

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6233625号
(P6233625)

(45) 発行日 平成29年11月22日(2017.11.22)

(24) 登録日 平成29年11月2日(2017.11.2)

(51) Int.Cl.		F I		
G 1 0 L 25/54	(2013.01)	G 1 0 L	25/54	
G 1 0 L 15/10	(2006.01)	G 1 0 L	15/10	5 0 0 Z
G 0 6 F 17/30	(2006.01)	G 0 6 F	17/30	3 5 0 C
		G 0 6 F	17/30	1 7 0 E

請求項の数 7 (全 20 頁)

(21) 出願番号	特願2013-37542 (P2013-37542)	(73) 特許権者	000002185
(22) 出願日	平成25年2月27日 (2013.2.27)		ソニー株式会社
(65) 公開番号	特開2014-115605 (P2014-115605A)		東京都港区港南1丁目7番1号
(43) 公開日	平成26年6月26日 (2014.6.26)	(74) 代理人	100121131
審査請求日	平成27年2月9日 (2015.2.9)		弁理士 西川 孝
(31) 優先権主張番号	特願2012-251809 (P2012-251809)	(74) 代理人	100082131
(32) 優先日	平成24年11月16日 (2012.11.16)		弁理士 稲本 義雄
(33) 優先権主張国	日本国(JP)	(72) 発明者	澁谷 崇
前置審査			東京都港区港南1丁目7番1号 ソニー株式会社社内
		(72) 発明者	安部 素嗣
			東京都港区港南1丁目7番1号 ソニー株式会社社内

最終頁に続く

(54) 【発明の名称】 音声処理装置および方法、並びにプログラム

(57) 【特許請求の範囲】

【請求項1】

同定対象となるコンテンツの入力信号のスペクトログラムを二次関数で近似した近似結果、および前記スペクトログラムのピーク周辺のスペクトルを二次関数で近似した近似結果に基づいて、距離関数を用いて各時間周波数領域における信号を平滑化フィルタによりフィルタリングすることで、正弦波らしさの時間平均量に基づく第1の音響特徴量を算出するとともに、前記各時間周波数領域における信号を一次微分フィルタによりフィルタリングすることで、正弦波らしさの時間変化量に基づく、前記第1の音響特徴量とは異なる第2の音響特徴量を算出する入力信号処理部と、

予め用意したコンテンツの参照信号に基づいて、前記第1の音響特徴量と前記第2の音響特徴量とを算出する参照信号処理部と、

前記入力信号の前記第1の音響特徴量および前記第2の音響特徴量と、前記参照信号の前記第1の音響特徴量および前記第2の音響特徴量とに基づいて、前記入力信号と前記参照信号の類似度を計算するマッチング処理部と

を備える音声処理装置。

【請求項2】

前記マッチング処理部は、前記入力信号と前記参照信号の前記第1の音響特徴量に基づいて、各時間周波数領域におけるコンテンツの信号らしさを示すマスクパターンを生成し、前記マスクパターン、前記第1の音響特徴量、および前記第2の音響特徴量に基づいて前記類似度を計算する

請求項 1 に記載の音声処理装置。

【請求項 3】

前記マッチング処理部は、前記入力信号の前記第 1 の音響特徴量と、前記参照信号の前記第 1 の音響特徴量との類似度をさらに算出し、前記マスクパターン、前記第 1 の音響特徴量の前記類似度、および前記第 2 の音響特徴量に基づいて、前記入力信号と前記参照信号の前記類似度を計算する

請求項 2 に記載の音声処理装置。

【請求項 4】

前記マッチング処理部は、前記第 1 の音響特徴量の前記類似度に対する前記入力信号の寄与率よりも、前記第 1 の音響特徴量の前記類似度に対する前記参照信号の寄与率をより大きくして、前記第 1 の音響特徴量の前記類似度を算出する

請求項 3 に記載の音声処理装置。

【請求項 5】

前記第 2 の音響特徴量は、前記入力信号または前記参照信号のスペクトログラムに基づいて算出され、時間軸および周波数軸において前記第 1 の音響特徴量と同じ粒度を有する

請求項 1 乃至請求項 4 の何れか一項に記載の音声処理装置。

【請求項 6】

同定対象となるコンテンツの入力信号のスペクトログラムを二次関数で近似した近似結果、および前記スペクトログラムのピーク周辺のスペクトルを二次関数で近似した近似結果に基づいて、距離関数を用いて各時間周波数領域における信号を平滑化フィルタによりフィルタリングすることで、正弦波らしさの時間平均量に基づく第 1 の音響特徴量を算出するとともに、前記各時間周波数領域における信号を一次微分フィルタによりフィルタリングすることで、正弦波らしさの時間変化量に基づく、前記第 1 の音響特徴量とは異なる第 2 の音響特徴量を算出し、

予め用意したコンテンツの参照信号に基づいて、前記第 1 の音響特徴量と前記第 2 の音響特徴量とを算出し、

前記入力信号の前記第 1 の音響特徴量および前記第 2 の音響特徴量と、前記参照信号の前記第 1 の音響特徴量および前記第 2 の音響特徴量とに基づいて、前記入力信号と前記参照信号の類似度を計算する

ステップを含む音声処理方法。

【請求項 7】

同定対象となるコンテンツの入力信号のスペクトログラムを二次関数で近似した近似結果、および前記スペクトログラムのピーク周辺のスペクトルを二次関数で近似した近似結果に基づいて、距離関数を用いて各時間周波数領域における信号を平滑化フィルタによりフィルタリングすることで、正弦波らしさの時間平均量に基づく第 1 の音響特徴量を算出するとともに、前記各時間周波数領域における信号を一次微分フィルタによりフィルタリングすることで、正弦波らしさの時間変化量に基づく、前記第 1 の音響特徴量とは異なる第 2 の音響特徴量を算出し、

予め用意したコンテンツの参照信号に基づいて、前記第 1 の音響特徴量と前記第 2 の音響特徴量とを算出し、

前記入力信号の前記第 1 の音響特徴量および前記第 2 の音響特徴量と、前記参照信号の前記第 1 の音響特徴量および前記第 2 の音響特徴量とに基づいて、前記入力信号と前記参照信号の類似度を計算する

ステップを含む処理をコンピュータに実行させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本技術は音声処理装置および方法、並びにプログラムに関し、特に、より高精度に任意のコンテンツの同定を行なうことができるようにした音声処理装置および方法、並びにプログラムに関する。

10

20

30

40

50

【背景技術】

【0002】

例えば、コンテンツを構成する音声の信号を参照信号として、任意の機器において参照信号に基づいて再生された音声を収音して得られた入力信号と、参照信号とに基づいて一致検索を行なえば、コンテンツを同定することができる。この場合、元の音源から出力された音声は、残響や雑音が混入されて収音されるので、入力信号に基づく音声は、コンテンツの音声に残響や雑音が重畳されたものとなる。

【0003】

このようなコンテンツ同定技術として、例えばCD (Compact Disc) 等のクリーンな音楽の信号を参照信号として、音楽以外の音が混入してしまった入力信号から、その背景音楽を同定する楽曲同定技術がある。

10

【0004】

楽曲同定技術では、クリーンな音楽の参照信号から抽出された音響特徴量と、入力信号から抽出された音響特徴量とのマッチング処理により、楽曲の同定が行なわれる。ここで、入力信号には雑音が混入していることが前提となっており、入力信号から得られる音響特徴量は雑音の影響を受けてしまう。

【0005】

そこで、例えばマッチング処理においては、マスクパターンが用いられる。このマスクパターンは音響特徴量を構成する要素のうち、信頼のおける要素を表現する情報である。マスクパターンが用いられるマッチング処理では、マスクパターンに基づいて、多次元の音響特徴量を構成する各要素が、信頼のおける要素と信頼のおけない要素に分けられて、信頼のおける要素のみが用いられてマッチングが行なわれる。

20

【0006】

このようにマスクパターンを利用する楽曲同定技術として、例えば時間周波数成分を持つ特徴行列に対して、所定の時間周波数領域をマスクする複数のマスクパターンを予め用意しておき、楽曲同定を行うものが提案されている(例えば、特許文献1参照)。

【0007】

この技術では、入力信号の特徴行列と、データベース内の楽曲の特徴行列、つまり参照信号の特徴行列に対して、予め用意した全てのマスクパターンを用いて算出した類似度のうちの最大値を入力信号と楽曲の類似度として楽曲同定が行なわれる。なお、この楽曲同定では、入力信号によらない固定の複数マスクパターンを記憶しておき、それらのマスクパターンが用いられてマッチング処理が行なわれる。

30

【先行技術文献】

【特許文献】

【0008】

【特許文献1】特開2009-276776号公報

【発明の概要】

【発明が解決しようとする課題】

【0009】

しかしながら、上述した技術では、コンテンツの同定は音楽の一致検索に特化されており、一般の任意のコンテンツ、例えば放送番組等のコンテンツについてコンテンツを同定することはできなかった。例えば、放送番組のコンテンツでは、音楽の流れていないシーンの音声信号を入力信号として検索したい場合もあるが、そのような場合には上述した技術によりコンテンツの同定を行なうことができない。

40

【0010】

また、上述した技術では音声の残響の影響が考慮されていないため、高精度なコンテンツ同定を実現することができない場合があった。すなわち、実使用環境では入力信号は残響の影響を受け、この残響は検索に悪影響を及ぼす。そのため、残響の強い環境では、コンテンツの一致検索の精度が低下してしまう。

【0011】

50

さらに、特許文献1に記載の技術では、固定のマスクパターンが使用されている。ところが、入力信号に乗っている混入雑音に関しては、どの時間にどのような特性の雑音に乗るかは予測不能であるから、入力信号に最適なマスクパターンを予め用意するのは困難である。したがって、予め用意したマスクパターンでは高精度にコンテンツを同定することができないことがある。

【0012】

本技術は、このような状況に鑑みてなされたものであり、より高精度に任意のコンテンツの同定を行なうことができるようにするものである。

【課題を解決するための手段】

【0013】

本技術の一側面の音声処理装置は、同定対象となるコンテンツの入力信号のスペクトログラムを二次関数で近似した近似結果、および前記スペクトログラムのピーク周辺のスペクトルを二次関数で近似した近似結果に基づいて、距離関数を用いて各時間周波数領域における信号を平滑化フィルタによりフィルタリングすることで、正弦波らしさの時間平均量に基づく第1の音響特徴量を算出するとともに、前記各時間周波数領域における信号を一次微分フィルタによりフィルタリングすることで、正弦波らしさの時間変化量に基づく、前記第1の音響特徴量とは異なる第2の音響特徴量を算出する入力信号処理部と、予め用意したコンテンツの参照信号に基づいて、前記第1の音響特徴量と前記第2の音響特徴量とを算出する参照信号処理部と、前記入力信号の前記第1の音響特徴量および前記第2の音響特徴量と、前記参照信号の前記第1の音響特徴量および前記第2の音響特徴量とに基づいて、前記入力信号と前記参照信号の類似度を計算するマッチング処理部とを備える。

【0014】

前記マッチング処理部には、前記入力信号と前記参照信号の前記第1の音響特徴量に基づいて、各時間周波数領域におけるコンテンツの信号らしさを示すマスクパターンを生成させ、前記マスクパターン、前記第1の音響特徴量、および前記第2の音響特徴量に基づいて前記類似度を計算させることができる。

【0015】

前記マッチング処理部には、前記入力信号の前記第1の音響特徴量と、前記参照信号の前記第1の音響特徴量との類似度をさらに算出させ、前記マスクパターン、前記第1の音響特徴量の前記類似度、および前記第2の音響特徴量に基づいて、前記入力信号と前記参照信号の前記類似度を計算させることができる。

【0016】

前記マッチング処理部には、前記第1の音響特徴量の前記類似度に対する前記入力信号の寄与率よりも、前記第1の音響特徴量の前記類似度に対する前記参照信号の寄与率をより大きくして、前記第1の音響特徴量の前記類似度を算出させることができる。

【0017】

前記第2の音響特徴量を、前記入力信号または前記参照信号のスペクトログラムに基づいて算出し、時間軸および周波数軸において前記第1の音響特徴量と同じ粒度を有するようにすることができる。

【0018】

本技術の一側面の音声処理方法またはプログラムは、同定対象となるコンテンツの入力信号のスペクトログラムを二次関数で近似した近似結果、および前記スペクトログラムのピーク周辺のスペクトルを二次関数で近似した近似結果に基づいて、距離関数を用いて各時間周波数領域における信号を平滑化フィルタによりフィルタリングすることで、正弦波らしさの時間平均量に基づく第1の音響特徴量を算出するとともに、前記各時間周波数領域における信号を一次微分フィルタによりフィルタリングすることで、正弦波らしさの時間変化量に基づく、前記第1の音響特徴量とは異なる第2の音響特徴量を算出し、予め用意したコンテンツの参照信号に基づいて、前記第1の音響特徴量と前記第2の音響特徴量とを算出し、前記入力信号の前記第1の音響特徴量および前記第2の音響特徴量と、前記

10

20

30

40

50

参照信号の前記第1の音響特徴量および前記第2の音響特徴量とに基づいて、前記入力信号と前記参照信号の類似度を計算するステップを含む。

【0019】

本技術の一側面においては、同定対象となるコンテンツの入力信号のスペクトログラムを二次関数で近似した近似結果、および前記スペクトログラムのピーク周辺のスペクトルを二次関数で近似した近似結果に基づいて、距離関数が用いられて各時間周波数領域における信号を平滑化フィルタによりフィルタリングすることで、正弦波らしさの時間平均量に基づく第1の音響特徴量が算出されるとともに、前記各時間周波数領域における信号を一次微分フィルタによりフィルタリングすることで、正弦波らしさの時間変化量に基づく、前記第1の音響特徴量とは異なる第2の音響特徴量が算出され、予め用意したコンテンツの参照信号に基づいて、前記第1の音響特徴量と前記第2の音響特徴量とが算出され、前記入力信号の前記第1の音響特徴量および前記第2の音響特徴量と、前記参照信号の前記第1の音響特徴量および前記第2の音響特徴量とに基づいて、前記入力信号と前記参照信号の類似度が計算される。

10

【発明の効果】

【0020】

本技術の一側面によれば、より高精度に任意のコンテンツの同定を行なうことができる。

【図面の簡単な説明】

【0021】

20

【図1】マスクパターンについて説明する図である。

【図2】音声処理装置の構成例を示す図である。

【図3】入力信号処理部の構成例を示す図である。

【図4】参照信号処理部の構成例を示す図である。

【図5】マッチング処理部の構成例を示す図である。

【図6】音響特徴量について説明する図である。

【図7】一致検索処理を説明するフローチャートである。

【図8】音響特徴量IA1の抽出処理を説明するフローチャートである。

【図9】音響特徴量IA2の抽出処理を説明するフローチャートである。

【図10】コンピュータの構成例を示す図である。

30

【発明を実施するための形態】

【0022】

以下、図面を参照して、本技術を適用した実施の形態について説明する。

【0023】

第1の実施の形態

本技術の特徴

本技術は、多機能携帯型電話機やタブレット型端末装置等の携帯端末装置の録音機能を使用して、ユーザが別の機器で視聴しているテレビ番組やラジオ番組、ストリーミング配信コンテンツ等の任意のコンテンツを同定するためのものである。

【0024】

40

処理対象となる音声をテレビジョン受像機やラジオ、パーソナルコンピュータ等の機器のスピーカから出力させ、その音声を携帯端末装置で録音させると、録音により得られる音声には、機器のスピーカから出力された音声は携帯端末装置までの空間を通過することによって、その空間での音声の残響も含まれてしまう。また、録音により得られる音声には、機器のスピーカから出力された音声以外の音声も混入する（以下、これを混入雑音と呼ぶこととする）。

【0025】

本技術では、これらの残響や混入雑音に頑健なコンテンツの一致検索を行うことが目的となる。より一般的には、所定の音源が空間を通過することによって残響や混入雑音が重畳された音源と、元の音源（ドライソース）との一致検索を行うことが本技術の目的とな

50

る。

【0026】

まず、本技術の特徴について説明する。本技術は、以下の5つの特徴を有している。

【0027】

特徴(1)

マスクパターンの生成には入力信号と参照信号について算出された、切り出された各々の時間周波数領域内での正弦波らしさの指標が用いられる

特徴(2)

正弦波らしさの指標は、微小時間内でのスペクトル形状の安定性によって定量化される

特徴(3)

上記の正弦波らしさは残響に頑健な指標とされる

特徴(4)

マスクパターンは入力信号の情報だけでなく、参照信号の情報も用いて生成される

特徴(5)

入力信号と参照信号との類似度算出の際に、それらを等価に扱うのではなく入力信号よりも参照信号に重きをおいた類似度計算が行なわれる

【0028】

例えば、本技術では図1に示す入力信号のスペクトログラムと参照信号のスペクトログラムが得られる。なお、図1において、縦軸および横軸は、それぞれ周波数および時間を示している。

【0029】

図1では、図中、右側には参照信号のスペクトログラムが示されており、図中、左側には入力信号のスペクトログラムが示されている。

【0030】

入力信号のスペクトログラム、つまり時間周波数領域において、実線で表される成分は参照信号にも含まれている音声成分を示しており、点線で表される成分は参照信号には含まれていない、混入雑音の成分を示している。

【0031】

本技術では、マスクパターンを生成することで、図中、斜線部分の領域である、信頼のおける時間周波数領域が特定され、その領域のみが用いられて入力信号と参照信号のマッチング処理が行なわれる。

【0032】

このような本技術によれば、以下の効果を得ることができる。

【0033】

効果(1)

音楽の流れているシーンだけでなく、音楽の流れていないシーンを使って、そのコンテンツの同定を行うことができる

効果(2)

残響のある空間においても、視聴番組等のコンテンツの同定が可能である

効果(3)

入力信号に元の参照信号に含まれていた音声以外の音声(混入雑音)が含まれていても、視聴番組等のコンテンツの同定が可能である

【0034】

音声処理装置の構成例

次に、本技術を適用した具体的な実施の形態について説明する。

【0035】

図2は、本技術を適用した音声処理装置の一実施の形態の構成例を示す図である。

【0036】

音声処理装置11は、入力信号処理部21、参照信号処理部22、およびマッチング処理部23から構成される。

10

20

30

40

50

【 0 0 3 7 】

音声処理装置 1 1 には、予め用意されたコンテンツの音声の参照信号と、所定の機器で再生された参照信号に基づく音声を、他の機器で録音（収音）することで得られた、同定対象のコンテンツの音声の入力信号とが入力される。例えば、入力信号は音声処理装置 1 1 による録音で得られた音声信号であってもよい。

【 0 0 3 8 】

なお、参照信号として、例えば複数のコンテンツの音声信号が入力される。また、音声処理装置 1 1 には、参照信号のコンテンツ属性データも入力される。ここで、コンテンツ属性データとは、コンテンツ名（番組名）や放送日時、出演者など、コンテンツに関するデータである。

10

【 0 0 3 9 】

入力信号処理部 2 1 は、供給された入力信号を解析して 2 種類の音響特徴量 I A 1 および音響特徴量 I A 2 を生成し、マッチング処理部 2 3 に供給する。

【 0 0 4 0 】

参照信号処理部 2 2 は、供給された、コンテンツの元音源である参照信号を解析して、音響特徴量 I A 1 および音響特徴量 I A 2 に対応する 2 種類の音響特徴量 R A 1 および音響特徴量 R A 2 を生成し、マッチング処理部 2 3 に供給する。

【 0 0 4 1 】

ここで、音響特徴量 I A 1 と音響特徴量 R A 1 は同じ特徴量（同じ種別の特徴量）であり、音響特徴量 I A 2 と音響特徴量 R A 2 も同じ特徴量である。以下、音響特徴量 I A 1 と音響特徴量 R A 1 を特に区別する必要のない場合、単に音響特徴量 A 1 とも称し、音響特徴量 I A 2 と音響特徴量 R A 2 を特に区別する必要のない場合、単に音響特徴量 A 2 とも称することとする。

20

【 0 0 4 2 】

マッチング処理部 2 3 は、入力信号処理部 2 1 からの音響特徴量 I A 1 および音響特徴量 I A 2 と、参照信号処理部 2 2 からの音響特徴量 R A 1 および音響特徴量 R A 2 とに基づいて、入力信号と参照信号とのマッチング処理を行い、コンテンツを同定する。また、マッチング処理部 2 3 は、供給されたコンテンツ属性データのうち、マッチング処理により同定されたコンテンツのコンテンツ属性データと、マッチング処理の結果とを出力する。

30

【 0 0 4 3 】

入力信号処理部の構成例

また、図 2 の入力信号処理部 2 1 は、より詳細には図 3 に示すように構成される。図 3 に示す入力信号処理部 2 1 は、入力信号切り出し部 5 1、時間周波数変換部 5 2、音響特徴量抽出部 5 3、および音響特徴量抽出部 5 4 から構成される。

【 0 0 4 4 】

入力信号切り出し部 5 1 は、供給された入力信号の所定時間長の区間を切り出して、時間周波数変換部 5 2 に供給する。時間周波数変換部 5 2 は、入力信号切り出し部 5 1 から供給された入力信号に対する時間周波数変換を行なって、入力信号を対数振幅スペクトログラムに変換し、音響特徴量抽出部 5 3 および音響特徴量抽出部 5 4 に供給する。

40

【 0 0 4 5 】

音響特徴量抽出部 5 3 は、時間周波数変換部 5 2 から供給された対数振幅スペクトログラムに基づいて音響特徴量 I A 1 を算出し、マッチング処理部 2 3 に供給する。音響特徴量抽出部 5 4 は、時間周波数変換部 5 2 から供給された対数振幅スペクトログラムに基づいて音響特徴量 I A 2 を算出し、マッチング処理部 2 3 に供給する。

【 0 0 4 6 】

ここで、音響特徴量 I A 1 と音響特徴量 I A 2 について説明する。

【 0 0 4 7 】

例えば、音響特徴量 I A 1 と音響特徴量 I A 2 は全て行列で表現され、2 つの軸は時間成分と周波数成分である。それぞれの行列は、以下のような特徴を有している。

50

【 0 0 4 8 】

すなわち、音響特徴量 I A 1 は、各々の時間周波数領域において、入力信号のその領域内での正弦波らしさを表現する特徴行列である。

【 0 0 4 9 】

また、音響特徴量 I A 2 は入力信号と参照信号のマッチングに用いる特徴量であり、信号の個性を表現する特徴行列である。但し、音響特徴量 I A 2 の時間軸および周波数軸の粒度は音響特徴量 I A 1 の時間軸および周波数軸の粒度と同じとされる。

【 0 0 5 0 】

さらに、入力信号処理部 2 1 が音響特徴量 I A 1 と音響特徴量 I A 2 を算出する具体的な処理について説明する。

10

【 0 0 5 1 】

まず、入力信号切り出し部 5 1 は、常に入力され続けている入力信号からある一定時間（例えば、5 秒間）の信号を切り出して、時間周波数変換部 5 2 に供給する。時間周波数変換部 5 2 は、切り出された入力信号を対数振幅スペクトログラム（以下、単にスペクトログラムとも称する）に変換する。

【 0 0 5 2 】

また、音響特徴量抽出部 5 3 は、スペクトログラムを、分割された時間周波数領域内のスペクトログラムの正弦波らしさを数値化した中間特徴量に変換する。

【 0 0 5 3 】

すなわち、正弦波らしさの数値化にはスペクトログラムの微小時間内の安定性が利用される。楽器音や人の音声は雑音と違い、微小時間（例えば、0.020 秒）内では周波数がほぼ一定の正弦波と見なすことができ、スペクトログラムの形状がほぼ一定となる。

20

【 0 0 5 4 】

音響特徴量抽出部 5 3 は、この性質を利用して、微小時間内でのスペクトログラムの安定性を周波数帯ごとに数値化し、これを正弦波らしさの指標とする。具体的には、音響特徴量抽出部 5 3 は、スペクトログラムの時間フレームごとにピーク検出処理を行ない、ピーク周辺の時間周波数領域について対数振幅スペクトログラムを次式（1）に示す双二次関数 $g(k, n)$ で近似する。

【 0 0 5 5 】

【 数 1 】

$$g(k, n) = \bar{a}k^2 + \bar{b}k + \bar{c} \quad \dots (1)$$

30

【 0 0 5 6 】

なお、式（1）において、 k はスペクトログラムの周波数 bin 番号を示しており、 n はスペクトログラムの時間フレーム番号を示している。また、対数振幅スペクトログラムの近似は、最小二乗法等の最適化手法によって行なわれる。

【 0 0 5 7 】

次に、音響特徴量抽出部 5 3 は検出されたピーク周辺の時間周波数領域の各時間フレームの対数振幅スペクトルを次式（2）に示す二次関数 $f_n(k)$ で近似する。

【 0 0 5 8 】

【 数 2 】

$$f_n(k) = a_n k^2 + b_n k + c_n \quad \dots (2)$$

40

【 0 0 5 9 】

この近似も最小二乗法等の最適化手法によって行なわれる。

【 0 0 6 0 】

さらに、音響特徴量抽出部 5 3 は、双二次関数 $g(k, n)$ および二次関数 $f_n(k)$ の 2 種類の関数への近似によって得られた係数を用いて、正弦波らしさを次式（3）により算出する。

【 0 0 6 1 】

50

【数3】

$$\eta(n, k) = 1 - \alpha \sqrt{\sum \{D_1(a_n, \bar{a}) + D_2(b_n, \bar{b})\}} \quad \dots (3)$$

【0062】

なお、式(3)において、 α は正の値を持つパラメータであり、 $D(x, y)$ 、つまり $D_1(x, y)$ および $D_2(x, y)$ は距離関数を示している。

【0063】

また、正弦波を時間周波数変換した場合には、二次関数の二次の係数の理論値 a が存在する。この理論値と算出された二次関数の二次の係数の近さを加味して、次式(4)により正弦波らしさが算出されるようにしてもよい。

【0064】

【数4】

$$\eta(n, k) = 1 - \alpha \sqrt{\sum \{D_1(a_n, \bar{a}) + D_2(b_n, \bar{b}) + D_3(a_n, a)\}} \quad \dots (4)$$

【0065】

なお、 $\eta(n, k)$ は各ピークにおける正弦波らしさを意味するので、 $\eta(n, k) < 0$ となった場合は $\eta(n, k) = 0$ とされる。これによって、 $\eta(n, k)$ は0から1の値を取ることとなる。

【0066】

また、ピークに該当しない周波数binについては $\eta(n, k) = 0$ とされ、当該時間フレームについて、各周波数binの正弦波らしさの情報を持ったベクトルが得られる。この正弦波らしさは残響にロバストな特徴量であるため、最終的に残響に頑健な検索を行うことが可能である。

【0067】

以上のようにして得られたベクトルを、時間フレームをずらしながら算出していき、得られたベクトルを時系列に並べて、時間軸方向にダウンサンプリングを行ったものが音響特徴量 $IA1$ である。ダウンサンプリングには平滑化フィルタ(ローパスフィルタ)が用いられる。フィルタリングされた値は各周波数における正弦波らしさの時間平均値を意味する。

【0068】

なお、得られた音響特徴量 $IA1$ の各要素について、量子化や対数関数、指数関数、シグモイド関数等の非線形処理が施されるようにしてもよい。

【0069】

また、音響特徴量抽出部54では、スペクトログラムが音響特徴量 $IA2$ に変換される。

【0070】

例えば、音響特徴量 $IA1$ と同様に算出された正弦波らしさの行列に対して時間軸方向に一次微分フィルタを施し、ダウンサンプリングしたものが音響特徴量 $IA2$ とされる。一次微分フィルタによってフィルタリングされた値は各周波数における正弦波らしさの時間変化量を意味する。

【0071】

なお、得られた音響特徴量 $IA2$ の各要素についても、量子化や対数関数、指数関数、シグモイド関数等の非線形処理が施されるようにしてもよい。さらに、音響特徴量 $IA2$ は信号の個性を表現するものであればよく、例えば、一定の時間区間内のスペクトルの時間平均を正規化したものなどを用いるようにしてもよい。

【0072】

参照信号処理部の構成例

また、図2の参照信号処理部22は、より詳細には図4に示すように構成される。図4に示す参照信号処理部22は、参照信号切り出し部81、時間周波数変換部82、音響特徴量抽出部83、および音響特徴量抽出部84から構成される。

10

20

30

40

50

【 0 0 7 3 】

参照信号切り出し部 8 1 は、供給された参照信号の所定時間長の区間を切り出して、時間周波数変換部 8 2 に供給する。時間周波数変換部 8 2 は、参照信号切り出し部 8 1 から供給された参照信号に対する時間周波数変換を行なって、参照信号を対数振幅スペクトログラムに変換し、音響特徴量抽出部 8 3 および音響特徴量抽出部 8 4 に供給する。

【 0 0 7 4 】

音響特徴量抽出部 8 3 は、時間周波数変換部 8 2 から供給された対数振幅スペクトログラムに基づいて音響特徴量 R A 1 を算出し、マッチング処理部 2 3 に供給する。音響特徴量抽出部 8 4 は、時間周波数変換部 8 2 から供給された対数振幅スペクトログラムに基づいて音響特徴量 R A 2 を算出し、マッチング処理部 2 3 に供給する。

10

【 0 0 7 5 】

なお、音響特徴量抽出部 8 3 と音響特徴量抽出部 8 4 は、音響特徴量抽出部 5 3 と音響特徴量抽出部 5 4 に対応しており、音響特徴量 I A 1 および音響特徴量 I A 2 と同じ時間軸・周波数軸粒度の音響特徴量 R A 1 および音響特徴量 R A 2 を出力する。

【 0 0 7 6 】

また、参照信号から抽出された音響特徴量 R A 1 および音響特徴量 R A 2 は、直接マッチング処理部 2 3 に供給されてもよいし、記憶装置に供給されてデータベースとして保存されてもよい。但し、音響特徴量 R A 1 および音響特徴量 R A 2 が記憶装置に供給される場合には、音響特徴量 R A 1 および音響特徴量 R A 2 は、参照信号のメタデータ（番組名、放送日時、出演者など）、つまりコンテンツ属性データとセットで保存される必要がある。

20

【 0 0 7 7 】

マッチング処理部の構成例

さらに、図 2 のマッチング処理部 2 3 は、より詳細には図 5 に示すように構成される。図 5 に示すマッチング処理部 2 3 は、マスクパターン生成部 1 1 1、類似度計算部 1 1 2、および比較統合部 1 1 3 から構成される。

【 0 0 7 8 】

マスクパターン生成部 1 1 1 は、音響特徴量抽出部 5 3 からの音響特徴量 I A 1 と、音響特徴量抽出部 8 3 からの音響特徴量 R A 1 とに基づいてマスクパターンを生成し、マスクパターンおよび音響特徴量 A 1 間の類似度を類似度計算部 1 1 2 に供給する。マスクパターンは、各時間周波数領域におけるコンテンツの信号らしさの信頼度、つまり信頼における時間周波数領域を示している。

30

【 0 0 7 9 】

類似度計算部 1 1 2 は、音響特徴量抽出部 5 4 からの音響特徴量 I A 2、音響特徴量抽出部 8 4 からの音響特徴量 R A 2、およびマスクパターン生成部 1 1 1 からのマスクパターンと類似度に基づいて、入力信号の参照信号との類似度を算出する。また、類似度計算部 1 1 2 は、算出した類似度とともに、供給されたコンテンツ属性データを比較統合部 1 1 3 に供給する。

【 0 0 8 0 】

比較統合部 1 1 3 は、類似度計算部 1 1 2 から供給された類似度に基づいて、参照信号のコンテンツと、入力信号に含まれるコンテンツが同一であるかを判定し、その判定結果とともにコンテンツ属性データを出力する。

40

【 0 0 8 1 】

マッチング処理部 2 3 では、参照信号と入力信号の類似度が計算される。例えば図 6 に示すように、一定時間（例えば、5 秒間）の入力信号に参照信号の断片が含まれている場合、一般に入力信号の音響特徴量 I A 1 および音響特徴量 I A 2 の行列は、参照信号の音響特徴量 R A 1 および音響特徴量 R A 2 に比べ時間方向の成分数が少ない。

【 0 0 8 2 】

そのため、参照信号の音響特徴量 R A 1 および音響特徴量 R A 2 の行列から入力信号の音響特徴量 I A 1 および音響特徴量 I A 2 の時間方向の長さと同じ長さの部分行列が切り

50

出されて類似度計算が行なわれる。この部分行列の切り出しについては、切り取り可能な部分行列すべてについて切り出しが行なわれる。切り出し処理はマスクパターン生成部 1 1 1 および類似度計算部 1 1 2 で行われる。

【 0 0 8 3 】

図 6 では、縦方向および横方向は、それぞれ周波数および時間を示している。また、矢印 Q 1 1 乃至矢印 Q 1 4 により示される長方形は、それぞれ参照信号の音響特徴量 R A 1、参照信号の音響特徴量 R A 2、入力信号の音響特徴量 I A 1、および入力信号の音響特徴量 I A 2 を表している。

【 0 0 8 4 】

この例では、参照信号から抽出された音響特徴量 R A 1 および音響特徴量 R A 2 は、入力信号から抽出された音響特徴量 I A 1 および音響特徴量 I A 2 よりも、10 図中、横方向、つまり時間方向に長くなっており、時間方向の成分が多いことが分かる。

【 0 0 8 5 】

そのため、音響特徴量 R A 1 や音響特徴量 R A 2 の一部分が切り出されて部分行列とされ、類似度計算に用いられる。

【 0 0 8 6 】

次に、マッチング処理部 2 3 で行なわれる具体的な処理について説明する。

【 0 0 8 7 】

マスクパターン生成部 1 1 1 は、入力信号の音響特徴量 I A 1 と参照信号の音響特徴量 R A 1 からマスクパターンを生成し、さらに音響特徴量 A 1 同士の類似度を算出する。20 マスクパターンは、音響特徴量 A 1 と同様に時間軸と周波数軸を持った二次元の行列で表現される。

【 0 0 8 8 】

例えば、入力信号の音響特徴量 I A 1 と参照信号の音響特徴量 R A 1 から、正弦波が存在しない時間周波数領域をマスクする行列がマスクパターンとして生成される。より具体的には、例えば次式 (5) を計算することによりマスクパターンが生成される。

【 0 0 8 9 】

【 数 5 】

$$W_{f(t+u)} = S_{fu}^{(1)} A_{f(t+u)}^{(1)} \dots (5)$$

30

【 0 0 9 0 】

なお、式 (5) において、 $W_{f(t+u)}$ はマスクパターンの行列要素を示しており、 $S_{fu}^{(1)}$ は入力信号の音響特徴量 I A 1 の行列要素を示しており、 $A_{f(t+u)}^{(1)}$ は参照信号の音響特徴量 R A 1 の部分行列の要素を示している。

【 0 0 9 1 】

また、 f は各行列の周波数成分を示しており、 u は各行列の時間成分を示しており、 t は部分行列の時間オフセットを示している。

【 0 0 9 2 】

このようにして算出されたマスクパターンは、後段の類似度計算部 1 1 2 において、時間周波数領域ごとの重みとして用いられる。40 つまり、マスクパターンの行列要素 $W_{f(t+u)}$ の値の大きい時間周波数領域に重きを置いた類似度が算出される。

【 0 0 9 3 】

音響特徴量 A 1 同士の類似度は、2 つの特徴量の近さを定量化した非負の指標で、例えば次式 (6) によって算出される。

【 0 0 9 4 】

【 数 6 】

$$R^{(1)}(t) = \frac{\sum S_{fu}^{(1)} A_{f(t+u)}^{(1)}}{\left(\sum S_{fu}^{(1)p}\right)^{1/p} \cdot \left(\sum A_{f(t+u)}^{(1)q}\right)^{1/q}} \dots (6)$$

50

【 0 0 9 5 】

なお、式(6)において、 $R^{(1)}(t)$ は $S^{(1)}_{f_u}$ と $A^{(1)}_{f(t+u)}$ の類似度を示している。また、 p および q は、入力信号の音響特徴量 $I A 1$ と参照信号の音響特徴量 $R A 1$ の類似度に対する寄与率を調整するパラメータである。すなわち、 p および q は $1/p + 1/q = 1$ を満たす1以上の値を持つ重み係数である。

【 0 0 9 6 】

例えば、 p を q より大きくすることで参照信号に含まれている音声に重きを置いた類似度計算が行なわれ、入力信号に参照信号とは関係のない混入雑音が含まれていても、その影響を軽減したマッチングを行うことができる。なお、音響特徴量同士の類似度として、上記の類似度の他に、二乗誤差や絶対誤差など2つの行列の相違に基づいて算出される値が用いられてもよい。

10

【 0 0 9 7 】

また、類似度計算部112は、入力信号の音響特徴量 $I A 2$ と参照信号の音響特徴量 $R A 2$ 、マスクパターン、および音響特徴量 $A 1$ 同士の類似度を用いて、最終的な類似度計算を行なう。

【 0 0 9 8 】

類似度計算部112により算出される類似度は、時間周波数領域における正弦波らしさの情報を持つマスクパターンを各時間周波数領域における信頼度と見なして重み付けして定量化した、入力信号の音響特徴量 $I A 2$ と参照信号の音響特徴量 $R A 2$ の近さの指標である。さらに、音響特徴量 $A 1$ 同士の類似度も加味されて、例えば次式(7)の計算により類似度 $R(t)$ が算出される。

20

【 0 0 9 9 】

【数7】

$$R(t) = \frac{\sum W_{f(t+u)} \exp\left(-\beta \left(S_{f_u}^{(2)} - A_{f(t+u)}^{(2)}\right)^2\right)}{\sum W_{f(t+u)}} R^{(1)}(t) \quad \dots (7)$$

【 0 1 0 0 】

なお、式(7)において、 $A^{(2)}_{f(t+u)}$ は参照信号の音響特徴量 $R A 2$ の部分行列を示しており、 $S^{(2)}_{f_u}$ は入力信号の音響特徴量 $I A 2$ の行列を示している。また、 β は正の値を持つパラメータである。

30

【 0 1 0 1 】

さらに、類似度計算には、上記の式(7)の計算以外に、二乗誤差や絶対誤差などの2つの行列(入力信号の音響特徴量 $I A 2$ と参照信号の音響特徴量 $R A 2$)の相違に基づいて計算される値が用いられてもよい。

【 0 1 0 2 】

比較統合部113では、類似度計算部112で算出された類似度から参照信号のコンテンツと、入力信号に含まれるコンテンツが同一のものであるかが判定される。

【 0 1 0 3 】

同一コンテンツであるかの判定方法は、複数の参照信号について得られた類似度のうち、所定の閾値を超えたもので最大の類似度を持つ参照信号が入力信号に含まれているコンテンツであると判定する方法とされる。なお、全ての参照信号の類似度が閾値を超えなかった場合は、参照信号の中に該当するコンテンツはないという判定がなされる。

40

【 0 1 0 4 】

また、ここで用いられる閾値は、常に固定の値とされてもよいし、入力信号と複数の参照信号から得られた複数の類似度から統計的に設定されるようにしてもよい。

【 0 1 0 5 】

一致検索処理の説明

ところで、音声処理装置11に入力信号と参照信号が供給され、コンテンツ同定が指示

50

されると、音声処理装置 11 は一致検索処理を行なってコンテンツの同定を行なう。以下、図 7 のフローチャートを参照して、音声処理装置 11 による一致検索処理について説明する。

【0106】

ステップ S 11 において、入力信号切り出し部 51 は、供給された入力信号を切り出して時間周波数変換部 52 に供給する。例えば、一定時間分の入力信号が切り出される。

【0107】

ステップ S 12 において、時間周波数変換部 52 は、入力信号切り出し部 51 から供給された入力信号に対して時間周波数変換を行い、入力信号を対数振幅スペクトrogramに変換し、音響特徴量抽出部 53 および音響特徴量抽出部 54 に供給する。

10

【0108】

ステップ S 13 において、音響特徴量抽出部 53 は音響特徴量 I A 1 の抽出処理を行なって、入力信号の音響特徴量 I A 1 を算出し、マッチング処理部 23 のマスクパターン生成部 111 に供給する。

【0109】

ここで、図 8 のフローチャートを参照して、ステップ S 13 の処理に対応する、音響特徴量抽出部 53 による音響特徴量 I A 1 の抽出処理について説明する。

【0110】

ステップ S 51 において、音響特徴量抽出部 53 は、時間周波数変換部 52 から供給された対数振幅スペクトrogramについて、時間フレームを選択する。

20

【0111】

ステップ S 52 において、音響特徴量抽出部 53 は、対数振幅スペクトrogramの選択した時間フレームについてピーク検出を行なう。

【0112】

ステップ S 53 において、音響特徴量抽出部 53 は、検出されたピーク周辺の時間周波数領域の対数振幅スペクトルを 2 種類の二次関数で近似する。例えば、対数振幅スペクトルが式 (1) および式 (2) に示した関数で近似される。

【0113】

ステップ S 54 において、音響特徴量抽出部 53 は、近似された二次関数の係数から正弦波らしさの指標へ変換し、保存する。例えば式 (3) の (n, k) が正弦波らしさの指標として算出される。

30

【0114】

ステップ S 55 において、音響特徴量抽出部 53 は入力信号の全時間フレームを処理したか否かを判定する。ステップ S 55 においてまだ全ての時間フレームを処理していないと判定された場合、処理はステップ S 51 に戻り、上述した処理が繰り返される。

【0115】

これに対して、ステップ S 55 において、全時間フレームを処理したと判定された場合、ステップ S 56 において、音響特徴量抽出部 53 は、保存した正弦波らしさの指標のベクトルを時系列順に並べて行列化する。

【0116】

ステップ S 57 において、音響特徴量抽出部 53 は、行列化された正弦波らしさの指標、つまり正弦波らしさの行列に対して時間軸方向にフィルタリングを施し、正弦波らしさの時間平均量を算出する。例えば、平滑化フィルタによりフィルタリングが行なわれる。

40

【0117】

ステップ S 58 において、音響特徴量抽出部 53 は、フィルタリングにより得られた正弦波らしさの時間平均量に対して時間軸方向にリサンプリングを行い、音響特徴量 I A 1 とする。音響特徴量抽出部 53 が、このようにして入力信号から抽出した音響特徴量 I A 1 をマスクパターン生成部 111 に供給すると、音響特徴量 I A 1 の抽出処理は終了し、その後、処理は図 7 のステップ S 14 へと進む。

【0118】

50

ステップS 14において、音響特徴量抽出部54は音響特徴量IA2の抽出処理を行なって、入力信号の音響特徴量IA2を算出し、マッチング処理部23の類似度計算部112に供給する。

【0119】

ここで、図9のフローチャートを参照して、ステップS 14の処理に対応する、音響特徴量抽出部54による音響特徴量IA2の抽出処理について説明する。なお、ステップS 91乃至ステップS 96の処理は、図8のステップS 51乃至ステップS 56の処理と同様であるので、その説明は省略する。

【0120】

ステップS 96の処理が行なわれると、正弦波らしさの行列が得られるので、ステップS 97において、音響特徴量抽出部54は、正弦波らしさの行列に対して時間軸方向にフィルタリングを施し、正弦波らしさの時間変化量を算出する。例えば、一次微分フィルタによりフィルタリングが行なわれる。

【0121】

ステップS 98において、音響特徴量抽出部54は、フィルタリングにより得られた正弦波らしさの時間変化量に対して時間軸方向にリサンプリングを行い、音響特徴量IA2とする。音響特徴量抽出部54が、このようにして入力信号から抽出した音響特徴量IA2を類似度計算部112に供給すると、音響特徴量IA2の抽出処理は終了し、その後、処理は図7のステップS 15へと進む。

【0122】

図7のフローチャートの説明に戻り、ステップS 15において、参照信号切り出し部81は、供給された参照信号を切り出して時間周波数変換部82に供給する。

【0123】

ステップS 16において、時間周波数変換部82は、参照信号切り出し部81から供給された参照信号に対して時間周波数変換を行い、参照信号を対数振幅スペクトログラムに変換し、音響特徴量抽出部83および音響特徴量抽出部84に供給する。

【0124】

ステップS 17において、音響特徴量抽出部83は音響特徴量RA1の抽出処理を行なって、参照信号の音響特徴量RA1を算出し、マッチング処理部23のマスクパターン生成部111に供給する。

【0125】

また、ステップS 18において、音響特徴量抽出部84は音響特徴量RA2の抽出処理を行なって、参照信号の音響特徴量RA2を算出し、マッチング処理部23の類似度計算部112に供給する。

【0126】

なお、これらのステップS 17およびステップS 18の処理は、ステップS 13およびステップS 14の処理と同様であるので、その説明は省略する。但し、ステップS 17およびステップS 18の処理では、処理対象の信号が入力信号ではなく参照信号とされる。

【0127】

ステップS 19において、マスクパターン生成部111は、音響特徴量抽出部53からの音響特徴量IA1と、音響特徴量抽出部83からの音響特徴量RA1とに基づいてマスクパターンを生成する。例えば、マスクパターン生成部111は式(5)の演算を行なうことでマスクパターンを生成する。

【0128】

ステップS 20において、マスクパターン生成部111は、音響特徴量A1の類似度を算出する。例えば、マスクパターン生成部111は式(6)により音響特徴量A1の類似度を算出する。そして、マスクパターン生成部111は、算出されたマスクパターンと、音響特徴量A1の類似度とを類似度計算部112に供給する。

【0129】

ステップS 21において類似度計算部112は、音響特徴量抽出部54からの音響特徴

10

20

30

40

50

量 I A 2、音響特徴量抽出部 8 4 からの音響特徴量 R A 2、およびマスクパターン生成部 1 1 1 からのマスクパターンと類似度に基づいて、入力信号と参照信号の最終的な類似度を算出する。

【 0 1 3 0 】

例えば、類似度計算部 1 1 2 は、式 (7) を計算することで、入力信号と参照信号、つまり入力信号のコンテンツと参照信号のコンテンツとの類似度を算出し、コンテンツ属性データとともに比較統合部 1 1 3 に供給する。

【 0 1 3 1 】

ステップ S 2 2 において、比較統合部 1 1 3 は、類似度計算部 1 1 2 から供給された類似度に基づいて、参照信号のコンテンツと、入力信号に含まれるコンテンツが同一コンテンツであるかの判定を行なう。

10

【 0 1 3 2 】

例えば、比較統合部 1 1 3 は、複数の参照信号について得られた類似度のうち、所定の閾値を超える最大のものを特定し、特定された類似度の参照信号のコンテンツが、入力信号のコンテンツであるとする。比較統合部 1 1 3 は、このようにして特定された入力信号のコンテンツのコンテンツ属性データと、コンテンツ同定の判定結果とを出力し、一致検索処理は終了する。

【 0 1 3 3 】

以上のようにして、音声処理装置 1 1 は、入力信号と参照信号から正弦波らしさを示す音響特徴量 A 1 を算出するとともに、音響特徴量 A 1 からマスクパターンを生成し、そのマスクパターンと、信号の個性を示す音響特徴量 A 2 とから類似度を算出する。

20

【 0 1 3 4 】

このように、入力信号から得られた音響特徴量 I A 1 と、参照信号から得られた音響特徴量 R A 1 とからマスクパターンを生成すれば、残響や混入雑音に頑健なものを得ることができる。これにより、より高精度に任意のコンテンツの同定を行なうことができる。

【 0 1 3 5 】

ところで、上述した一連の処理は、ハードウェアにより実行することもできるし、ソフトウェアにより実行することもできる。一連の処理をソフトウェアにより実行する場合には、そのソフトウェアを構成するプログラムが、コンピュータにインストールされる。ここで、コンピュータには、専用のハードウェアに組み込まれているコンピュータや、各種のプログラムをインストールすることで、各種の機能を実行することが可能な、例えば汎用のパーソナルコンピュータなどが含まれる。

30

【 0 1 3 6 】

図 1 0 は、上述した一連の処理をプログラムにより実行するコンピュータのハードウェアの構成例を示すブロック図である。

【 0 1 3 7 】

コンピュータにおいて、CPU (Central Processing Unit) 7 0 1 , ROM (Read Only Memory) 7 0 2 , RAM (Random Access Memory) 7 0 3 は、バス 7 0 4 により相互に接続されている。

【 0 1 3 8 】

バス 7 0 4 には、さらに、入出力インターフェース 7 0 5 が接続されている。入出力インターフェース 7 0 5 には、入力部 7 0 6、出力部 7 0 7、記録部 7 0 8、通信部 7 0 9、及びドライブ 7 1 0 が接続されている。

40

【 0 1 3 9 】

入力部 7 0 6 は、キーボード、マウス、マイクロホン、撮像素子などよりなる。出力部 7 0 7 は、ディスプレイ、スピーカなどよりなる。記録部 7 0 8 は、ハードディスクや不揮発性のメモリなどよりなる。通信部 7 0 9 は、ネットワークインターフェースなどよりなる。ドライブ 7 1 0 は、磁気ディスク、光ディスク、光磁気ディスク、又は半導体メモリなどのリムーバブルメディア 7 1 1 を駆動する。

【 0 1 4 0 】

50

以上のように構成されるコンピュータでは、CPU 701が、例えば、記録部708に記録されているプログラムを、入出力インターフェース705及びバス704を介して、RAM703にロードして実行することにより、上述した一連の処理が行われる。

【0141】

コンピュータ(CPU701)が実行するプログラムは、例えば、パッケージメディア等としてのリムーバブルメディア711に記録して提供することができる。また、プログラムは、ローカルエリアネットワーク、インターネット、デジタル衛星放送といった、有線または無線の伝送媒体を介して提供することができる。

【0142】

コンピュータでは、プログラムは、リムーバブルメディア711をドライブ710に装着することにより、入出力インターフェース705を介して、記録部708にインストールすることができる。また、プログラムは、有線または無線の伝送媒体を介して、通信部709で受信し、記録部708にインストールすることができる。その他、プログラムは、ROM702や記録部708に、あらかじめインストールしておくことができる。

【0143】

なお、コンピュータが実行するプログラムは、本明細書で説明する順序に沿って時系列に処理が行われるプログラムであっても良いし、並列に、あるいは呼び出しが行われたとき等の必要なタイミングで処理が行われるプログラムであっても良い。

【0144】

また、本技術の実施の形態は、上述した実施の形態に限定されるものではなく、本技術の要旨を逸脱しない範囲において種々の変更が可能である。

【0145】

例えば、本技術は、1つの機能をネットワークを介して複数の装置で分担、共同して処理するクラウドコンピューティングの構成をとることができる。

【0146】

また、上述のフローチャートで説明した各ステップは、1つの装置で実行する他、複数の装置で分担して実行することができる。

【0147】

さらに、1つのステップに複数の処理が含まれる場合には、その1つのステップに含まれる複数の処理は、1つの装置で実行する他、複数の装置で分担して実行することができる。

【0148】

さらに、本技術は、以下の構成とすることも可能である。

【0149】

[1]

同定対象となるコンテンツの入力信号に基づいて、各時間周波数領域における信号の正弦波らしさを示す第1の音響特徴量と、前記第1の音響特徴量とは異なる第2の音響特徴量とを算出する入力信号処理部と、

予め用意したコンテンツの参照信号に基づいて、前記第1の音響特徴量と前記第2の音響特徴量とを算出する参照信号処理部と、

前記入力信号の前記第1の音響特徴量および前記第2の音響特徴量と、前記参照信号の前記第1の音響特徴量および前記第2の音響特徴量とに基づいて、前記入力信号と前記参照信号の類似度を計算するマッチング処理部と

を備える音声処理装置。

[2]

前記マッチング処理部は、前記入力信号と前記参照信号の前記第1の音響特徴量に基づいて、各時間周波数領域におけるコンテンツの信号らしさを示すマスクパターンを生成し、前記マスクパターン、前記第1の音響特徴量、および前記第2の音響特徴量に基づいて前記類似度を計算する

[1]に記載の音声処理装置。

10

20

30

40

50

[3]

前記マッチング処理部は、前記入力信号の前記第 1 の音響特徴量と、前記参照信号の前記第 1 の音響特徴量との類似度をさらに算出し、前記マスクパターン、前記第 1 の音響特徴量の前記類似度、および前記第 2 の音響特徴量に基づいて、前記入力信号と前記参照信号の前記類似度を計算する

[2] に記載の音声処理装置。

[4]

前記マッチング処理部は、前記第 1 の音響特徴量の前記類似度に対する前記入力信号の寄与率よりも、前記第 1 の音響特徴量の前記類似度に対する前記参照信号の寄与率をより大きくして、前記第 1 の音響特徴量の前記類似度を算出する

10

[3] に記載の音声処理装置。

[5]

前記第 2 の音響特徴量は、前記入力信号または前記参照信号のスペクトログラムに基づいて算出され、時間軸および周波数軸において前記第 1 の音響特徴量と同じ粒度を有する

[1] 乃至 [4] の何れかに記載の音声処理装置。

【符号の説明】

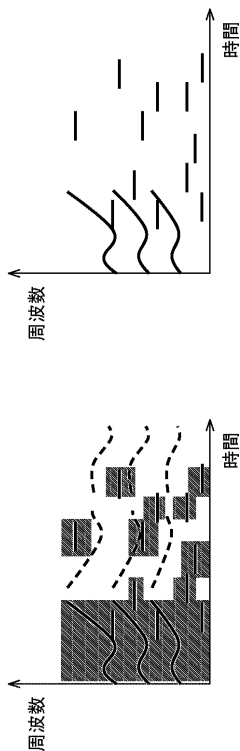
【 0 1 5 0 】

1 1 音声処理装置, 2 1 入力信号処理部, 2 2 参照信号処理部, 2 3 マッチング処理部, 5 3 音響特徴量抽出部, 5 4 音響特徴量抽出部, 8 3 音響特徴量抽出部, 8 4 音響特徴量抽出部, 1 1 1 マスクパターン生成部, 1 1 2 類似度計算部, 1 1 3 比較統合部

20

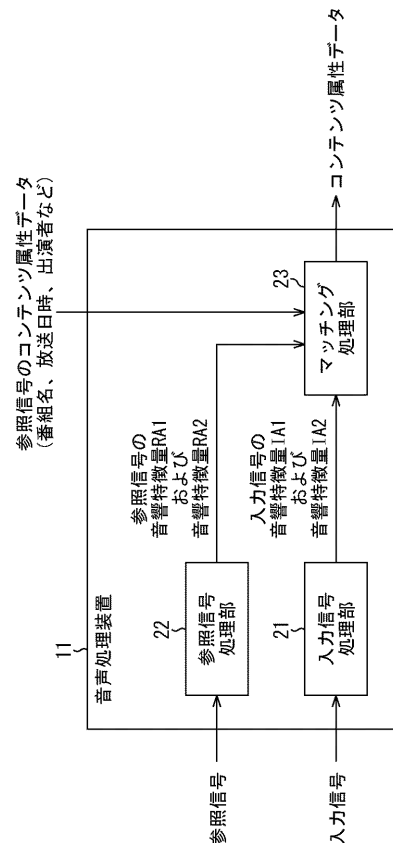
【図 1】

図1

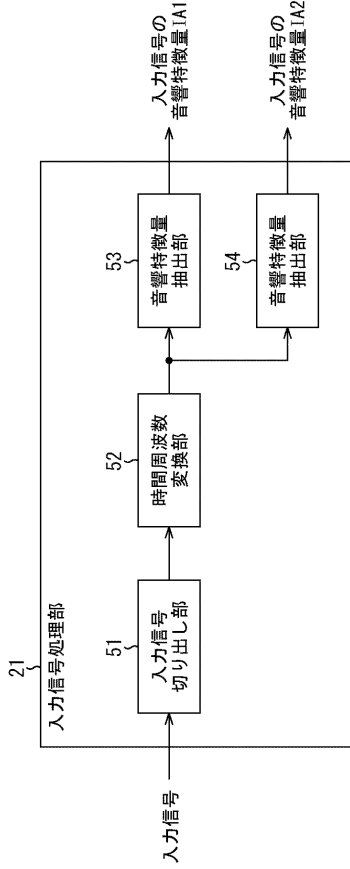


【図 2】

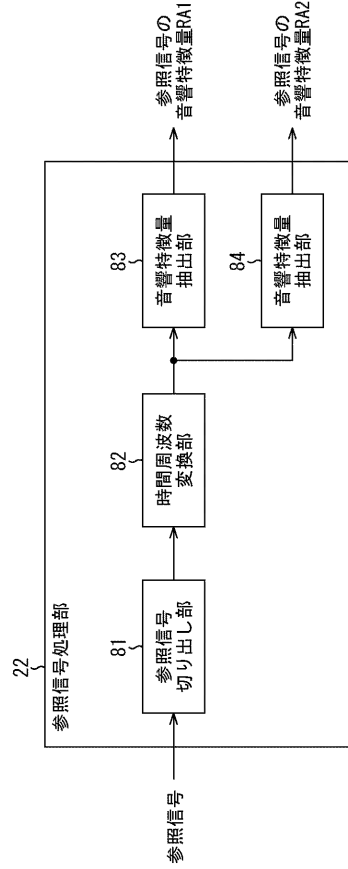
図2



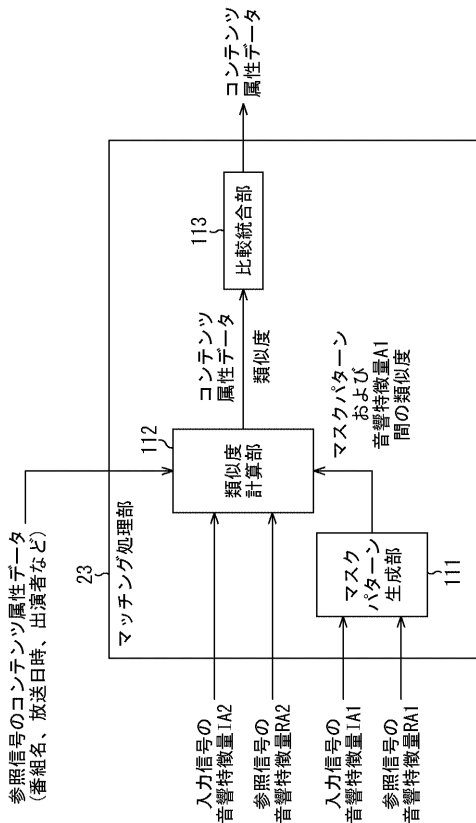
【図3】
図3



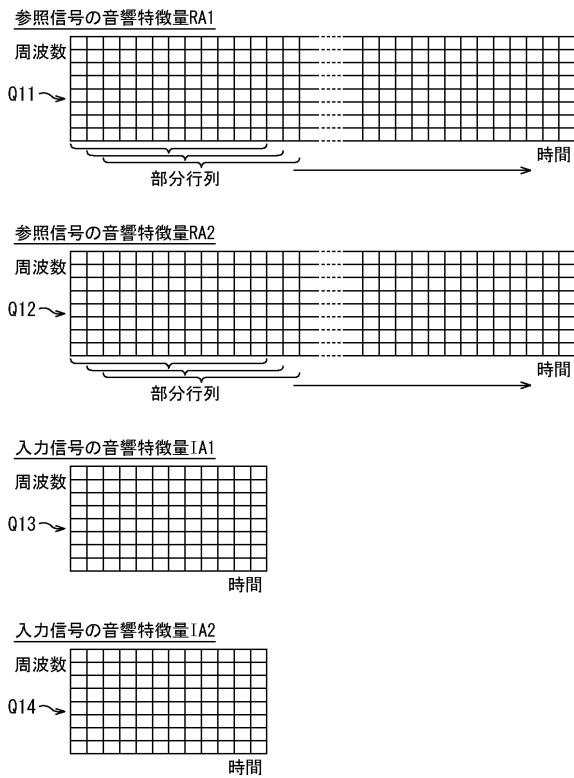
【図4】
図4



【図5】
図5

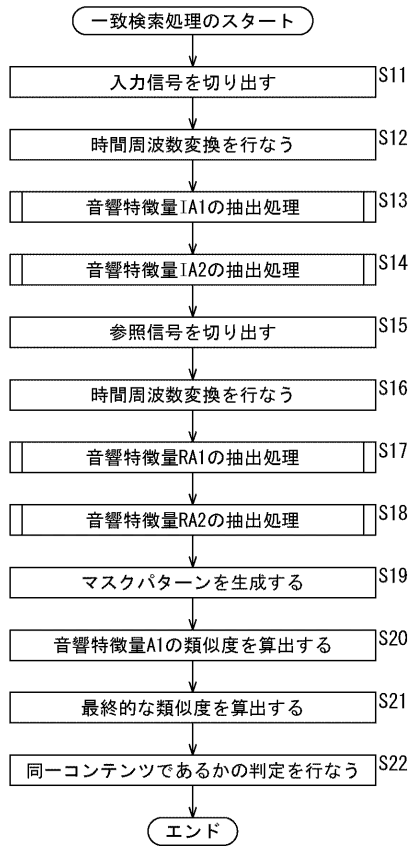


【図6】
図6



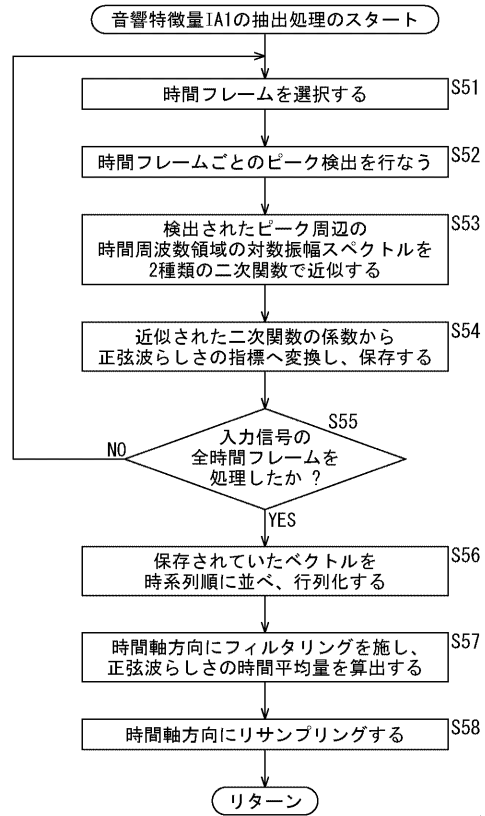
【図7】

図7



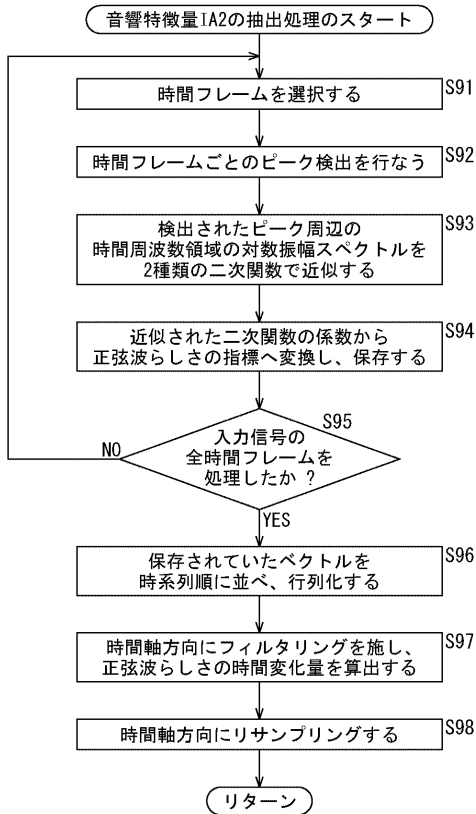
【図8】

図8



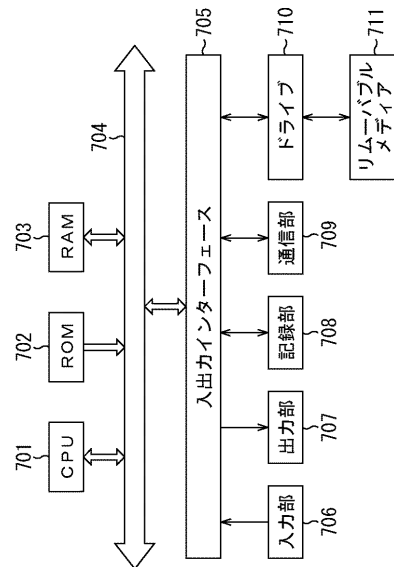
【図9】

図9



【図10】

図10



フロントページの続き

(72)発明者 西口 正之
東京都港区港南1丁目7番1号 ソニー株式会社内

審査官 菊池 智紀

(56)参考文献 特開2012-226080(JP,A)
特開2012-098360(JP,A)
澁谷崇 他, "トーン構造記述子を用いた高速背景音楽検索", 情報処理学会研究報告, 2011年7月29日, Vol.2011-MUS-91, No.17, pp.1-8

(58)調査した分野(Int.Cl., DB名)
G10L 25/00 - 25/93,
15/00 - 15/34
G06F 17/30