



US006968337B2

(12) **United States Patent**
Wold

(10) **Patent No.:** **US 6,968,337 B2**
(45) **Date of Patent:** **Nov. 22, 2005**

(54) **METHOD AND APPARATUS FOR IDENTIFYING AN UNKNOWN WORK**

(Continued)

FOREIGN PATENT DOCUMENTS

- (75) Inventor: **Erling H. Wold**, El Cerrito, CA (US)
- (73) Assignee: **Audible Magic Corporation**, Los Gatos, CA (US)

EP	0402210 A *	12/1990	G06F/11/08
EP	0517405 A2 *	12/1992	G07C/9/00
WO	WO0123981 A1 *	4/2001	G06F/1/00
WO	WO 02/15035 A2	2/2002	G06F/17/00

(*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 458 days.

OTHER PUBLICATIONS

L. Baum et al., A Maximization Technique Occuring in the Statistical Analysis of Probabilistic Functions of Markov Chains, *The Annals of Mathematical Statistics*, vol. 41, No. 1 pp. 164-171, 1970 (no month).

(21) Appl. No.: **10/192,783**

(22) Filed: **Jul. 9, 2002**

(Continued)

(65) **Prior Publication Data**

US 2003/0023852 A1 Jan. 30, 2003

Primary Examiner—Frantz Coby
(74) *Attorney, Agent, or Firm*—Sierra Patent Group, Ltd.

Related U.S. Application Data

(60) Provisional application No. 60/304,647, filed on Jul. 10, 2001.

(51) **Int. Cl.**⁷ **G06F 7/00**

(52) **U.S. Cl.** **707/100; 707/101; 707/102; 707/103 R; 707/104.1**

(58) **Field of Search** 707/100, 101, 707/102, 103 R, 104.1; 705/26; 386/83; 725/22; 380/54, 28

ABSTRACT

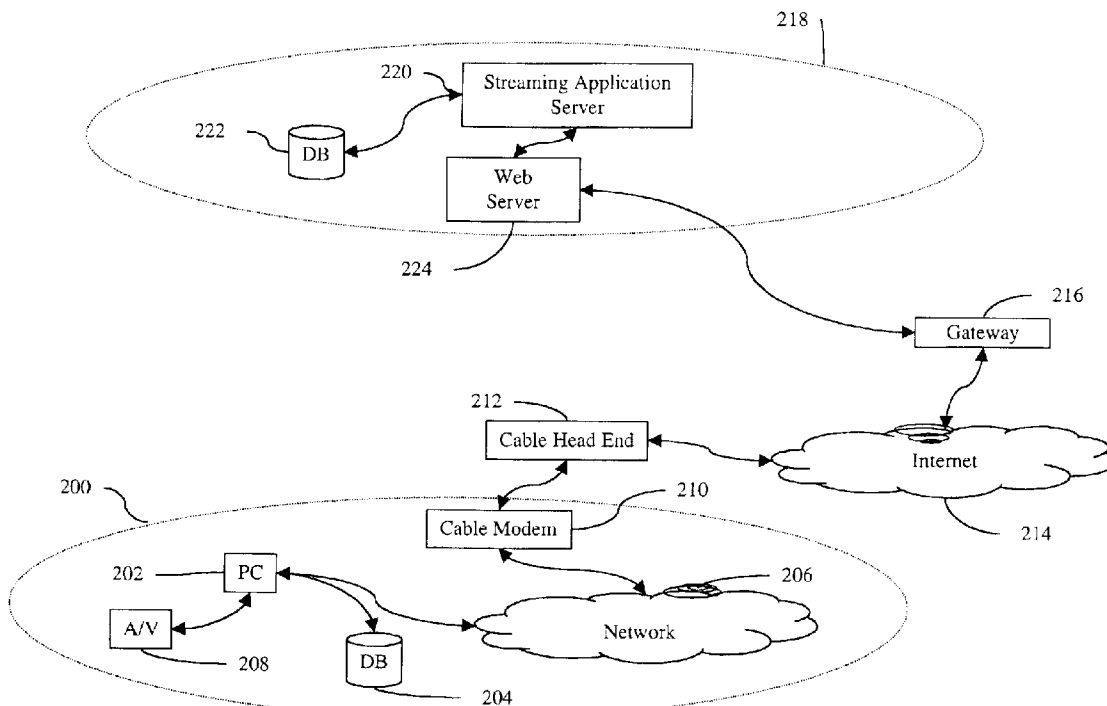
A system for determining an identity of a received work. The system receives audio data for an unknown work. The audio data is divided into segments. The system generates a signature of the unknown work from each of the segments. Reduced dimension signatures are then generated at least a portion of the signatures. The reduced dimension signatures are then compared to reduced dimensions signatures of known works that are stored in a database. A list of candidates of known works is generated from the comparison. The signatures of the unknown works are then compared to the signatures of the known works in the list of candidates. The unknown work is then identified as the known work having signatures matching within a threshold.

(56) **References Cited**

U.S. PATENT DOCUMENTS

3,919,479 A	11/1975	Moon et al.	179/1 SB
4,230,990 A	10/1980	Lert, Jr. et al.	455/67
4,450,531 A	5/1984	Kenyon et al.	364/604

26 Claims, 9 Drawing Sheets



U.S. PATENT DOCUMENTS

4,677,455	A	6/1987	Okajima	357/38
4,677,466	A	6/1987	Lert, Jr. et al.	358/84
4,739,398	A	4/1988	Thomas et al.	358/84
4,843,562	A	6/1989	Kenyon et al.	364/487
4,918,730	A	4/1990	Schulze	381/43
5,210,820	A	5/1993	Kenyon	395/2
5,283,819	A	2/1994	Glick et al.	379/90
5,437,050	A	7/1995	Lamb et al.	455/2
5,504,518	A *	4/1996	Ellis et al.	725/22
5,581,658	A	12/1996	O'Hagan et al.	395/22
5,613,004	A *	3/1997	Cooperman et al.	380/28
5,710,916	A	1/1998	Barbara et al.	395/609
5,918,223	A	6/1999	Blum et al.	707/1
5,930,369	A *	7/1999	Cox et al.	380/54
5,949,885	A *	9/1999	Leighton	380/54
6,006,256	A	12/1999	Zdepski et al.	709/217
6,011,758	A	1/2000	Dockes et al.	369/30
6,026,439	A	2/2000	Chowdhury et al.	709/233
6,044,402	A	3/2000	Jacobson et al.	709/225
6,118,450	A	9/2000	Proehl et al.	345/349
6,192,340	B1	2/2001	Abecassis	704/270
6,243,615	B1	6/2001	Neway et al.	700/108
6,253,193	B1	6/2001	Ginter et al.	705/57
6,253,337	B1	6/2001	Maloney et al.	714/38
6,422,061	B1	7/2002	Sunshine et al.	73/29.01
6,771,885	B1 *	8/2004	Agnihotri et al.	386/83
2002/0082999	A1	6/2002	Lee et al.	705/51
2002/0087885	A1	7/2002	Peled et al.	713/201

2002/0198789	A1 *	12/2002	Waldman	705/26
--------------	------	---------	---------	--------

OTHER PUBLICATIONS

A. P. Dempster et al. "Maximum Likelihood from Incomplete Data via the EM Algorithm", *Journal of the Royal Statistical Society, Series B (Methodological)*, vol. 39, Issue 1, pp. 1–38, 1977 (no month).

D. Reynolds et al., "Robust Text-Independent Speaker Identification Using Gaussian Mixture Speaker Models", *IEEE Transactions on Speech and Audio Processing*, vol. 3, No. 1, pp. 72–83, Jan. 1995.

B. Pellom et al., "Fast Likelihood Computation Techniques in Nearest-Neighbor Based search for Continuous Speech Recognition", *IEEE Signal Processing Letters*, vol. 8, No. * pp. 221–224, Aug. 2001.

J. Haitisma et al., "Robust Audio hashing for Content Identification", *CBMI 2001, Second International Workshop on Content Based Multimedia and Indexing*, Sep. 19–21, 2001, Brescia, Italy, Sep. 19–21, 2001.

PacketHound Tech Specs, palisadesys.com/products/packethound/tech_specs/prod_Phtechspecs.shtml, 2002 (no month).

"How Does PacketHound Work?", palisadesys.com/products/packethound/how_does_itwork/prod_Phhow.shtml, 2002 (no month).

* cited by examiner

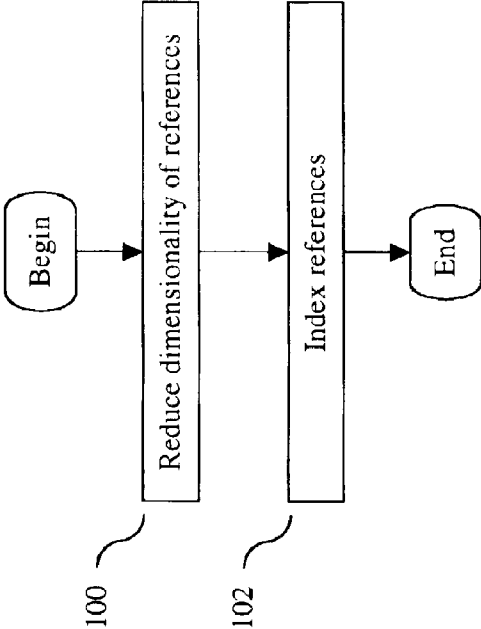


FIG. 1A
Present Invention

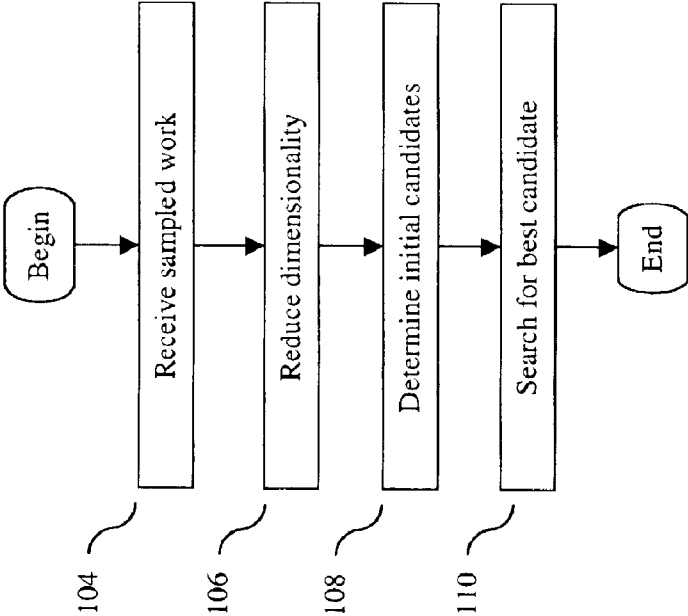


FIG. 1B
Present Invention

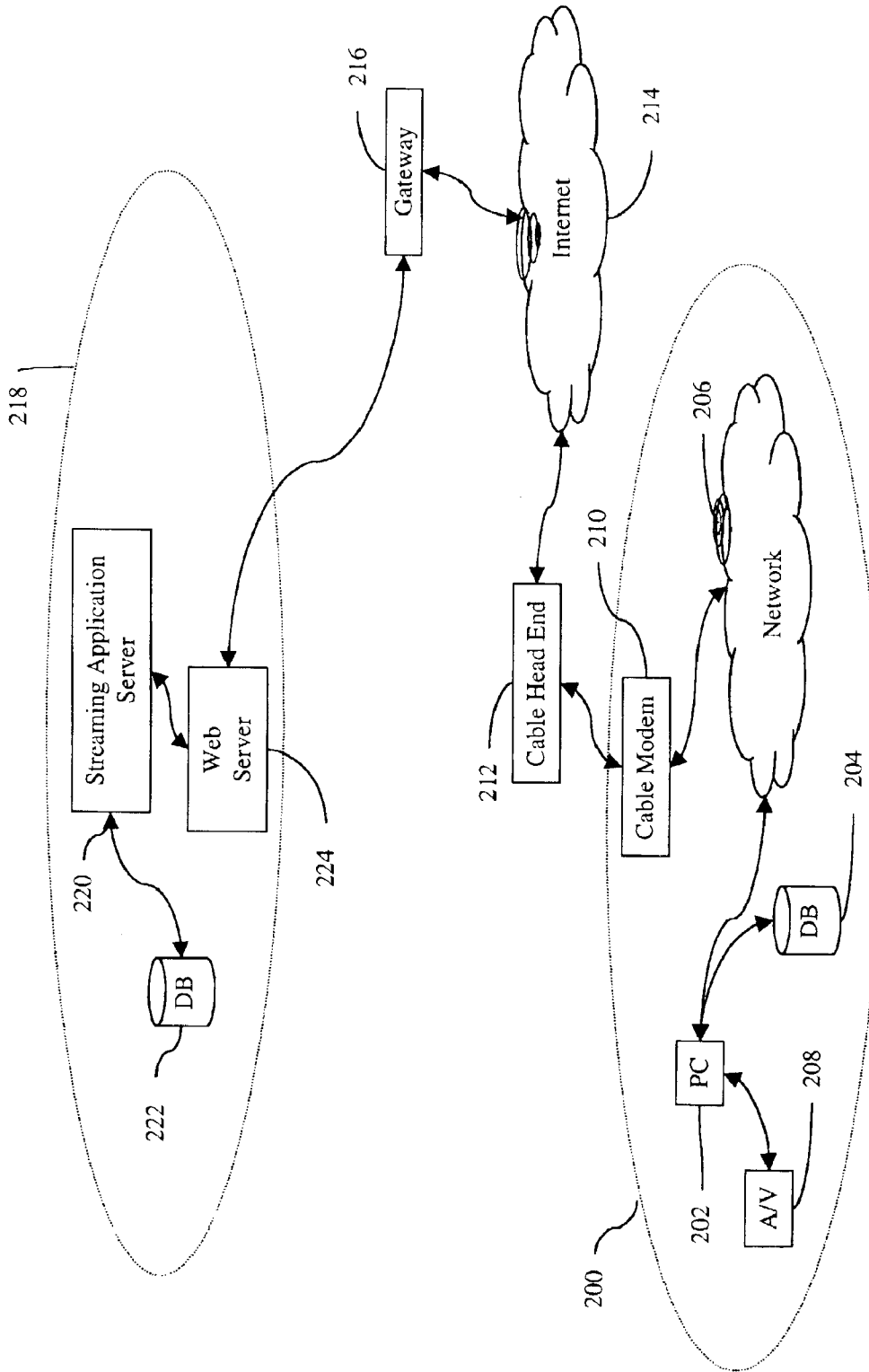


FIG. 2
Present Invention

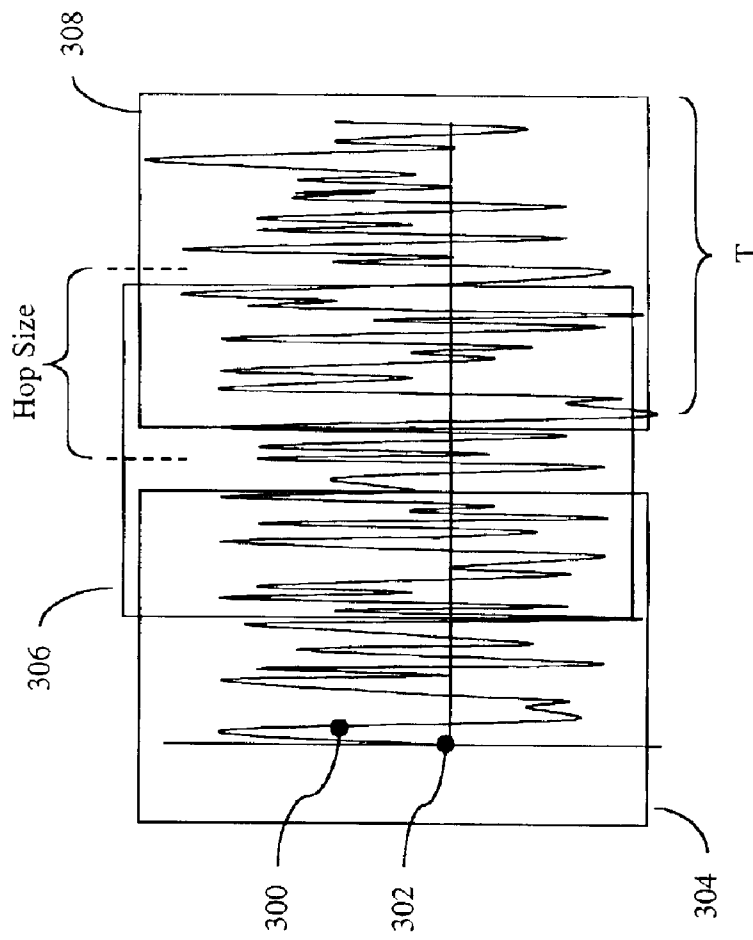


FIG. 3
Present Invention

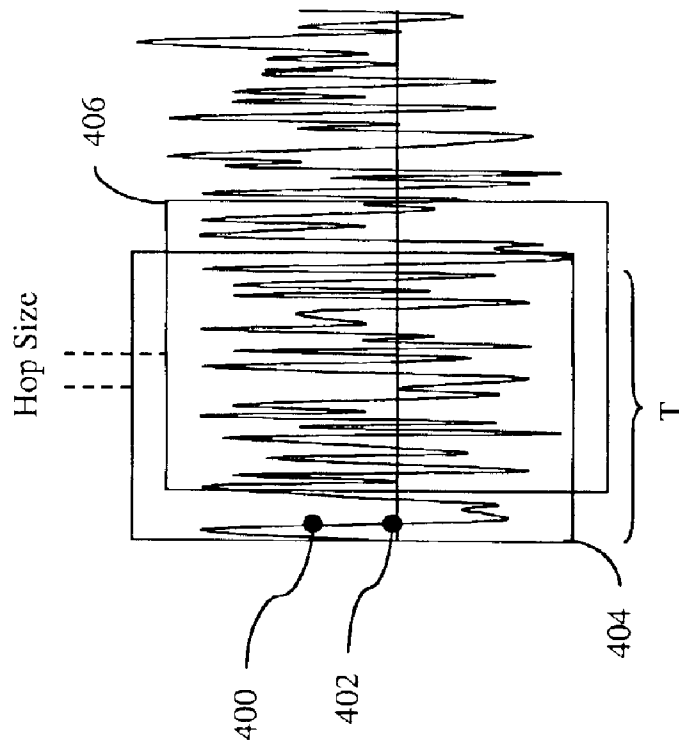


FIG. 4
Present Invention

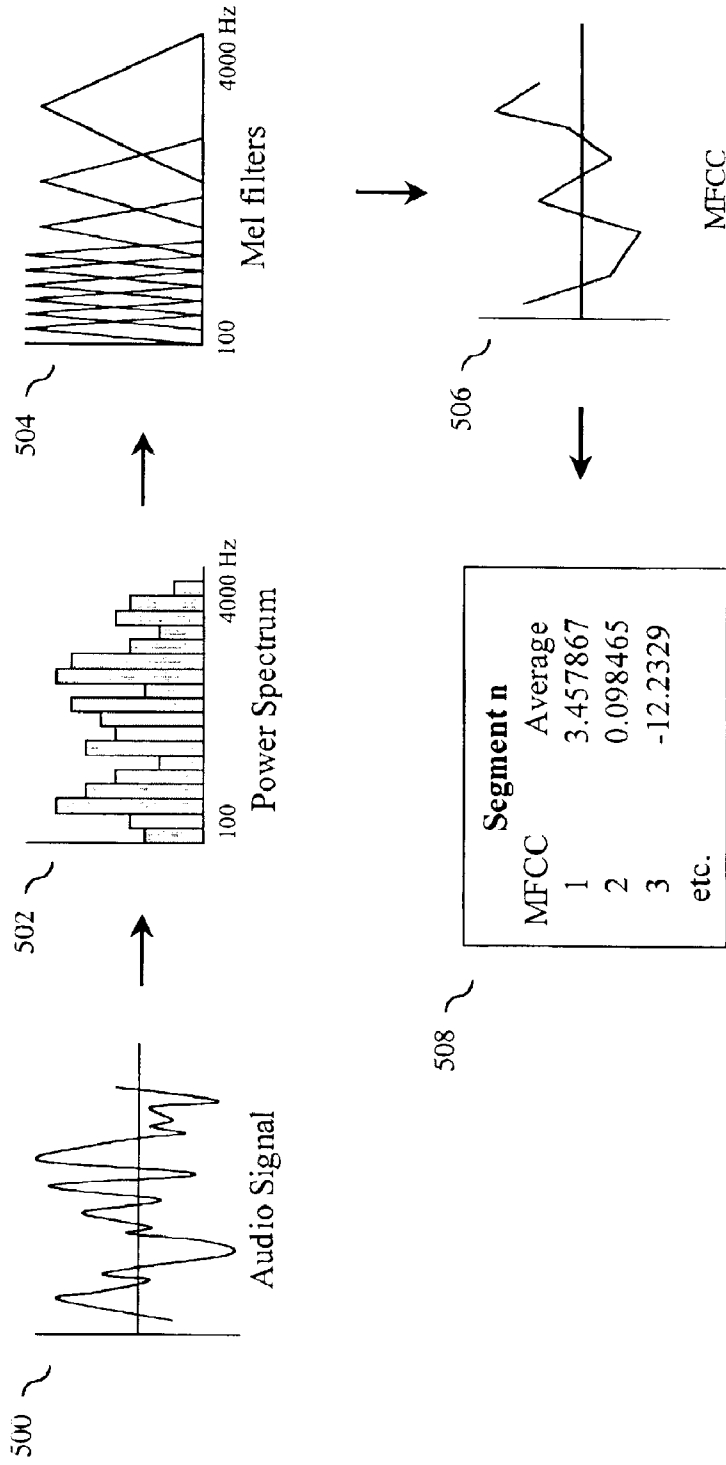


FIG. 5
Present Invention

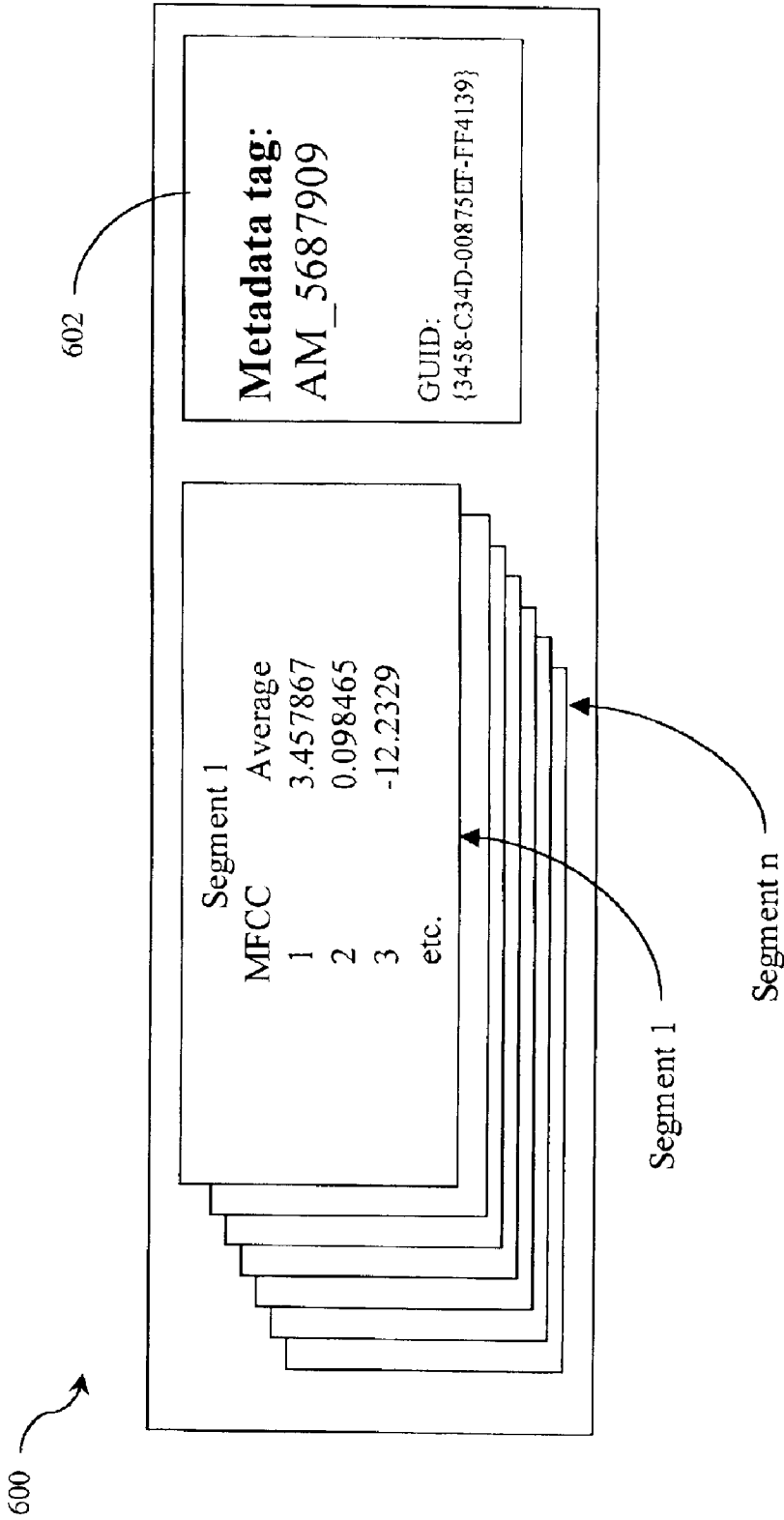


FIG. 6
Present Invention

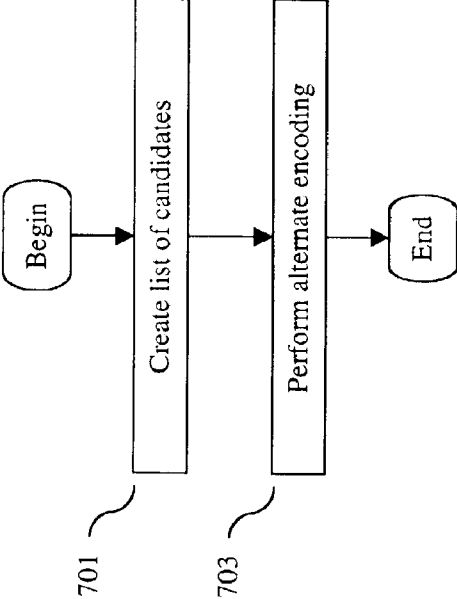


FIG. 7A
Present Invention

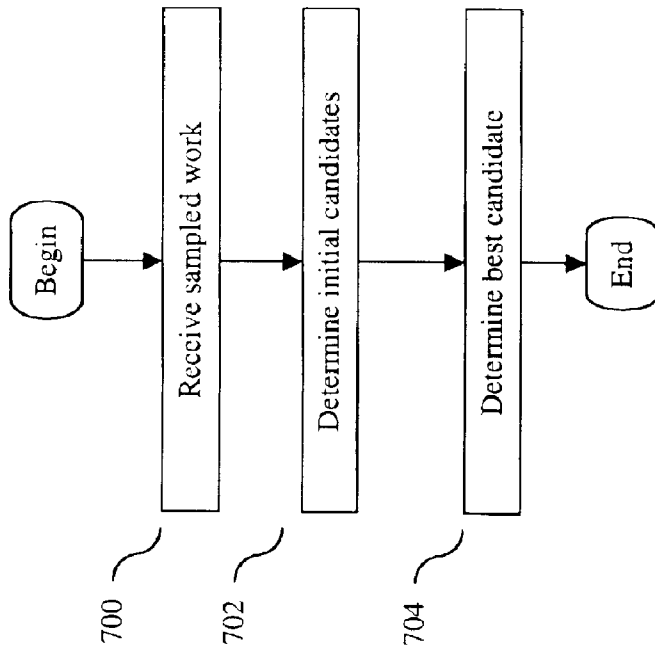


FIG. 7B
Present Invention

METHOD AND APPARATUS FOR IDENTIFYING AN UNKNOWN WORK

PRIORITY CLAIM

This application claims the benefit of U.S. Provisional Application Ser. No. 60/304,647, filed Jul. 10, 2001.

BACKGROUND OF THE INVENTION

1. Field of the Invention

The present invention relates to data communications. In particular, the present invention relates to a novel method and apparatus for identifying an unknown work.

BACKGROUND

2. The Prior Art

Digital audio technology has greatly changed the landscape of music and entertainment. Rapid increases in computing power coupled with decreases in cost have made it possible for individuals to generate finished products having a quality once available only in a major studio. One consequence of modern technology is that legacy media storage standards, such as reel-to-reel tapes, are being rapidly replaced by digital storage media, such as the Digital Versatile Disk (DVD), and Digital Audio Tape (DAT). Additionally, with higher capacity hard drives standard on most personal computers, home users may now store digital files such as audio or video tracks on their home computers.

Furthermore, the Internet has generated much excitement, particularly among those who see the Internet as an opportunity to develop new avenues for artistic expression and communication. The Internet has become a virtual gallery, where artists may post their works on a Web page. Once posted, the works may be viewed by anyone having access to the Internet.

One application of the Internet that has received considerable attention is the ability to transmit recorded music over the Internet. Once music has been digitally encoded, the audio may be both downloaded by users for play, or broadcast ("streamed") over the Internet. When audio is streamed, it may be listened to by Internet users in a manner much like traditional radio stations.

Given the widespread use of digital media, digital audio files, or digital video files containing audio information, may need to be identified. The need for identification of digital files may arise in a variety of situations. For example, an artist may wish to verify royalty payments or generate their own Arbitron®-like ratings by identifying how often their works are being streamed or downloaded. Additionally, users may wish to identify a particular work. The prior art has made efforts to create methods for identifying digital audio works.

However, systems of the prior art suffer from certain disadvantages. One area of difficulty arises when a large number of reference signatures must be compared to an unknown audio recording.

The simplest method for comparing an incoming audio signature (which could be from a file on the Internet, a recording of a radio or Internet radio broadcast, a recording from a cell phone, etc) to a database of reference signatures for the purpose of identification is to simply compare the incoming signature to every element of the database. However, since it may not be known where the reference signatures might have occurred inside the incoming signature, this comparison must be done at many time

locations within the incoming signature. Each individual signature-to-signature comparison at each point in time may also be done in a "brute-force" manner using techniques known in the art; essentially computing the full Euclidean distance between the entire signatures' feature vectors. A match can then be declared when one of these comparisons yields a score or distance that is above or below some threshold, respectively.

However, when an audio signature or fingerprint contains a large number of features such a brute-force search becomes too expensive computationally for real-world databases which typically have several hundred thousand to several million signatures.

Many researchers have worked on methods for multi-dimensional indexing, although the greatest effort has gone into geographical (2-dimensional) or spatial (3-dimensional) data. Typically, all of these methods order the elements of the database based on their proximity to each other.

For example, the elements of the database can be clustered into hyper-spheres or hyper-rectangles, or the space can be organized into a tree form by using partitioning planes. However, when the number of dimensions is large (on the order of 15 or more), it can be shown mathematically that more-or-less uniformly distributed points in the space all become approximately equidistant from each other. Thus, it becomes impossible to cluster the data in a meaningful way, and comparisons can become both lengthy and inaccurate.

Hence, there exists a need to provide a means for data comparison which overcomes the disadvantages of the prior art.

BRIEF DESCRIPTION OF THE INVENTION

A method and apparatus for identifying an unknown work is disclosed. In one aspect, a method may include the acts of providing a reference database having a reduced dimensionality containing signatures of sampled works; receiving a sampled work; producing a signature from the work; and reducing the dimensionality of the signature.

BRIEF DESCRIPTION OF THE DRAWING FIGURES

FIG. 1A is a flowchart of a method according to the present invention.

FIG. 1B is a flowchart of another method according to the present invention.

FIG. 2 is a diagram of a system suitable for use with the present invention.

FIG. 3 is a diagram of segmenting according to the present invention.

FIG. 4 is a detailed diagram of segmenting according to the present invention showing hop size.

FIG. 5 is a graphical flowchart showing the creating of a segment feature vector according to the present invention.

FIG. 6 is a diagram of a signature according to the present invention.

FIG. 7A is a flowchart of a method for preparing a reference database according to the present invention.

FIG. 7B is a flowchart of method for identifying an unknown work according to the present invention.

DETAILED DESCRIPTION OF PREFERRED EMBODIMENTS

Persons of ordinary skill in the art will realize that the following description of the present invention is illustrative

only and not in any way limiting. Other embodiments of the invention will readily suggest themselves to such skilled persons having the benefit of this disclosure.

It is contemplated that the present invention may be embodied in various computer and machine-readable data structures. Furthermore, it is contemplated that data structures embodying the present invention will be transmitted across computer and machine-readable media, and through communications systems by use of standard protocols such as those used to enable the Internet and other computer networking standards.

The invention further relates to machine-readable media on which are stored embodiments of the present invention. It is contemplated that any media suitable for storing instructions related to the present invention is within the scope of the present invention. By way of example, such media may take the form of magnetic, optical, or semiconductor media.

The present invention may be described through the use of flowcharts. Often, a single instance of an embodiment of the present invention will be shown. As is appreciated by those of ordinary skill in the art, however, the protocols, processes, and procedures described herein may be repeated continuously or as often as necessary to satisfy the needs described herein. Accordingly, the representation of the present invention through the use of flowcharts should not be used to limit the scope of the present invention.

The present invention may also be described through the use of web pages in which embodiments of the present invention may be viewed and manipulated. It is contemplated that such web pages may be programmed with web page creation programs using languages standard in the art such as HTML or XML. It is also contemplated that the web pages described herein may be viewed and manipulated with web browsers running on operating systems standard in the art, such as the Microsoft Windows® and Macintosh® versions of Internet Explorer® and Netscape®. Furthermore, it is contemplated that the functions performed by the various web pages described herein may be implemented through the use of standard programming languages such as Java® or similar languages.

The present invention will first be described in general overview. Then, each element will be described in further detail below.

Referring now to FIG. 1A, a flowchart is shown which provides a general overview of the present invention as related to the preparation of a database of reference signatures. Two overall acts are performed to prepare a reference database in accordance with the present invention: in act 100, the present invention reduces the dimensionality of reference signatures; and the reference database is indexed in act 102.

Referring now to FIG. 1B, a flowchart is shown which provides a general overview of the present invention as related to the identification of an unknown signature in accordance with the present invention. In act 104, a sampled work is received. In act 106, the present invention reduces the dimensionality of the received work. In act 108, the present invention determines initial candidates. In act 110, the present invention searches for the best candidate.

Prior to presenting a detailed overview of each act of FIGS. 1A and 1B, some background will first be presented.

Structural Embodiment of the Present Invention

Referring now to FIG. 2, a diagram of a system suitable for use with the present invention is shown. FIG. 2 includes a client system 200. It is contemplated that client system 200

may comprise a personal computer 202 including hardware and software standard in the art to run an operating system such as Microsoft Windows®, MAC OS® Palm OS, UNIX, or other operating systems standard in the art. Client system 200 may further include a database 204 for storing and retrieving embodiments of the present invention. It is contemplated that database 204 may comprise hardware and software standard in the art and may be operatively coupled to PC 202. Database 204 may also be used to store and retrieve the works and segments utilized by the present invention.

Client system 200 may further include an audio/video (A/V) input device 208. A/V device 208 is operatively coupled to PC 202 and is configured to provide works to the present invention which may be stored in traditional audio or video formats. It is contemplated that A/V device 208 may comprise hardware and software standard in the art configured to receive and sample audio works (including video containing audio information), and provide the sampled works to the present invention as digital audio files. Typically, the A/V input device 208 would supply raw audio samples in a format such as 16-bit stereo PCM format. A/V input device 208 provides an example of means for receiving a sampled work.

It is contemplated that sampled works may be obtained over the Internet, also. Typically, streaming media over the Internet is provided by a provider, such as provider 218 of FIG. 2. Provider 218 includes a streaming application server 220, configured to retrieve works from database 222 and stream the works in a format standard in the art, such as Real®, Windows Media®, or QuickTime®. The server then provides the streamed works to a web server 224, which then provides the streamed work to the Internet 214 through a gateway 216. Internet 214 may be any packet-based network standard in the art, such as IP, Frame Relay, or ATM.

To reach the provider 218, the present invention may utilize a cable or DSL head end 212 standard in the art operatively, which is coupled to a cable modem or DSL modem 210 which is in turn coupled to the system's network 206. The network 206 may be any network standard in the art, such as a LAN provided by a PC 202 configured to run software standard in the art.

It is contemplated that the sampled work received by system 200 may contain audio information from a variety of sources known in the art, including, without limitation, radio, the audio portion of a television broadcast, Internet radio, the audio portion of an Internet video program or channel, streaming audio from a network audio server, audio delivered to personal digital assistants over cellular or wireless communication systems, or cable and satellite broadcasts.

Additionally, it is contemplated that the present invention may be configured to receive and compare segments coming from a variety of sources either stored or in real-time. For example, it is contemplated that the present invention may compare a real-time streaming work coming from streaming server 218 or A/V device 208 with a reference segment stored in database 204.

Segmenting Background

It is contemplated that a wide variety of sampled works may be utilized in the present invention. However, the inventors have found the present invention especially useful with segmented works. An overview of a segmented work will now be provided.

FIG. 3 shows a diagram showing the segmenting of a work according to the present invention. FIG. 3 includes

audio information **300** displayed along a time axis **302**. FIG. **3** further includes a plurality of segments **304**, **306**, and **308** taken of audio information **300** over some segment size **T**.

In an exemplary non-limiting embodiment of the present invention, instantaneous values of a variety of acoustic features are computed at a low level, preferably about 100 times a second. In particular, 10 MFCCs (cepstral coefficients) are computed. It is contemplated that any number of MFCCs may be computed. Preferably, 5–20 MFCCs are computed, however, as many as 30 MFCCs may be computed, depending on the need for accuracy versus speed.

Segment-level features are disclosed U.S. Pat. No. 5,918,223 to Blum, et al., which is assigned to the assignee of the current disclosure and incorporated by reference as though fully set forth herein. In an exemplary non-limiting embodiment of the present invention, the segment-level acoustical features comprise statistical measures as disclosed in the '223 patent of low-level features calculated over the length of each segment. The data structure may store other book-keeping information as well (segment size, hop size, item ID, UPC, etc). As can be seen by inspection of FIG. **3**, the segments **304**, **306**, and **308** may overlap in time. This amount of overlap may be represented by measuring the time between the center point of adjacent segments. This amount of time is referred to herein as the hop size of the segments, and is so designated in FIG. **3**. By way of example, if the segment length **T** of a given segment is one second, and adjacent segments overlap by 50%, the hop size would be 0.5 second.

The hop size may be set during the development of the software. Additionally, the hop sizes of the reference database and the real-time signatures may be predetermined to facilitate compatibility. For example, the reference signatures in the reference database may be precomputed with a fixed hop and segment size, and thus the client applications should conform to this segment size and have a hop size which integrally divides the reference signature hop size. It is contemplated that one may experiment with a variety of segment sizes in order to balance the tradeoff of accuracy with speed of computation for a given application.

The inventors have found that by carefully choosing the hop size of the segments, the accuracy of the identification process may be significantly increased. Additionally, the inventors have found that the accuracy of the identification process may be increased if the hop size of reference segments and the hop size of segments obtained in real-time are each chosen independently. The importance of the hop size of segments may be illustrated by examining the process for segmenting pre-recorded works and real-time works separately.

Reference Signatures

Prior to attempting to identify a given work, a reference database of signatures must be created. When building a reference database, a segment length having a period of less than three seconds is preferred. In an exemplary non-limiting embodiment of the present invention, the segment lengths have a period ranging from 0.5 seconds to 3 seconds. For a reference database, the inventors have found that a hop size of approximately 50% to 100% of the segment size is preferred.

It is contemplated that the reference signatures may be stored on a database such as database **204** as described above. Database **204** and the discussion herein provide an example of means for providing a plurality of reference signatures each having a segment size and a hop size.

Unknown Signatures

The choice of the hop size is important for the signatures of the audio to be identified, hereafter referred to as “unknown audio.”

FIG. **4** shows a detailed diagram of the segmentation of unknown audio according to the present invention. FIG. **4** includes unknown audio information **400** displayed along a time axis **402**. FIG. **4** further includes segments **404** and **406** taken of audio information **400** over some segment length **T**. In an exemplary non-limiting embodiment of the present invention, the segment length of unknown audio segments is chosen to range from 0.5 to 3 seconds.

As can be seen by inspection of FIG. **4**, the hop size of unknown audio segments is chosen to be smaller than that of reference segments. In an exemplary non-limiting embodiment of the present invention, the hop size of unknown audio segments is less than 50% of the segment size. In yet another exemplary non-limiting embodiment of the present invention, the unknown audio hop size may be 0.1 seconds.

The inventors have found such a small hop size advantageous for the following reasons. The ultimate purpose of generating unknown audio segments is to analyze and compare them with the reference segments in the database to look for matches. The inventors have found at least two major reasons why an unknown audio recording would not match its counterpart in the database. One is that the broadcast channel does not produce a perfect copy of the original. For example, the work may be edited or processed or the announcer may talk over part of the work. The other reason is that larger segment boundaries may not line up in time with the original segment boundaries of the target recordings.

The inventors have found that by choosing a smaller hop size, some of the segments will ultimately have time boundaries that line up with the original segments, notwithstanding the problems listed above. The segments that line up with a “clean” segment of the work may then be used to make an accurate comparison while those that do not so line up may be ignored. The inventors have found that a hop size of 0.1 seconds seems to be the maximum that would solve this time shifting problem.

As mentioned above, once a work has been segmented, the individual segments are then analyzed to produce a segment feature vector. FIG. **5** is a diagram showing an overview of how the segment feature vectors may be created using the methods described in U.S. Pat. No. 5,918,223 to Blum, et al. It is contemplated that a variety of analysis methods may be useful in the present invention, and many different features may be used to make up the feature vector. The inventors have found that the pitch, brightness, bandwidth, and loudness features of the '223 patent to be useful in the present invention. Additionally, spectral features may be used analyzed, such as the energy in various spectral bands. The inventors have found that the cepstral features (MFCCs) are very robust (more invariant) given the distortions typically introduced during broadcast, such as EQ, multi-band compression/limiting, and audio data compression techniques such as MP3 encoding/decoding, etc.

In act **500**, the audio segment is sampled to produce a segment. In act **502**, the sampled segment is then analyzed using Fourier Transform techniques to transform the signal into the frequency domain. In act **504**, mel frequency filters are applied to the transformed signal to extract the significant audible characteristics of the spectrum. In act **506**, a Discrete Cosine Transform is applied which converts the signal into mel frequency cepstral coefficients (MFCCs).

Finally, in act 508, the MFCCs are then averaged over a predetermined period. In an exemplary non-limiting embodiment of the present invention, this period is approximately one second. Additionally, other characteristics may be computed at this time, such as brightness or loudness. A segment feature vector is then produced which contains a list containing at least the 10 MFCCs corresponding average.

The disclosure of FIGS. 3, 4, and 5 provide examples of means for creating a signature of a sampled work having a segment size and a hop size.

FIG. 6 is a diagram showing a complete signature 600 according to the present invention. Signature 600 includes a plurality of segment feature vectors 1 through n generated as shown and described above. Signature 600 may also include an identification portion containing a unique ID. It is contemplated that the identification portion may contain a unique identifier provided by the RIAA (Recording Industry Association of America) or some other audio authority or cataloging agency. The identification portion may also contain information such as the UPC (Universal Product Code) of the various products that contain the audio corresponding to this signature. Additionally, it is contemplated that the signature 600 may also contain information pertaining to the characteristics of the file itself, such as the hop size, segment size, number of segments, etc., which may be useful for storing and indexing.

Signature 600 may then be stored in a database and used for comparisons.

The following computer code in the C programming language provides an example of a database structure in memory according to the present invention:

```
typedef struct
{
    float hopSize;           /* hop size */
    float segmentSize;      /* segment size */
    MFSignature* signatures; /* array of signatures */
} MFDatabase;
```

The following provides an example of the structure of a segment according to the present invention:

```
typedef struct
{
    char* id;                /* unique ID for this audio clip */
    long numSegments;       /* number of segments */
    float* features;        /* feature array */
    long size;              /* size of per-segment feature vector */
    float hopSize;
    float segmentSize;
} MFSignature;
```

The discussion of FIG. 6 provides an example of means for storing segments and signatures according to the present invention.

A more detailed description of the operation of the present invention will now be provided.

Referring now to FIG. 7A, a flowchart showing one aspect of a method according to the present invention is presented.

Reference Database Preparation

Prior to the identification of an unknown sample, a database of reference signatures is prepared in accordance with the present invention.

In an exemplary non-limiting embodiment of the present invention, a reference signature may comprise an audio

signature derived from a segmentation of the original audio work as described above. In a presently preferred embodiment, reference signatures have 20 non-overlapping segments, where each segment is one second in duration, with one-second spacing from center to center, as described above. Each of these segments is represented by 10 Mel filtered cepstral coefficients (MFCCs), resulting in a feature vector of 200 dimensions. Since indexing a vector space of this dimensionality is not practical, the number of dimensions used for the initial search for possible candidates is reduced according to the present invention.

Reducing the Dimensionality

FIG. 7A is a flowchart of dimension reduction according to the present invention. The number of dimensions used for the initial search for possible candidates is reduced, resulting in what the inventors refer to as a subspace. By having the present invention search a subspace at the outset, the efficiency of the search may be greatly increased.

Referring now to FIG. 7A, the present invention accomplishes two tasks to develop this subspace: (1) the present invention uses less than the total number of segments in the reference signatures in act 701; and (2) the present invention performs a principal components analysis to reduce the dimensionality in act 703.

Using Less Segments to Perform an Initial Search

The inventors empirically have found that using data from two consecutive segments (i.e., a two-second portion of the signature) to search for approximately 500 candidates is a good tradeoff between computation complexity and accuracy. The number of candidates can be altered for different applications where either speed or accuracy is more or less important.

For example, the present invention may be configured to extract a predetermined percentage of candidates. In an exemplary non-limiting embodiment of the present invention, a list of candidates may comprise 2% of the size of the reference signature database when using 2 segments for the initial search. In another exemplary non-limiting embodiment of the present invention, a list of candidates may be those reference signatures whose distances based on the initial 2 segment search are below a certain threshold.

As will be appreciated by those of ordinary skill in the art, the dimension reduction of the present invention may be used to perform initial search using fewer segments for data other than MFCC-based feature vectors. It is contemplated that any feature-based vector set may be used in the present invention.

Furthermore, the segments used in the initial search do not have to be the same size as the segments used for the final search. Since it may be better to use as few dimensions as possible in the initial search for candidates, a smaller segment size is advantageous here. The full segment size can then be used in the final search. In an exemplary non-limiting embodiment of the current invention, the initial search may use the higher-order MFCCs (since these are the most robust)—this is a simple way to reduce the dimensionality.

In the next section, we will discuss another, more sophisticated, method for reducing the segment size for the initial candidate search.

Perform Alternate Encoding

The second step is to use an alternate encoding of the MFCC data which has the same information but with fewer features.

To accomplish this, the present invention first performs an eigenanalysis of N candidates to determine the principal components of the MFCCs for our typical audio data. In an

exemplary non-limiting embodiment of the present invention, the present invention examines 25,000 audio signatures of 20 segments each—each taken from a different recording, which gives provides 500,000 sets of MFCCs. The inventors have found that this is enough to be a good statistical sample of the feature vectors.

As is appreciated by those of ordinary skill in the art, the number examined in the present invention may be adjusted to provide a good statistical sample of different kinds of music. For example, 100 or a 1000 segments may be satisfactory.

Next, a Karhunen-Loève transformation is derived. Each set of 10 MFCCs becomes a column of a matrix A . We then compute $A^T A$ and find the 10 eigenvalues and eigenvectors of this matrix. Sorting the eigenvectors by eigenvalue (largest eigenvalue first) results in a list of orthogonal basis vectors that are the principal components of the segment data. For a database of typical music recordings, 95% of the information in the MFCCs is contained in the first 7 components of this new basis.

As is known by those having ordinary skill in the art, the Karhunen-Loève transformation is represented by the matrix that has the all 10 of the above eigenvectors as its rows. This transformation is applied to all the segments of all the reference signatures in the database as well as to all the segments of any signatures that are to be identified. This allows approximate distances to be computed by using the first few components of the transformed segment MFCC vectors for a small tradeoff in accuracy. Most importantly, it reduces the initial search dimension to 14 (7 components times 2 segments), which can be indexed with reasonable efficiency.

As will be appreciated by those of ordinary skill in the art, dimension reduction according to the present invention may be utilized to examine subspaces for feature sets other than MFCCs. The dimension reduction of the present invention may be applied to any set of features since such sets comprise vectors of floating point numbers. For example, given a feature vector comprising spectral coefficients and loudness, one could still apply the KL-process of the present invention to yield a smaller and more easily searched feature vector.

Furthermore, the transform of the present invention may be applied to each segment separately. For example, prior art identification methods may use a single 30-second segment of sound over which they compute an average feature vector. Of course, the accuracy of such methods are much lower, but the process of the present invention may work for such features as well. Moreover, such prior art methods may be used as an initial search.

The dimension reduction aspect of the present invention provides significant efficiency gains over prior art methods. For example, in a “brute force” method, the signature of the incoming sampled work is tested against every reference signature in the database. This is time-consuming because the comparison of any two signatures is a 200-dimensional comparison and because there are a lot of reference signatures in the database. Either alone are not unsatisfactory, but both together takes a long time. The present invention solves the first problem by searching only a subspace, i.e., using less than all 200 dimensions in the comparison.

In addition to the raw speedup given by searching a subspace, the reduced dimensionality also allows one to practically index the database of reference signatures. As mentioned above, it is impractical to index a 200-dimensional database, but 14 is practical.

The present disclosure thus provides for several manners in which the dimensionality may be reduced:

- (1) searching for the top N candidates over a subspace;
- (2) searching for the top N candidates using less than the total number of segments from the reference signature;
- (3) searching for the top N candidates by projecting the reference signatures and signature of the work to be identified onto a subspace; and
- (4) searching for the top N candidates by projecting the reference signatures and signature of the work to be identified onto a subspace, where the subspace is determined by a Karhunen-Loève transformation.

The preparation of the reference database may occur at any time. For example, the results of the preparation may occur each time the server is started up. Additionally, the results could be saved and reused from then on, or the results may be prepared once and used over again. It may need to be recomputed whenever a new reference signature is added to the database.

Computing the Index

The present invention may also compute an index of the reference signatures. As is appreciated by those having ordinary skill in the art, many indexing strategies are available for use in the present invention. Examples include the k-d tree, the SS-tree, the R-tree, the SR-tree, and so on. Any look-up method known in the art may be used in the present disclosure. Common to all indexing strategies is that the multidimensional space is broken into a hierarchy of regions which are then structured into a tree. As one progress down the tree during the search process, the regions become smaller and have fewer elements. All of these trees have tradeoffs that affect the performance under different conditions, e.g., whether the entire tree fits into memory, whether the data is highly clustered, and so on.

In an exemplary non-limiting embodiment of the present invention, a binary k-d tree indexing method is utilized. This is a technique well-known in the art, but a brief overview is given here. At the top level, the method looks to see which dimension has the greatest extent, and generates a hyperplane perpendicular to this dimension that splits the data into two regions at its median. This yields two subspaces on either side of the plane. This process is continued by recursion on the data in each of these subspaces until each of the subspaces has a single element.

After the reference database has been prepared, the present invention may be used to identify an unknown work. Such a process will now be shown and described.

Identification of an Unknown Work

Referring now to FIG. 7B, a flowchart of a method for identifying an unknown work is shown. In act 700, the present invention receives a sampled work. In act 702, the present invention determines a set of initial candidates. Finally, in act 704, the present invention determines the best candidate. Each act will now be described in more detail.

Receiving a Sampled Work

Beginning with act 700, a sampled work is provided to the present invention. It is contemplated that the work will be provided to the present invention as a digital audio stream. It should be understood that if the audio is in analog form, it may be digitized in any manner standard in the art.

Indexed Lookup.

In act 702, the present invention determines the initial candidates. In a preferred embodiment, the present invention uses the index created above to perform an indexed candidate search.

An index created in accordance with the present invention may be used to do the N nearest neighbor search required to find the initial candidates.

Candidate Search.

Once a set of N nearest neighbors is determined, the closest candidate may then be determined in act 704. In an exemplary non-limiting embodiment of the present invention, a brute-force search method may be used to determine which candidate is the closest to the target signature. In another preferred embodiment, the present invention may compare the distance of this best candidate to a predetermined threshold to determine whether there is a match.

There are a number of techniques that may be applied to the candidate search stage which make it much faster. In one aspect, these techniques may be used in a straightforward brute-force search that did not make use of any of the steps previously described above. That is, one could do a brute-force search directly on the reference signature database without going through the index search of step 702, for example. Since there is some overhead in doing step 702, direct brute-force search may be faster for some applications, especially those that need only a small reference database, e.g., generating a playlist for a radio station that plays music from a small set of possibilities.

Speedups of Brute-force Search.

Any reference signature that is close to the real-time signature has to be reasonably close to it for every segment in the signature. Therefore, in one aspect, several intermediate thresholds are tested as the distance is computed and the computation is exited if any of these thresholds are exceeded. In a further aspect, each single segment-to-segment distance is computed as the sum of the squared differences of the MFCCs for the two corresponding segments. Given the current computation of the MFCCs, average segment-to-segment distances for matches are about approximately 2.0. In an exemplary non-limiting embodiment of the invention, we exit the computation and set the distance to infinity if any single segment-to-segment distance is greater than 20. In further aspects, the computation is exited if any two segment-to-segment distances are greater than 15, or if any four segment-to-segment distances are greater than 10. It should be clear to anyone skilled in the art that other thresholds for other combinations of intermediate distances could easily be implemented and set using empirical tests.

Since any match will also be close to a match at a small time-offset, we may initially compute the distances at multiples of the hop size. If any of these distances are below a certain threshold, we compute the distances for hops near it. In an exemplary non-limiting embodiment of the invention, we compute distances for every third hop. If the distance is below 8.0, we compute the distances for the neighboring hops. It should be clear to anyone skilled in the art that other thresholds for other hop-skippings could easily be implemented and set using simple empirical tests.

While embodiments and applications of this invention have been shown and described, it would be apparent to those skilled in the art that many more modifications than mentioned above are possible without departing from the inventive concepts herein. For example, the teachings of the present disclosure may be used to identify a variety of sampled works, including, but not limited to, images, video and general time-based media. The invention, therefore, is not to be restricted except in the spirit of the appended claims.

What is claimed is:

1. A method for determining an identity of a received work comprising:

receiving audio data for an unknown work;

dividing said audio data into a plurality of segments;

generating a plurality of signatures of said unknown work wherein each signature is generated from one of said plurality of segments;

generating a plurality of reduced dimension signatures of said unknown works wherein each of said plurality of reduced dimension signatures is generated from one of at least a portion of said plurality of signatures;

comparing said plurality of reduced dimension signatures to at least one reduced dimension signature for each one of a portion of a plurality of known works having a record stored in a works database wherein each record includes an identification of said work and at least one signature of said work;

determining a list of candidates from said plurality of known works responsive to said comparisons;

comparing said plurality of signatures to at least a portion of said plurality of signatures of each of said plurality of known works in said list of candidates;

determining one of said plurality of known works in said list of candidates that matches said unknown work from said comparison; and

identifying said unknown work as said one of said plurality of known works that matches said unknown work.

2. The method of claim 1, wherein said step of generating said plurality of signatures comprises:

calculating a plurality of mel frequency cepstral coefficients (MFCCs) for each of a plurality segments to generate each of said plurality of signatures.

3. The method of claim 1, wherein said step of generating said plurality of signatures comprises:

calculating each of said plurality of signatures using a plurality of acoustical features from one of said plurality of segments selected from a group consisting of loudness, pitch, brightness, bandwidth, spectrum and mel frequency cepstral coefficients (MFCCs).

4. The method of claim 1 further comprising:

indexing said records of said plurality of known references in said database.

5. The method of claim 4 wherein said indexing is performed by using an indexing strategy chosen from the group consisting of: the k-d tree, the SS-tree, the R-tree, and the SR-tree.

6. The method of claim 1, wherein each of said plurality of segments has a segment size of approximately 0.5 to 3 seconds.

7. The method of claim 1, wherein each of said plurality of segments has a segment size of approximately 1 second.

8. The method of claim 1, wherein a hop size between consecutive ones of said plurality of segments is less than 50% of a segment size of each of said plurality of segments.

9. The method of claim 1, wherein a hop size between consecutive ones of said plurality of segments is approximately 0.1 seconds.

10. The method of claim 1, wherein said step of generating said plurality of signatures comprises:

calculating an average for each of a plurality of acoustical features selected from a group consisting of: loudness, pitch, brightness, bandwidth, spectral features, and mod frequency cepstral coefficients.

11. The method of claim 1, wherein said step of generating said plurality of reduced dimension signatures comprises:

projecting features of each said portion of said plurality of signatures in a Karhunen-Loeve basis.

13

12. The method of claim 1 wherein said step of generating said plurality of reduced dimension signatures comprises:

selecting predetermined one of a plurality of mid frequency cepstral coefficients from each one of said portion of said plurality of signatures.

13. The method of claim 1, wherein said list of candidates includes a predetermined number of said plurality of known works.

14. An apparatus for determining an identity of a received work comprising:

means for receiving audio data for an unknown work; means for dividing said audio data into a plurality of segments;

means for generating a plurality of signatures of said unknown work wherein each signature is generated from one of said plurality of segments;

means for generating a plurality of reduced dimension signatures of said unknown works wherein each of said plurality of reduced dimension signatures is generated from one of at least a portion of said plurality of signatures;

means for comparing said plurality of reduced dimension signatures to at least one reduced dimension signature for each one of a portion of a plurality of known works having a record stored in a works database wherein each record includes a identification of said work and at least one signature of said work;

means for determining a list of candidates from said plurality of known works responsive to said comparisons;

means for comparing said plurality of signatures to at least a portion of said plurality of signatures of each of said plurality of known works in said list of candidates; and

means for determining one of said plurality of known works in said list of candidates that matches said unknown work from said comparison.

15. The apparatus of claim 14, wherein said means for generating said plurality of signatures comprises:

means for calculating a plurality of mel frequency cepstral coefficients (MFCCs) for each of a plurality segments to generate each of said plurality of signatures.

16. The apparatus of claim 14, wherein said means for generating said plurality of signatures comprises:

14

means for calculating each of said plurality of signatures using a plurality of acoustical features from one of said plurality of segments selected from a group consisting of loudness, pitch, brightness, bandwidth, spectrum and mel frequency cepstral coefficients (MFCCs).

17. The apparatus of claim 14 further comprising: means for indexing said records of said plurality of known references in said database.

18. The apparatus of claim 17 wherein said indexing is performed by using an indexing strategy chosen from the group consisting of: the k-d tree, the SS-tree, the R-tree, and the SR-tree.

19. The apparatus of claim 14, wherein each of said plurality of segments has a segment size of approximately 0.5 to 3 seconds.

20. The apparatus of claim 14, wherein each of said plurality of segments has a segment size of approximately 1 second.

21. The apparatus of claim 14, wherein a hop size between consecutive ones of said plurality of segments is less than 50% of a segment size of each of said plurality of segments.

22. The apparatus of claim 14, wherein a hop size between consecutive ones of said plurality of segments is approximately 0.1 seconds.

23. The apparatus of claim 14, wherein said means for generating said plurality of signatures comprises:

means for calculating an average for each of a plurality of acoustical features selected from a group consisting of: loudness, pitch, brightness, bandwidth, spectral features, and mod frequency cepstral coefficients.

24. The apparatus of claim 14, wherein said means for generating said plurality of reduced dimension signatures comprises:

means for projecting features of each said portion of said plurality of signatures in a Karhunen-Loeve basis.

25. The apparatus of claim 14 wherein said means for generating said plurality of reduced dimension signatures comprises:

means for selecting predetermined ones of a plurality of mid frequency cepstral coefficients from each one of said portion of said plurality of signatures.

26. The apparatus of claim 14, wherein said list of candidates includes a predetermined number of said plurality of known works.

* * * * *