

Statistical Characterisation of MP3 Encoders for Steganalysis

Rainer Böhme
Technische Universität Dresden
Institute for System Architecture
01062 Dresden, Germany
rainer.boehme@inf.tu-dresden.de

Andreas Westfeld
Technische Universität Dresden
Institute for System Architecture
01062 Dresden, Germany
westfeld@inf.tu-dresden.de

ABSTRACT

This paper outlines a strategy to discriminate different ISO/MPEG 1 Audio Layer-3 (MP3) encoding programs by statistical particularities of the compressed audio streams. We use Bayesian logic to deduce the most probable encoder on the basis of a feature vector that can be extracted from arbitrary MP3 files. All appropriate features used for the classification are discussed and example results for sets of test data from 20 different codecs are given. Possible applications include advances in information hiding, increases in the reliability of steganographic attacks, and inferences about the origin of MP3 files for forensic purpose. We demonstrate that a pre-classification of MP3 encoders reduces the false alarm rate for a steganographic detection method. Implications for the generalisability of the proposed scheme to other file formats are addressed.

Categories and Subject Descriptors

D.2.11 [Software Architectures]: Information Hiding

General Terms

Security

Keywords

Steganalysis, MP3 Encoder Classification, Digital Forensics

1. INTRODUCTION

The invention of the ISO/MPEG 1 Audio Layer-3 (MP3) audio compression algorithm [5, 12] is probably one of the most remarkable and far-reaching developments in the area of digital media processing. The MP3 format enables compression rates of about 1/10 of the size of uncompressed digital audio while degrading the audible quality only marginally. Together with the moderate complexity of the compression

algorithm—software implementations of MP3 coders/decoders (codecs) with acceptable performance even on low budget home computers soon became available—the format simplified the interchange of music and resulted in worldwide popularity for its users and sleepless nights for the music industry. The popularity of the format fostered demand for encoding tools and opened a market for a variety of programs for different needs. Today we count hundreds of MP3 encoder front-ends based on several dozens of encoding engines ranging from proof of concept hacks to targeted products either tuned for high speed, or optimised to costly and flexible tools for professional studio requirements.

Given these facts, the MP3 format became an interesting carrier for steganographically hidden data. Steganography, which is somewhat related to cryptography, aims to conceal the very existence of a confidential message by hiding it imperceptibly within other, less suspicious data [19]. MP3 is a promising carrier format for steganography in three ways. At first, the popularity of the format is an advantage, because exchanging common and widely used types of data is less conspicuous to an observer. For example, sharing an MP3 file over the Internet is a completely common task and doing so is a plausible form of communication. Second, MP3 files are typically between 2 and 4 megabytes (MB) in size and thus are larger than other common formats (e.g., text documents or photographs as e-mail attachments). All forms of information hiding suffer from a small proportion of payload compared to the total amount of information, necessary to cover the message. So, larger file sizes simplify the handling of medium-sized payloads (e.g., a text message or a photograph). The inconveniences that come with splitting up messages over different carriers can be almost avoided for MP3 files. Third, the nature of the lossy MP3 compression itself makes it attractive for steganographic use. The information loss that is a concomitant of the encoding process creates a certain amount of unpredictability that can be exploited to carry hidden information securely.

Compared to the suitability of MP3 files for steganography, the amount of known steganographic tools for this format is still quite limited. **MP3Stego** [20] is based on the **8hz-mp3** encoder [1] and hides message bits in the parity of block length. Although this procedure is limited to a very low capacity, it is (under certain conditions, see below) detectable [23]. The attack is based on the analysis of statistical properties, i.e., the variance of block lengths in the MP3 stream. **Stego-Lame** [22] pursues another approach and embeds into uncompressed *Pulse Code Modulation* (PCM) au-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

MM&Sec'04, September 20-21, 2004, Magdeburg, Germany.
Copyright 2004 ACM 1-58113-854-7/04/0009 ...\$5.00.

dio data. The amount of information is so small and the embedding procedure so carefully selected, that a subsequent lossy MP3 compression does not erase the hidden information. This tool is still in an experimental stage. An appropriate attack is delivered in the same bundle. In addition to these publicly known stego-tools we expect some more being used in the wild. Although the complexity of MP3 compression exceeds those of typical steganographic tools (e.g., LSB image embedding), the availability of commented source codes for MP3 encoders facilitates the composition of derivatives with steganographic extensions. Hence, advances in the detection of steganographic data in MP3 files are relevant.

The experience with the existing attack against **MP3Stego** shows that the detection procedure can distinguish MP3 files with and without steganographic content quite reliably if they are encoded with either **MP3Stego** or its underlying encoding engine [1]. However, files from other encoders tend to have similar statistical properties as steganograms from **MP3Stego** and thus are identified as false positives. Hence, the reliability of the detection algorithm heavily depends on the prior knowledge about the encoder of a particular file. While this situation might be sufficient for an academic attack or proof of concept, it is definitely not optimal for real world applications. In the fieldwork, we usually cannot expect any prior knowledge about the source of an arbitrary MP3 file. We therefore present a procedure to determine the encoder of MP3 files on the basis of statistical features that are typical for a certain implementation of the MP3 format specification. The insertion of a preclassification of MP3 encoders allows a steganalyst to run the appropriate detection algorithm for the determined encoder and thus dramatically decrease the amount of false positives. Thus it is believed that statistical classification of MP3 encoders can increase the reliability of detection procedures.

The rest of this paper is organised as follows. In the next section we briefly review the relevant particularities of the MP3 format that are analysed for the extraction of statistical features. The features themselves are explained in Section 3. Experimental results that back the performance of the proposed scheme are presented in Section 4, before we discuss further applications and possible generalisations to other file formats in Section 5.

2. ANALYSIS OF MP3 SPECIFICATION

The purpose of this section is not to repeat the architecture and specification of MP3 compression [2, 12], but to give a brief overview of those principles that are relevant as features for our proposed statistical classification. Hence, we focus on the latitudes that are left in the ISO specification, which leave space for different implementations. It is the vaguely defined particularities that finally lead to different output streams for the same input data.

2.1 Principles of MP3 Compression

The developers of MP3 audio compression included several techniques to maximise the relationship between perceived audio quality and storage volume. In contrast to previous schemes, they designed a two-track approach. On the *first track*, the audio information is split up into 32 equally spaced frequency sub-bands. These components are separately mapped into the time domain with a *Modified Discrete Cosine Transformation* (MDCT). The following quan-

tisation step reduces the precision of the MDCT coefficients. As a last step, a lossless entropy encoding of the quantised coefficients leads to the compact representation of MP3 audio data. The *second track* is very important for the performance of MP3 encoding, because it is used as a control track. Also starting from the PCM input data, a 1024-point Fourier transformation is used to fit the local frequency spectrum as input to a psycho-acoustic model. This model emulates the particularities of human auditory perception and derives appropriate masking functions for the input signal. The model controls the choice of block types and quantisation factors in the first track. Hence, this two-track approach adaptively finds an optimal trade-off between data reduction and audible degradation for a given input signal.

Regarding the underlying data format, an MP3 stream consists of a series of *frames*. Synchronisation tags separate frames from other information sharing the same transmission or storage stream (e.g., video frames). For a given bit rate, all MP3 frames have a fixed compressed size and represent a fixed amount of 1152 PCM samples. Usually, an MP3 frame contains 32 bits of header information, an optional 16 bit *Cyclic Redundancy Check* (CRC) checksum, and two *granules* of compressed audio data. Each granule can be subdivided into one (mono) or two (stereo) *blocks*. Since the actual block size depends on the amount of information that is required to describe the input signal, it may vary between frames. To match the floating block sizes with the fixed frame sizes without wasting bandwidth, the MP3 standard introduces a so-called *reservoir* mechanism. Frames that do not use their full capacity are filled up (partly) with block data of subsequent frames. This method assures that local highly dynamic sections in the input stream can be stored with over-average precision, while less demanding sections allocate under-average space. However, the extent of reservoir usage is limited in order to decrease the interdependencies between more distant frames and to facilitate resynchronisation in the middle of a stream.

2.2 Level of Analysis and Related Work

In order to perform a statistical characterisation of MP3 encoders we have to find differences in the encoding process. These differences may have multiple causes. At the first glance, all loosely defined parameters in the specification are subject to different interpretations. However, the standard precisely describes a large set of critical parameters including the exact coefficients for the filter bank and threshold values for the psycho-acoustic model. Nevertheless, some implementations seem to vary or fine tune these parameters. In addition, performance evaluations may have led to sloppy implementations of the standard, such as shortcuts in the inner quantisation loop or the choice of non-optimal Huffman tables. Also, a number of parameter defaults for meta information are up to the implementor (e.g., the *Serial Copy Management System* (SCMS) flags, also known as *protection bit* [9]). All these variations together cause particular features in the output stream that are indications of a specific encoder and therefore are subject to a detailed analysis.

To structure the occurrences of implementation specific particularities in the MP3 encoding process, we will subdivide the process into three layers as shown in Table 1. The *transformation layer* includes all “passive” operations that directly affect the audio data, namely the filter bank,

Table 1: Structure of MP3 encoding process

Functionality	Points for analysis
<i>Transformation Layer</i>	
- Filter bank	- Frequency range
- MDCT transform	- Filter noise
- FFT transform	- Audible artefacts
<i>Modelling Layer</i>	
- Quantisation loop	- Size control
- Model computation	- Model decisions
- Table selection	- Capability usage
<i>Bitstream Layer</i>	
- Auxiliary data	- Surface information
- Frame header bits	- SCMS protection bit
- Checksums	- SCMS original bit
- Stream formatting	

and the MDCT and *Fast Fourier* (FFT) time to frequency transformations, respectively. In this layer, variations in the filter coefficients or in the precision of the floating point operations may cause measurable features such as typical frequency ranges or additional noise components.

We define all “active” components of the compression algorithm as part of the *modelling layer*. These sub-processes are less close to the underlying audio data and mainly perform the trade-off between size and quality of the compressed data. In this layer, encoder differences basically occur in three ways:

1. Calculation of size control quantities, e. g., whether net or gross file sizes are used as reference for the bit rate control.
2. Model decisions: Different threshold values lead to different marginal distributions of control parameters over the data stream.
3. Capability usage: Some encoders do not support all compression modes specified in the MP3 standard.

The uppermost layer, which we call *bit stream layer*, handles the formatting of already compressed MP3 frames into a valid bit stream. These operations include the composition of frame headers, the optional calculation of CRC checksums for error detection, and the insertion of meta data. For instance, quasi-standardised ID3 tags [13] contain information about the names of artists, interprets, and publishers of audio files. Optional VBR (*variable bit rate*) headers store additional data evaluated by some MP3 players to display valid progress bars and enable efficient skipping within MP3 files with variable bit rate.¹ The existence of a certain kind

¹As MP3 has been specified for *constant bit rates* (CBR) the majority of MP3 files are encoded as CBR with one of the predefined rates. However, some encoding programs optionally encode each frame with a different bit rate (out of the predefined rates), thus enabling *variable bit rate* (VBR) streams with MP3.

of meta information and its default values may be used as indicator for the encoding program.

EncSpot [4], the only tool for MP3 encoder detection we know, relies on the deterministic surface parameters of the bit stream layer. As these parameters are easily accessible, it is also simple to erase or change their values and therefore trick this kind of encoder detection. Therefore we decided to use statistical features related with deeper structures of the encoder and thus are more difficult to manipulate. Our initial experiments with parameters of the transformation layer showed that those tend to be dependent on the type of audio data actually encoded. For example, it is impossible to measure encoder characteristics, such as the upper frequency bound, if the encoded audio material does not use the full range. Also, artefacts occur at typical envelopes or frequency changes that do not appear similarly in all kinds of music. Hence, we decided to focus our level of analysis to the modelling layer, which promises to deliver the most robust features in terms of source data independency and difficulty of manipulation.

2.3 Terminology and Procedure

To precisely describe the nature of the features we introduce some formal notations. We denote a medium m as m_0 for the source (i. e., uncompressed) representation and as $m_i = e_i(m_0)$ if it is encoded with encoding program e_i . e_i is element of the set of n encoders $E = \{e_1, e_2, \dots, e_n\}$. We write the set of all files encoded using e_i as $M_i = e_i(M_0)$, where M_0 is the set of all uncompressed source media.

The function $f(m)$ extracts a symbolic feature x from m . The vector of k different features

$$\mathbf{x} = \mathbf{f}(m) = (f_1(m), f_2(m), \dots, f_k(m))$$

is called feature vector. The components of the feature vector \mathbf{x} are selected to be as similar as possible for different media $m \in M_i$ encoded with the same encoder e_i , and also as dissimilar as possible for all encoded media $\bar{m} \in \{e_j(m_0) | j \neq i\}$ that are derived from m_0 by encoding it with other encoders. Therefore the information about the characteristics of the encoding program is consolidated in the value of \mathbf{x} .

Classifiers are algorithms which automatically classify an object, i. e., assign it according to its features to one of several predefined classes. As the literature contains multiple options, the choice of a specific algorithm for our purpose was determined by the conditions given in our application. *Fisher Linear Discriminant* (FLD) methods and *Support Vector Machines* (SVM) have already been successfully applied for steganalysis [16, 6]. These methods perform well for numeric (i. e., continuous) features, but are less suitable for symbolic features. Hence, we chose to apply a classifier which is based on Bayesian logic [15]. As we will show in Section 4, we get notable results with the simple *Naïve Bayes Classifier* (NBC) [3].²

We use a classifier c to establish the relation between a specific instantiation of $\mathbf{x} = \mathbf{f}(m_i)$ and the encoding program e_i that was used to create m_i . If we do not have any knowledge about the encoder, we can only derive probabilistic evidence about this assignment. For a given medium m

²These results are coherent with the findings from a comprehensive evaluation of different classifiers: Compared to a set of complex classification models, the simple NBC performed equal or superior for many realistic decision problems [14].

a classifier tries to compute the conditional probabilities

$$P(e_i|\mathbf{f}(m)) = P(e_i|x_1 = f_1(m), x_2 = f_2(m), \dots, x_k = f_k(m)),$$

with $1 \leq i \leq n$, and then selects the most probable encoder e_i , so that

$$P(e_i|\mathbf{f}(m)) > P(e_j|\mathbf{f}(m)), \forall e_j \in E \setminus \{e_i\}, i = c(\mathbf{f}(m)).$$

The classifier's performance depends on its parameterisation, which can be induced from data. Therefore we assemble a training set

$$T = \{(i, e_i(m)) | 1 \leq i \leq n \wedge m \in M_0\}.$$

Each element of T contains a compressed representation of medium m and a reference to the known encoding program. We note a classifier trained with the training set T as c_T . The encoder prediction of a specific instantiation of \mathbf{x} , and of an underlying medium m will be denoted as $c_T(\mathbf{x})$ and $c_T(\mathbf{f}(m))$, respectively. To evaluate the quality of the classification, we regard the proportion p of correctly classified cases when the classifier is run on elements of a test set S , which is composed similarly to the training set T :

$$p(c, S) = \frac{|\{(i, m_i) \in S | i = c(\mathbf{f}(m_i))\}|}{|S|}$$

As a weak form of reliability evaluation, the same training set T can be reclassified, thus $c_T(\mathbf{f}(m_i))$ with $(i, m_i) \in T$. A somewhat stronger measure can be achieved for disjoint test and training sets, so that $S \cap T = \emptyset$.

3. DESCRIPTION OF FEATURES

As a result of iterative comparisons and analyses of MP3 encoder differences, we discovered a set of 10 features in the *modelling layer*. For a structured presentation, the features are assigned to categories, which will be discussed separately in the following subsections.

3.1 Calculation of Size Control Quantities

Distinct encoders seem to differ in the way the target bit rate is calculated, as we discovered measurable differences in the effective bit rate. According to the MP3 standard, each block can be encoded with one of 14 predefined bit rates.³ However, because of the difficulty to reach an exact compressed size, these act just as guiding numbers. Some encoders treat these rates as an upper limit, others as an average. Also, the encoders differ in the scope of frames that are evaluated as control parameters for the compression loop. If broader scopes are considered, or fixed headers at the beginning of MP3 files are also reflected in the quantisation loop, then the effective bit rate varies with the file length and converges to a target value with an increasing number of frames.

These phenomena are depicted in Figure 1 for four selected encoders on the basis of files with a nominal bit rate of 128 kbps. The curves are drawn according to a least square estimate with a linear and a hyperbolic term over measured data points.⁴ The effective bit rates β_{eff} of **8hz-mp3** and **mp3comp** depend on the number of frames p , while there is

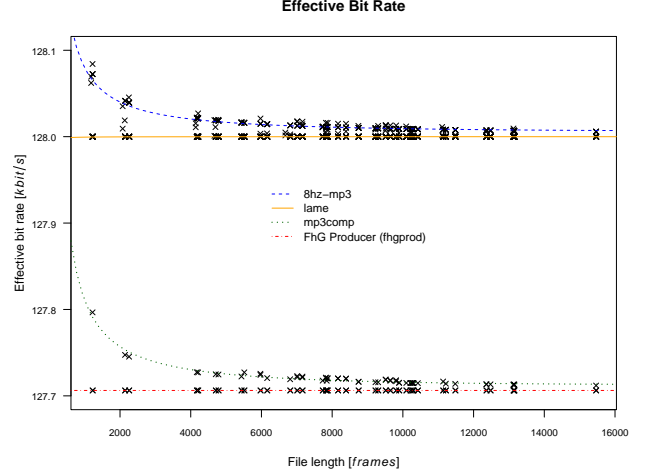


Figure 1: Relation between effective bit rate and file length for selected encoders

no influence for files encoded with **lame** or **fhgprod**. We calculate the effective bit rate as

$$\beta_{\text{eff}} = \frac{[(\text{filesize}) - (\text{junkbytes}) - (\text{meta information})] \cdot 8 \cdot \varphi}{1152 \cdot p},$$

with $\varphi = 44.1$ kHz as sampling frequency. Even for large files we observe a measurable difference in the marginal β_{eff} between all four encoders. To derive a bit rate independent feature from this observation, we calculate a criteria ϱ_1 as ratio between the effective bit rate β_{eff} and the nominal bit rate β_{nom} :

$$\varrho_1 = \frac{\beta_{\text{eff}}}{\beta_{\text{nom}}}, \quad \text{with} \quad \beta_{\text{nom}} = \frac{1}{p} \sum_{i=1}^p \beta_{\text{nom}}^{(i)},$$

where $\beta_{\text{nom}}^{(i)}$ is the nominal bit rate given in the header of the i -th frame. To map this ratio to a symbolic feature x_1 , we define the extraction function f_1 as follows:

$$f_1(m) = \begin{cases} 0 & \text{for } \varrho_1 < 1 - 1 \cdot 10^{-4} \\ 1 & \text{for } 1 - 1 \cdot 10^{-4} \leq \varrho_1 \leq 1 \\ 2 & \text{for } 1 < \varrho_1 \leq 1 + 5 \cdot 10^{-6} \\ 3 & \text{else.} \end{cases}$$

The number of levels and the exact boundaries for this feature, as well as for the following ones, are determined by an iterative process of comparing a set of test audio files. We report the functions which lead to the best experimental results, even though further optimisation is still possible.

In Section 2.1, we mentioned that an MP3 stream consists of a sequence of *frames*. Again, two *granules* share a frame of fixed size. The quantisation loop adjusts the size of the granules separately according to two criteria:

1. Size: The granule must fit into the available space.
2. Quality: Signal noise shall remain imperceptible.

For some encoders, e.g., **shine**, we observed a slight bias for quality over size. As the ‘hard’ space limit counts on

³Bit rates for Layer-3 in kbps: 32, 40, 48, 56, 64, 80, 96, 112, 128, 160, 192, 224, 256, 320

⁴ R^2 values range between 0.83 and 0.97.

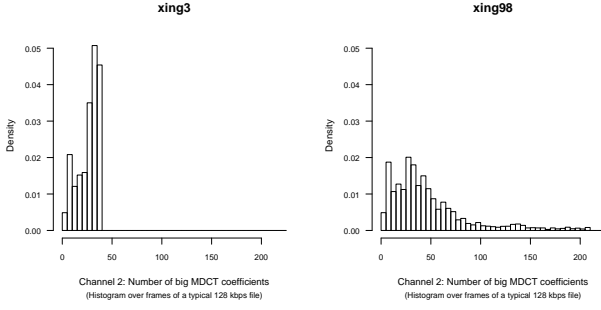


Figure 2: Comparison of size control in stereo files encoded with xing3 and xing98

both granules together, the first granules $g_1^{(i)}$ of all frames ($1 \leq i \leq p$) tend to get bigger than the second ones $g_2^{(i)}$. Hence, we measure the proportion of frames in the file where the length of the first granule $\text{len}(g_1)$ dominates the second one $\text{len}(g_2)$:

$$\varrho_2 = \frac{1}{p} \sum_{i=1}^p G(i), \quad \text{with}$$

$$G(x) = \begin{cases} 1 & \text{for } \text{len}(g_1^{(x)}) > \text{len}(g_2^{(x)}) \\ 0 & \text{else.} \end{cases}$$

Again, we define a mapping function, now for feature x_2 :

$$f_2(m) = \begin{cases} 0 & \text{for } \varrho_2 < 0.50 \\ 1 & \text{for } 0.50 \leq \varrho_2 < 0.55 \\ 2 & \text{for } 0.55 \leq \varrho_2 < 0.70 \\ 3 & \text{else.} \end{cases}$$

The next feature makes use of characteristics of the reservoir mechanism. We found that the abruptness of the rise in reservoir usage between silent and dynamic parts in the audio stream differs between some encoders. Other encoders even do not use the reservoir at all. As the vast majority of audio files start with a tiny silence, we derive the feature x_3 from the amount of bytes shared between the first and the second frame $v_{(1,2)}$:

$$f_3(m) = \begin{cases} 0 & \text{for } v_{(i,i+1)} = 0 \quad \forall i: 1 \leq i < p \\ 1 & \text{for } v_{(1,2)} > 300 \\ 2 & \text{else.} \end{cases}$$

The function $f_3(m)$ is zero if the reservoir is not used in the whole file. The values 1 and 2 identify hard and soft reservoir usage, respectively.

The last feature in this category is less theoretically based and our evaluations show that it has little impact on the classification result, except for a better separation between two versions of the Xing encoder, namely **xing98** and **xing3**. However, we report it for the sake of completeness. We observed that **xing3** uses a different size control mechanism for the second block of every granule of stereo files. The differences are clearly visible in the histogram of lengths of *big value* MDCT coefficients (see Figure 2). Following the ISO/ MPEG 1 Audio Layer-3 terminology [12], *big values* are the partition of spectral coefficients with absolute values

greater than 1. This partition holds the most energy of the transformed audio signal and thus the average number of big values is a valid indicator for the extent of size reduction in the quantisation loop. To derive a continuous feature from the different spread of histogram values in the stereo channel, we measure the entropy from the histogram with the approximation given in [17]:

$$H \approx - \sum_{j=1}^{d_{\max}} d_j \log d_j + \log \Delta x,$$

with d_j denoting the density of occurrences in the j -th bin and Δx as bin size. Since Δx is constant for all encoders, we use a simplified function to calculate feature x_4 :

$$f_4(m) = - \sum_{j=1}^{60} d_j \log d_j$$

Note that in contrast to previous features, $f_4(m)$ is a continuous feature that is modelled by the classifier as a normal distributed random variable with mean $\mu_i^{(4)}$ and standard deviation $\sigma_i^{(4)}$ for the i -th encoder e_i . However, as this feature evaluates the characteristics of the second channel in stereo data, it is not applicable to mono files; hence, we cannot discriminate between **xing3** and **xing98** for mono files.

3.2 Model Decision

The psycho-acoustic model is a second source for distinguishing features. Differences in the computation of control parameters or modifications in the choice of threshold values lead to typical marginal distributions of measurable parameters.

The binary value *preflag* causes an additional amplification of high frequencies and is individually set for each compressed block b_i ($1 \leq i \leq q$, with q as number of blocks in a file). Concerning the treatment of this parameter, the ISO/ MPEG 1 Audio Layer-3 standard explicitly leaves latitude:

“The condition to switch on the preemphasis is up to the implementation.” [12, p. 110]

To derive an operable feature we calculate the proportion of blocks with preflag set

$$\varrho_5 = \frac{1}{q} \sum_{i=1}^q \text{preflag}(b_i)$$

and map it into disjoint regions for the symbolic feature x_5 :⁵

$$f_5(m) = \begin{cases} 0 & \text{for } \varrho_5 = 0.00 \\ 1 & \text{for } 0.00 < \varrho_5 \leq 0.01 \\ 2 & \text{for } 0.01 < \varrho_5 \leq 0.05 \\ 3 & \text{for } 0.05 < \varrho_5 \leq 0.10 \\ 4 & \text{for } 0.10 < \varrho_5 \leq 0.21 \\ 5 & \text{for } 0.21 < \varrho_5 \leq 0.35 \\ 6 & \text{for } 0.35 < \varrho_5 \leq 0.62 \\ 7 & \text{for } 0.62 < \varrho_5 \leq 0.77 \\ 8 & \text{else.} \end{cases}$$

⁵Our experiments show that the symbolic interpretation of x_5 leads to better classification results than a treatment as continuous feature with assumed normal distribution.

The MP3 audio format offers different block types, which allow an optimal trade-off for sections requiring higher time resolution at the cost of frequency resolution and vice versa. The majority of blocks are encoded with block type 0, the *long block* with lower time and higher frequency resolution. Block type 2 defines a *short block*, which offers less coefficients to be stored for three different points in time. Two more block types are specified to perform smooth shifts between the above mentioned types. Hence, the standard defines a graph of valid block transitions between two adjacent blocks b_i and b_{i+1} , as shown in Figure 3.

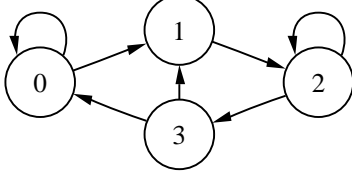


Figure 3: Valid MP3 block type transitions

An evaluation of block type transitions of MP3 files from different encoders uncovers two interesting details: First, some encoders (**shine**, all **xing***) do not use short blocks at all and thus always encode with block type 0. Second, other encoders (**lame**, **gogo**, and **plugger**) include specific “illegal” transitions, mainly at the beginning of a file. As these transitions are rarely observable from other encoders, they identify the encoder reliably. Hence, we construct the extraction function for feature x_6 as follows:⁶

$$f_6(m) = \begin{cases} 0 & \text{for } \text{type}(b_i) = 0 \quad \forall i: 1 \leq i \leq q \\ 1 & \text{for } \text{type}(b_1) = 0 \quad \wedge \quad \text{type}(b_2) = 2 \\ 2 & \text{for } \text{type}(b_1) = 2 \quad \wedge \quad \text{type}(b_2) = 3 \\ 3 & \text{for } |\{b_i | \text{type}(b_i) = 2\}| = \\ & |\{b_i | \text{type}(b_i) = 3\}| = 1 \\ 4 & \text{else.} \end{cases}$$

We have no other explication for these strange transitions than assuming that they are intended to leave a kind of encoder fingerprint in the output data. It is up to a deeper analysis of these particularities in the source code to reveal further evidence.

3.3 Capability Usage

The third category of features exploits the fact that some encoders do not implement all functions specified in the MP3 standard. We call this category *capability usage* and clearly separate these capabilities from surface parameters, such as header flags, because the latter can easily be changed without touching the compressed data.

The *Scale Factor Selection Information* (SCFSI) is a parameter that allows an encoder to reuse scale factors for subsequent parts of the stream if they do not change over time. However, only few encoders use this compression method, namely **lame**, **gogo**, and **xingac21** (“AudioCatalyst”). We

define a feature x_7 reflecting the use of SCFSI:

$$f_7(m) = \begin{cases} 0 & \text{for } \text{scfsi}(b_i) = 0 \quad \forall i: 1 \leq i \leq q \\ 1 & \text{else.} \end{cases}$$

Although MP3 frames have a fixed length, the amount of information used to describe the respective audio signal may vary. We refer to this quantity as *effective frame length* $\text{len}_{\text{eff}}(i)$. According to the MP3 standard, the effective frame length has no constraints to match a multiple of bytes, words or quad-words. However, we observed that some encoders (**8hz-mp3**, **bladeenc**, **m3ec**, **plugger**, **shine**, **soloh**) adjust all effective frame lengths to byte boundaries, while others do not. We use this characteristics as feature x_8 :

$$f_8(m) = \begin{cases} 0 & \text{for } \text{len}_{\text{eff}}(i) = 0 \bmod 8 \quad \forall i: 1 \leq i \leq p \\ 1 & \text{else.} \end{cases}$$

After the quantisation, the MDCT coefficients are further compressed by a Huffman style entropy coder. In contrast to the method proposed by Huffman [8], the tables are not computed from the marginal symbol distribution. In order to avoid the transmission of marginal distributions or table data, the developers of MP3 standardised a set of 28 predefined Huffman tables that were empirically optimised for the most probable cases in audio compression. In the very rare case of longer code words an escape mechanism allows storage of uncompressed values. An MP3 encoder chooses the most suitable table separately and independently for each of the three *regions* of the *big value* MDCT coefficients. As there is no efficient method to perform an optimal table selection, some encoders increase performance by using heuristics to quickly select a suitable table, rather than the optimal one. From a comparison of table usage frequencies, we found two noteworthy characteristics: First, all Xing encoders seem to avoid strictly using table number 0 for region 2.⁷ Second, only a few encoders (**m3ec**, **mp3enc31**, **uzura**) use table 23 for the regions 1 and 2. We exploit these observations as additional information for our classification:

$$f_9(m) = \begin{cases} 0 & \text{for } \text{table}_2(b_i) \neq 0 \quad \forall i: 1 \leq i \leq q \\ 1 & \text{for } \exists (b_i, j): \text{table}_j(b_i) = 23, \\ & 1 \leq i \leq q, j = 1, 2 \\ 2 & \text{else.} \end{cases}$$

Also, **shine** uses only a subset of the defined tables. However, as we can already identify this rarely used encoder with several other features, we refrain from adjusting this feature for the detection of **shine**.

3.4 Miscellaneous

Since our last feature does not fit in any of the above categories, we decided to explain it separately. Independent from whether the reservoir mechanism is used or not, there may be a couple of bytes unused and filled up to meet the fixed frame length. These so-called *stuffing bits* can be set to any arbitrary values. For a closer examination of these values, we composed histograms of the byte values in the stuffing areas. While most encoders set all stuffing bits to zero, we still found some exceptions and mapped them into a symbolic feature x_{10} :

⁶For simplicity we give the relations for mono files. Stereo files work similar if blocks are evaluated in pairs. The given definition is not disjoint, hence the values are assigned by the first matching condition.

⁷According to the standard, we count the regions from zero.

$$f_{10}(m) = \begin{cases} 0 & \text{for stuffing with zeros} \\ 1 & \text{for no stuffing at all} \\ 2 & \text{for stuffing with 0x55 or 0xaa} \\ 3 & \text{for stuffing with "GOGO" (0x47 and 0x4f)} \\ 4 & \text{else.} \end{cases}$$

The enumeration of features in this section is a subset of particularities we took into account and from which we selected the most promising ones. The selection is far comprehensive, so it is still feasible to find further differentiating features. Such features may be necessary to reliably separate new encoders, or encoders that were not included in our initial analysis.

4. EXPERIMENTAL RESULTS

For our experimental work, we used the R Statistical Framework [10, 21] and implemented an extension for statistical analyses of MP3 files on the basis of the open source MP3 player `mpg123` [7]. All results are based on an MP3 database of about 2,400 files encoded with 20 different encoders (cf. Table 4 in the appendix). The audio data was selected from different sources to make the measurements independent from specific types of music or speech. We included tracks from a re-mastered CD of Grammy Nominees, from a compilation of Blues Brothers (including some live recordings), further piano music by Chopin, as well as *Sound Quality Assessment Material* (SQAM) files with speech and instrumental sounds. All source files were read from CD recordings and stored as PCM wave files with 44.1 kHz, 16 bit, stereo.

If provided, we encoded every source audio with three constant bit rates that we believe are the most widely used rates for MP3 files, namely 112 kbps, 128 kbps, and 192 kbps. Additional MP3 files with variable bit rates, with two quality settings each, were generated by the encoders `iTunes`, `lame`, and `xingac21`.

4.1 Validity for Known Data

To measure the performance of our proposed method, we implemented a Naïve Bayes Classifier (NBC) [3] for fixed feature vectors of both symbolic and continuous features.

In the first experiment, we trained the classifier c_{T_1} with a training set T_1 of about 2,400 cases. For each case, we extract a feature vector $\mathbf{f}(m_i)$ from a file encoded with a defined encoder e_i and use these tuples to induce classification parameters for c_{T_1} . To evaluate the performance of c_{T_1} we use the same feature vectors, because $S_1 = T_1$, as input to the classifier and compare the predicted encoders to the known true values. In this experiment we reach a hit rate of $p(c_{T_1}, T_1) = 96.2\%$. As a measure of confidence, we calculate the average a-posteriori probability over the predicted encoders $\bar{P}_{\max} = 96.1\%$. The classifier calculates the a-posteriori probability $\max_i P(e_i | \mathbf{f}(m))$ for each file on the basis of the feature vector.

Table 3 summarises the features proposed in Section 3. We use a jack-knife method to empirically evaluate the importance of each feature for the classification result. Therefore the training and classification procedure is repeated several times, while excluding individual features one by one. The resulting increase of misses in the classification table is a measure for the importance of a feature. According to these values, the effective bit rate seems to be the most important feature, followed by the method of reservoir usage.

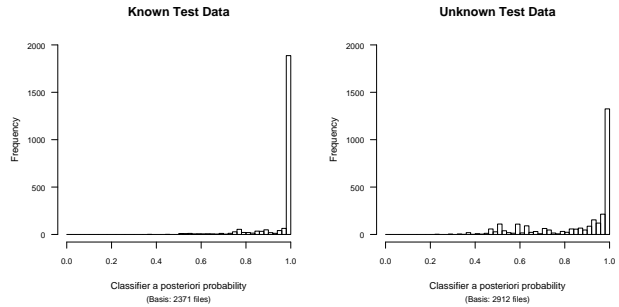


Figure 4: Comparison of classification confidence between self generated test data (left) and data collected in the wild (right).

A closer look at the results shows that the main sources for classification errors occur between tightly related encoding engines, such as the DOS and UNIX versions of Fraunhofer’s `l3enc`, and between two subsequent versions of Xing encoders (`xing3` and `xing99`). Also, `soloh` produces false classifications as `8hz-mp3`, especially for source files from the Blues Brothers CD. To explore these misclassifications, we debugged the `soloh` binary and found references to an early version of the `8hz-mp3` encoder. Hence, the similarity in statistical features may reveal insights about the “intellectual origin” of certain encoders.

To support our results and reduce the risk of tautological finding, we repeat the experiment with a split-half method. We trained the classifier c_{T_2} with a sub sample T_2 of the first training set T_1 . All other elements from T_1 are used in the test set S_2 , so that $T_2 \cap S_2 = \emptyset$. The results of this second experiment are shown in Table 5 (in the appendix). We found an overall hit rate of $p(c_{T_2}, S_2) = 94.9\%$ and an average a-posteriori probability of 95.9%. As both quality measures differ only marginally from the first experiment (−1.3 and −0.2 percentage points, respectively), we conclude that the proposed method can also reliably identify the encoders of unknown MP3 files.

4.2 Reliability for Unknown Data

We used classifier c_{T_1} to determine the encoders of a random sample of 3,000 MP3 files drawn from different sources, with a total amount of more than 19,000 MP3 files. The overall average a-posteriori probability is 86.9%. This is about 10 percentage points below the values for known data. We still consider this a good value because we are aware that our training set certainly does not include all available encoders. In addition, some well separable encoders in our assembled test database, such as `uzura` and `shine`, have not been identified in the mass of unknown files.

Figure 4 shows the distribution of a-posteriori probabilities for known and unknown test data.⁸ The average a-posteriori probability for the most frequent encoders is shown separately in Table 2.

Encoders e_* not included in the training procedure may lead to misclassifications in either of two ways: If the feature vector $\mathbf{f}(m_*)$ is similar to one of the trained encoders then

⁸The latter has been reduced to bit rates between 112 kbps and 192 kbps, keeping 2912 files.

Table 2: A-posteriori probabilities for frequently identified encoders

Rank	Encoder	Share	Confidence $P_{\max}(e \mathbf{x})$		
			μ	σ	n
1.	fhgprod	17.3 %	0.87	0.14	381
2.	xingac21	11.6 %	0.99	0.04	256
3.	l3enc272	10.7 %	0.95	0.08	235
4.	soundjam	10.0 %	0.97	0.11	220
5.	bladeenc	9.0 %	0.83	0.16	198
6.	mp3comp	8.8 %	0.88	0.16	193
7.	lame ^a	8.6 %	0.56	0.11	190

^aThe low confidence may be due to different versions of `lame`; our training data has been encoded with the recent V 3.93.

we face a misclassification without noticing it. In this case, the classifier reports a high a-posteriori probability, also interpreted as confidence measure, and the known features are “blind” towards the differences between the two encoders. To overcome this problem, one has to search for new features between the existing and the new encoders. In the second case, the feature vector $\mathbf{f}(m_*)$ is dissimilar from the typical values of the trained encoders. Then the classifier reports a low a-posteriori probability signifying the difficulties in assigning the actual feature vector to one of the trained classes. This case is more favourable because the classification problem is identified. New encoders can be added by retraining the parameters of the classifier with an extended set.

4.3 Application for Steganalysis

To demonstrate the advances in steganalysis due to pre-classification, we assembled a test set of about 500 pristine MP3 files from different encoders together with 369 steganograms from `MP3stego` [20]. If we run the attack against `MP3stego` [23] directly on the test set, we clearly identify all 369 steganograms but face an additional 377 false positives (75.4 %). Using the proposed method as pre-classifier to filter all files from other encoders but `8hz-mp3` removes all false alarms, while still 312 steganograms are reliably detected. The miss rate of 15 % can further be reduced by using a specially trained classifier for this purpose. Only in combination with source classification does the detection method have sufficient discriminative power to be suitable for a large scale search for steganograms in MP3 files.

5. DISCUSSION AND CONCLUSION

In this paper, a method is presented to determine the encoder of ISO/ MPEG 1 Audio Layer-3 data on the basis of statistical features extracted from the data. We explained a set of 10 features that were used with a *Naïve Bayes Classifier* to discriminate between 20 different MP3 encoders. The results show, that the proposed method is quite reliable for special purpose test data as well as for a sample of arbitrary MP3 files. However, the proposed scheme is far from being the ultimate solution and it needs further refinement for real world applications.

5.1 Limitations and Future Directions

The first obstacle is the relatively narrow range of supported bit rates. In order to keep the test database operable, we decided to concentrate on the most widely used bit rates. Moreover, we tried to keep the features independent from the bit rate. This approach appears to have been effective, as we do not have any problems when classifying variable bit rate (VBR) files despite never explicitly designing a feature for VBR data. However, as some encoders change the stereo model for different bit rates—especially for more extreme settings—further analyses of the robustness of the features against bit rate changes may increase the reliability of the classification.

As already stated, MP3 files support different stereo modes and most encoders offer a variety of options to fine tune the encoding result. Since the test database always uses the (most likely) default settings and the presented features do not care about other encoding modes, sophisticated encoding parameters may cause false classifications. Hence, the influence of stereo modes and other encoding options is subject to further research.

In addition, some of the present features rely on file parameters (e.g., total file size) or precisely evaluate the beginning of a track (e.g., the initial silence). These features will fail if only fragments of a stream shall be classified.

Regarding the composition of encoders in the training set, we mainly cover open source encoders and the most widely used encoders from Fraunhofer and Xing. The versions we researched were not systematically selected. Even if we are quite confident that additional software encoders can be added with moderate effort, we still have not examined the characteristics of hardware encoders which, for example, are installed in portable digital audio recorders. The typical optimisations that are necessary to implement the MP3 encoding algorithm in DSP hardware might cause features of a different kind than those we exploit to differentiate software encoders.

To complete the list of open research questions, we refer to possible interactions between statistical features used for source classification and audio watermarking algorithms.

5.2 Transferability to Other Formats

The results on MP3 files show that encoder detection is feasible and has useful applications for steganalysis and related areas. Hence, it might be an interesting question as to whether the approach can be generalised—certainly with adapted features—to other data formats.

Obviously, the MP3 format is a good candidate for encoder detection for two reasons: First, the popularity of the format, and thus the demand for encoders, developed a market for a couple of parallel developments in the late 1990s. Second, the inclusion of a psycho-acoustic model simplifies the task of feature discovery, because it leverages small numerical differences in the signal decomposition to measurable statistics, such as block type frequencies. From this point of view, MPEG 2 audio or MPEG 4 video seem to be promising formats for similar research. Other formats, for example the popular JPEG image compression scheme, might be quite harder to classify. This format is less complicated—at least in the way it is used in the overwhelming majority of cases—and the *Independent JPEG Group* (IJG) offers a standard implementation that is included in many applications [11].

However, judging from our experience with MP3, we are confident that similar methods can be constructed for most complex standards that leave latitude for implementations. Assuming that latitude increases with complexity, we can even be quite optimistic for future formats. Some discoveries we made, for example the block type signature of open source encoders, back our optimism: As long as programmers leave identifying traces by even violating the standards, whether unintentional or motivated for one's ego, classification will be feasible. Nevertheless, it is likely to always remain as an iterative analytical task, which is difficult to automate.

5.3 Related Applications

Apart from the advances in steganalytic reliability, the proposed method may have applications in two further ways. From an academic point of view, the insights gained from the analysis of inter-encoder differences in MP3 files can be used to construct new steganographic algorithms. If we know the parameters that are treated differently by different encoders, we can consider them as indeterministic and modify them to carry steganographic messages. Also, the design of watermarking algorithms, which are robust against MP3 compression, gains from further knowledge about encoder differences.

Last but not least, a more practical application for tools derived from this approach is digital forensics. Knowledge about the encoder of a suspicious file may lead to inferences about a possible creator. However, we must note that it is still possible to fool any of the presented features, at least if some effort is spent. The output of any of these classifiers is always a probabilistic guess and must not be considered as outright proof.

6. ACKNOWLEDGEMENTS

The work on this paper was supported by the Air Force Office of Scientific Research under the research grant number FA8655-03-1-3A46. The U. S. Government is authorised to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation there on. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the Air Force Office of Scientific Research, or the U. S. Government.

The travel to the ACM Multimedia and Security Workshop has been supported in part by the European Commission through the IST Programme under contract IST-2002-507932 ECRYPT.

The authors thank Stefan Köpsell for his helpful comments.

7. REFERENCES

- [1] 8hz-mp3, 8Hz Productions, <http://www.8hz.com/mp3/>, 1998.
- [2] Brandenburg, K., Stoll, G.: ISO-MPEG-1 Audio: A Generic Standard for Coding of High-Quality Digital Audio. *Journal of the Audio Engineering Society* 42 (10), 1994, pp. 780–794.
- [3] Duda, R. O., Hart, P. E.: Pattern Classification and Scene Analysis. Wiley, New York, 1973.
- [4] EncSpot - An MP3 Analyzer, <http://www.guerillasoft.nstep.com/EncSpot2>.
- [5] Fraunhofer Institut Integrierte Schaltungen (IIS), <http://www.iis.fraunhofer.de>.
- [6] Fridrich, J.: Feature-based Steganalysis for JPEG Images and its Implications for Future Design of Steganographic Schemes. Paper presented at the Sixth Workshop on Information Hiding, Toronto, Canada, May 2004.
- [7] Hipp, M.: MPG123 – Fast MP3 Player for Linux and UNIX Systems, <http://www.mpg123.de>, 2001.
- [8] Huffman, D.: A Method for the Construction of Minimum Redundancy Codes. *Proc. of the IRE* 40, 1962, pp. 1098–1101.

- [9] IEC 958, Digital Audio Interface, International Standard, 1990.
- [10] Ihaka, R., Gentleman, R.: R – A Language for Data Analysis and Graphics. *Journal of Computational Graphics and Statistics* 5 (3), 1996, pp. 299–314.
- [11] Independent JPEG Group, <http://www.ijg.org>.
- [12] ISO/IEC 13818-3, Information Technology. Generic Coding of Moving Pictures and Associated Audio: Audio. International Standard, 1994.
- [13] Nilsson, M.: ID3v2 – The Audience is Informed, <http://www.id3.org>, 1998.
- [14] Langley, P., Iba, W., Thompson, K.: An Analysis of Bayesian Classifiers. *Proc. of the 10th Conference on Artificial Intelligence*, MIT Press, 1992, pp. 223–228.
- [15] Lindley, D. V.: Bayesian Statistics – A Review. Society for Industrial and Applied Mathematics, 1995.
- [16] Lyu, S., Farid, H.: Detecting Hidden Messages Using Higher-Order Statistics and Support Vector Machines. In Petitcolas, F. A. P. (Ed.): Information Hiding. Fifth International Workshop, LNCS 2578, Springer-Verlag, Berlin, Heidelberg, 2003, pp. 340–354.
- [17] Moddemeijer, R.: On Estimation of Entropy and Mutual Information of Continuous Distributions. *Signal Processing* 16 (3), 1989, pp. 233–246.
- [18] Pearl, L.: Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference (2nd ed.). Morgan Kaufman, New York, 1992.
- [19] Petitcolas, F. A. P., Anderson, R. J., Kuhn, M. G.: Information Hiding – A Survey. *Proc. of the IEEE* 87 (7), 1999, pp. 1062–1078.
- [20] Petitcolas, F. A. P.: MP3Stego, <http://www.cl.cam.ac.uk/~fapp2/steganography/mp3stego/>, 2002.
- [21] The R Project for Statistical Computing, <http://www.r-project.org>.
- [22] <http://sourceforge.net/projects/stego-lame>.
- [23] Westfeld, A.: Detecting Low Embedding Rates. In F. A. P. Petitcolas (Ed.): Information Hiding. Fifth International Workshop, LNCS 2578, Springer-Verlag, Berlin, Heidelberg, 2003, pp. 324–339.

APPENDIX

Table 3: Overview of features used for classification

No.	Description	Levels	Importance ^a
Size control features			
x_1	Effective bit rate ratio	4	8.35
x_2	Granule size balance	4	0.08
x_3	Reservoir usage ramp	3	5.01
x_4	Entropy of big values	cont.	2.15
Model decision features			
x_5	Preflag ratio	9	1.73
x_6	Block type transitions	5	1.56
Capability usage features			
x_7	SCFSI usage	2	0.50
x_8	Frame length alignment	2	0.92
x_9	Huffman table selection	3	0.63
Miscellaneous features			
x_{10}	Stuffing byte values	5	0.88

^aThe importance is measured with a jack-knife method: The column shows the additional overall classification error in percentage points if the feature is left out. Hence, higher values indicate higher importance of a feature.

Table 4: List of Analysed MP3 Encoders

Mnemonic	Name	Publisher	Version	Year
8hz-mp3	8HZ-MP3 Encoder	8Hz Productions	02b	1998
bladeenc	BladeEnc	Tord Jansson	0.94.2	2001
fastenc	FastEnc	Fraunhofer IIS	1.02	2000
fhgprod	Fraunhofer MP3 Producer	Opticom	2.1	1998
gogo	gogo301 petit	Herumi and Pen	3.01	2001
iTunes	Apple iTunes	Apple Computer Inc.	4.1-52	2003
l3enc272	l3enc (Linux)	Fraunhofer IIS	2.72	1997
l3encdos	l3enc (MS-DOS)	Fraunhofer IIS	2.60	1996
lame	LAME Ain't an MP3 Encoder	Mike Cheng et al.	3.93	2003
m3ec	M3E Command Line Version	N/A	0.98b	2000
mp3comp	MP3 Compressor	MP3hC	0.9f	1997
mp3enc31	mp3enc (Demo)	Fraunhofer IIS	3.1	1998
plugger	Plugger	Alberto Demichelis	0.4	1998
shine	Shine	Gabriel Bouvigne	0.1.4	2001
soloh	SoloH MPEG Encoder	N/A	0.07a	1998
soundjam	SoundJam (Macintosh)	Casady and Greene	2.5.1	2000
uzura	Uzura 3	N/A (Fortran code)	1.0	2002
xing3	Xing MP3 Encoder	Xing Technology Corp.	3.0-32	1997
xing98	Xing MP3 Encoder (x3enc)	Xing Technology Corp.	1.02	1998
xingac21	AudioCatalyst	Xing Technology Corp.	2.10	1999

Note: All trademarks are the property of their respective owners.

Table 5: Classifier performance on disjoint test data ($S \cap T = \emptyset$)

	True Encoder																			
	8hz-mp3	bladeenc	fastenc	fhgprod	gogo	iTunes	l3enc272	l3encdos	lame	m3ec	mp3comp	mp3enc31	plugger	shine	soloh	soundjam	uzura	xing3	xing98	xingac21
% of files classified as ...																				
8hz-mp3	95	-	-	-	-	-	-	-	-	-	-	-	-	-	23	-	-	-	-	-
bladeenc	-	100	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
fastenc	-	-	100	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
fhgprod	-	-	-	94	-	-	-	-	-	-	38	-	-	-	-	-	-	-	-	-
gogo	-	-	-	-	100	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
iTunes	-	-	-	-	-	100	-	-	-	-	-	2	-	-	-	-	-	-	-	-
l3enc272	-	-	-	-	-	-	84	-	-	-	-	-	-	-	-	-	-	-	-	-
l3encdos	-	-	-	-	-	-	16	100	-	-	-	-	-	-	-	-	-	-	-	-
lame	-	-	-	-	-	-	-	-	100	-	-	-	-	-	-	-	-	-	-	-
m3ec	-	-	-	-	-	-	-	-	-	100	-	-	-	-	3	-	-	-	-	-
mp3comp	-	-	-	6	-	-	-	-	-	-	62	-	-	-	-	-	-	-	-	-
mp3enc31	-	-	-	-	-	-	-	-	-	-	-	95	-	-	-	-	-	-	-	-
plugger	-	-	-	-	-	-	-	-	-	-	-	-	100	-	-	-	-	-	-	-
shine	-	-	-	-	-	-	-	-	-	-	-	-	-	100	-	-	-	-	-	-
soloh	5	-	-	-	-	-	-	-	-	-	-	-	-	-	74	-	-	-	-	-
soundjam	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100	-	-	-	-
uzura	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100	-	-	-
xing3	-	-	-	-	-	-	-	-	-	-	-	3	-	-	-	-	-	85	13	-
xing98	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	15	87	-
xingac21	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100

Basis: $|T| = |S| \approx 1200$ files, $n = 20$ encoders, $k = 10$ features, overall error rate: 5.1 %, average classification confidence $\bar{P}_{\max} = 0.959$