta_avg	0.57
mfcc_f1_avg	0.53
mfcc_f1_delta_avg	0.3

Table 4.14: Weights by performance for various descriptors.

## 4.4 Examples

Some examples of potential duplicate pairs were picked from our collection, so the properties of the various methods in different situations can be examined. Only a subset of the averaging methods are presented here for the sake of readability. Only the top 8 best performing descriptors are considered. For reference, the weights used are found in table 4.14.

(a) Scores		(b) Combined scores		
Descriptor	Score	Method	Score	SD
fdmf_0	0.63	Truncated mean	0.8307	0.1894
fdmf_1	0.81	Weighted arithmetic mean	0.7790	0.5147
fdmf_2	0.88	Arithmetic mean	0.8054	0.1620
libfooid	0.85	Root-mean square	0.8216	0.1628
mfcc_avg	0.98	Bayesian	0.9999	-
mfcc_delta_avg	0.84	Median	0.8436	0.1665
mfcc_f1_avg	0.98	Weighted RMS	0.8313	0.2272
mfcc_f1_delta_avg	0.48			

Table 4.15: Scores for “Dido - Do You Have A Little Time” in MP3 and WMA encodings. The system correctly scores this as a good match.

## 4.5 Improving search speed

The analyzed system is not in any way optimized for quick searches. In fact, the calculations required to perform a search that matches all files against each other takes several hours on a fairly up to date computer. An indexing scheme is needed.

A very simplistic approach is to use any one of the scalar descriptors. A small script (see B.4.5) was created to assess the viability of these descriptors automatically. For each descriptor the range of the distances between verified duplicate pairs was found. Only fingerprints within the range  $[v - \delta, v + \delta]$  were examined as candidate matches, where  $v$  is the scalar descriptor value for the query file, and  $\delta$  is two times the average distance between verified duplicate pair fingerprints for that descriptor.

(a) Scores		(b) Combined scores		
Descriptor	Score	Method	Score	SD
fdmf_0	0	Truncated mean	0.4158	0.2909
fdmf_1	0.10	Weighted arithmetic mean	0.4201	0.7753
fdmf_2	0	Arithmetic mean	0.4248	0.2480
libfooid	0.65	Root-mean square	0.5358	0.2172
mfcc_avg	0.63	Bayesian	0.5035	-
mfcc_delta_avg	0.81	Median	0.6379	0.2243
mfcc_f1_avg	0.41	Weighted RMS	0.5260	0.3024
mfcc_f1_delta_avg	0.81			

Table 4.16: Scores for two separate recordings of “Billie Holiday - Solitude”. The songs are so different they were not considered a match when verified by a human listener. Most combinations score this as a low-probability match, yet still sees the similarity in the songs.

(a) Scores		(b) Combined scores		
Descriptor	Score	Method	Score	SD
fdmf_0	0	Truncated mean	0.21051	0.3429
fdmf_1	0	Weighted arithmetic mean	0.2705	0.5862
fdmf_2	0	Arithmetic mean	0.2932	0.2460
libfooid	0	Root-mean square	0.4353	0.1701
mfcc_avg	0.53	Bayesian	0.5425	-
mfcc_delta_avg	0.73	Median	0.6315	0.1362
mfcc_f1_avg	0.30	Weighted RMS	0.4030	0.1618
mfcc_f1_delta_avg	0.77			

Table 4.17: Scores for “David Bowie - Ricochet” and “Vamp - Svin på skog”. Despite some similarities reported by the mfcc-series of descriptors, fdmf and libfooid correctly sees no match in these two songs. Yet, the both the median and the naïve Bayes classifier reports higher scores than wanted.

The length descriptor is the clear winner, as it results in both the lowest number of candidate files on average, and finds a higher percentage of the verified duplicate pairs than any other descriptor.

When trying multiple combinations for descriptors, one expects a sharp drop in the number of combinations that has to be tried, and an increase in the number of verified matches that was removed. An attempt was made were combinations of one or more scalar descriptors were used as an index, to see that no combination of descriptors produced any kind of “magic index”.

Descriptor	Avg. # of comp.	% found	Avg. distance
Centroid	3538	84.71%	370
Track length	2487	87.71%	18.2
Mean/square ratio	3595	86.27%	0.017
Steepness	3465	84.72%	125
Zerocrossings	3463	85.7%	9100
Zerocrossings w/ guard band	3520	85.37%	5049

Table 4.18: Performance characteristics of the six scalar descriptors if they were to be used as an indexing scheme, when searching for all the verified duplicate pair files.

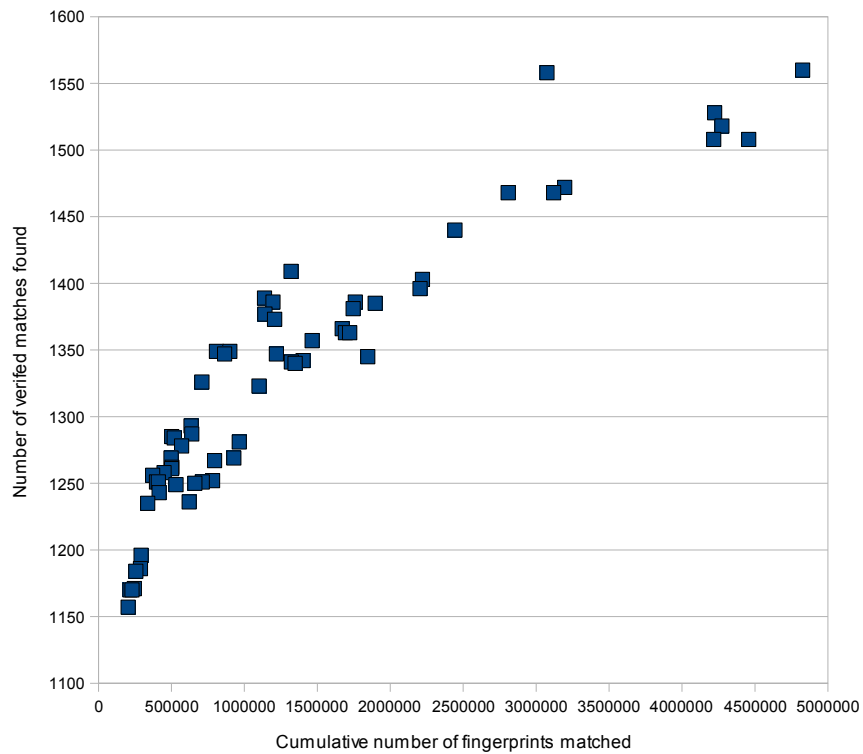


Figure 4.8: The number of verified duplicate pairs that were removed versus the number of files that would be included in the search when sequentially searching for all duplicate pairs. An optimum solution would be ranked in the upper left corner.



# Chapter 5

## Conclusion

Large music collections are now more common than ever before. Yet, search technology for music is still in its infancy. Audio fingerprinting is one method that allows searching for music.

In this thesis several audio fingerprinting solutions are combined into a single solution to determine if such a combination can yield better results than any of the solutions can separately. The solution is used to find duplicate music files in a personal collection.

A real-life music collection with 9536 music files was used. Audio fingerprints for all files were created using three software suites, namely fdmf, libFooID and Gunderson's master thesis. By combining the likelihood from multiple system that two files match, recognition rates were increased.

### 5.1 Results

This thesis finds that using different averaging methods to produce a more stable audio fingerprint works satisfactorily, and can be used to effectively combine different systems.

Among the methods of combination that was tested, the weighted root-mean square most effectively ranked the results in a satisfying manner. It was notably better than the other approaches tried. The WRMS produced 61% more correct matches than the original FDMF solution, and 49% more correct matches than libFooID.

The individual descriptors can be ranked according to performance in the following order, from best to worst: libfooid, fdmf\_0, mfcc\_avg, mfcc\_delta\_avg, mfcc\_f1\_avg, fdmf\_2, fdmf\_1 and mfcc\_f1\_delta\_avg. The remaining descriptors all ranked considerably worse than the ones listed here, and were found to be too untrustworthy to be included in calculations.

### 5.2 Evaluation

The implementation of this thesis went smoothly. However, there are certain aspects that could be attempted for improving the proposed system.

The most important improvement is likely to be a strictly controlled music collection. Constructing a music collection from original CDs, and encoding the

audio with lossless and lossy encodings, and introducing a certain amount of noise into the recordings would be helpful. With such a collection it would also be possible to create precision vs. recall curves for the system.

Writing the glue code using PHP and MySQL seems to have been a mistake, as it was chosen because of the author's familiarity with the language. A better choice would have been to use C and a binary database system, as putting the data into MySQL proved to be slow. Data processing in PHP is very slow compared to native languages, such as C.

An approach that tries to use averages on all possible combinations of descriptors could potentially improve the solution even further. This approach could be used to determine whether the results from individual descriptors have high degrees of overlap, and could potentially exclude some descriptors from the recommended set.

### 5.3 Further work

The combination of software used in this thesis has not been optimized for speed, but for being malleable enough for conducting the tests required to provide better results. As such, speed has not been a priority, and the system therefore is very slow, even on a relatively small collection such as the one we tested on.

Most of the fingerprinting engines are written in C or C++, and it will be relatively easy to combine these into a modular system that allows plugging in new fingerprinting engines easily. By writing the fingerprint comparer in C or C++, instead of an interpreted language such as PHP, the system could easily be made to search for a fingerprint in an already fingerprinted collection in  $O(n)$  time.

# Bibliography

- [1] Michael Fink; Michele Covell; Shumeet Baluja. Social- and interactive-television applications based on real-time ambient-audio identification. 2007. Accessed on the 24th of March 2008.
- [2] S. Baluja and M. Covell. Audio fingerprinting: Combining computer vision & data stream processing. *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*, 2:II-213–II-216, 15-20 April 2007.
- [3] R. H. Blessing. Outlier treatment in data merging. *Journal of Applied Crystallography*, 30(4):421–426, 1997.
- [4] Mark Bocko. Music file compressed 1,000 times smaller than mp3, 2008. Accessed on April 2nd 2008.
- [5] K. Brandenburg. MP3 and AAC explained. *Proceedings of AES 17th International Conference, XP008004053*, pages 99–110.
- [6] K. Brandenburg, J. Herre, J.D. Johnston, Y. Mahieux, and E. Schroeder. ASPEC: Adaptive spectral entropy coding of high quality music signals. *Proc. 90th Conv. Aud. Eng. Soc*, 1991.
- [7] The Echo Nest Corporation. The echo nest. Accessed on the 21st of May 2008.
- [8] M. Covell and S. Baluja. Known-audio detection using waveprint: Spectrogram fingerprinting by wavelet hashing. *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*, 1:I-237–I-240, 15-20 April 2007.
- [9] J. Stephen Downie. Music information retrieval (chapter 7). *Annual Review of Information Science and Technology 37*, pages 295–340, 2003.
- [10] T. Eriksson and H.G. Kang. Pitch quantization in low bit-rate speech coding. *Acoustics, Speech, and Signal Processing, 1999. ICASSP'99. Proceedings., 1999 IEEE International Conference on*, 1, 1999.
- [11] Xiph.org foundation. Vorbis i specification, 2004. Accessed on the 21st of May 2008.
- [12] Gracenote. Gracenote: Musicid. Accessed on the 21st of May 2008.
- [13] Paul Graham. *Hackers & Painters*. O'Reilly, 2004. ISBN 059600662-4.

- [14] Amara Graps. An introduction to wavelets. *IEEE Comput. Sci. Eng.*, 2(2):50–61, 1995.
- [15] Steinar H. Gunderson. Musical descriptors: An assessment of psychoacoustical models in the presence of lossy compression. Master’s thesis, 2007.
- [16] Jaap Haitsma, Ton Kalker, and Job Oostveen. Robust audio hashing for content identification. In *International Workshop on Content-Based Multimedia Indexing (CBMI’01)*, Brescia, Italy, September 2001. Accessed on the 1st of February 2008.
- [17] Richard Jones. Audio fingerprinting for clean metadata, 2007. Accessed on the 21st of May 2008.
- [18] Yan Ke, D. Hoiem, and R. Sukthankar. Computer vision for music identification. *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, 1:597–604 vol. 1, 20-25 June 2005.
- [19] Vegard Andreas Larsen. Determining audio fingerprint boundaries. Norwegian University of Science and Technology, 2007.
- [20] Alison Latham. *The Oxford Companion to Music*. Oxford University Press, 2002. ISBN 0198662122.
- [21] D.D. Lewis. Naive (Bayes) at forty: The independence assumption in information retrieval. *Proceedings of ECML-98, 10th European Conference on Machine Learning*, pages 4–15, 1998.
- [22] Ricardo Miotto and Nicola Orio. A methodology for the segmentation and identification of music works. In *ISMIR 2007 - Proceedings of the 8th International Conference on Music Information Retrieval, Vienna, Austria, September 23-27*, pages 273–278. Österreichische Computer Gesellschaft, 2007.
- [23] H. Nyquist. Certain topics in telegraph transmission theory. *Proceedings of the IEEE*, 90(2):280–305, 2002.
- [24] Justis og politidepartementet. *Lov om patenter*. 2007. ISBN 82-504-1193-5. Accessed on the 4th of February 2008.
- [25] IFLA Study Group on the Functional Requirements for Bibliographic Records. *Functional Requirements for Bibliographic Records: Final Report*. UBCIM Publications-New Series. K.G.Saur, München. Accessed on the 21st of May 2008.
- [26] L. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web, 1998.
- [27] Gian-Carlo Pascutto. foosic - the living music database, 2006. Accessed on the 3rd of February 2008.
- [28] Michaël Betser; Patrice Collen; Jean-Bernard Rault. Audio identification using sinusoidal modeling and application to jingle detection. In *ISMIR*

- 2007 - *Proceedings of the 8th International Conference on Music Information Retrieval, Vienna, Austria, September 23-27*, pages 139–142. Österreichische Computer Gesellschaft, 2007. Accessed on the 30th of March 2008.
- [29] Trent D. Stephens Rod R. Seeley, Philip Tate. *Essentials of anatomy and physiology*. McGraw-Hill, 6th edition edition, 2006. ISBN 007110746-0.
- [30] Kurt Rosenfeld. Merging music collections without redundancy, October 2005. Accessed on the 27th of April 2008.
- [31] Kurt Rosenfeld. Learning to recognize duplicate music files, May 2007. Accessed on the 27th of April 2008.
- [32] J.G. Shao-Jen Lim; Harris. Analog implementation of ratio spectrum computation. *Circuits and Systems, 1998. ISCAS '98. Proceedings of the 1998 IEEE International Symposium on*, 1:277–280 vol.1, 31 May-3 Jun 1998.
- [33] Joel Spolsky. *Joel on Software: And on Diverse and Occasionally Related Matters That Will Prove of Interest to Software Developers, Designers, and Managers, and to Those Who, Whether by Good Fortune or Ill Luck, Work with Them in Some Capacity*. Apress, 1 edition, 2004. ISBN 1590593898. 125-130 pp.
- [34] Hartmut Traunmuller. Analytical expressions for the tonotopic sensory scale. *The Journal of the Acoustical Society of America*, 88(1):97–100, 1990.
- [35] University of Jyväskylä. *A Matlab toolbox for music feature extraction from audio*, September 2007.
- [36] J.M. Valin and C. Montgomery. Improved noise weighting in CELP coding of speech-applying the Vorbis psychoacoustic model to Speex. *Audio Engineering Society Convention*, 2006.
- [37] Eric W Weisstein. Metric, 2004. Accessed on the 21st of May 2008.
- [38] Eric W Weisstein. Wave equation, 2004. Accessed on the 21st of May 2008.



# Index

- AAC, 10
- advertisement, 13
- amplitude, 5, 16
- architecture, 26, 28
- arithmetic mean, 31, 50, 55
- asymptotic, 15
- ATRAC, 10
- audio fingerprint, 11, 12
  
- band energy, 17
- Bark scale, 18
- Bayes, 32, 55
- bibliographic facet, 11
- bibliographic records, 22
- bin, 15
- bit rate, 10, 19
- BPM, 14
  
- central tendency, 29
- centroid, 17–19, 46, 47, 49
- chromatic scale, 8
- chunk metric, 17
- cluster, 23
- cochlea, 7
- commutative, 29
- compactness, 12
- compression, 10
- computer vision, 16
- copyright, 16
- correlation coefficient, 18
- cycle, 5
  
- decoding, 26
- destructiveness, 12
- diatonic scale, 8
- duplicate, 23
  
- edit distance, 17, 18
- energy spectrum, 17
- equivalence, 21, 26
- Euclidean distance, 20, 27
- evaluation, 39
  
- expressions, 22
  
- F1CC, 18, 20, 49
- facet, 11
- fair use, 16
- false positive, 39
- fdmf, 16, 17, 26, 28, 29, 41, 45
- feature extraction, 11
- file-sharing, 13
- filter bank, 17
- fingerprint, 11
- FLAC, 10
- floor-1 cepstral coefficients, 18, 20, 49
- Fourier transform, 15
- FRBR, 22
- frequency, 5, 16
- frequency spectrum, 14, 18
- FT, 15
  
- genre, 11
- geometric mean, 34
- GNU C, 17
- Google, 13, 16
- GPL, 16
- Gracenote, 13
- granularity, 12
- guard band, 19, 48
- Gunderson's descriptors, 18, 27, 29
  
- Hamming distance, 17, 18, 26
- harmonic mean, 34
- harmonics, 8, 11
- hash, 12
- hidden Markov models, 16
  
- intellectual property, 16
- IP, 16
- item, 22
  
- last.fm, 16
- length, 18, 19, 46
- libFingerprint, 16

- libFooID, 16, 18, 27, 29, 45, 47
- lossless, 21
- lossless compression, 10
- lossy compression, 10, 20
- lyric, 11
  
- Mahalanobis distance, 20
- manifestation, 22
- Matlab, 17
- MD5, 12, 28
- mean/square ratio, 18, 19, 47
- median, 32, 55
- mel frequency cepstral coefficients, 18, 19, 48
- mel scale, 19
- metadata, 19
- MFCC, 17–19, 48, 49
- MIDI, 9
- MIRtoolbox, 17
- MP3, 10, 19, 21, 41
- music collection, 21
- MySQL, 28
  
- naïve Bayes, 32, 55
- noise, 10
- normalize, 18
- Nyquist-Shannon sampling theorem, 9
  
- octave, 8
- Ogg Vorbis, 10, 20, 21
- Organ of Corti, 7
- ossicle, 7
- overtone, 8
  
- patent, 16
- peak amplitude, 5
- perceived frequency, 11
- period, 5
- Perl, 17
- phase, 7, 16
- pitch, 8, 11, 17
- play time, 18
- power spectrum, 17
- precision, 39
- probability, 33
  
- QBH, 11
- query by humming, 11
  
- rate of zero crossings, 18, 19, 48
- ratio spectrum, 17
  
- recall, 39
- reliability, 12
- resampling, 18
- rhythm, 14
- RMS, 31, 50
- robustness, 12
- root mean square, 31
- root-mean square, 50
  
- sample, 8
- sampling theorem, 9
- Schmitt triggering, 18, 19
- sheet music, 9
- sinusoidal model, 16
- skewness, 17
- software patent, 16
- song length, 18, 19, 46
- SonyEricsson, 13
- sound pressure level, 7
- spatial placement, 12
- spectral, 14
- spectral centroid, 18
- speech recognition, 10
- standard deviation, 31
- statistical measure, 31
- steepness, 18, 19, 48, 49
- STFT, 15, 16
  
- tempo, 11
- threshold, 18, 26, 27
- timbre, 11
- time-based representation, 14
- token, 23
- TrackID, 13
- transcript, 10
- Traunmüller's formula, 18
- Truncated mean, 32
- truncated mean, 31, 55
- twist spectrum, 17
- tympanic membrane, 7
  
- variable bit rate, 19
- verification, 25
- volume, 7
  
- WAM, 32, 50
- waveform, 14
- wavelet, 14–16
- wavelet transform, 15
- Waveprint, 16



weighted arithmetic mean, 32, 50  
weighted root-mean square, 32, 50  
weights, 34  
Windows Media Audio, 10  
Windows Media Audio Lossless, 10  
WMA, 10, 21, 41  
WMA lossless, 10  
work, 22  
WRMS, 32, 50  
  
zero crossings rate, 18, 19, 48



# Appendix A

## Software used

**GNU GCC** GNU's Compiler Collection (GCC) is one of the most widely used compilers for C code. Available from <http://www.gnu.org/>.

**LyX** A free software document processor that uses *L<sup>A</sup>T<sub>E</sub>X*. Available from <http://www.lyx.org/>.

**Eclipse** Programming environment, originally designed for Java development. Available from <http://www.eclipse.org/>.

**PHP** Scripting language originally designed for creating dynamic web pages. Available from <http://www.php.net/>.

**PHPEclipse** Plugin for Eclipse that simplifies development in PHP. Available from <http://www.phpeclipse.de/>.

**phpMyAdmin** MySQL administration interface written in PHP. Available from <http://www.phpmyadmin.net/>.

**MySQL** Open source DBMS. Available from <http://www.mysql.com/>.

**Ubuntu** Free GNU/Linux-based operating system . Available from <http://www.ubuntu.org/>.



# Appendix B

## Source code

The source code presented in this appendix is licensed under the same license as the software it modifies.

### B.1 fdmf

The original source code and the modifications made to fdmf are licensed under the GPL.

#### B.1.1 db.h

Database connectivity headers.

```
1 #ifndef DB_H_
2 #define DB_H_
3
4 #include <stdio.h>
5 #include <fcntl.h>
6 #include <unistd.h>
7 #include <sys/types.h>
8 #include <sys/stat.h>
9 #include <pwd.h>
10 #include <string.h>
11 #include <gdbm.h>
12 #include <assert.h>
13 #include </usr/include/mysql/mysql.h>
14
15 int db_setup(void);
16 void db_close(void);
17 int db_find_file(char *, int);
18 int db_insert_file(char *, int, int);
19 void db_insert_fingerprint(int, int, char *, int);
20 void db_add_score(int, int, int, int, int);
21
22 #endif /*DB_H_*/
```

#### B.1.2 db.c

Database connectivity and functions for inserting data.

```
1 #include "vector_pairs.h"
2
3 #define SQL_INSERT_SCORE "INSERT INTO tblScore (ixFile1, ixFile2, ixType
, dblScore) VALUES (?, ?, ?, 1 - (? / ?))"
4 #define SQL_SELECT_FINGERPRINT "SELECT ixFile FROM tblFingerprint WHERE
ixFile = ? AND ixType = ?"
```

```

5 #define SQL_INSERT_FINGERPRINT "INSERT INTO tblFingerprint (ixFile ,
   ixType, sFingerprint) VALUES (?, ?, ?)"
6 #define SQL_SELECT_FILE "SELECT ixFile FROM tblFile WHERE md5Path = MD5
   (?)"
7 #define SQL_INSERT_FILE "INSERT INTO tblFile(sPath, md5Path, cbFile)
   VALUES (?, MD5(?), ?)"
8
9 MYSQL *db_conn;
10 MYSQL_STMT *sqlInsertFile;
11 MYSQL_STMT *sqlSelectFingerprint;
12 MYSQL_STMT *sqlInsertFingerprint;
13 MYSQL_STMT *sqlSelectFile;
14 MYSQL_STMT *sqlInsertScore;
15
16 /*
17  * Lets create a local cache of every ixFile value we get from
18  * MySQL, so we won't have to query for it every time.
19  */
20 int *file_cache;
21 int file_cache_ready = 0;
22
23 int db_setup()
24 {
25     char *server      = "localhost";
26     char *user        = "master";
27     char *password    = "7aBPK8huJQz3TNQS";
28     char *database    = "master";
29
30     db_conn = mysql_init(NULL);
31     if (!mysql_real_connect(db_conn, server, user, password, database,
32         0, NULL, 0))
33     {
34         fprintf(stderr, "%s\n", mysql_error(db_conn));
35         exit(0);
36     }
37     printf("database connected\n");
38
39     sqlInsertScore = mysql_stmt_init(db_conn);
40     mysql_stmt_prepare(sqlInsertScore, SQL_INSERT_SCORE, strlen(
41         SQL_INSERT_SCORE));
42
43     sqlSelectFingerprint = mysql_stmt_init(db_conn);
44     mysql_stmt_prepare(sqlSelectFingerprint, SQL_SELECT_FINGERPRINT,
45         strlen(SQL_SELECT_FINGERPRINT));
46
47     sqlInsertFingerprint = mysql_stmt_init(db_conn);
48     mysql_stmt_prepare(sqlInsertFingerprint, SQL_INSERT_FINGERPRINT,
49         strlen(SQL_INSERT_FINGERPRINT));
50
51     sqlSelectFile = mysql_stmt_init(db_conn);
52     mysql_stmt_prepare(sqlSelectFile, SQL_SELECT_FILE, strlen(
53         SQL_SELECT_FILE));
54
55     sqlInsertFile = mysql_stmt_init(db_conn);
56     mysql_stmt_prepare(sqlInsertFile, SQL_INSERT_FILE, strlen(
57         SQL_INSERT_FILE));
58
59     return(0);
60 }
61
62 void db_close()
63 {
64     mysql_close(db_conn);
65
66     mysql_stmt_close(sqlInsertScore);
67     mysql_stmt_close(sqlSelectFingerprint);
68     mysql_stmt_close(sqlInsertFingerprint);
69     mysql_stmt_close(sqlSelectFile);
70     mysql_stmt_close(sqlInsertFile);
71 }
72
73 void init_file_cache(int size)
74 {
75     int i;

```

```

70     if (!file_cache_ready)
71     {
72         file_cache = (int*) malloc(sizeof(int) * size);
73         for (i = 0; i < size; i++)
74             file_cache[i] = -1;
75         file_cache_ready = 1;
76     }
77 }
78
79 int db_find_file(char *file, int cFileNo)
80 {
81     int ixFile = -1;
82     MYSQL_BIND bind[1];
83     unsigned long length;
84     my_bool is_null;
85     my_bool error;
86
87     if (file_cache[cFileNo] != -1)
88         return file_cache[cFileNo];
89
90     memset(bind, 0, sizeof(bind));
91
92     length = strlen(file);
93
94     bind[0].buffer_type = MYSQL_TYPE_STRING;
95     bind[0].buffer = file;
96     bind[0].buffer_length = strlen(file);
97     bind[0].is_null = 0;
98     bind[0].length = &length;
99
100    if (mysql_stmt_bind_param(sqlSelectFile, bind))
101    {
102        fprintf(stderr, "mysql_stmt_bind_param() failed in db_find_file
103                ()\n");
104        fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFile));
105        exit(0);
106    }
107    if (mysql_stmt_execute(sqlSelectFile))
108    {
109        fprintf(stderr, "mysql_stmt_execute() failed in db_find_file()\n
110                ");
111        fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFile));
112        exit(0);
113    }
114
115    memset(bind, 0, sizeof(bind));
116
117    bind[0].buffer_type = MYSQL_TYPE_LONG;
118    bind[0].buffer = (char *)&ixFile;
119    bind[0].buffer_length = 4;
120    bind[0].length = &length;
121    bind[0].is_null = &is_null;
122    bind[0].error = &error;
123
124    if (mysql_stmt_bind_result(sqlSelectFile, bind))
125    {
126        fprintf(stderr, "mysql_stmt_bind_result() failed in db_find_file
127                ()\n");
128        fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFile));
129        exit(0);
130    }
131
132    if (mysql_stmt_store_result(sqlSelectFile))
133    {
134        fprintf(stderr, "mysql_stmt_store_result() failed in
135                db_find_file()\n");
136        fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFile));
137        exit(0);
138    }
139
140    if (mysql_stmt_num_rows(sqlSelectFile) == 0)
141        return -1;
142
143    if (mysql_stmt_fetch(sqlSelectFile))

```

```

140     {
141         fprintf(stderr, "mysql_stmt_fetch() failed in db_find_file()\n");
142         ;
143         fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFile));
144         exit(0);
145     }
146     if (ixFile != -1)
147         file_cache[cFileNo] = ixFile;
148     /* fprintf(stderr, "%s %d %ld %d\n", file, ixFile, length, error);
149        */
150     return ixFile;
151 }
152 /* returns ixFile, -1 if error */
153 int db_insert_file(char *file, int cFileNo, int cFiles)
154 {
155     struct stat fileBuf;
156     int ixFile = -1;
157     long cbFile = 0;
158     MYSQL_BIND bind[3];
159     unsigned long str_len;
160
161     init_file_cache(cFiles);
162
163     ixFile = db_find_file(file, cFileNo);
164     if (ixFile != -1)
165         return ixFile;
166
167     memset(bind, 0, sizeof(bind));
168
169     stat(file, &fileBuf);
170     cbFile = fileBuf.st_size;
171
172     str_len = strlen(file);
173
174     bind[0].buffer_type = MYSQL_TYPE_STRING;
175     bind[0].buffer = (char *)file;
176     bind[0].buffer_length = strlen(file);
177     bind[0].length = &str_len;
178     bind[0].is_null = 0;
179
180     // yes, we're duping info, let mysql handle MD5 hashing
181     bind[1].buffer_type = MYSQL_TYPE_STRING;
182     bind[1].buffer = (char *)file;
183     bind[1].buffer_length = strlen(file);
184     bind[1].length = &str_len;
185     bind[1].is_null = 0;
186
187     bind[2].buffer = (char *)&cbFile;
188     bind[2].buffer_type = MYSQL_TYPE_LONG;
189     bind[2].is_null = 0;
190     bind[2].length = 0;
191
192     if (mysql_stmt_bind_param(sqlInsertFile, bind))
193     {
194         fprintf(stderr, "mysql_stmt_bind_param() failed in
195             db_insert_file()\n");
196         fprintf(stderr, "%s\n", mysql_stmt_error(sqlInsertFile));
197         exit(0);
198     }
199
200     if (mysql_stmt_execute(sqlInsertFile))
201     {
202         fprintf(stderr, "mysql_stmt_execute failed in db_insert_file()!\n");
203         fprintf(stderr, "%s\n", mysql_stmt_error(sqlInsertFile));
204         exit(0);
205     }
206
207     /* return db_find_file(file); */
208     return mysql_insert_id(db_conn);
209 }

```



```

210 int db_has_fingerprint(int ixFile, int ixType)
211 {
212     int ixReturn = -1;
213     MYSQL_BIND bind[2];
214     unsigned long length;
215     my_bool is_null;
216     my_bool error;
217
218     memset(bind, 0, sizeof(bind));
219
220     bind[0].buffer_type = MYSQL_TYPE_LONG;
221     bind[0].buffer = (char *)&ixFile;
222     bind[0].is_null = 0;
223     bind[0].length = 0;
224
225     bind[0].buffer_type = MYSQL_TYPE_LONG;
226     bind[0].buffer = (char *)&ixType;
227     bind[0].is_null = 0;
228     bind[0].length = 0;
229
230     if (mysql_stmt_bind_param(sqlSelectFingerprint, bind))
231     {
232         fprintf(stderr, "mysql_stmt_bind_param() failed in
233             db_has_fingerprint()\n");
234         fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFingerprint));
235         exit(0);
236     }
237     if (mysql_stmt_execute(sqlSelectFingerprint))
238     {
239         fprintf(stderr, "mysql_stmt_execute() failed in
240             db_has_fingerprint()\n");
241         fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFingerprint));
242         exit(0);
243     }
244     memset(bind, 0, sizeof(bind));
245
246     bind[0].buffer_type = MYSQL_TYPE_LONG;
247     bind[0].buffer = (char *)&ixReturn;
248     bind[0].buffer_length = 4;
249     bind[0].length = &length;
250     bind[0].is_null = &is_null;
251     bind[0].error = &error;
252
253     if (mysql_stmt_bind_result(sqlSelectFingerprint, bind))
254     {
255         fprintf(stderr, "mysql_stmt_bind_result() failed in
256             db_has_fingerprint()\n");
257         fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFingerprint));
258         exit(0);
259     }
260     if (mysql_stmt_store_result(sqlSelectFingerprint))
261     {
262         fprintf(stderr, "mysql_stmt_store_result() failed in
263             db_has_fingerprint()\n");
264         fprintf(stderr, "%s\n", mysql_stmt_error(sqlSelectFingerprint));
265         exit(0);
266     }
267     if (mysql_stmt_num_rows(sqlSelectFingerprint) == 0)
268         return 0;
269     else
270         return 1;
271 }
272 void db_insert_fingerprint(int ixFile, int ixType, char *fingerprint,
273     int fingerprint_length)
274 {
275     char *sqlInsert;

```

```

275     char *sFingerprint;
276
277     if (db_has_fingerprint(ixFile, ixType))
278         return;
279
280     sqlInsert = (char *) malloc(2048);
281     sFingerprint = (char *) malloc(2048);
282     mysql_real_escape_string(db_conn, sFingerprint, fingerprint,
283                             fingerprint_length);
284
285     sprintf(sqlInsert, "INSERT INTO tblFingerprint (ixFile, ixType,
286                  sFingerprint) VALUES (%d, %d, '%s')", ixFile, ixType,
287                  sFingerprint);
288     mysql_query(db_conn, sqlInsert);
289 }
290
291 void db_add_score(int ixFile1, int ixFile2, int ixType, int bDistance,
292                 int bThreshold)
293 {
294     MYSQL_BIND bind[5];
295     int i;
296
297     memset(bind, 0, sizeof(bind));
298
299     for (i = 0; i < 5; i++)
300     {
301         bind[i].buffer_type = MYSQL_TYPE_LONG;
302         bind[i].is_null = 0;
303         bind[i].length = 0;
304     }
305
306     if (ixFile1 < ixFile2)
307     {
308         bind[0].buffer = (char *)&ixFile1;
309         bind[1].buffer = (char *)&ixFile2;
310     }
311     else
312     {
313         bind[0].buffer = (char *)&ixFile2;
314         bind[1].buffer = (char *)&ixFile1;
315     }
316     bind[2].buffer = (char *)&ixType;
317     bind[3].buffer = (char *)&bDistance;
318     bind[4].buffer = (char *)&bThreshold;
319
320     if (mysql_stmt_bind_param(sqlInsertScore, bind))
321     {
322         fprintf(stderr, "mysql_stmt_bind_param failed in db_add_score()
323                 !\n");
324         fprintf(stderr, "%s\n", mysql_stmt_error(sqlInsertScore));
325         exit(0);
326     }
327
328     if (mysql_stmt_execute(sqlInsertScore))
329     {
330         /* failing this insert is perfectly legitimate, as this comparison
331            may already
332            * exist somewhere in the database, so we won't do any error
333            reporting.
334            */
335         /*
336         fprintf(stderr, "mysql_stmt_execute failed in db_add_score()!\n
337                ");
338         fprintf(stderr, "%s\n", mysql_stmt_error(sqlInsertScore));
339         exit();
340         */
341     }
342 }

```

### B.1.3 run\_tests.c

Modifications for inserting music files into the database, along with potential duplicate pairs.

```

1 #include "vector_pairs.h"
2
3 int run_tests(char *dblk, int *hashes, char **names,
4             int vecs, int *thres, char delim) {
5     /* compare all pairs of vectors for possible matches */
6     int i, j, k, bitcount[256], type1_err=0, type2_err=0;
7 #ifdef INSERT_INTO_DB
8     int ixFile = 0;
9     int ixFile2 = 0;
10    char *temp;
11 #endif
12    char terminator;
13    struct stat stat_results;
14    if (delim == '\\0')
15    {
16        terminator = '\\0';
17    }
18    else
19    {
20        terminator = '\\n';
21    }
22    setup_bitcount_tbl(bitcount);
23    for (i = 0; i < vecs; i++)
24    {
25        /* code for mysql database insertion */
26
27 #ifdef INSERT_INTO_DB
28        ixFile = db_insert_file(names[i], i, vecs);
29        temp = (char *) malloc(VECTOR_BYTES);
30        strncpy(temp, dblk + i * MULTI_VEC_LEN + 0 * VECTOR_BYTES,
31              VECTOR_BYTES);
32        db_insert_fingerprint(ixFile, 1, temp, VECTOR_BYTES); /* fdmf_0
33        */
34        strncpy(temp, dblk + i * MULTI_VEC_LEN + 1 * VECTOR_BYTES,
35              VECTOR_BYTES);
36        db_insert_fingerprint(ixFile, 2, temp, VECTOR_BYTES); /* fdmf_1
37        */
38        strncpy(temp, dblk + i * MULTI_VEC_LEN + 2 * VECTOR_BYTES,
39              VECTOR_BYTES);
40        db_insert_fingerprint(ixFile, 3, temp, VECTOR_BYTES); /* fdmf_2
41        */
42        free(temp);
43 #endif
44
45    for (j = i + 1; j < vecs; j++) {
46 #ifdef INSERT_INTO_DB
47        ixFile2 = db_insert_file(names[j], j, vecs);
48 #endif
49
50        int distance[NUM_TESTS], score, test;
51        char *ptr_a = dblk + i * MULTI_VEC_LEN;
52        char *ptr_b = dblk + j * MULTI_VEC_LEN;
53        int basenames_match, contents_similar;
54        /* calculate distance between two vectors for each test*/
55        for (score = test = 0; test < NUM_TESTS; test++) {
56            for (distance[test] = k = 0; distance[test] <= thres[
57                test] && k < VECTOR_BYTES; k++) {
58                unsigned char c = *(ptr_a + k) ^ *(ptr_b + k);
59                distance[test] += bitcount[c];
60            }
61            ptr_a += VECTOR_BYTES;
62            ptr_b += VECTOR_BYTES;
63
64            if (distance[test] < thres[test])
65            {
66                score++;
67 #ifdef INSERT_INTO_DB
68                db_add_score(ixFile, ixFile2, test+1, distance[test]
69                    , thres[test]);
70 #endif
71            }
72        }
73    }
74 }

```

```

62     }
63     } /* post: score contains the number of tests that passed */
64     contents_similar = (score >= SCORE_THRESHOLD) ? 1 : 0;
65     basenames_match = (hashes[i] == hashes[j]) ? 2 : 0;
66 #ifndef IGNORE_GHOST_FILES
67     /* If either file of a pair does not exist, ignore the pair
68     */
69     if (contents_similar + basenames_match) {
70         if (stat(names[i], &stat_results)) continue;
71         if (stat(names[j], &stat_results)) continue;
72     }
73 #endif
74     switch (contents_similar + basenames_match) {
75     case 0: /* different contents, different basenames */
76         /* do nothing */
77         break;
78     case 1: /* similar contents, different basenames */
79         printf("%s%c%s%c", names[i], delim, names[j],
80             terminator);
81         type1_err++;
82         break;
83     case 2: /* different contents, same basenames */
84         /* printf("type 2: %s%c%s%c", names[i], delim,
85             names[j], terminator); */
86         type2_err++;
87         break;
88     case 3: /* similar contents, same basenames */
89         printf("%s%c%s%c", names[i], delim, names[j],
90             terminator);
91         break;
92     }
93     if (contents_similar + basenames_match == 1 ||
94         contents_similar + basenames_match == 3)
95     {
96         printf("%d%c%d%c%d%c", distance[0], delim, distance[1],
97             delim, distance[2], delim);
98     }
99     }
100     fprintf(stderr, "%d\t%d\t%d\t", thres[0], thres[1], thres[2]);
101     fprintf(stderr, "%d\t%d\n", type1_err, type2_err);
102     return(0);
103 }

```

### B.1.4 fdmf

Software that runs fdmf and fooid on all the files in a specific folder. Modified to run fooid.

```

1  #!/usr/bin/perl -w
2  # Kurt Rosenfeld 2004, 2005, 2006
3  # GPL
4  use strict;
5  use GDBM_File;
6  use Digest::MD5 qw(md5 md5_hex);
7  use Encode qw(encode_utf8);
8  use File::Copy;
9  use File::Basename;
10 use Cwd;
11
12 my $db_file = glob '~/.fdmf';
13 my $musicdir = $ARGV[0] or die "RTFM";
14 my $NUM_BANDS = 4;
15
16 my $SONICREDUCER = find_sr() or die "Can't find good sonic reducer";
17 # $SONICREDUCER = "cat /mnt/ramdisk/audio.raw | ".$SONICREDUCER;
18
19 die "Unable to run spline. Is it in your path?" if !spline /dev/null;
20
21 # We cache the table of filenames and their envelope spectra.

```

```

22 if (stat $db_file) { # keep the old db file just in case
23     File::Copy::copy "$db_file", "$db_file.old" or die "copy: $!";
24 }
25
26 tie my %DB, 'GDBM_File', glob("$db_file"), &GDBM_WRCREAT|&GDBM_SYNC,
    0640;
27 # If your perl/GDBM can't do SYNC mode, comment the line above and use
    this:
28 #tie my %DB, 'GDBM_File', glob("$db_file"), &GDBM_WRCREAT, 0640;
29
30 # Recurse through the directory, searching for mp3s
31 my @filelist = recurse_from_dir($musicdir);
32
33 # Compare that filelist to the already-cached files in the database.
34 my @uncached = grep((not $DB{"$_"}), @filelist);
35
36 my $filecount = 0;
37
38 print STDERR "Found ",$#filelist+1," files , ",$#uncached+1," uncached.\n
    ";
39
40 # this loop will add any uncached files to the database
41 FILE: for my $file (@uncached) {
42     open(FILE, $file) or die "Can't open $file: $!";
43     my $file_md5 = Digest::MD5->new->addfile(*FILE)->digest;
44     print STDERR $filecount++, "/", $#uncached, " waiting on $file\n";
45     # If an identical file is already indexed, reuse its summary.
46     foreach my $k (keys %DB) {
47         my $db_file_md5 = substr $DB{$k}, 100, 16;
48         if ($file_md5 eq $db_file_md5) {
49             print STDERR "Identical files (using cached summary):\n";
50             print STDERR "\t$file\n\t\tAND\n\t\t$k\n";
51             my $cached_summary = substr $DB{$k}, 0, 96;
52             $DB{$file} = $cached_summary . base_hash($file) . $file_md5;
53             next FILE;
54         }
55     }
56     my $summary;
57     next FILE if sonic_reduce($file, \$summary);
58     $DB{$file} = $summary . base_hash($file) . $file_md5;
59 } # post: %DB is complete
60
61 untie %DB;
62
63 #####
64 ##### THE END #####
65 #####
66
67 sub base_hash {
68     substr md5(encode_utf8(basename(shift))), 0, 4;
69 }
70
71
72 sub recurse_from_dir {
73     my ($dir) = @_;
74     my @filelist;
75     my $file_ptrn = '^[^.]'; # ignore dot files
76
77     unless (opendir(PARSEDIR, "$dir")) {
78         printf STDERR "skipping $dir because $!\n";
79         return;
80     }
81
82     foreach my $file (sort(readdir(PARSEDIR))) {
83         my $fullname = $dir . "/" . $file;
84         $fullname =~ s/\/+\/\//g;
85         if (-d $fullname && $file !~ /\^\.\.*\/) {
86             push @filelist, recurse_from_dir($fullname);
87         }
88         elsif ($file =~ ^/$file_ptrn/i && -r $fullname) {
89             $filelist[+ $#filelist] = $fullname;
90         }
91     }
92     closedir(PARSEDIR);

```

```

93     return @filelist;
94 }
95
96
97 sub sonic_reduce {
98     # ARGUMENT: filename of music file
99     # RETURN by REFERENCE: 768-bit string for fdmf database
100    # RETURN VALUE: 0 for success, nonzero for nonsuccess.
101    my $f = shift;
102    my $summary_ref = shift;
103    my @SR;
104    my @DECODE_CMD;
105    return -1 if decode_cmd($f, \@DECODE_CMD);
106
107    pipe PIPE1_OUT, PIPE1_IN;
108
109    my $pid1 = fork;
110    if ($pid1 == 0) { # we are child 1 (the decoder)
111        close(PIPE1_OUT) or die "$!";
112        open(STDOUT, ">&PIPE1_IN") or die "$!";
113        exec @DECODE_CMD or die "$!";
114    }
115    else { # we are still the parent (having had one child of two)
116        close(PIPE1_IN) or die; # child 1 will write on this
117        pipe PIPE2_OUT, PIPE2_IN;
118        # we will wait for the decode to be done
119        # the original fdmf did not do this, as it got input from
120        # stdin, that was produced in real time by the decoder
121        # we now need all that input for later software
122        waitpid $pid1, 0;
123        my $pid2 = fork;
124        if ($pid2 == 0) { # we are child 2 (the sonic_reducer process)
125            close(PIPE2_OUT) or die "$!";
126            open(STDIN, "<&PIPE1_OUT") or die "$!";
127            open(STDOUT, ">&PIPE2_IN") or die "$!";
128            exec $SONICREDUCER." /mnt/ramdisk/audio.raw" or die "$!";
129        }
130        else { # we are still the parent (having had both children)
131            close(PIPE2_IN) or die "$!";
132            @SR = <PIPE2_OUT>;
133        }
134        # waitpid $pid1, 0;
135        waitpid $pid2, 0;
136    }
137
138    if ($#SR != 767) {
139        print STDERR "sonic_reduce had trouble with $f. Corrupt audio
140            file?\n";
141        return -2;
142    }
143    my @e = @SR[0 .. 255]; # energy spectrum summary
144    my @r = @SR[256 .. 511]; # ratio (high/low) spectrum summary
145    my @t = @SR[512 .. 767]; # twist (odd/even ratio) spectrum summary
146    my $j = join("", quantize(@e), quantize(@r), quantize(@t));
147    $$summary_ref = pack("b*", $j);
148
149    system ("/home/vegard/Masteroppgave/Kode/fdmf/fooid", "/mnt/ramdisk/
150        audio.raw", $f);
151    my $md5 = md5_hex($f);
152    system ("/home/vegard/Masteroppgave/Kode/fdmf/descriptor /mnt/
153        ramdisk/audio.raw > /home/vegard/Masteroppgave/Kode/results/
154        descriptors/$md5");
155    system ("echo \"$md5\" >> /home/vegard/Masteroppgave/Kode/results/
156        descriptors/filelist.txt");
157
158    return 0;
159 }
160
161 sub decode_cmd {
162     # This routine looks at the filename extension of a music file.
163     # It returns the shell command for decoding it.
164     # These also work, if you don't want to use mplayer:
165     # @ $cmd_ref = ("mpg123", "-s", "-q", "$f");
166     # @ $cmd_ref = ("ogg123", "-d", "raw", "-f", "-", "$f");

```

```

162     my $f = shift;
163     my $cmd_ref = shift;
164     my $filename_extension = lc substr $f, -3;
165     if ($filename_extension =~ /mp3|ogg|m4a|wma|wav|\\.ra/) {
166         # @cmd_ref = ("mplayer", "-nortc", "-ao", "pcm",
167         #           "-aofile", "/dev/stdout", "$f");
168         # @cmd_ref = ("mplayer", "-ao", "pcm:file=/dev/stdout", "$f");
169         @cmd_ref = ("mplayer", "-really-quiet", "-msglevel", "all=1",
170                   "-nojoystick", "-nolirc", "-ao", "pcm:fast:file=/mnt/ramdisk",
171                   "/audio.raw", "$f");
172         return 0;
173     }
174     else {
175         print STDERR "Filename $f doesn't have an extension that we
176             handle.\n";
177         return -1;
178     }
179 }
180
181 sub find_sr {
182     my $sr_path;
183     $sr_path = dirname($0) . "/sonic_reducer"; # look in same dir as
184         fdmf
185     return $sr_path if -x $sr_path;
186     $sr_path = getcwd() . "/sonic_reducer"; # look in the current
187         directory
188     return $sr_path if -x $sr_path;
189     return 0; # we failed to find sonic_reducer
190 }
191
192 sub quantize {
193     # one bit quantize with median value as threshold
194     my @s;
195     my $median = median(@_);
196     foreach my $i (0 .. $#_) {
197         $s[$i] = ($_[ $i ] > $median) ? 1 : 0;
198     }
199     return @s;
200 }
201
202 sub median {
203     my $size;
204     my $median;
205     my @sorted = sort {$a <=> $b} @_;
206     $size = scalar @sorted;
207     if ($size % 2){
208         $median = $sorted[( $size - 1 ) / 2];
209     }
210     else {
211         $median = ($sorted[ $size / 2 ] + $sorted[( $size / 2 ) - 1]) / 2;
212     }
213 }

```

## B.2 libFooID

libFooID also uses db.h and db.c from fdmf. The source code is available under the GPL.

### B.2.1 foo.h

Headers.

```

1 #ifndef FOO_H
2 #define FOO_H
3
4 #include <assert.h>
5 #include <stdlib.h>
6 #include <string.h>

```

```

7 #include <stdio.h>
8 #include <math.h>
9 #include <limits.h>
10 #include <fooid.h>
11 #include <sys/stat.h>
12 #include "db.h"
13
14 #endif /*FOO_H*/

```

### B.2.2 foo.c

Console tool that reads WAVE files from disk, generates fingerprints (using libFooID), and writes them to the database.

```

1 #include "foo.h"
2
3 #define SAMPLE_RATE 44100
4 #define CHANNELS 2
5 #define N_BLOCKS 1024
6
7 t_fooid *fooid;
8
9 int fsize(char *name)
10 {
11     struct stat stbuf;
12
13     if (stat(name, &stbuf) == -1)
14     {
15         fprintf(stderr, "fsize: can't access %s\n", name);
16         return 0;
17     }
18     return stbuf.st_size;
19 }
20
21 int main(int argc, char *argv[])
22 {
23     FILE *fp;
24     signed short *data;
25     int fp_more = 1;
26     int fp_data_size;
27     unsigned long fp_fingerprint_size;
28     unsigned char *fp_fingerprint;
29     int file_size;
30     int audio_length;
31     int whereami, ixFile;
32
33     if (argc == 1 || argc > 3)
34     {
35         fprintf(stderr, "Usage:\n");
36         fprintf(stderr, "\t%s <raw-audiofile> [<original-filename>]:\n",
37             argv[0]);
38         exit(1);
39     }
40     else
41     {
42         if ((fp = fopen(argv[1], "rb")) == NULL)
43         {
44             fprintf(stderr, "Can't open %s\n", argv[1]);
45             exit(1);
46         }
47         else
48         {
49             data = (signed short *) malloc(sizeof(signed short) *
50                 N_BLOCKS);
51             fooid = fp_init(SAMPLE_RATE, CHANNELS);
52             if (ferror(fp))
53             {
54                 fprintf(stderr, "Error %d!\n", ferror(fp));
55             }
56             /*

```



```

56         * Lets skip the header (we know we get 44100 Hz, 16 bit, 2
57         channel)
58     */
59     fseek(fp, 44, SEEK_SET);
60     whereami = 44;
61     while ((!feof(fp)) && (fp_more == 1) && !ferror(fp))
62     {
63         fp_data_size = fread(data, sizeof(signed short),
64                               N_BLOCKS, fp);
65         whereami += fp_data_size;
66         fp_more = fp_feed_short(fooid, data, fp_data_size);
67     }
68     fclose(fp);
69     fp_fingerprint_size = fp_getsize(fooid);
70     file_size = fsize(argv[1]);
71     /*
72     * The * 2 is because it is 16 bits, while the file size is
73     * in 8 bits
74     */
75     audio_length = (file_size / (float) (SAMPLE_RATE * CHANNELS
76     * 2)) * 100;
77
78     fp_fingerprint = (unsigned char *)malloc(sizeof(unsigned
79     char) * fp_fingerprint_size);
80     if (fp_calculate(fooid, audio_length, fp_fingerprint) < 0)
81     {
82         fprintf(stderr, "fp_calculate failed!\n");
83     }
84     if (argc == 3)
85     {
86         db_setup();
87         ixFile = db_insert_file(argv[2], 0, 1);
88         db_insert_fingerprint(ixFile, 4, (char *) fp_fingerprint
89         , fp_fingerprint_size);
90         db_close();
91     }
92     fp_free(fooid);
93     free(data);
94     free(fp_fingerprint);
95 }
96
97 return 0;
98 }
99
100 }
101
102 }

```

### B.2.3 Database.java

Database connectivity for Java.

```

1  import java.sql.*;
2
3  public class Database {
4      public Connection conn;
5      public Database()
6      {
7          try {
8              Class.forName("com.mysql.jdbc.Driver").newInstance();
9              conn = DriverManager.getConnection(
10                 "jdbc:mysql://localhost/master?" +
11                 "user=master&password=didyouthinkyouwouldfindithere");
12          } catch (Exception ex) {
13              // handle any errors
14              System.out.println("SQLException: " + ex.getMessage());

```



```

56             insertScore.execute();
57         }
58     }
59     System.out.println("Finished " + rgFiles.get(i).ixFile +
60         "!");
61 }
62 }
63 catch (Exception e) {
64     e.printStackTrace();
65 }
66 }
67
68 public FingerPrint() {
69     r = new int[87][16];
70     dom = new int[87];
71 }
72
73 public void read(DataInputStream fds) throws java.io.IOException,
74     Exception {
75     byte buff[] = new byte[424];
76     int read = fds.read(buff);
77
78     ByteArrayInputStream ds = new ByteArrayInputStream(buff);
79
80     version = ServerUtil.readLEShort(ds);
81
82     if (version != 0 || read != 424) {
83         throw new Exception("Fingerprint version not supported.\n");
84     }
85
86     length = ServerUtil.readLEInt(ds);
87     avg_fit = ServerUtil.readLEShort(ds);
88     avg_dom = ServerUtil.readLEShort(ds);
89
90     for (int i = 0; i < 87; i++) {
91         for (int j = 0; j < 16; j += 4) {
92             int temp = ds.read();
93
94             r[i][j + 0] = (temp >>> 6) & 0x3;
95             r[i][j + 1] = (temp >>> 4) & 0x3;
96             r[i][j + 2] = (temp >>> 2) & 0x3;
97             r[i][j + 3] = temp & 0x3;
98         }
99     }
100
101     int dompos = 0;
102
103     for (int i = 0; i < (66 / 3); i++) {
104         int temp = ds.read();
105         int temp2 = ds.read();
106         int temp3 = ds.read();
107
108         // 3 * 8 bits = 24 bits = 4 * 6 bit doms
109         dom[dompos++] = ((temp >>> 2) & 0x3F);
110         dom[dompos++] = (((temp & 0x03) << 4) | ((temp2 >>> 4) & 0
111             x0F));
112         dom[dompos++] = (((temp2 & 0x0F) << 2) | ((temp3 >>> 6) & 0
113             x3));
114         if (dompos < 87) {
115             dom[dompos++] = (temp3 & 0x3F);
116         }
117     }
118
119     ds.close();
120 }
121
122 public void readFrom(String s) throws java.io.FileNotFoundException,
123     java.io.IOException, Exception
124     {
125     read(new DataInputStream(new FileInputStream(s)));
126 }
127
128 public void displaySummary() {

```

```

125     System.out.println("\tLength: " + length);
126     System.out.println("\tAVG DOM: " + avg_dom);
127     System.out.println("\tAVG FIT: " + avg_fit);
128 }
129
130 public void displayFull() {
131     displaySummary();
132
133     for (int i = 0; i < 87; i++) {
134         StringBuffer line = new StringBuffer(2 * 16 + 1);
135
136         for (int j = 0; j < 16; j++) {
137             int tr = r[i][j];
138             switch (tr) {
139                 case 0:
140                     line.append("0 ");
141                     break;
142                 case 1:
143                     line.append("1 ");
144                     break;
145                 case 2:
146                     line.append("2 ");
147                     break;
148                 case 3:
149                     line.append("3 ");
150                     break;
151                 default:
152                     line.append("X ");
153                     break;
154             }
155         }
156
157         System.out.println(line + " DOM: " + dom[i]);
158     }
159 }
160
161 public float Match(FingerPrint fp) {
162     if (quickMatch(fp)) {
163         return fullMatch(fp);
164     } else {
165         return 0.0f;
166     }
167 }
168
169 public boolean quickMatch(FingerPrint fp) {
170     /*
171      * length within 30 seconds
172      */
173     if (length + 100 * 30 < fp.length) {
174         return false;
175     }
176     if (length - 100 * 30 > fp.length) {
177         return false;
178     }
179
180     /*
181      * average FIT within 0.4
182      */
183     if (avg_fit + 400 < fp.avg_fit) {
184         return false;
185     }
186     if (avg_fit - 400 > fp.avg_fit) {
187         return false;
188     }
189
190     /*
191      * average DOM within 6 units
192      */
193     if (avg_dom + 600 < fp.avg_dom) {
194         return false;
195     }
196     if (avg_dom - 600 > fp.avg_dom) {
197         return false;
198     }

```

```

199     return true;
200 }
201
202
203 // XXX: early exit here
204 public float fullMatch(FingerPrint fp) {
205     /*
206         determine max sensible frame
207     */
208     int maxframe = Math.round((length * 0.9765625f) / 100.0f);
209     maxframe = Math.min(maxframe, 87);
210
211     /*
212         set up 'flaw' counters
213     */
214     int[] rf = new int[4];
215     int[] df = new int[64];
216     int tdf, trf;
217
218     rf[0] = 0;
219     rf[1] = 0;
220     rf[2] = 0;
221     rf[3] = 0;
222     tdf = 0;
223     trf = 0;
224
225     for (int i = 0; i < 64; i++) {
226         df[i] = 0;
227     }
228
229     for (int f = 0; f < maxframe; f++) {
230         for (int b = 0; b < 16; b++) {
231             int rdiff = Math.abs(r[f][b] - fp.r[f][b]);
232             rf[rdiff]++;
233             trf += rdiff;
234         }
235         int ddiff = Math.abs(dom[f] - fp.dom[f]);
236         df[ddiff]++;
237         tdf += ddiff;
238     }
239
240     /*
241         DOM flaws are linear, penalty points = how far we're off
242
243         FIT flaws are nonlinear
244         1-off: 1 penalty point
245         2-off: 3 penalty points
246         3-off: 9 penalty points
247     */
248
249     final int maxrflaw = 9 * 16 * maxframe;
250     final int maxdflaw = (63 * maxframe) / 4;
251
252     int w_trf = rf[1] + rf[2] * 4 + rf[3] * 9;
253     int w_tdf = tdf / 4;
254
255     int totalflaws = w_trf + w_tdf;
256     final int maxflaws = maxrflaw + maxdflaw;
257
258     /*
259         percentage is ratio of our
260         flaws to max theoretical flaws
261     */
262     float perc = ((float) (totalflaws)) / ((float) (maxflaws));
263
264     /*
265         we expect a random track to get
266         about halfway there, so confidence
267         of 50% -> 0%, and a match to get
268         nowhere, so confidence of 0% -> 100%
269     */
270     float conf = ((1.0f - perc) - 0.5f) * 2.0f;
271
272     /*
273         limit to sane values

```

```

273     */
274     conf = Math.min(conf, 1.0f);
275     conf = Math.max(conf, 0.0f);
276
277     return conf;
278 }
279
280 public int getVersion() {
281     return version;
282 }
283
284 public int getAvg_dom() {
285     return avg_dom;
286 }
287
288 public void setAvg_dom(int avg_dom) {
289     this.avg_dom = avg_dom;
290 }
291
292 public int getAvg_fit() {
293     return avg_fit;
294 }
295
296 public void setAvg_fit(int avg_fit) {
297     this.avg_fit = avg_fit;
298 }
299
300 public int getLength() {
301     return length;
302 }
303
304 public void setLength(int length) {
305     this.length = length;
306 }
307 }

```

## B.3 Gunderson's descriptors

The software is licensed under the GPL.

### B.3.1 wave.h

Headers for reading WAVE files from disk.

```

1 #ifndef _WAVE_H
2 #define _WAVE_H 1
3
4 #include "common.h"
5
6 void read_wave_file(Sample &s, char *filename);
7
8 #endif /* !defined(_WAVE_H) */

```

### B.3.2 wave.cpp

Functions for reading WAVE files from disk.

```

1 #include <stdio.h>
2 #include <assert.h>
3
4 #include "wave.h"
5
6 #define READ_BLOCKS 1024
7 #define WAVE_HEADER_SIZE 44
8
9 void read_wave_file(Sample &s, char *filename)
10 {

```

```

11     unsigned file_pos = 0;
12     long file_len;
13     long data_size;
14     signed short *data;
15     FILE *fp;
16
17     fp = fopen(filename, "rb");
18     assert(fp != NULL);
19     fseek(fp, 0, SEEK_END);
20     file_len = ftell(fp);
21     assert(file_len > FEATURE_SAMPLES * 2 + WAVE_HEADER_SIZE);
22
23
24     long samples = (file_len - WAVE_HEADER_SIZE) / (2 * 2);
25     s.samples.reserve(FEATURE_SAMPLES);
26
27     fseek(fp, 44, SEEK_SET);
28     data = (signed short *)malloc(sizeof(signed short) * READ_BLOCKS);
29     while ((!feof(fp)) && file_pos < FEATURE_SAMPLES * 2)
30     {
31         data_size = fread(data, sizeof(signed short), READ_BLOCKS, fp);
32         // data_size - 1 because we need to be able to read at least 2
33         // samples
34         for (int i = 0; i < data_size - 1; i += 2)
35         {
36             signed short l, r;
37             l = data[i];
38             r = data[i + 1];
39
40             #if 0
41                 /*
42                  * There's a lot of codec delay in MP3. This skips
43                  * the first 1105 samples (LAME's estimate of the
44                  * delay for our test MP3s) to make lining up plots
45                  * easier.
46                  */
47                 static unsigned delayed = 0;
48                 if (++delayed < 1105)
49                     continue;
50             #endif
51             s.samples.push_back(0.5 * (double(r) + double(l)));
52         }
53         file_pos += data_size;
54     }
55     s.length = (samples / 44100.0);
56     free(data);
57 }

```

### B.3.3 read\_descriptors.php

Reads the output from Gunderson's descriptors into the database.

```

1 <?php
2 /**
3  * This file reads a list of files produced by Gunderson's descriptors,
4  * and inserts all the fingerprints into tblFingerprint.
5  *
6  * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
7  */
8 require_once('include/include.php');
9
10 // where do we find our files?
11 define('PATH_RESULTS', '/home/vegard/Masteroppgave/Kode/results/
12     descriptors/');
13
14 $oDB = CDatabase::Get();
15
16 // Let's have a lookup list of the index of a fingerprint type.
17 // E.g. $rgixType['centroid'] == 6
18 $sql = "SELECT sType, ixType FROM tblType";

```

```

18 $rgixType = $oDB->extended->GetAssoc($sql, null, null, null,
    MDB2_FETCHMODE_ASSOC, false);
19
20 // Read the list of files and remove excessive fat
21 $rgFiles = file(PATH_RESULTS.'filelist.txt');
22 $rgFiles = array_map('trim', $rgFiles);
23
24 // Our queries
25 $sqlInsert = "INSERT IGNORE INTO tblFingerprint (ixFile, ixType,
    sFingerprint) VALUES (?, ?, ?)";
26 $sqlFind = "SELECT ixFile FROM tblFile WHERE md5Path LIKE ?";
27
28 $i = 0;
29 echo "Starting...\n";
30 // Loop through files
31 foreach ($rgFiles as $md5Path)
32 {
33     $rgDescriptors = file(PATH_RESULTS.$md5Path);
34     $rgDescriptors = array_map('trim', $rgDescriptors);
35     $ixFile = $oDB->extended->GetOne($sqlFind, null, array($md5Path),
        array('text'));
36     foreach ($rgDescriptors as $line)
37     {
38         // Line structure as .ini files
39         list($sDescriptor, $sValues) = explode('=', $line);
40
41         /*
42          * Some of the fingerprints have very large values, and this
43          * is something Gunderson does in his Voronoi program. It should
44          * have been done when the fingerprints were generated.
45          *
46          * We are taking the Y-th root of any moments named *_momentY.
47          */
48         if (preg_match("/^(_moment)(\d)$/", $sDescriptor, $matches))
49         {
50             $rgValues = explode(',', $sValues);
51             foreach ($rgValues as $k => $dbl)
52             {
53                 if ($dbl < 0.0)
54                     $rgValues[$k] = -pow(-$dbl, 1.0 / $matches[2]);
55                 else
56                     $rgValues[$k] = pow($dbl, 1.0 / $matches[2]);
57             }
58             $sValues = implode(',', $rgValues);
59         }
60         // Do we know what kind of descriptor this is?
61         // The fingerprint files also have a bunch of non-necessary data.
62         if (isset($rgixType[$sDescriptor]))
63         {
64             $ixType = $rgixType[$sDescriptor];
65             // Insert data into tblFingerprint
66             $oDB->extended->execParam($sqlInsert, array($ixFile, $ixType,
                $sValues), array('integer', 'integer', 'text'));
67         }
68     }
69     $i++;
70     if ($i % 100 == 0)
71         echo "\t[" . date('H:i:s') . "] Processed $i files\n";
72 }
73 echo "Done!\n";
74
75 ?>

```

### B.3.4 process\_descriptor\_scores.php

Compares the fingerprints and generates the scores for potential duplicate pairs.

```

1 <?php
2 /**
3  * This file compares the fingerprints produced by Gunderson's
    descriptors using

```



```

4  * Euclidean distance, and generates a score for each fingerprint
   * comparison.
5  *
6  * Outputs SQL insert statements to descriptor_inserts.sql that should
7  * processed by MySQL at a later stage.
8  *
9  * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
10 */
11 require_once('include/include.php');
12 ob_end_flush();
13
14 $rgMax = array(
15     6 => 50.0,           // centroid
16     7 => 5.0,           // time
17     8 => 0.005,         // msratio
18     9 => 25.0,          // sleepness
19     10 => 2000.0,       // zero crossings
20     11 => 1000.0,       // zero crossing_guard
21     12 => 40.0,         // mfcc_avg
22     13 => 3.0,          // mfcc_delta_avg
23     18 => 20.0,         // mfcc_f1_avg
24     19 => 5.0,          // mfcc_f1_delta_avg
25
26     14 => 50,           // mfcc_delta_moment2
27     15 => 12500,        // mfcc_delta_moment3
28     16 => 6000000,      // mfcc_delta_moment4
29     17 => 4350000000,   // mfcc_delta_moment5
30
31     20 => 1500,         // mfcc_f1_delta_moment2
32     21 => 1250000,      // mfcc_f1_delta_moment3
33     22 => 1175000000,   // mfcc_f1_delta_moment4
34     23 => 1410000000000, // mfcc_f1_delta_moment5
35
36     24 => 6000,         // mfcc_delta_moment2
37     25 => 6100000,      // mfcc_delta_moment3
38     26 => 10000000000,  // mfcc_delta_moment4
39     27 => 12800000000000, // mfcc_delta_moment5
40
41     28 => 2000,         // mfcc_moment2
42     29 => 1700000,      // mfcc_moment3
43     30 => 2900000000,   // mfcc_moment4
44     31 => 4100000000000, // mfcc_moment5
45 );
46
47 $oDB = CDatabase::Get();
48
49 // we can ignore the _delta_moments of we need to
50 // $rgTypes = array(14, 15, 16, 17, 20, 21, 22, 23, 24, 25, 26, 27, 28,
51 // 29, 30, 31);
52 $rgTypes = array_keys($rgMax);
53
54 foreach ($rgTypes as $ixType)
55 {
56     if (!isset($rgMax[$ixType]))
57         continue;
58
59     $fp = fopen("inserts3.sql", "w");
60
61     fwrite($fp, "INSERT DELAYED IGNORE INTO tblScore (ixFile1, ixFile2,
62         ixType, dblScore) VALUES ");
63     $fCommaFirst = false;
64
65     $sqlSelectFingerprints = "SELECT ixFile, sFingerprint FROM
66         tblFingerprint WHERE ixType = ? ORDER BY ixFile";
67     $cbBefore = memory_get_usage(true);
68     foreach ($rgTypes as $ixType)
69     {
70         // So we don't have to do sqrt() a few million times, we cache
71         // the maximum Euclidean distance squared
72         $dblMaxDistance[$ixType] = $rgMax[$ixType] * $rgMax[$ixType];
73         $rgoFingerprints[$ixType] = $oDB->extended->GetAll(
74             $sqlSelectFingerprints, null, array($ixType), array('
75                 integer'), MDB2_FETCHMODE_ORDERED);

```

```

72     foreach ($rgoFingerprints[$SixType] as $k => $f)
73     {
74         $rgoFingerprints[$SixType][$k][1] = explode(' ',
75             $rgoFingerprints[$SixType][$k][1]);
76     }
77     $cbAfter = memory_get_usage(true);
78
79     echo (( $cbAfter - $cbBefore ) / 1024) . "kB memory in use before starting
80         \n";
81
82     $cFingerprints = count($rgoFingerprints[$rgTypes[0]]);
83     $cInserts = 0;
84     for ($i = 0; $i < $cFingerprints; $i++)
85     {
86         echo "\ti=$i\n";
87         for ($j = $i + 1; $j < $cFingerprints; $j++)
88         {
89             foreach ($rgTypes as $SixType)
90             {
91                 $dblSum = 0;
92                 /* Calculate the Euclidean squared.
93                  * This could have been in a method
94                  * call, but as we are inside our
95                  * third loop already, we'll keep the
96                  * overhead to a minimum.
97                  */
98                 foreach ($rgoFingerprints[$SixType][$i][1] as $k =>
99                     $dblValue1)
100                 {
101                     $dblDiff = $dblValue1 - $rgoFingerprints[$SixType][$j]
102                         [1][$k];
103                     $dblSum += $dblDiff * $dblDiff;
104                 }
105                 // Is the Euclidean distance acceptably small?
106                 if ($dblSum < $dblMaxDistance[$SixType])
107                 {
108                     // Only take the sqrt() if necessary
109                     $dblEucl = sqrt($dblSum);
110                     $ixFile1 = $rgoFingerprints[$SixType][$i][0];
111                     $ixFile2 = $rgoFingerprints[$SixType][$j][0];
112                     $dblScore = 1.0 - ($dblEucl / $rgMax[$SixType]);
113                     if ($fCommaFirst)
114                         $sOut .= ",\n";
115                     else
116                         $fCommaFirst = true;
117
118                     $cInserts++;
119                     $sOut .= "($ixFile1 , $ixFile2 , $SixType , $dblScore)";
120                 }
121             }
122         }
123     }
124     // Let's have a maximum of 100 value groups in each INSERT.
125     if ($cInserts >= 100)
126     {
127         fwrite($fp, $sOut);
128         $sOut = '';
129         fwrite($fp, ";\n\n");
130         fwrite($fp, "INSERT DELAYED IGNORE INTO tblScore ($ixFile1 ,
131             $ixFile2 , $ixType , dblScore) VALUES ");
132         $fCommaFirst = false;
133         echo "Wrote $cInserts inserts to file.\n";
134         $cInserts = 0;
135     }
136 }
137
138 // If we have something more to write, let's do it now.
139 fwrite($fp, $sOut);
140 $sOut = '';
141 fwrite($fp, ";\n\n");
142 fclose($fp);
143 echo "Done with everything.\n";
144 }

```

141 ?&gt;

## B.4 Combining descriptors

Software listed here is licensed under the Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States License.

### B.4.1 duplicate\_functions.php

This script contains functions to check if two file names are likely to be the same song.

```

1 <?php
2 /**
3  * This script contains a variety of functions for comparing two file
4  * names against each other, and determining if they could be the same
5  * song.
6  *
7  * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
8  */
9 function tokenize_filename($s)
10 {
11     // remove anything before a space-hyphen-space sequence
12     // (often artist or album name)
13     $s = preg_replace("/.* - /i", '', $s);
14     $s = preg_replace("/[-\s\d]+/", ' ', strtolower($s));
15     $rgTokens = split(' ', $s);
16     foreach ($rgTokens as $key => $value)
17     {
18         if (strlen($value) <= 2)
19             unset($rgTokens[$key]);
20     }
21     return $rgTokens;
22 }
23
24 function filename($s)
25 {
26     return pathinfo($s, PATHINFO_FILENAME);
27 }
28
29 function artistname($s)
30 {
31     $sDir = pathinfo($s, PATHINFO_DIRNAME);
32     $rgDirs = explode('/', $sDir);
33     return strtolower($rgDirs[3]);
34 }
35
36 function compare_filename($s1, $s2, $rgTokens1 = NULL, $rgTokens2 = NULL
37 )
38 {
39     if ($rgTokens1 == NULL)
40         $rgTokens1 = tokenize_filename(filename($s1));
41     if ($rgTokens2 == NULL)
42         $rgTokens2 = tokenize_filename(filename($s2));
43     $cHits = 0;
44     $rgAllTokens = array_unique(array_merge($rgTokens1, $rgTokens2));
45     foreach ($rgAllTokens as $token)
46     {
47         if (in_array($token, $rgTokens1) && in_array($token, $rgTokens2)
48             )
49             $cHits++;
50     }
51     $cTokens = max(count($rgTokens1), count($rgTokens2));
52     if ($cTokens < 1) return false;
53     return ((artistname($s1) == artistname($s2) &&
54         ($cHits / (float) $cTokens) > 0.6));

```

```
54 }
55 ?>
```

### B.4.2 estimate\_filename\_accuracy.php

This script attempts to match 500 000 random combinations of file names in the database, and prints out the ones that match. See 3.1.3.

```
1 <?php
2 /**
3  * This script picks 500 000 random files to compare, and checks
4  * if they have file names that match. The file names that match
5  * are printed to stdout, and was verified by a human.
6  *
7  * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
8  */
9 require_once('include/include.php');
10 $oDB = CDatabase::Get();
11 $oSmarty = CSmarty::Get();
12
13 include('duplicate_finder.php');
14
15 $sqlSelectFile = "SELECT ixFile, sPath FROM tblFile";
16 $rgoFiles = $oDB->extended->GetAssoc($sqlSelectFile, null, null, null,
17     MDB2_FETCHMODE_ASSOC, false);
18
19 $rgoTokens = array();
20 foreach ($rgoFiles as $ixFile => $sPath)
21 {
22     $rgoTokens[$ixFile] = tokenize_filename(filename($sPath));
23 }
24
25 $cNumMatch = 0;
26 for ($i = 0; $i < 500000; $i++)
27 {
28     $x = rand(0, count($rgoFiles));
29     $y = rand(0, count($rgoFiles));
30     while ($x == $y || !isset($rgoFiles[$x]) || !isset($rgoFiles[$y]))
31     {
32         $x = rand(0, count($rgoFiles));
33         $y = rand(0, count($rgoFiles));
34     }
35     if (compare_filename($rgoFiles[$x], $rgoFiles[$y], $rgoTokens[$x],
36         $rgoTokens[$y]))
37     {
38         $cNumMatch++;
39         echo "<br/>\n".str_replace("/mnt/media/", '', $rgoFiles[$x])."<br/>\n";
40         echo str_replace("/mnt/media/", '', $rgoFiles[$y])."<br/>\n";
41     }
42 }
43 echo $cNumMatch."\n";
44 ?>
```

### B.4.3 estimate\_duplicate\_count.php

This script finds the number of potential duplicate pairs in the music collection by looking at their file names. See 3.1.3.

```
1 <?php
2 /**
3  * This script finds the number of potential duplicates
4  * when only looking at the file names.
5  *
6  * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
7  */
8 require_once('include/include.php');
9 $oDB = CDatabase::Get();
```

```

10
11 include('duplicate_functions.php');
12
13 $sqlSelectFile = "SELECT ixFile, sPath FROM tblFile ORDER BY ixFile ASC"
14 ;
15 $rgoFiles = $oDB->extended->GetAssoc($sqlSelectFile, null, null, null,
16     MDB2_FETCHMODE_ASSOC, false);
17
18 /*
19  * Pre-process the file names into tokens, since we
20  * will be doing about 45M checks.
21  */
22 $rgoTokens = array();
23 $ixFileMax = 0;
24 foreach ($rgoFiles as $ixFile => $sPath)
25 {
26     $rgoTokens[$ixFile] = tokenize_filename(filename($sPath));
27     $ixFileMax = $ixFile;
28 }
29
30 /*
31  * Start counting duplicates.
32  */
33 $cDupCount = 0;
34 $rgoMatch = array();
35 for ($i = 0; $i <= $ixFileMax; $i++)
36 {
37     // Progress indicator
38     // echo "$i\t$cDupCount\n";
39
40     if (!isset($rgoFiles[$i])) continue;
41
42     for ($j = $i + 1; $j <= $ixFileMax; $j++)
43     {
44         if (!isset($rgoFiles[$j])) continue;
45
46         if (compare_filename($rgoFiles[$i], $rgoFiles[$j],
47             $rgoTokens[$i], $rgoTokens[$j]))
48         {
49             $cDupCount++;
50         }
51     }
52 }
53 echo $cDupCount;
54 ?>

```

#### B.4.4 calculate\_summed\_score.php

This script combines the scores from individual descriptors into combined scores.

```

1 <?php
2 /**
3  * This script combines match scores for several descriptors using the
4  * methods outlined in the thesis. It needs a lot of memory, probably
5  * around 512MB or so.
6  *
7  * Inserts the results directly into MySQL.
8  *
9  * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
10 */
11 require_once('include/include.php');
12 $oDB = CDatabase::Get();
13
14 define('C_AM', 1); // arithmetic mean
15 define('C_RMS', 2); // root-mean square
16 define('C_WAM_SUM', 3); // weighted arithmetic mean by sum
17 define('C_TM', 4); // truncated mean
18 define('C_WAM_COUNT', 5); // weighted arithmetic mean by count
19 define('C_WAM_MAN1', 6); // weighted arithmetic mean by manual weights 1
20 define('C_WAM_PERF', 7); // weighted arithmetic mean by performance
21 define('C_WAM_AVG', 8); // weighted arithmetic mean by average
22 define('C_WAM_MAN2', 9); // weighted arithmetic mean by manual weights 2

```

```

23 define('C_WAM_MAN3', 10); // weighted arithmetic mean by manual weights 2
24 define('C_BAYES', 11); // naÃve Bayes
25 define('C_WRMS', 18); // weighted root mean square
26 /*
27  * Methods operating on a limited set of descriptors.
28  */
29 define('C_BAYES_CHOSEN', 12); // naÃve Bayes
30 define('C_AM_CHOSEN', 14); // arithmetic mean
31 define('C_RMS_CHOSEN', 15); // root mean square
32 define('C_TM_CHOSEN', 16); // truncated mean
33 define('C_MEDIAN_CHOSEN', 17); // median
34 define('C_WRMS_CHOSEN', 19); // weighted RMS
35
36 define('BAYES_PROB', 0.8);
37
38 define('CUTOFF_SCORE', 0.2);
39
40 define('FILES_IN_PASS', 25);
41
42
43 /**
44  * Root mean square
45  */
46 function calculate_rms($rgValues, $cCount)
47 {
48     $dblSum = 0;
49     foreach ($rgValues as $dblValue)
50         $dblSum += ($dblValue * $dblValue);
51     return sqrt($dblSum / (float) $cCount);
52 }
53
54 /**
55  * Arithmetic mean
56  */
57 function calculate_am($rgValues, $cCount)
58 {
59     return array_sum($rgValues) / (float) $cCount;
60 }
61
62 /**
63  * Weighted arithmetic mean
64  */
65 function calculate_wam($rgValues, $rgWeight, $dblSumWeights)
66 {
67     if ($dblSumWeights == 0) return 0;
68     $dblSum = 0;
69     foreach ($rgValues as $ix => $dbValue)
70     {
71         $dblSum += ($dbValue * $rgWeight[$ix]);
72     }
73     return $dblSum / $dblSumWeights;
74 }
75
76 /**
77  * Weighted root mean square
78  */
79 function calculate_wrms($rgValues, $rgWeight, $dblSumWeights)
80 {
81     $dblSum = 0;
82     foreach ($rgValues as $ix => $dblValue)
83         $dblSum += $rgWeight[$ix] * ($dblValue * $dblValue);
84     return sqrt($dblSum / (float) $dblSumWeights);
85 }
86
87 /**
88  * Median
89  */
90 function calculate_median($rgValues, $cCount)
91 {
92     array_pad($rgValues, $cCount, 0);
93     sort($rgValues);
94     while (count($rgValues) > 2)
95     {

```

```

96     array_pop($rgValues);
97     array_shift($rgValues);
98 }
99 if (count($rgValues) == 1)
100     return array_pop($rgValues);
101 else
102     return array_sum($rgValues) / count($rgValues);
103 }
104
105 /**
106  * Truncated arithmetic mean.
107  */
108 function calculate_truncated_mean($rgValues, $cCount)
109 {
110     // make sure 0-values are counted
111     array_pad($rgValues, $cCount, 0);
112     // sort
113     sort($rgValues);
114     // remove last and first
115     array_pop($rgValues);
116     array_shift($rgValues);
117     return calculate_am($rgValues, $cCount - 2);
118 }
119
120 /**
121  * Naive Bayes
122  */
123 function calculate_bayes($rgValues)
124 {
125     $dblProd = 1; $dblInvProd = 1;
126     foreach ($rgValues as $dblValue)
127     {
128         $dblValue = max(0.01, min(0.99, $dblValue));
129         $dblProd *= $dblValue;
130         $dblInvProd *= (1 - $dblValue);
131     }
132     $dblProd /= BAYES_PROB;
133     $dblInvProd /= (1 - BAYES_PROB);
134     $dblSum = $dblProd + $dblInvProd;
135     if ($dblSum != 0)
136         return $dblProd / $dblSum;
137     else
138         return 0;
139 }
140
141 /**
142  * Standard deviation
143  */
144 function calculate_sd($dblMean, $rgValues, $cCount)
145 {
146     $dblSum = 0;
147     foreach ($rgValues as $ix => $dblValue)
148     {
149         $dblSum += ($dblValue - $dblMean) * ($dblValue - $dblMean);
150     }
151     return sqrt($dblSum / (float) $cCount);
152 }
153
154 /**
155  * Which weights do we need?
156  */
157 $rgWAMs = array(C_WAM_SUM, C_WAM_COUNT, C_WAM_AVG, C_WAM_MAN1,
158               C_WAM_PERF, C_WAM_MAN2, C_WAM_MAN3);
159
160 /**
161  * Let's figure out the weights for the weighted arithmetic means.
162  */
163 $sqlSelectWeight = "SELECT ixType, dblWeightSum, dblWeightCount,
164                   dblWeightAvg, ".
165                   "dblWeightMan1, dblWeightMan2, dblWeightMan3, dblWeightPerformance ".
166                   "FROM tblType";
167 $rW =& $oDB->query($sqlSelectWeight);
168 $rgWeight = array();

```

```

167 $cCount = 0;
168 while ($r = $rW->fetchRow())
169 {
170     $rgWeight[C_WAM_SUM][ $r->ixType] = $r->dblWeightSum;
171     $rgWeight[C_WAM_COUNT][ $r->ixType] = $r->dblWeightCount;
172     $rgWeight[C_WAM_MAN1][ $r->ixType] = $r->dblWeightMan1;
173     $rgWeight[C_WAM_MAN2][ $r->ixType] = $r->dblWeightMan2;
174     $rgWeight[C_WAM_MAN3][ $r->ixType] = $r->dblWeightMan3;
175     $rgWeight[C_WAM_PERF][ $r->ixType] = $r->dblWeightPerformance;
176     $rgWeight[C_WAM_AVG][ $r->ixType] = $r->dblWeightAvg;
177     $cCount++;
178 }
179 /*
180  * Build an array of weight sums, indexed by which set of weights we are
181  * using
182  * This is used later on for the weighted arithmetic mean, we'll just
183  * cache
184  * this for later use.
185  */
186 foreach ($rgWAMs as $ixWAM)
187 {
188     $rgSumWeights[$ixWAM] = array_sum($rgWeight[$ixWAM]);
189 }
190 $rW->free();
191 /*
192  * Determine the number of files. (Or more accurately, the range of
193  * ixFile)
194  */
195 $cNumFiles = $oDB->extended->GetOne("SELECT MAX(ix File) FROM tblFile");
196 /*
197  * This section allows us to automatically start multiple subprocesses
198  * of this script, to use more than one processor, and automatically
199  * dividing the work load.
200  * Examples:
201  * $ php calculate_summed_score.php 4
202  * Runs 4 subprocesses. Output in output.txt.
203  * $ php calculate_summed_score.php
204  * Runs in 1 subprocess. Output to stdout.
205  * $ php calculate_summed_score.php 40 50 "one"
206  * Runs one process, calculating for files 40-90, using the
207  * process identifier "one".
208  */
209 $ixRow = 0;
210 $sProcess = '';
211 if ($argc == 4 && is_numeric($argv[1]) && is_numeric($argv[2]))
212 {
213     // use parameters as range to use
214     $ixRow = (int) $argv[1];
215     $cNumFiles = (int) $argv[2];
216     $sProcess = $argv[3];
217 }
218 else if ($argc == 2 && is_numeric($argv[1]))
219 {
220     // we are being asked to use multiple processes
221     $cNumProc = (int) $argv[1];
222
223     $cFilesPerProc = (int) ($cNumFiles / $cNumProc);
224     // make sure we are a multiple of FILES_IN_PASS
225     $cFilesPerProc += FILES_IN_PASS - ($cFilesPerProc % FILES_IN_PASS);
226
227     exec('rm output.txt');
228
229     $cmd = "php calculate_summed_score.php %d %d '%s' >> output.txt &";
230     $n = 1;
231     for ($i = 0; $i < $cNumFiles; $i += $cFilesPerProc)
232     {
233         $imax = min($i + ($cFilesPerProc - 1), $cNumFiles);
234         echo "Launching subprocess ($i -> $imax)...\n";
235         exec(sprintf($cmd, $i, $imax, $n));
236         $n++;
237     }

```



```

238     die();
239 }
240 }
241
242 echo "[${sProcess}] Starting to sum (${cNumFiles} files in total)... \n";
243
244 /*
245  * Start calculating everything.
246  */
247
248 $rgSelectedTypes = array(1,2,3,4,12,13,18,19);
249 $cSelectedCount = count($rgSelectedTypes);
250 while ($ixRow <= $cNumFiles)
251 {
252     $sqlSelectScores = "SELECT ixFile1, ixFile2, ixType, dblScore "
253     "FROM tblScore WHERE ixFile1 BETWEEN $ixRow AND "
254     "($ixRow + FILES_IN_PASS - 1)." ORDER BY ixFile1, ixFile2";
255     $res =& $oDB->query($sqlSelectScores);
256
257     while ($row = $res->fetchRow())
258     {
259         $sKey = $row->ixFile1.'-'. $row->ixFile2;
260         if ($sPrevKey != $sKey && !empty($sPrevKey))
261         {
262             // remember that $row has changed, so we can't look at those
263             // values
264             list($ixFile1, $ixFile2) = split('-', $sPrevKey);
265
266             /*
267              * First run everything on the set of all descriptors
268              */
269             $dblRMS = calculate_rms($rgValues, $cCount);
270             if ($dblRMS > CUTOFF_SCORE)
271             {
272                 $dblRMS_SD = calculate_sd($dblRMS, $rgValues, $cCount);
273                 $rgInsert[] = '('.$ixFile1.', '.$ixFile2.', '.C_RMS.', '
274                 .$dblRMS.', '.$dblRMS_SD.')';
275             }
276
277             $dblWRMS = calculate_wrms($rgValues, $rgWeight[C_WAM_MAN2],
278                 $rgSumWeights[C_WAM_MAN2]);
279             if ($dblWRMS > CUTOFF_SCORE)
280             {
281                 $dblWRMS_SD = calculate_sd($dblWRMS, $rgValues,
282                     $rgSumWeights[C_WAM_MAN2]);
283                 $rgInsert[] = '('.$ixFile1.', '.$ixFile2.', '.C_WRMS.', '
284                 .$dblWRMS.', '.$dblWRMS_SD.')';
285             }
286
287             $dblAM = calculate_am($rgValues, $cCount);
288             if ($dblAM > CUTOFF_SCORE)
289             {
290                 $dblAM_SD = calculate_sd($dblAM, $rgValues, $cCount);
291                 $rgInsert[] = '('.$ixFile1.', '.$ixFile2.', '.C_AM.', '
292                 .$dblAM.', '.$dblAM_SD.')';
293             }
294
295             foreach ($rgWAMs as $ixWAM)
296             {
297                 $dblWAM = calculate_wam($rgValues, $rgWeight[$ixWAM],
298                     $rgSumWeights[$ixWAM]);
299                 if ($dblWAM > CUTOFF_SCORE)
300                 {
301                     $dblWAM_SD = calculate_sd($dblWAM, $rgValues,
302                         $rgSumWeights[$ixWAM]);
303                     $rgInsert[] = '('.$ixFile1.', '.$ixFile2.', '.$ixWAM
304                     .', '.$dblWAM.', '.$dblWAM_SD.')';
305                 }
306             }
307
308             $dblITM = calculate_truncated_mean($rgValues, $cCount);
309             if ($dblITM > CUTOFF_SCORE)
310             {

```

```

302         $dblTM_SD = calculate_sd($dblTM, $rgValues, $cCount - 2)
303         ;
304         $rgInsert [] = '('.$ixFile1.', '.$ixFile2.', '.C_TM.', '.
305         $dblTM.', '.$dblTM_SD.')';
306     }
307     $dblBayes = calculate_bayes($rgValues)
308     if ($dblBayes > CUTOFF_SCORE)
309     {
310         $rgInsert [] = '('.$ixFile1.', '.$ixFile2.', '.C_BAYES.',
311         '.$dblBayes.',0)';
312     }
313     /*
314     * Lets build a separate value array for the set of chosen
315     * descriptors.
316     * This allows using the same functions to calculate summed
317     * scores.
318     */
319     $rgSelectedValues = array();
320     foreach ($rgSelectedTypes as $ixType)
321     {
322         $rgSelectedValues[$ixType] = $rgValues[$ixType];
323     }
324     /*
325     * Calculations for the chosen descriptors.
326     */
327     $dblMed = calculate_median($rgSelectedValues,
328     $cSelectedCount);
329     if ($dblMed > CUTOFF_SCORE)
330     {
331         $dblMed_SD = calculate_sd($dblMed, $rgSelectedValues,
332         $cSelectedCount);
333         $rgInsert [] = '('.$ixFile1.', '.$ixFile2.', '.
334         C_MEDIAN_CHOSEN.', '.$dblMed.', '.$dblMed_SD.')';
335     }
336     $dblRMS = calculate_rms($rgSelectedValues, $cSelectedCount);
337     if ($dblRMS > CUTOFF_SCORE)
338     {
339         $dblRMS_SD = calculate_sd($dblRMS, $rgSelectedValues,
340         $cSelectedCount);
341         $rgInsert [] = '('.$ixFile1.', '.$ixFile2.', '.
342         C_RMS_CHOSEN.', '.$dblRMS.', '.$dblRMS_SD.')';
343     }
344     /*
345     * Why is the WRMS here, yet not the WAM? Because the WAM
346     * was calculated using
347     * all sets of weights, but we only chose one set of weights
348     * for the WRMS.
349     * This means that when limiting the set of descriptors, yet
350     * using the same weights
351     * we effectively get the third set of manual weights. Weird
352     * , yet correct.
353     */
354     $dblWRMS = calculate_wrms($rgSelectedValues, $rgWeight[
355     C_WAM_MAN2], $rgSumWeights[C_WAM_MAN2]);
356     if ($dblWRMS > CUTOFF_SCORE)
357     {
358         $dblWRMS_SD = calculate_sd($dblWRMS, $rgSelectedValues,
359         $rgSumWeights[C_WAM_MAN2]);
360         $rgInsert [] = '('.$ixFile1.', '.$ixFile2.', '.
361         C_WRMS_CHOSEN.', '.$dblWRMS.', '.$dblWRMS_SD.')';
362     }
363     $dblAM = calculate_am($rgSelectedValues, $cSelectedCount);
364     if ($dblAM > CUTOFF_SCORE)
365     {
366         $dblAM_SD = calculate_sd($dblAM, $rgSelectedValues,
367         $cSelectedCount);
368         $rgInsert [] = '('.$ixFile1.', '.$ixFile2.', '.
369         C_AM_CHOSEN.', '.$dblAM.', '.$dblAM_SD.')';

```

```

357     }
358
359     $dblTM = calculate_truncated_mean($rgSelectedValues,
360                                     $cSelectedCount);
361     if ($dblTM > CUTOFF_SCORE)
362     {
363         $dblTM_SD = calculate_sd($dblTM, $rgSelectedValues,
364                                 $cSelectedCount - 2);
365         $rgInsert [] = '('.$ixFile1.', '.$ixFile2.', '.
366                       C_TM_CHOSEN.', '.$dblTM.', '.$dblTM_SD.')';
367     }
368
369     $dblBayes = calculate_bayes($rgSelectedValues)
370     if ($dblBayes > CUTOFF_SCORE)
371     {
372         $rgInsert [] = '('.$ixFile1.', '.$ixFile2.', '.
373                       C_BAYES_CHOSEN.', '.$dblBayes.', 0)';
374     }
375
376     /*
377     * Reset the value array for the next run.
378     */
379     unset($rgValues);
380     $rgValues = array();
381 }
382 $rgValues[$row->ixType] = $row->dblScore;
383 $sPrevKey = $sKey;
384 }
385 $res->free();
386
387 /*
388 * Insert all our cumulated results. Uses INSERT DELAYED since we
389 * don't need
390 * to inspect this right away.
391 */
392 $sql = "INSERT DELAYED INTO tblSummedScore (ixFile1, ixFile2,
393       ixScoreMethod, dblScore, dblSD) VALUES ";
394 $sql.= implode(' ', $rgInsert);
395 $oDB->query($sql);
396
397 $i += count($rgInsert);
398 $rgInsert = array();
399 $ixRow += FILES_IN_PASS;
400 echo "\t[".date('H:i:s')."] [$$Process] Processed files starting
401       from $ixRow ($i cumulative database inserts).\n";
402 }
403
404 echo "[$$Process] Done!\n";
405
406 echo "[$$Process] Entered $i scores into tblSummedScore!\n";
407 echo "[$$Process] Finished!\n";
408 ?>

```

### B.4.5 find\_index.php

This script finds the average distance within verified duplicate pairs for the simple descriptors song length, centroid, steepness, mean/square ratio and the rate of zero crossings. These are all scalar descriptors that can be indexed by. See 4.5 for motivation.

```

1 <?php
2 /**
3  * This script finds the average — and the minimum and maximum —
4  * distance for verified duplicates with the simple descriptors
5  * song length, centroid, msratio, steepness and zerocrossings.
6  * It also lists the percentage of verified duplicates that will
7  * be found when searching within 2 times the average distance
8  * for a given descriptor.
9  *
10 * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>

```

```

11  */
12  require_once('include/include.php');
13  $oDB = CDatabase::Get();
14
15  /*
16  * Get a list of descriptor names.
17  */
18  $sql = "SELECT ixType, sType FROM tblType";
19  $rgsTypes = $oDB->extended->GetAssoc($sql, null, null, null,
20      MDB2_FETCHMODE_ASSOC, false);
21
22  /*
23  * Find all the verified matches.
24  */
25  $sql = "SELECT * FROM tblVerifiedMatch";
26  $rgoVerified = $oDB->extended->GetAll($sql);
27  $rgVerified = array();
28  foreach ($rgoVerified as $o)
29  {
30      $rgVerified[$o->ixFile1][$o->ixFile2] = true;
31      $rgVerified[$o->ixFile2][$o->ixFile1] = true;
32  }
33  $rgoVerified = null;
34  unset($rgoVerified);
35
36  function array_average($rg)
37  {
38      return (array_sum($rg) / (float) count($rg));
39  }
40
41  /*
42  * The descriptors to look at.
43  */
44  $rgixType = array(6, 7, 8, 9, 10, 11);
45  foreach ($rgixType as $ixType)
46  {
47      /*
48      * Find all the fingerprints for the verified matches.
49      */
50      $sql =
51          "SELECT f1.sFingerprint AS fp1, f2.sFingerprint AS fp2, ixFile1
52          "FROM tblVerifiedMatch v, tblFingerprint f1, tblFingerprint f2 "
53          "WHERE f1.ixFile = v.ixFile1 AND f2.ixFile = v.ixFile2 AND "
54          "f1.ixType = ? AND f2.ixType = f1.ixType";
55      $rgoFingerprint = $oDB->extended->GetAll($sql, null,
56          array($ixType), array('integer'));
57
58      /*
59      * Find the distance for each verified match.
60      * This cannot be done in MySQL, because the column
61      * that stores the fingerprints is in binary format,
62      * to accomodate storing of all types of fingerprints,
63      * and MySQL wouldn't know what to do with a
64      * subtraction of binary data.
65      */
66      $rgDiff = array();
67      foreach ($rgoFingerprint as $oF)
68      {
69          $rgDiff[] = abs($oF->fp1 - $oF->fp2);
70      }
71      sort($rgDiff);
72      $avg = array_average($rgDiff);
73      $min = $rgDiff[0];
74      $max = round($rgDiff[count($rgDiff)-1], 2);
75
76      /*
77      * Now lets figure out how many fingerprints we would
78      * need to search (on average) when limiting our search
79      * to fingerprints within +/- 2 * $avg, and how many of
80      * the results we find are verified.
81      */

```

```

82     $sql = "SELECT ixFile , sFingerprint FROM tblFingerprint WHERE ixType
           = ?";
83     $rgFingerprint = $oDB->extended->GetAll($sql , null ,
84         array($ixType) , array('integer'));
85     $rgixSearchFor = array_keys($rgoFingerprint);
86     $rgValues = array();
87     $rgNumVerified = array();
88     $delta = $avg * 2;
89     foreach ($rgixSearchFor as $ix)
90     {
91         $dblValue = $rgoFingerprint[$ix]->fp1;
92         $ixFile = $rgoFingerprint[$ix]->ixFile1;
93         $c = 0;
94         $cNumVerified = 0;
95         foreach ($rgFingerprint as $oFingerprint)
96         {
97             $dblFingerprint = (float) $oFingerprint->sFingerprint;
98             if ($dblFingerprint < ($dblValue + $delta) &&
99                 $dblFingerprint > ($dblValue - $delta))
100            {
101                $c++;
102                $ixFile2 = (float) $oFingerprint->ixFile;
103                if (isset($rgVerified[$ixFile][$ixFile2]))
104                    $cNumVerified++;
105            }
106        }
107        $rgValues[] = $c;
108        $rgNumVerified[] = $cNumVerified / count($rgVerified[$ixFile]);
109    }
110    $c = round(array_average($rgValues) , 0);
111    $dblFindRate = round(array_average($rgNumVerified) , 4);
112
113    echo "avg: $avg\tmin: $min\tmax: $max\tsearch: $c\t"
114        ."find rate: $dblFindRate\t\t".$rgsTypes[$ixType]."\n";
115
116    ?>

```

### B.4.6 find\_num\_correct.php

This script finds the number of verified duplicate pairs in partitions of 100 duplicates, for all the individual descriptors and all the combined solutions.

```

1  <?php
2  /**
3   * This script finds the number of correct duplicates ,
4   * divided into partitions of 100 duplicates , for each
5   * individual descriptor and each combined solution .
6   *
7   * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
8   */
9  require_once('include/include.php');
10 $cNumRows = 1000;
11
12 /*
13  * HTML table headers
14  */
15 ?>
16 <table border="1">
17 <tr>
18 <th>Type/method</th>
19 <th>Name</th>
20 <?php
21 for ($i = 1; $i <= 10; $i++)
22 {
23     echo "<th>".(( $i -1)*100+1) . "-" . ( $i *100) . "</th>\n";
24 }
25 ?>
26 <th>%</th>
27 <th>Errors in first <?= $cNumRows ?></th>
28 <th>First error at #</th>
29 </tr>

```

```

30 <?php
31 $oDB = CDatabase::Get();
32
33 $cNumFetchRows = $cNumRows;
34
35
36 /*
37  * Get a list of the verified duplicates.
38  */
39 $sql = "SELECT * FROM tblVerifiedMatch";
40 $rgoVerified = $oDB->extended->GetAll($sql);
41 $rgVerified = array();
42 foreach ($rgoVerified as $o)
43 {
44     $rgVerified[$o->ixFile1][$o->ixFile2] = true;
45 }
46 $rgoVerified = null;
47 unset($rgoVerified);
48
49 function display($ixMethod, $sTitle, $rgoScores)
50 {
51     global $rgVerified;
52     global $cNumRows;
53 ?>
54 <tr>
55 <td><?= $ixMethod ?></td>
56 <td><?= $sTitle ?></td>
57 <?php
58     $rgoScoresVerified = array();
59     $rgoScoresWrong = array();
60     $i = 0;
61     $cWrong = 0;
62     $ixFirstWrong = 0;
63     $g = 0; // counter
64     $cWrongInPartition = 0;
65     for ($z = 0; $z < $cNumRows; $z++)
66     {
67         $o = $rgoScores[$z];
68         if (isset($rgoScores[$z]) &&
69             $rgVerified[$o->ixFile1][$o->ixFile2])
70         {
71             $rgoScoresVerified[] = $o;
72         }
73         else
74         {
75             $rgoScoresWrong[] = $o;
76             if (empty($ixFirstWrong))
77                 $ixFirstWrong = $i;
78             if ($i < $cNumRows)
79             {
80                 $cWrong++;
81                 $cWrongInPartition++;
82             }
83         }
84     }
85     unset($rgoScores[$k]);
86     $i++;
87     if ($i % 100 == 0)
88     {
89         $g++;
90         echo "<td>$cWrongInPartition</td>\n";
91         $cWrongInPartition = 0;
92     }
93 }
94 ?>
95 <td><?= $cWrong / $cNumRows ?></td>
96 <td><?= $cWrong ?></td>
97 <td><?= $ixFirstWrong ?></td>
98 </tr>
99 <?php
100 }
101 }
102
103 /*

```

```

104  * A list of the names of all the scoring methods.
105  */
106  $sqlSelectTitles = "SELECT ixScoreMethod, sTitle FROM tblScoreMethods";
107  $rgsTitle = $oDB->extended->GetAssoc($sqlSelectTitles, null, null, null,
108      MDB2_FETCHMODE_ASSOC, false);
109  // Get scores for all the scoring methods
110  $rgixScoreMethods = array_keys($rgsTitle);
111
112  /*
113  * Now let's examine them.
114  */
115  $sql = "SELECT dblScore, dblSD, ixFile1, ixFile2 ".
116      "FROM tblSummedScore ".
117      "WHERE ixScoreMethod = ? ".
118      "ORDER BY dblScore DESC, dblSD ASC ".
119      "LIMIT 0, $cNumFetchRows";
120
121  foreach ($rgixScoreMethods as $ixScoreMethod)
122  {
123      $sTitle = $rgsTitle[$ixScoreMethod];
124      $rgoScores = $oDB->extended->GetAll($sql, null,
125          array($ixScoreMethod), array('integer'));
126      display('C'.$ixScoreMethod, $sTitle, $rgoScores);
127  }
128
129  /*
130  * Let's look at the individual descriptors
131  */
132
133  $sqlSelectTypes = "SELECT ixType, sType FROM tblType";
134  $rgoTypes = $oDB->extended->GetAssoc($sqlSelectTypes, null, null, null,
135      MDB2_FETCHMODE_ASSOC, false);
136
137  $sql = "SELECT ixFile1, ixFile2, dblScore ".
138      "FROM tblScore WHERE ixType = ? ".
139      "ORDER BY dblScore DESC LIMIT 0, $cNumFetchRows";
140
141  foreach ($rgoTypes as $ixType => $sType)
142  {
143      $rgoScores = $oDB->extended->GetAll($sql, null,
144          array($ixType), array('integer'));
145      display('I'.$ixType, $sType, $rgoScores);
146  }
147
148  ?>
149  </table>

```

### B.4.7 display\_duplicates.php

This script displays the potential duplicate pairs, and allows for basic filtering based on whether the file names match, or if the duplicate pair have been verified.

```

1  <?php
2  /**
3   * This script display potential duplicates in the database,
4   * and highlights the duplicates that match by file name,
5   * or have been verified to be correct.
6   *
7   * Options:
8   *   fVerified=1           — show only verified duplicates
9   *   fVerified=0         — show only non-verified duplicates
10  *   fMatch=1             — show only matching file names
11  *   fMatch=0            — show only non-matching file names
12  *   ixPage=<page>       — skip <page> * 1000 results
13  *
14  * If fVerified and fMatch is not specified, no restrictions
15  * are placed on the duplicates status.
16  *
17  * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
18  */

```

```

19 require_once('include/include.php');
20 $oDB = CDatabase::Get();
21
22 $cResultsPerPage = 1000;
23 $ixPage = 0;
24 $ixStart = 0;
25 $ixScoreMethod = 1;
26
27 require_once('duplicate_functions.php');
28
29 /**
30  * Get a list of the verified duplicates.
31  */
32 $sql = "SELECT * FROM tblVerifiedMatch";
33 $rgoVerified = $oDB->extended->GetAll($sql);
34 $rgoVerified = array();
35 foreach ($rgoVerified as $o)
36 {
37     $rgoVerified[$o->ixFile1][$o->ixFile2] = true;
38 }
39 $rgoVerified = null;
40 unset($rgoVerified);
41
42 $sqlSelectFile = "SELECT ixFile, sPath FROM tblFile";
43 $rgoFiles = $oDB->extended->GetAssoc($sqlSelectFile, null, null, null,
44     MDB2_FETCHMODE_ASSOC, false);
45 $rgoFiles = str_replace('/home/vegard/Music/', '', $rgoFiles);
46 $rgoFiles = str_replace('/mnt/media/', '', $rgoFiles);
47
48
49 /**
50  * Now lets get all the scores from the database
51  */
52
53 // Lower bound for values extracted.
54 $dblScoreLimit = 0.05;
55 if (isset($_GET['ixType']))
56 {
57     $ixType = (int) $_GET['ixType'];
58     $sqlSelectScores = "SELECT ixFile1, ixFile2, dblScore ".
59         "FROM tblScore WHERE ".
60         "dblScore > $dblScoreLimit AND ixType = $ixType ".
61         "ORDER BY dblScore DESC LIMIT ?, ?";
62 }
63 elseif (isset($_GET['ixScoreMethod']))
64 {
65     $ixScoreMethod = (int) $_GET['ixScoreMethod'];
66     $sqlSelectScores = "SELECT ixFile1, ixFile2, dblScore ".
67         "FROM tblSummedScore ".
68         "WHERE dblScore > $dblScoreLimit ".
69         "AND ixScoreMethod = $ixScoreMethod ".
70         "ORDER BY dblScore DESC LIMIT ?, ?";
71 }
72 else
73 {
74     /**
75      * No type specified, display a menu of choices.
76      */
77     echo '<h1>Combined methods</h1>';
78     $sql = "SELECT ixScoreMethod, sTitle FROM tblScoreMethods";
79     $r = $oDB->extended->GetAll($sql);
80     foreach ($r as $row)
81     {
82         echo '<a href="?ixScoreMethod='.$row->ixScoreMethod.'">'.
83             $row->sTitle.'</a><br/>';
84     }
85
86     echo '<h1>Individual descriptors</h1>';
87     $sql = "SELECT ixType, sType FROM tblType";
88     $r = $oDB->extended->GetAll($sql);
89     foreach ($r as $row)
90     {
91         echo '<a href="?ixType='.$row->ixType.'">'.
92             $row->sType.'</a><br/>';

```



```

93     }
94     die();
95 }
96
97 /**
98  * Allow parameters to change what is fetched.
99  */
100 if ($_GET['ixPage'])
101 {
102     $ixPage = (int) $_GET['ixPage'];
103     if ($ixPage < 0) $ixPage = 0;
104     $ixStart = $cResultsPerPage * $ixPage;
105 }
106
107 /**
108  * Go fetch!
109  */
110 $res = $oDB->extended->GetAll($sqlSelectScores, null,
111     array($ixStart, $cResultsPerPage),
112     array('integer', 'integer'));
113
114 /*
115  * Output the HTML page headers
116  */
117 ?>
118 <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN"
119     "http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd">
120 <html xmlns="http://www.w3.org/1999/xhtml" xml:lang="no" lang="no">
121     <head>
122         <title>Master thesis</title>
123         <link rel="stylesheet" type="text/css" href="css/main.css" />
124     </head>
125     <body>
126     <form method="post" action="save_verified.php">
127     <input type="hidden" name="redirect_to"
128         value="display_unverified_scores.php<?php
129     echo (isset($ixType) ? "?ixType=$ixType" : "?ixScoreMethod=
130         $ixScoreMethod");
131     ?>" />
132     <input type="submit" value="Save verified" />
133     <table>
134     <tr>
135     <th><abbr title="Match number">M#</abbr></th>
136     <th><abbr title="Do file names match?">M?</abbr></th>
137     <th><abbr title="Is this match verified">V?</abbr></th>
138     <th><abbr title="Check to verify">?</abbr></th>
139     <th><abbr title="File numbers (ix File)">F#</abbr></th>
140     <th>Files</th>
141     <th>Score</th>
142     </tr>
143     <?
144     $j = $ixStart;
145     foreach ($res as $row)
146     {
147
148         $fVerified = isset($rgVerified[$row->ixFile1][$row->ixFile2]);
149
150         $fCompares = compare_filename(
151             $rgoFiles[$row->ixFile1],
152             $rgoFiles[$row->ixFile2]);
153
154         if ((!isset($_GET['fVerified']) ||
155             ($_GET['fVerified'] == 1 && $fVerified) ||
156             ($_GET['fVerified'] == 0 && !$fVerified))
157             &&
158             (!isset($_GET['fMatch']) ||
159             ($_GET['fMatch'] == 1 && $fCompares) ||
160             ($_GET['fMatch'] == 0 && !$fCompares)))
161         {
162             /*
163              * Output ugly HTML table data.
164              */
165             echo '<tr class="'.($j % 2 == 0 ? 'colored' : '').'.">';

```

```

166     echo '<td rowspan="2">'.($j+1).'

```

#### B.4.8 save\_verified.php

This script saves verified duplicate pairs to the database.

```

1 <?php
2 /**
3  * Simple support script that saves two files as a verified
4  * duplicate.
5  *
6  * @author Vegard Andreas Larsen <vegarl@stud.ntnu.no>
7  */
8 require_once('include/include.php');
9
10 $oDB = CDatabase::Get();
11
12 $st = $oDB->prepare(
13     'INSERT IGNORE INTO tblVerifiedMatch (ixFile1, ixFile2) VALUES (?,
14         ?)',
15     array('integer', 'integer'), MDB2_PREPARE_MANIP);
16
17 /*
18  * We get multiple values as an array of strings.
19  */
20 if (isset($_POST['match']))
21 {
22     $rgMatch = $_POST['match'];
23     foreach ($rgMatch as $value)
24     {
25         // The strings are splittable by a hyphen
26         list($ixFile1, $ixFile2) = split('-', $value);
27         $st->execute(array($ixFile1, $ixFile2));
28     }
29 }
30
31 /*
32  * Someone might want us to send the user back afterwards.
33  */
34 if (isset($_POST['redirect_to']))
35 {
36     header('Location: '.$_POST['redirect_to']);
37     die();
38 }
39 ?>

```

## B.5 Licenses

### B.5.1 The GNU General Public License

Version 2, June 1991

Copyright © 1989, 1991 Free Software Foundation, Inc.

51 Franklin Street, Fifth Floor, Boston, MA 02110-1301, USA

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

#### B.5.1.1 Preamble

The licenses for most software are designed to take away your freedom to share and change it. By contrast, the GNU General Public License is intended to guarantee your freedom to share and change free software—to make sure the software is free for all its users. This General Public License applies to most of the Free Software Foundation's software and to any other program whose authors commit to using it. (Some other Free Software Foundation software is covered by the GNU Library General Public License instead.) You can apply it to your programs, too.

When we speak of free software, we are referring to freedom, not price. Our General Public Licenses are designed to make sure that you have the freedom to distribute copies of free software (and charge for this service if you wish), that you receive source code or can get it if you want it, that you can change the software or use pieces of it in new free programs; and that you know you can do these things.

To protect your rights, we need to make restrictions that forbid anyone to deny you these rights or to ask you to surrender the rights. These restrictions translate to certain responsibilities for you if you distribute copies of the software, or if you modify it.

For example, if you distribute copies of such a program, whether gratis or for a fee, you must give the recipients all the rights that you have. You must make sure that they, too, receive or can get the source code. And you must show them these terms so they know their rights.

We protect your rights with two steps: (1) copyright the software, and (2) offer you this license which gives you legal permission to copy, distribute and/or modify the software.

Also, for each author's protection and ours, we want to make certain that everyone understands that there is no warranty for this free software. If the software is modified by someone else and passed on, we want its recipients to know that what they have is not the original, so that any problems introduced by others will not reflect on the original authors' reputations.

Finally, any free program is threatened constantly by software patents. We wish to avoid the danger that redistributors of a free program will individually obtain patent licenses, in effect making the program proprietary. To prevent this, we have made it clear that any patent must be licensed for everyone's free use or not licensed at all.

The precise terms and conditions for copying, distribution and modification follow.

### B.5.1.2 Terms and Conditions For Copying, Distribution and Modification

0. This License applies to any program or other work which contains a notice placed by the copyright holder saying it may be distributed under the terms of this General Public License. The “Program”, below, refers to any such program or work, and a “work based on the Program” means either the Program or any derivative work under copyright law: that is to say, a work containing the Program or a portion of it, either verbatim or with modifications and/or translated into another language. (Hereinafter, translation is included without limitation in the term “modification”.) Each licensee is addressed as “you”.

Activities other than copying, distribution and modification are not covered by this License; they are outside its scope. The act of running the Program is not restricted, and the output from the Program is covered only if its contents constitute a work based on the Program (independent of having been made by running the Program). Whether that is true depends on what the Program does.

1. You may copy and distribute verbatim copies of the Program’s source code as you receive it, in any medium, provided that you conspicuously and appropriately publish on each copy an appropriate copyright notice and disclaimer of warranty; keep intact all the notices that refer to this License and to the absence of any warranty; and give any other recipients of the Program a copy of this License along with the Program.

You may charge a fee for the physical act of transferring a copy, and you may at your option offer warranty protection in exchange for a fee.

2. You may modify your copy or copies of the Program or any portion of it, thus forming a work based on the Program, and copy and distribute such modifications or work under the terms of Section 1 above, provided that you also meet all of these conditions:
  - (a) You must cause the modified files to carry prominent notices stating that you changed the files and the date of any change.
  - (b) You must cause any work that you distribute or publish, that in whole or in part contains or is derived from the Program or any part thereof, to be licensed as a whole at no charge to all third parties under the terms of this License.
  - (c) If the modified program normally reads commands interactively when run, you must cause it, when started running for such interactive use in the most ordinary way, to print or display an announcement including an appropriate copyright notice and a notice that there is no warranty (or else, saying that you provide a warranty) and that users may redistribute the program under these conditions, and telling the user how to view a copy of this License. (Exception: if the Program itself is interactive but does not normally print such an announcement, your work based on the Program is not required to print an announcement.)

These requirements apply to the modified work as a whole. If identifiable sections of that work are not derived from the Program, and can be reasonably considered independent and separate works in themselves, then this License, and its terms, do not apply to those sections when you distribute them as separate works. But when you distribute the same sections as part of a whole which is a work based on the Program, the distribution of the whole must be on the terms of this License, whose permissions for other licensees extend to the entire whole, and thus to each and every part regardless of who wrote it.

Thus, it is not the intent of this section to claim rights or contest your rights to work written entirely by you; rather, the intent is to exercise the right to control the distribution of derivative or collective works based on the Program.

In addition, mere aggregation of another work not based on the Program with the Program (or with a work based on the Program) on a volume of a storage or distribution medium does not bring the other work under the scope of this License.

3. You may copy and distribute the Program (or a work based on it, under Section 2) in object code or executable form under the terms of Sections 1 and 2 above provided that you also do one of the following:
  - (a) Accompany it with the complete corresponding machine-readable source code, which must be distributed under the terms of Sections 1 and 2 above on a medium customarily used for software interchange; or,
  - (b) Accompany it with a written offer, valid for at least three years, to give any third party, for a charge no more than your cost of physically performing source distribution, a complete machine-readable copy of the corresponding source code, to be distributed under the terms of Sections 1 and 2 above on a medium customarily used for software interchange; or,
  - (c) Accompany it with the information you received as to the offer to distribute corresponding source code. (This alternative is allowed only for noncommercial distribution and only if you received the program in object code or executable form with such an offer, in accord with Subsection b above.)

The source code for a work means the preferred form of the work for making modifications to it. For an executable work, complete source code means all the source code for all modules it contains, plus any associated interface definition files, plus the scripts used to control compilation and installation of the executable. However, as a special exception, the source code distributed need not include anything that is normally distributed (in either source or binary form) with the major components (compiler, kernel, and so on) of the operating system on which the executable runs, unless that component itself accompanies the executable.

If distribution of executable or object code is made by offering access to copy from a designated place, then offering equivalent access to copy the

source code from the same place counts as distribution of the source code, even though third parties are not compelled to copy the source along with the object code.

4. You may not copy, modify, sublicense, or distribute the Program except as expressly provided under this License. Any attempt otherwise to copy, modify, sublicense or distribute the Program is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.
5. You are not required to accept this License, since you have not signed it. However, nothing else grants you permission to modify or distribute the Program or its derivative works. These actions are prohibited by law if you do not accept this License. Therefore, by modifying or distributing the Program (or any work based on the Program), you indicate your acceptance of this License to do so, and all its terms and conditions for copying, distributing or modifying the Program or works based on it.
6. Each time you redistribute the Program (or any work based on the Program), the recipient automatically receives a license from the original licensor to copy, distribute or modify the Program subject to these terms and conditions. You may not impose any further restrictions on the recipients' exercise of the rights granted herein. You are not responsible for enforcing compliance by third parties to this License.
7. If, as a consequence of a court judgment or allegation of patent infringement or for any other reason (not limited to patent issues), conditions are imposed on you (whether by court order, agreement or otherwise) that contradict the conditions of this License, they do not excuse you from the conditions of this License. If you cannot distribute so as to satisfy simultaneously your obligations under this License and any other pertinent obligations, then as a consequence you may not distribute the Program at all. For example, if a patent license would not permit royalty-free redistribution of the Program by all those who receive copies directly or indirectly through you, then the only way you could satisfy both it and this License would be to refrain entirely from distribution of the Program.

If any portion of this section is held invalid or unenforceable under any particular circumstance, the balance of the section is intended to apply and the section as a whole is intended to apply in other circumstances.

It is not the purpose of this section to induce you to infringe any patents or other property right claims or to contest validity of any such claims; this section has the sole purpose of protecting the integrity of the free software distribution system, which is implemented by public license practices. Many people have made generous contributions to the wide range of software distributed through that system in reliance on consistent application of that system; it is up to the author/donor to decide if he or she is willing to distribute software through any other system and a licensee cannot impose that choice.

This section is intended to make thoroughly clear what is believed to be a consequence of the rest of this License.

8. If the distribution and/or use of the Program is restricted in certain countries either by patents or by copyrighted interfaces, the original copyright holder who places the Program under this License may add an explicit geographical distribution limitation excluding those countries, so that distribution is permitted only in or among countries not thus excluded. In such case, this License incorporates the limitation as if written in the body of this License.
9. The Free Software Foundation may publish revised and/or new versions of the General Public License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns.

Each version is given a distinguishing version number. If the Program specifies a version number of this License which applies to it and “any later version”, you have the option of following the terms and conditions either of that version or of any later version published by the Free Software Foundation. If the Program does not specify a version number of this License, you may choose any version ever published by the Free Software Foundation.

10. If you wish to incorporate parts of the Program into other free programs whose distribution conditions are different, write to the author to ask for permission. For software which is copyrighted by the Free Software Foundation, write to the Free Software Foundation; we sometimes make exceptions for this. Our decision will be guided by the two goals of preserving the free status of all derivatives of our free software and of promoting the sharing and reuse of software generally.

## NO WARRANTY

11. BECAUSE THE PROGRAM IS LICENSED FREE OF CHARGE, THERE IS NO WARRANTY FOR THE PROGRAM, TO THE EXTENT PERMITTED BY APPLICABLE LAW. EXCEPT WHEN OTHERWISE STATED IN WRITING THE COPYRIGHT HOLDERS AND/OR OTHER PARTIES PROVIDE THE PROGRAM “AS IS” WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE ENTIRE RISK AS TO THE QUALITY AND PERFORMANCE OF THE PROGRAM IS WITH YOU. SHOULD THE PROGRAM PROVE DEFECTIVE, YOU ASSUME THE COST OF ALL NECESSARY SERVICING, REPAIR OR CORRECTION.
12. IN NO EVENT UNLESS REQUIRED BY APPLICABLE LAW OR AGREED TO IN WRITING WILL ANY COPYRIGHT HOLDER, OR ANY OTHER PARTY WHO MAY MODIFY AND/OR REDISTRIBUTE THE PROGRAM AS PERMITTED ABOVE, BE LIABLE TO YOU FOR DAMAGES, INCLUDING ANY GENERAL, SPECIAL, INCIDENTAL OR CONSEQUENTIAL DAMAGES ARISING OUT OF THE USE OR INABILITY TO USE THE PROGRAM (INCLUDING BUT NOT LIMITED TO LOSS OF DATA OR DATA BEING RENDERED INACCURATE OR LOSSES SUSTAINED BY YOU OR THIRD PARTIES OR A FAILURE OF THE PROGRAM TO OPERATE WITH ANY OTHER PROGRAMS), EVEN IF SUCH

HOLDER OR OTHER PARTY HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

## END OF TERMS AND CONDITIONS

### Appendix: How to Apply These Terms to Your New Programs

If you develop a new program, and you want it to be of the greatest possible use to the public, the best way to achieve this is to make it free software which everyone can redistribute and change under these terms.

To do so, attach the following notices to the program. It is safest to attach them to the start of each source file to most effectively convey the exclusion of warranty; and each file should have at least the “copyright” line and a pointer to where the full notice is found.

```
one line to give the program's name and a brief idea of what it does.
Copyright (C) yyyy name of author
```

```
This program is free software; you can redistribute it and/or modify
it under the terms of the GNU General Public License as published
by the Free Software Foundation; either version 2 of the License, or
(at your option) any later version.
```

```
This program is distributed in the hope that it will be useful, but
WITHOUT ANY WARRANTY; without even the implied warranty
of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU
General Public License for more details.
```

```
You should have received a copy of the GNU General Public License
along with this program; if not, write to the Free Software Founda-
tion, Inc., 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301,
USA.
```

Also add information on how to contact you by electronic and paper mail.

If the program is interactive, make it output a short notice like this when it starts in an interactive mode:

```
Gnomovision version 69, Copyright (C) yyyy name of author
Gnomovision comes with ABSOLUTELY NO WARRANTY; for de-
tails type 'show w'.
This is free software, and you are welcome to redistribute it under
certain conditions; type 'show c' for details.
```

The hypothetical commands `show w` and `show c` should show the appropriate parts of the General Public License. Of course, the commands you use may be called something other than `show w` and `show c`; they could even be mouse-clicks or menu items—whatever suits your program.

You should also get your employer (if you work as a programmer) or your school, if any, to sign a “copyright disclaimer” for the program, if necessary. Here is a sample; alter the names:



Yoyodyne, Inc., hereby disclaims all copyright interest in the program 'Gnomovision' (which makes passes at compilers) written by James Hacker.

signature of Ty Coon, 1 April 1989  
Ty Coon, President of Vice

This General Public License does not permit incorporating your program into proprietary programs. If your program is a subroutine library, you may consider it more useful to permit linking proprietary applications with the library. If this is what you want to do, use the GNU Library General Public License instead of this License.

## B.5.2 Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States

### B.5.2.1 Short text

You are free:

**to Share** to copy, distribute, display and perform the work

**to Remix** to make derivative works

Under the following conditions:

**Attribution.** You must attribute the work in the manner specified by the author or licensor (but not in any way that suggests that they endorse you or your use of the work).

**Noncommercial.** You may not use this work for commercial purposes.

**Share Alike.** If you alter, transform, or build upon this work, you may distribute the resulting work only under the same or similar license to this one.

- For any reuse or distribution, you must make clear to others the license terms of this work. The best way to do this is with a link to this web page.
- Any of the above conditions can be waived if you get permission from the copyright holder.
- Apart from the remix rights granted under this license, nothing in this license impairs or restricts the author's moral rights.

### B.5.2.2 Full license

THE WORK (AS DEFINED BELOW) IS PROVIDED UNDER THE TERMS OF THIS CREATIVE COMMONS PUBLIC LICENSE ("CCPL" OR "LICENSE"). THE WORK IS PROTECTED BY COPYRIGHT AND/OR OTHER

APPLICABLE LAW. ANY USE OF THE WORK OTHER THAN AS AUTHORIZED UNDER THIS LICENSE OR COPYRIGHT LAW IS PROHIBITED.

BY EXERCISING ANY RIGHTS TO THE WORK PROVIDED HERE, YOU ACCEPT AND AGREE TO BE BOUND BY THE TERMS OF THIS LICENSE. TO THE EXTENT THIS LICENSE MAY BE CONSIDERED TO BE A CONTRACT, THE LICENSOR GRANTS YOU THE RIGHTS CONTAINED HERE IN CONSIDERATION OF YOUR ACCEPTANCE OF SUCH TERMS AND CONDITIONS.

#### 1. Definitions

- (a) "Collective Work" means a work, such as a periodical issue, anthology or encyclopedia, in which the Work in its entirety in unmodified form, along with one or more other contributions, constituting separate and independent works in themselves, are assembled into a collective whole. A work that constitutes a Collective Work will not be considered a Derivative Work (as defined below) for the purposes of this License.
  - (b) "Derivative Work" means a work based upon the Work or upon the Work and other pre-existing works, such as a translation, musical arrangement, dramatization, fictionalization, motion picture version, sound recording, art reproduction, abridgment, condensation, or any other form in which the Work may be recast, transformed, or adapted, except that a work that constitutes a Collective Work will not be considered a Derivative Work for the purpose of this License. For the avoidance of doubt, where the Work is a musical composition or sound recording, the synchronization of the Work in timed-relation with a moving image ("synching") will be considered a Derivative Work for the purpose of this License.
  - (c) "Licensor" means the individual, individuals, entity or entities that offer(s) the Work under the terms of this License.
  - (d) "Original Author" means the individual, individuals, entity or entities who created the Work.
  - (e) "Work" means the copyrightable work of authorship offered under the terms of this License.
  - (f) "You" means an individual or entity exercising rights under this License who has not previously violated the terms of this License with respect to the Work, or who has received express permission from the Licensor to exercise rights under this License despite a previous violation.
  - (g) "License Elements" means the following high-level license attributes as selected by Licensor and indicated in the title of this License: Attribution, Noncommercial, ShareAlike.
2. Fair Use Rights. Nothing in this license is intended to reduce, limit, or restrict any rights arising from fair use, first sale or other limitations on the exclusive rights of the copyright owner under copyright law or other applicable laws.

3. License Grant. Subject to the terms and conditions of this License, Licensor hereby grants You a worldwide, royalty-free, non-exclusive, perpetual (for the duration of the applicable copyright) license to exercise the rights in the Work as stated below:
  - (a) to reproduce the Work, to incorporate the Work into one or more Collective Works, and to reproduce the Work as incorporated in the Collective Works;
  - (b) to create and reproduce Derivative Works provided that any such Derivative Work, including any translation in any medium, takes reasonable steps to clearly label, demarcate or otherwise identify that changes were made to the original Work. For example, a translation could be marked "The original work was translated from English to Spanish," or a modification could indicate "The original work has been modified.";
  - (c) to distribute copies or phonorecords of, display publicly, perform publicly, and perform publicly by means of a digital audio transmission the Work including as incorporated in Collective Works;
  - (d) to distribute copies or phonorecords of, display publicly, perform publicly, and perform publicly by means of a digital audio transmission Derivative Works;

The above rights may be exercised in all media and formats whether now known or hereafter devised. The above rights include the right to make such modifications as are technically necessary to exercise the rights in other media and formats. All rights not expressly granted by Licensor are hereby reserved, including but not limited to the rights set forth in Sections 4(e) and 4(f).

4. Restrictions. The license granted in Section 3 above is expressly made subject to and limited by the following restrictions:
  - (a) You may distribute, publicly display, publicly perform, or publicly digitally perform the Work only under the terms of this License, and You must include a copy of, or the Uniform Resource Identifier for, this License with every copy or phonorecord of the Work You distribute, publicly display, publicly perform, or publicly digitally perform. You may not offer or impose any terms on the Work that restrict the terms of this License or the ability of a recipient of the Work to exercise the rights granted to that recipient under the terms of the License. You may not sublicense the Work. You must keep intact all notices that refer to this License and to the disclaimer of warranties. When You distribute, publicly display, publicly perform, or publicly digitally perform the Work, You may not impose any technological measures on the Work that restrict the ability of a recipient of the Work from You to exercise the rights granted to that recipient under the terms of the License. This Section 4(a) applies to the Work as incorporated in a Collective Work, but this does not require the Collective Work apart from the Work itself to be made subject to the terms of this License. If You create a Collective Work,

upon notice from any Licensor You must, to the extent practicable, remove from the Collective Work any credit as required by Section 4(d), as requested. If You create a Derivative Work, upon notice from any Licensor You must, to the extent practicable, remove from the Derivative Work any credit as required by Section 4(d), as requested.

- (b) You may distribute, publicly display, publicly perform, or publicly digitally perform a Derivative Work only under: (i) the terms of this License; (ii) a later version of this License with the same License Elements as this License; or, (iii) either the unported Creative Commons license or a Creative Commons license for another jurisdiction (either this or a later license version) that contains the same License Elements as this License (e.g. Attribution-NonCommercial-ShareAlike 3.0 (Unported)) ("the Applicable License"). You must include a copy of, or the Uniform Resource Identifier for, the Applicable License with every copy or phonorecord of each Derivative Work You distribute, publicly display, publicly perform, or publicly digitally perform. You may not offer or impose any terms on the Derivative Works that restrict the terms of the Applicable License or the ability of a recipient of the Work to exercise the rights granted to that recipient under the terms of the Applicable License. You must keep intact all notices that refer to the Applicable License and to the disclaimer of warranties. When You distribute, publicly display, publicly perform, or publicly digitally perform the Derivative Work, You may not impose any technological measures on the Derivative Work that restrict the ability of a recipient of the Derivative Work from You to exercise the rights granted to that recipient under the terms of the Applicable License. This Section 4(b) applies to the Derivative Work as incorporated in a Collective Work, but this does not require the Collective Work apart from the Derivative Work itself to be made subject to the terms of the Applicable License.
- (c) You may not exercise any of the rights granted to You in Section 3 above in any manner that is primarily intended for or directed toward commercial advantage or private monetary compensation. The exchange of the Work for other copyrighted works by means of digital file-sharing or otherwise shall not be considered to be intended for or directed toward commercial advantage or private monetary compensation, provided there is no payment of any monetary compensation in connection with the exchange of copyrighted works.
- (d) If You distribute, publicly display, publicly perform, or publicly digitally perform the Work (as defined in Section 1 above) or any Derivative Works (as defined in Section 1 above) or Collective Works (as defined in Section 1 above), You must, unless a request has been made pursuant to Section 4(a), keep intact all copyright notices for the Work and provide, reasonable to the medium or means You are utilizing: (i) the name of the Original Author (or pseudonym, if applicable) if supplied, and/or (ii) if the Original Author and/or Licensor designate another party or parties (e.g. a sponsor institute, publishing entity, journal) for attribution ("Attribution Parties") in Licensor's copyright notice, terms of service or by other reasonable means,

the name of such party or parties; the title of the Work if supplied; to the extent reasonably practicable, the Uniform Resource Identifier, if any, that Licensor specifies to be associated with the Work, unless such URI does not refer to the copyright notice or licensing information for the Work; and, consistent with Section 3(b) in the case of a Derivative Work, a credit identifying the use of the Work in the Derivative Work (e.g., "French translation of the Work by Original Author," or "Screenplay based on original Work by Original Author"). The credit required by this Section 4(d) may be implemented in any reasonable manner; provided, however, that in the case of a Derivative Work or Collective Work, at a minimum such credit will appear, if a credit for all contributing authors of the Derivative Work or Collective Work appears, then as part of these credits and in a manner at least as prominent as the credits for the other contributing authors. For the avoidance of doubt, You may only use the credit required by this Section for the purpose of attribution in the manner set out above and, by exercising Your rights under this License, You may not implicitly or explicitly assert or imply any connection with, sponsorship or endorsement by the Original Author, Licensor and/or Attribution Parties, as appropriate, of You or Your use of the Work, without the separate, express prior written permission of the Original Author, Licensor and/or Attribution Parties.

- (e) For the avoidance of doubt, where the Work is a musical composition:
  - i. Performance Royalties Under Blanket Licenses. Licensor reserves the exclusive right to collect whether individually or, in the event that Licensor is a member of a performance rights society (e.g. ASCAP, BMI, SESAC), via that society, royalties for the public performance or public digital performance (e.g. webcast) of the Work if that performance is primarily intended for or directed toward commercial advantage or private monetary compensation.
  - ii. Mechanical Rights and Statutory Royalties. Licensor reserves the exclusive right to collect, whether individually or via a music rights agency or designated agent (e.g. Harry Fox Agency), royalties for any phonorecord You create from the Work ("cover version") and distribute, subject to the compulsory license created by 17 USC Section 115 of the US Copyright Act (or the equivalent in other jurisdictions), if Your distribution of such cover version is primarily intended for or directed toward commercial advantage or private monetary compensation.
- (f) Webcasting Rights and Statutory Royalties. For the avoidance of doubt, where the Work is a sound recording, Licensor reserves the exclusive right to collect, whether individually or via a performance-rights society (e.g. SoundExchange), royalties for the public digital performance (e.g. webcast) of the Work, subject to the compulsory license created by 17 USC Section 114 of the US Copyright Act (or the equivalent in other jurisdictions), if Your public digital performance is primarily intended for or directed toward commercial advantage or private monetary compensation.

5. Representations, Warranties and Disclaimer  
UNLESS OTHERWISE MUTUALLY AGREED TO BY THE PARTIES IN WRITING, LICENSOR OFFERS THE WORK AS-IS AND ONLY TO THE EXTENT OF ANY RIGHTS HELD IN THE LICENSED WORK BY THE LICENSOR. THE LICENSOR MAKES NO REPRESENTATIONS OR WARRANTIES OF ANY KIND CONCERNING THE WORK, EXPRESS, IMPLIED, STATUTORY OR OTHERWISE, INCLUDING, WITHOUT LIMITATION, WARRANTIES OF TITLE, MARKETABILITY, MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, NONINFRINGEMENT, OR THE ABSENCE OF LATENT OR OTHER DEFECTS, ACCURACY, OR THE PRESENCE OF ABSENCE OF ERRORS, WHETHER OR NOT DISCOVERABLE. SOME JURISDICTIONS DO NOT ALLOW THE EXCLUSION OF IMPLIED WARRANTIES, SO SUCH EXCLUSION MAY NOT APPLY TO YOU.
6. Limitation on Liability.  
EXCEPT TO THE EXTENT REQUIRED BY APPLICABLE LAW, IN NO EVENT WILL LICENSOR BE LIABLE TO YOU ON ANY LEGAL THEORY FOR ANY SPECIAL, INCIDENTAL, CONSEQUENTIAL, PUNITIVE OR EXEMPLARY DAMAGES ARISING OUT OF THIS LICENSE OR THE USE OF THE WORK, EVEN IF LICENSOR HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.
7. Termination
  - (a) This License and the rights granted hereunder will terminate automatically upon any breach by You of the terms of this License. Individuals or entities who have received Derivative Works (as defined in Section 1 above) or Collective Works (as defined in Section 1 above) from You under this License, however, will not have their licenses terminated provided such individuals or entities remain in full compliance with those licenses. Sections 1, 2, 5, 6, 7, and 8 will survive any termination of this License.
  - (b) Subject to the above terms and conditions, the license granted here is perpetual (for the duration of the applicable copyright in the Work). Notwithstanding the above, Licensor reserves the right to release the Work under different license terms or to stop distributing the Work at any time; provided, however that any such election will not serve to withdraw this License (or any other license that has been, or is required to be, granted under the terms of this License), and this License will continue in full force and effect unless terminated as stated above.
8. Miscellaneous
  - (a) Each time You distribute or publicly digitally perform the Work (as defined in Section 1 above) or a Collective Work (as defined in Section 1 above), the Licensor offers to the recipient a license to the Work on the same terms and conditions as the license granted to You under this License.

- (b) Each time You distribute or publicly digitally perform a Derivative Work, Licensor offers to the recipient a license to the original Work on the same terms and conditions as the license granted to You under this License.
- (c) If any provision of this License is invalid or unenforceable under applicable law, it shall not affect the validity or enforceability of the remainder of the terms of this License, and without further action by the parties to this agreement, such provision shall be reformed to the minimum extent necessary to make such provision valid and enforceable.
- (d) No term or provision of this License shall be deemed waived and no breach consented to unless such waiver or consent shall be in writing and signed by the party to be charged with such waiver or consent.
- (e) This License constitutes the entire agreement between the parties with respect to the Work licensed here. There are no understandings, agreements or representations with respect to the Work not specified here. Licensor shall not be bound by any additional provisions that may appear in any communication from You. This License may not be modified without the mutual written agreement of the Licensor and You.





# Appendix C

## Music collection

Due to space constraints only the albums are listed in this appendix. A complete list of songs can be found in the accompanying database. All of the songs were in a folder structure: <artist>/<album>/<song>. This list was generated by parsing the <artist>/<album>-structure, and cleaning the list manually.

### A

- ACDC
  - ACDC - High Voltage
  - Back In Black
  - Dirty Deeds Done Dirt Cheap
  - For Those About To Rock
  - High Voltage
  - Highway to Hell
  - If You Want Blood You've Got It
  - Let There Be Rock
  - Powerage
- ATB
  - Addicted To Music
  - Dedicated
  - No Silence
- Aerosmith
  - Aerosmith 1973
  - Get Your Wings 1974
  - Toys In The Attic 1975
  - Rocks 1976
  - Draw the Line 1978
  - Night In The Ruts 1980
  - Rock In A Hard Place 1982
  - Done With Mirrors 1985
- Permanent Vacation 1987
- Pump 1989
- Get A Grip 1993
- Nine Lives 1998
- Just Push Play 2001
- Al Di Meola
  - 1976 - Elegant Gypsy
  - 1976 - Land Of The Midnight Sun
  - 1977 - Casino
  - 1979 - Splendido Hotel
  - 1981 - Electric Rendevous
  - 1983 - Scenario
  - 1987 - Tirami Su
  - 1988 - Kiss My Axe
  - 1990 - World Sinfonia
  - 1993 - World Sinfonia - Heart Of Immigrant
  - 1994 - Orange And Blue
- Al Stewart
  - Year of the Cat (remastered)
- Alanis Morissette
  - Jagged Little Pill
  - Jagged Little Pill Acoustic
  - Supposed Former Infatuation Junkie
- Allan Edwall

- Den lilla bäcken
  - Edwalls Blandning 1979-84
  - Mina Visor 1 - Färdknäpp
  - Mina Visor 1 - Grovdoppa
  - Mina Visor 3 - Aftonro
  - Mina Visor 2 - Gnällspik
  - Mina Visor 2 - Vetahuteri
  - Ramsor om Dom och Oss
  - Alphaville
    - First Harvest 1984-1992
  - Amalia
    - Com Que Voz
  - Amy Winehouse
    - Back To Black
    - Frank
  - Annie Lennox
    - Bare
  - Apoptygma Berzerk
    - 7
    - APBL2000 (Live 2000 Version)
    - Black
    - Harmonizer
    - Kathy's Song (6-Track Maxi-Single)
    - The Apocalyptic Manifesto
    - Welcome To Earth
    - You And Me Against The World
  - Arctic Monkeys
    - Favourite Worst Nightmare
  - Astor Piazzolla
    - Armagedon (1977)
    - In Concert (1983)
    - La Camorra (1989)
    - Libertango (1981)
    - Tanguedia de Amor (1984)
    - The Rough Dancer And The Cyclical Night (1989)
    - Astor Piazzola & Gary Burton 'The New Tango' (1986)
    - Tango Sensations (1994)
  - Astor Piazzalla & David Tononbaum - El Porteno (1994)
  - Tango - Zero Hour (1986)
  - Eight Seasons (1996)
  - Musicues De Films (Tango, Henri IV) (1986)
  - Gidon Kremer - Hommage a Piazzolla (1996)
  - The Lausanne Concert (1989)
  - El Nuevo Tango de Buenos Aires (1989)
  - El Tango (1997)
  - Bandoneon Sinfonico (1990)
  - Maria De Buenos Aires Tango Operita CD1 (1998)
  - Ballet Tango (1992)
  - Maria De Buenos Aires Tango Operita CD2 (1998)
  - Concierto De Nacar (1997)
  - Tres tangos - concierto (1999)
  - Live At The 'Bouffes Du Nord' (1998)
- B**
- BT
    - Emotional Technology
  - Beatles
    - Let It Be... Naked
    - Revolver & Magical Mystery Tour
    - White Album (disc 1)
    - White Album (disc 2)
  - Beethoven
    - Violin Sonatas
  - Ben Webster
    - Jazz Ballads CD2
    - Jazz Ballads CD1
  - Billie Holiday
    - 1949-52 Radio & TV Broadcasts
    - 1953-56 Radio and TV Broadcasts
    - A Fine Romance
    - Anthology 1944-1959 Disc 1
    - Anthology 1944-1959 Disc 2
    - At Monterey

- Billie Holiday as Stratford '57
- Billie Holiday at Storyville
- Billie's Best
- Billie's Blues
- Broadcast Performances Volume 4
- Broadcast Performances Volume 3 (1956 - 58)
- Compact Jazz: Billie Holiday
- Control Booth Series, Vol. 1 1940-1941
- Control Booth Series, Vol. 2
- Greatest Hits
- I Loves You Porgy
- Jazz at the Philharmonic
- Lady Day & Prez 1937-1941
- Lady Day - The Best Of Billie Holiday (Disc 1)
- Lady Day - The Best Of Billie Holiday (Disc 2)
- Lady Day - The Storyville Concerts (Disc 2 of 2)
- Lady Day - The Storyville Concerts (Disk 1 of 2)
- Lady Day: The Complete Billie Holiday on Columbia (1933-1944) - Disc 1 through 10
- Lady Live!
- Lady Sings the Blues
- Lady in Satin
- My Man
- New Orleans
- Remixed Hits
- Solitude
- Songs for Distingué Lovers
- Storyville - Masters of Jazz
- The Best of Billie Holiday
- The Commodore Master Takes
- The Complete 1951 Storyville Club Sessions
- The Complete Billie Holiday On Verve, 1945-1959
- The Complete Billie Holiday on Verve 1945-1959 (1 through 10)
- The Complete Commodore Recordings (1 of 2)
- The Complete Commodore Recordings (2 of 2)
- The Complete Decca Recordings (1 of 2)
- The Complete Decca Recordings (2 of 2)
- The Complete Decca Recordings (US Release)
- The Complete Verve Studio Master Takes (Disc 1 through 6)
- The Quintessential Billie Holiday Volume 9 (1940-1942)
- The Quintessential Billie Holiday, Vol. 1 through 9
- The Sensitive Billie Holiday 1940-1949
- Vdisc Musical Contribution For Our Armed Forces Overseas
- Volume 10: 1940-1941
- Billy Idol
  - Greatest Hits
  - Rebel Yell (Expanded Edition)
- Björk
  - Greatest Hits
- Björk Guðmundsdóttir & Trió Guðmundar Ingólfssonar
  - Gling-Gló
- Black Eyed Peas
  - Elephunk
  - Monkey Business
- Bowie, David
  - Reality Bonus CD
- Brian Eno
  - Ambient 1
  - Another Day On Earth
  - Another Green World
  - Before and after Science
  - Nerve Net
- Brian Eno & David Byrne
  - My Life In The Bush Of Ghosts (Remastered)
  - My Life in the Bush of Ghosts
- Brian Eno & John Cale
  - Wrong Way Up
- Bryan Ferry
  - As Time Goes By

## C

- Charlie Haden
  - Land of the Sun
- Craig David
  - Born To Do It
- Crash Test Dummies
  - God Shuffled His Feet
- Creedence Clearwater Revival
  - Chronicle- 24-Karat Gold Disc
  - The Concert
- Tonight (1984)
- Diamond Dogs (1974)
- Never Let Me Down (1987)
- David Live
- Young Americans (1975)
- 1.Outside
- Singles
- Stage
- Heathen (Bonus Disc)
- Heathen (Disc 1)
- Reality
- Reality (CD 1)

## D

- DJ Bobo
  - Just For You
  - The Ultimate Megamix 99
- Dalbello
  - whomanfoursays
- Damien Rice
  - O
- Dane Cook
  - Retaliation
- David Bowie
  - David Bowie (1969)
  - Station To Station (1976)
  - Low (1977)
  - Space Oddity (1969)
  - Heroes (1977)
  - The Man Who Sold The World (1970)
  - Hunky Dory (1971)
  - Lodger (1979)
  - Scary Monsters (1980)
  - The Rise and Fall of Ziggy Stardust (1972)
  - Aladdin Sane (1973)
  - Let's Dance (1983)
  - Pin Ups (1973)
- David Byrne
  - David Byrne
  - Feelings
  - Grown Backwards
  - Look Into The Eyeball
  - Rei Momo
  - Uh-Oh
- Dido
  - Life For Rent
  - No Angel
- Don McLean
  - American Pie
- Doors
  - Waiting For The Sun
- Dr Hook & The Medicine Show
  - The Very Best Of
- Dune
  - Dune
  - Forever
- deLillos
  - Festen er ikke over... det er kake igjen Disc 1
  - Festen er ikke over... det er kake igjen Disc 2
  - Før Var Det Morsomt Med Sne
  - Ikke gå
  - Kast Alle Papirene
  - Kjerringvik-demoen del 1

- Kjerringvik-demoen del 2
  - Mere Disc 1
  - Mere Disc 2
  - Midt i begynnelsen
  - Suser Avgårde
  - Varme Mennesker
  - deLillos'85
    - Suser videre
    - evig forelsket da
- E**
- Echo And The Bunnymen
    - Heaven Up Here
  - Eiffel 65
    - Europop
  - Enya
    - A Day Without Rain
    - Amarantine
    - Paint The Sky With Stars
    - Shepherd Moons
    - The Celts
    - The Memory Of Trees
    - Watermark
  - Erik Satie
    - Gymnopedies
- F**
- Faithless
    - Forever Faithless- The Greatest Hits
  - Fragma
    - Embrace
    - Toca
  - Franz Ferdinand
    - Franz Ferdinand
    - Franz Ferdinand (Special Edition)
    - Remixes
    - This Ffire
    - You Could Have It So Much Better
- G**
- Gary Numan
    - Documents
    - Replicas
    - Telekon
    - The Pleasure Principle
    - Tubeway Army
    - Warriors
  - Gina G
    - Fresh
  - Glenn Miller
    - The Glenn Miller Story
  - Green Day
    - International Superhits!
  - Grinderman
    - Grinderman
  - Gwen Stefani
    - Love Angel Music Baby
    - The Sweet Escape
  - Gåte
    - Jygri
- H**
- Hampton The Hampster
    - The Hampsterdance Song
  - Hugh Cornwell
    - Beyond Elysian Fields
    - Footprints In The Desert
    - Guilty
    - Hi Fi

## I

- Ian Dury
  - Reasons To Be Cheerful - The Best Of Ian Dury (Disc1)
  - Reasons To Be Cheerful - The Best Of Ian Dury (Disc2)
- Ibrahim Ferrer
  - Buena Vista Social Club Presents Ibrahim Ferrer
- Iggy Pop
  - Blah-Blah-Blah
  - Brick By Brick
  - Lust For Life
  - New Values
  - Pop Music
- Imperiet
  - Alltid Rött Alltid Rätt: en samling 1983-88
- Infernal
  - From Paris To Berlin
- iio
  - Poetica

## J

- James Blunt
  - All The Lost Souls (Bonus Track)
  - Back To Bedlam (Edited)
- Jan Johansson
  - Folkvisor
  - Jazz på svenska
- Japan
  - Gentlemen Take Polaroids
  - Tin Drum
- Jerry Harrison
  - Casual Gods
- Jethro Tull

- Aqualung

- John Cale
  - 5 tracks
  - Black Acetat
  - Fragments Of A Rainy Season
  - Hobosapiens
  - Island Years (Disc 1)
  - Island Years (Disc 2)
  - Walking On Locusts
- John Legend
  - Get Lifted
- Johnny Cash
  - American III - Solitary Man
  - American IV The Man Comes Around
  - American V: A Hundred Highways
  - Unchained
  - Unearthed Volume 4: My Mother's Hymn Book

## K

- K.C. & the Sunshine Band
  - Shake Your Booty
- Kaizers Orchestra
  - Død manns tango
  - Evig Pint
  - Maestro
  - Mann Mot Mann [EP]
  - Ompa Til Du Dør
- Karin Krog
  - Where you at?
- Keith Jarret
  - 1977 - Byablue
  - 1980 - Sacred Hymns
  - 1980 - The Celestial Hawk
  - 1981 - Concerts
  - 1983 - Changes
  - 1983 - Standards Vol 1
  - 1983 - Standards Vol 2
  - 1987 - Changeless

- 1990 - Paris Concert
  - 1991 - Bye Bye Blackbird
  - 1991 - Vienna Concert
  - 1999 - The Melody At Night With You
  - Keith Jarrett
    - My Song
    - Personal Mountains
  - Ketil Bjornstad and David Darling
    - The River
  - Kevin Bloody Wilson
    - 20 Years Of Kev
    - Born again piss tank
    - Kalgoorlie Love Songs
    - Kev's back (The Return of the Yobbo)
    - Kev's Kristmas
    - Let's call him ... Kev!
    - Let Loose Live In The Outback
    - My Australian roots
    - The second kumin of Kev
    - The worst of Kevin Bloody Wilson
    - Youre average Australian yobbo
  - Kevin Coyne
    - Marjory Razorblade
    - Pointing the Finger
    - Sign Of The Times
  - Kraftwerk
    - Autobahn
    - The Man - Machine
  - Kylie Minogue
    - Body Language
- L**
- Laila Dalseth
    - one of a kind - CD1
    - one of a kind - CD2
  - Lars Winnerbäck
    - Daugava
  - Led Zeppelin
    - Houses Of The Holy
    - Led Zeppelin (Remaster 1994)
    - Led Zeppelin II
    - Led Zeppelin IV
  - Linkin Park
    - Dirt Off Your Shoulder-Lying From You- MTV Ultimate Mash-Ups Presents Collision Course (Parental Advisory)
    - Hybrid Theory (Bonus Tracks)
    - Live In Texas
    - Meteora (Bonus Tracks)
    - Minutes To Midnight (Parental Advisory)
    - Reanimation (Bonus Tracks)
  - Lisa Ekdahl
    - En Samling Sångers
  - Lloyd Cole
    - Antidepressant
    - Bad Vibes
    - Don't Get Weird On Me Baby
    - Love Story
    - Rattlesnakes
  - Lomsk
    - Amerikabrevet
  - Lou Reed
    - New York
    - Songs For Drella
    - The Very Best of Lou Reed

## M

- Madison Avenue
  - The Polyester Embassy
- Madonna
  - American Life (Edited)
  - GHV2
  - Get Together (Maxi-Single)
  - Hung Up (DJ Version)
  - I'm Going To Tell You A Secret (Live)
  - Jump (6-Track Maxi-Single)
  - Like A Virgin
  - Like A Virgin (Remastered-Bonus Tracks)
  - Music [Import Box Set]
  - Ray Of Light
  - Sorry (CD Maxi-Single)
  - The Confessions Tour (Live) (Parental Advisory)
- Madrugada
  - Grit
  - Industrial Silence
  - The Deep End
- Magga Stína
  - syngur Megas
- Mari Boine
  - Gula Gula (1989)
  - Goaskinviellja - Eagle Brother (1993)
  - Leahkastin (Unfolding) (1994)
  - Eallin - Live (1996)
  - Balvvoslatjna (Room of Worship) (1998)
  - Winter In Moscow (2001)
  - Eight Seasons (2001)
- Mariah Carey
  - #1's
- Mark Lanegan
  - Bubbegum
  - Whiskey for the Holy Ghost
- Michael Andrews
  - Donnie Darko (Score)
- Michael Jackson
  - Essential Michael Jackson
  - HIStory- Past, Present And Future, Book I
  - Number Ones
- Mike Scott
  - Bring 'em all in
  - Still Burning
- Miles Davis
  - Porgy and Bess
  - So What
  - Tutu
- Modern Talking
  - The Final Album- The Ultimate Best Of Modern Talking
- Morrissey
  - Bona Drag
- Mr. President
  - Space Gate
- Muse (UK)
  - Black Holes And Revelations



**N**

- N-Trance
  - Happy Hour
- Natacha Atlas
  - Something dangerous
- Neil Young
  - Everybody's Rockin'
  - Freedom
  - Rust Never Sleeps
  - Sleeps With Angels
- Nena
  - 99 Luftballons
- New Order
  - Get Ready
  - International
  - Republic
  - Technique
  - Waiting For The Sirens' Call
- New York Philharmonic
  - Adagio For Strings-Violin Concerto-In Praise Of Shahn (Expanded Edition)
- Nick Cave & The Bad Seeds
  - Abattoir Blues
  - Acoustic Versions of Songs from "Tender Prey"
  - As I Sat Sadly By Her Side (CD Single)
  - B-Sides & Rarities - Volume I
  - B-Sides & Rarities - Volume II
  - B-Sides & Rarities - Volume III
  - Do You Love Me?
  - Henry's Dream
  - Here Comes the Sun
  - Kicking Against the Pricks
  - Let Love In
  - Love Letter
  - Murder Ballads
  - No More Shall We Part
  - Nocturama

- Tender Prey
- The Boatman's Call
- The Good Son
- The Lyre Of Orpheus
- The Mercy Seat
- The Ship Song
- What A Wonderful World
- Where The Wild Roses Grow

- No Doubt
  - No Doubt
  - Return Of Saturn
  - Rock Steady
  - The Beacon Street Collection
  - The singles 1992 - 2003
  - Tragic Kingdom

**O**

- Oscar Peterson
  - Jazz Ballads 8 - Disc 1
  - Jazz Ballads 8 - Disc 2
  - Night Train
- Oslo Kammerkor - Sondre Bratland - Berit Opheim
  - Dâm

**P**

- Patti Smith
  - Trampin'
- Paul Oakenfold
  - Bunkka
- Paul Simon
  - The Rhythm of the Saints
- Paul Van Dyk
  - Global
  - Reflections
- Paul Weller
  - Illumination
  - Stanley Road
  - Wild Wood

- Pepito Ross
  - Vamos A La Playa
- Perssons Pack
  - Diamanter
  - Diamanter (cd 2)
  - Kanoner och små, små saker
  - Kärlek och dynamit
  - Nyårsafton i New York
  - Sekunder i Sverige
  - Svenska hjärtan
  - Äkta Hjärtan
- Pet Shop Boys
  - Actually (Remastered)
  - Popart- The Hits
- Peter Gabriel
  - Hit (CD 1)
  - Hit (CD 2)
  - III
  - Up
- Pink Floyd
  - Dark Side of the Moon
  - The Wall (Disc 1)
  - The Wall (Disc 2)
  - Wish You Were Here
- Pixies
  - Surfer Rosa & Come On Pilgrim

## Q

- Queen
  - A Day At The Races
  - A Night At The Opera
  - Jazz
  - News of the World

## R

- R.E.M.
  - In Time- The Best Of R.E.M. 1988-2003 Rarities And B-Sides
  - What's The Frequency, Kenneth-
- R.E.M.
  - In Time 1988-2003
  - Lifes Rich Pageant
- Raga Rockers
  - Forbudte Følelser
- Rage Against the Machine
  - Rage Against The Machine
- Ramones
  - Anthology (Disc 2)
  - Ramones - Anthology - Hey Ho Let's Go 1
- Rednex
  - Cotton Eye Joe (Sex & Violins)
- Robbie Williams
  - Escapology
- Robert Miles
  - Dreamland
- Roxy Music
  - Flesh + Blood
- Rufus Wainwright
  - Poses
  - Rufus Wainwright
  - Want One
  - Want two
- Rush
  - Rush 1974
  - Caress of Steel 1975
  - Fly By Night 1975
  - All The World's A Stage 1976
  - 2112 1976

- A Farewell To Kings 1977
  - Hemispheres 1978
  - Permanent Waves 1980
  - Exit...Stage Left 1981
  - Moving Pictures 1981
  - Signals 1982
  - Grace Under Pressure 1984
  - Hold Your Fire
  - Retrospective, Vol. 1 (1974-1980)
  - Ry Cooder
    - Paris Texas
- S**
- Scooter
    - 24 Carat Gold
  - Seu Jorge
    - The Life Aquatic Studio Sessions
  - Sex Pistols
    - Never Mind The Bollocks
  - Sharon Jones and the Dap-Kings
    - 100 Days, 100 Nights
  - Slade
    - Nobody's fool
    - Slade in flame
    - Whatever Happened To Slade
  - Smiths
    - Singles
  - Some Like It Hot
    - Some Like It Hot
  - Sophie Ellis Bextor
    - Murder On The Dancefloor
  - Stefan Sundström
    - Sundström spelar Allan
  - Stephen Lynch
    - A Little Bit Special
    - Superhero
    - The Craig Machine
  - System of a Down
    - Toxicity
- T**
- Talking Heads
    - Naked
    - Remain In Light
  - Teenage Fanclub
    - Bandwagonesque
    - Grand Prix
  - Television
    - Adventure
    - Marquee Moon (2003 Remaster)
  - The Aller Vørste
    - Disniland I De Tusen Hjem
    - Materialtrettthet
    - The Aller Vørste
  - The Clash
    - London Calling
  - The Countdown Quartet
    - Hits Of The 80's
  - The Cure
    - Disintegration
    - Galore (The Singles 1987-1997)
    - Kiss Me Kiss Me Kiss Me (Delux Edition - CD1)
    - Staring At The Sea The Singles
    - The Head On The Door
    - The Top
    - Wish
  - The Doors
    - Strange Days
    - Waiting for the Sun
  - The Dukes of Stratospear
    - Chips from the Chocolate Fireball

- The Jam
  - All Mod Cons
  - Setting Sons
  - Sound Affects
  - The Gift
- The Killers
  - Hot Fuss
  - Sam's Town
- The Lime Spiders
  - Nine Miles High
- The Prodigy
  - Music For The Jilted Generation
  - Out Of Space
  - The Fat Of The Land (Parental Advisory)
- The Raconteurs
  - Broken Boy Soldiers
- The Real McCoy
  - Platinum & Gold Collection-The Best Of Real McCoy
- The Rolling Stones
  - A Bigger Bang
  - Bridges To Babylon
  - Emotional Rescue
  - Exile On Main Street
  - Forty Licks (CD ONE)
  - Forty Licks (disc 2)
  - Goats Head Soup
  - It's Only Rock 'N Roll
  - Some Girls
  - Steel Wheels
  - Sticky Fingers
  - Tattoo You
  - The Very Best 1962 - 1975
  - Under Cover
  - Voodoo Lounge
- The Stranglers
  - 10
  - Aural Sculpture
- Dreamtime
- Feline - Extended Edition (incl. 6 Bonus Tracks)
- La Folie (Remastered)
- No More Heroes [Bonus Tracks]
- Rattus Norvegicus
- Sweet Smell Of Success Best Of The Epic Years
- Sweet Smell of Success
- The Best of the Epic Years
- The Raven
- The Streets
  - A Grand Don't Come For Free
  - Original Pirate Material
  - The Hardest Way To Make An Easy Living
- The Style Council
  - Our Favourite Shop
  - The Sound Of
- The The
  - 45 RPM
- The Triffids
  - Born Sandy Devotional
  - Born Sandy Devotional (original)
  - Calenture
  - Calenture (bonus disc)
  - Calenture (original)
  - In The Pines (original)
  - In The Pines [2007 Remastered & Expanded]
  - The Black Swan
- The Waterboys
  - Book Of Lightning
  - Dream Harder
  - Fisherman's Blues
  - Room To Roam
  - This Is The Sea [1 of 2]
  - This Is The Sea: Additional Recordings [2 of 2]
  - Universal Hall
- Tiesto
  - Elements Of Life
  - Parade Of The Athletes

- Tom Robinson Band
  - Power in the Darkness
  - TRB TWO
- Tom Waits
  - Bounced Checks
  - One From The Heart OST (Tom Waits and Crystal Gayle)
  - Swordfishtrombones 1983
  - Rain Dogs 1985
  - Franks Wild Years 1987
  - Bone Machine 1992
  - Night On Earth (Soundtrack) 1992
  - The Black Rider 1993
  - Mule Variations 1999
  - Alice (2002)
  - Blood Money (2002)
  - Big Time 1988 (Part2)
  - Orphans: Bastards (d3)
  - Orphans: Bawlers (d2)
  - Orphans: Brawlers (d1)
  - Tales From The Underground (1994) Volume 1 through 6
  - Tom Waits (Gavin Bryars with Tom Waits) - (1993) Jesus' Blood Ne
  - Tom Waits - 'Alice' (The Original Demos)
  - Tom Waits - (1977) Everytime I Hear T., Live in Hamburg, Germany
  - Tom Waits - (1977) Invitation To The Blues, April 26, Germany
  - Tom Waits - (1979) Cold Beer On A Hot Night
  - Tom Waits - (1979) Cold Beer and Warm Women (live Sydney)
  - Tom Waits - (1979) Fast Women & Slow Horses (live) (128)
  - Tom Waits - (1979) On Broadway
  - Tom Waits - (1982) Shadow of intolerance, Ontario
  - Tom Waits - (1998) Dead Man Walking, March 29, Los Angeles
  - Tom Waits - (1999) Hold On (EP)
  - Tom Waits - (1999) VH1 Storytellers
  - Tom Waits - (2000) May 26, Warsaw, Sala Kongresawa
  - Tom Waits - (2000) May 30, Paris, Le Grans Rex
  - Tom Waits - (2004) Amsterdamdamned, November 21, Amsterdam
  - Tom Waits - Miscellaneous
  - Tom Waits Live
  - Tom Waits Not - Released Studio Records
- Toni Braxton
  - Ultimate Toni Braxton
- Touch & Go
  - I Find You Very Attractive
- Toy Dolls
  - We're Mad (the anthology) - disc 2
  - we're mad! the anthology (Disc.1)
- Travis
  - 12 Memories
  - Good Feeling
  - The Boy With No Name
  - The Invisible Band (Bonus Tracks)
  - The Man Who
- Two to Tango
  - Two To Tango (1)
  - Two To Tango (2)

## U

- U2
  - 1980 Boy
  - 1981 October
  - 1983 Under A Blood Red Sky
  - 1983 War
  - 1984 The Unforgettable Fire
  - 1985 Wide Awake In America
  - 1987 The Joshua Tree
  - 1988 Rattle And Hum
  - 1991 Achtung Baby
  - 1993 Zooropa
  - 1997 Pop
  - 1998 The Best Of 1980-1990 & B-sides Disc 1
  - 1998 The Best Of 1980-1990 & B-sides Disc 2
  - 2000 All That You Can't Leave Behind
  - 2002 The Best And The B-Sides Of 1990-2000 CD1
  - 2002 The Best And The B-Sides Of 1990-2000 CD2
  - Achtung Baby
  - How To Dismantle An Atomic Bomb
  - The Best Of 1980 - 1990
  - The Best Of 1990-2000 (A Sides US Version)
  - Who's Gonna Ride Your Wild Horses
- U96
  - Heaven
  - Replugged
- Ultravox
  - 77 Ha-ha-ha
  - 77 Ultravox
- Perssons Pack
  - Nyårsafton i New York

## V

- VNV Nation
  - Empires
  - Genesis.2
  - Matter + Form
- Vamp
  - 13 humler
  - En annen sol
  - Flua på veggen
  - Godmorgen, søster
  - Horisonter
  - Månemannen
  - Vamp i full symfoni med kringkastingsorkesteret
  - siste stikk
- Various
  - Buffy The Vampire Slayer - The Album
  - Donnie Darko (Original Soundtrack)
  - I'm Your Fan (Tribute to Leonard Cohen)
  - McMusic 21
  - Suuret suomalaiset tangot : Satumaa
  - Taube 1
  - Taube 2
  - Until the End of the World
  - Viva Cuba!
  - Done Again (In The Style Of The Beatles)- The Beatles, Vol.8
- Vladimir Vissotsky
  - 1
  - 2
- Vladimir Vissotsky & Marina Vlady
  - Vlady Vissotsky
- Vømmøl Spellmannslag
  - Vømlingen
  - Vømmølmusikken

**W**

- William Orbit
  - Pieces In A Modern Style
- Wolfgang Amadeus Mozart
  - 4 Hornkonzerte - Concertos for Horn and Orchestra
- World mix
  - Deep Forest

**X**

- XTC
  - Apple Venus Volume 1
  - English Settlement [Remastered 2001]
  - Homegrown
  - Homespun The Apple Venus Volume One Home Demos
  - Mummer (remastered)
  - Nonsuch
  - Oranges & Lemons
  - Skylarking
  - The Big Express (remastered)
  - The Compact XTC The Singles 1978-1985
  - Wasp Star [Apple Venus Volume 2]