

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4580297号
(P4580297)

(45) 発行日 平成22年11月10日(2010.11.10)

(24) 登録日 平成22年9月3日(2010.9.3)

| | |
|-------------------------|-----------------------|
| (51) Int.Cl. | F I |
| G 1 O L 21/04 (2006.01) | G 1 O L 21/04 1 1 O Z |
| G 1 O L 11/00 (2006.01) | G 1 O L 11/00 4 O 2 L |

請求項の数 8 (全 24 頁)

| | | | |
|-----------|------------------------------|-----------|---------------------|
| (21) 出願番号 | 特願2005-204211 (P2005-204211) | (73) 特許権者 | 000005821 |
| (22) 出願日 | 平成17年7月13日(2005.7.13) | | パナソニック株式会社 |
| (65) 公開番号 | 特開2007-25039 (P2007-25039A) | | 大阪府門真市大字門真1006番地 |
| (43) 公開日 | 平成19年2月1日(2007.2.1) | (74) 代理人 | 100098291 |
| 審査請求日 | 平成20年5月29日(2008.5.29) | | 弁理士 小笠原 史朗 |
| | | (72) 発明者 | 三崎 正之 |
| | | | 大阪府門真市大字門真1006番地 松下 |
| | | | 電器産業株式会社内 |
| | | (72) 発明者 | 正木 芽衣子 |
| | | | 大阪府門真市大字門真1006番地 松下 |
| | | | 電器産業株式会社内 |
| | | (72) 発明者 | 河村 岳 |
| | | | 大阪府門真市大字門真1006番地 松下 |
| | | | 電器産業株式会社内 |

最終頁に続く

(54) 【発明の名称】 音声再生装置、音声録音再生装置、およびそれらの方法、記録媒体、集積回路

(57) 【特許請求の範囲】

【請求項 1】

入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための音声再生装置であって、

前記音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別部と、

前記音声区間および前記非音声区間に基づき、算出用フレーム長に対して当該音声区間が含まれる比率を示す音声含有率と、前記算出用フレーム長の音声含有率の平均値および標準偏差である統計値を少なくとも算出する音声情報算出部と、

前記音声速度変換処理の比率が1以上の速度比を基準値として予め設定し、前記統計値および前記算出用フレーム長の音声含有率に応じて、前記基準値から前記音声区間の速度比を算出する速度比算出部とを備え、

前記音声情報算出部は、それぞれ時間長が異なる前記算出用フレーム長を複数設定してそれぞれ前記音声含有率および前記統計値を算出し、

前記速度比算出部は、前記算出用フレーム長における音声含有率の前記平均値に対する差分値および前記標準偏差によって前記算出用フレーム長毎に得られる値を用いて算出される係数を前記速度比の基準値に乗ずることで、前記算出用フレーム長の音声含有率が相対的に高いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より小さな値に変更し、前記算出用フレーム長の音声含有率が相対的に低いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より大きな値に変更し、音声含有

10

20

率の時間的な変動に対して適応的に前記音声区間の速度比を調整する速度制御を行う音声再生装置。

【請求項 2】

前記速度比算出部は、ユーザの操作に応じて前記短縮された再生時間を設定し、音声区間における前記変更した速度比に基づいて、前記音声信号の再生時間が前記設定された再生時間となるように前記非音声区間の速度比を算出することを特徴とする、請求項 1 に記載の音声再生装置。

【請求項 3】

前記速度比算出部は、前記設定された再生時間内において前記非音声区間の速度比を一定に算出することを特徴とする、請求項 2 に記載の音声再生装置。

10

【請求項 4】

前記入力される音声信号のうち、少なくとも前記算出用フレーム長分の音声信号を含むように当該音声信号を順次更新しながら記録するバッファと、

前記バッファに記録された音声信号に対して音声速度変換処理を行って出力する速度変換部とを、さらに備え、

前記判別部は、前記バッファに記録された前記算出用フレーム長の音声信号に対して前記音声区間と前記非音声区間とを判別し、

前記音声情報算出部は、これまでに算出した統計値を単位時間毎に順次更新し、

前記速度比算出部は、前記単位時間ごとに更新される前記統計値および当該更新時の前記算出用フレーム長に設定された音声含有率に応じて前記音声区間の速度比を算出し、

20

前記速度変換部は、前記バッファで順次更新される音声信号に対して、前記単位時間ごとに算出された前記音声区間の速度比を用いて順次音声速度変換処理を行うことを特徴とする、請求項 1 に記載の音声再生装置。

【請求項 5】

入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための音声再生方法であって、

前記音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別ステップと、

前記音声区間および前記非音声区間に基づき、算出用フレーム長に対して当該音声区間が含まれる比率を示す音声含有率と、前記算出用フレーム長の音声含有率の平均値および標準偏差である統計値を少なくとも算出する音声情報算出ステップと、

30

前記音声速度変換処理の比率が 1 以上の速度比を基準値として予め設定し、前記統計値および前記算出用フレーム長の音声含有率に応じて、前記基準値から前記音声区間の速度比を算出する速度比算出ステップとを備え、

前記音声情報算出ステップは、それぞれ時間長が異なる前記算出用フレーム長を複数設定してそれぞれ前記音声含有率および前記統計値を算出し、

前記速度比算出ステップは、前記算出用フレーム長における音声含有率の前記平均値に対する差分値および前記標準偏差によって前記算出用フレーム長毎に得られる値を用いて算出される係数を前記速度比の基準値に乗ずることで、前記算出用フレーム長の音声含有率が相対的に高いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より小さな値に変更し、前記算出用フレーム長の音声含有率が相対的に低いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より大きな値に変更し、音声含有率の時間的な変動に対して適応的に前記音声区間の速度比を調整する速度制御を行う音声再生方法。

40

【請求項 6】

入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するためのコンピュータで実行される音声再生プログラムを記録した当該コンピュータで読み取り可能な記録媒体であって、

前記コンピュータに、

前記音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する

50

判別ステップと、

前記音声区間および前記非音声区間に基づき、算出用フレーム長に対して当該音声区間が含まれる比率を示す音声含有率と、前記算出用フレーム長の音声含有率の平均値および標準偏差である統計値を少なくとも算出する音声情報算出ステップと、

前記音声速度変換処理の比率が1以上の速度比を基準値として予め設定し、前記統計値および前記算出用フレーム長の音声含有率に応じて、前記基準値から前記音声区間の速度比を算出する速度比算出ステップとを含み、

前記音声情報算出ステップは、それぞれ時間長が異なる前記算出用フレーム長を複数設定してそれぞれ前記音声含有率および前記統計値を算出し、

前記速度比算出ステップは、前記算出用フレーム長における音声含有率の前記平均値に対する差分値および前記標準偏差によって前記算出用フレーム長毎に得られる値を用いて算出される係数を前記速度比の基準値に乗ずることで、前記算出用フレーム長の音声含有率が相対的に高いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より小さな値に変更し、前記算出用フレーム長の音声含有率が相対的に低いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より大きな値に変更し、音声含有率の時間的な変動に対して適応的に前記音声区間の速度比を調整する速度制御を行うプログラムを記録した、コンピュータに読み取り可能な記録媒体。

【請求項7】

入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための集積回路であって、

前記音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別部と、

前記音声区間および前記非音声区間に基づき、算出用フレーム長に対して当該音声区間が含まれる比率を示す音声含有率と、前記算出用フレーム長の音声含有率の平均値および標準偏差である統計値を少なくとも算出する音声情報算出部と、

前記音声速度変換処理の比率が1以上の速度比を基準値として予め設定し、前記統計値および前記算出用フレーム長の音声含有率に応じて、前記基準値から前記音声区間の速度比を算出する速度比算出部とを備え、

前記音声情報算出部は、それぞれ時間長が異なる前記算出用フレーム長を複数設定してそれぞれ前記音声含有率および前記統計値を算出し、

前記速度比算出部は、前記算出用フレーム長における音声含有率の前記平均値に対する差分値および前記標準偏差によって前記算出用フレーム長毎に得られる値を用いて算出される係数を前記速度比の基準値に乗ずることで、前記算出用フレーム長の音声含有率が相対的に高いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より小さな値に変更し、前記算出用フレーム長の音声含有率が相対的に低いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より大きな値に変更し、音声含有率の時間的な変動に対して適応的に前記音声区間の速度比を調整する速度制御を行う集積回路。

【請求項8】

入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための音声録音再生装置であって、

前記入力される音声信号を記録する情報記録部と、

前記情報記録部に記録される前の音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別部と、

前記音声区間および前記非音声区間に基づき、算出用フレーム長に対して当該音声区間が含まれる比率を示す音声含有率と、前記算出用フレーム長の音声含有率の平均値および標準偏差である統計値を少なくとも算出する音声情報算出部と、

前記音声速度変換処理の比率が1以上の速度比を基準値として予め設定し、前記統計値および前記算出用フレーム長の音声含有率に応じて、前記基準値から前記音声区間の速度比を算出する速度比算出部とを備え、

前記音声情報算出部は、それぞれ時間長が異なる前記算出用フレーム長を複数設定してそれぞれ前記音声含有率および前記統計値を算出し、

前記速度比算出部は、前記算出用フレーム長における音声含有率の前記平均値に対する差分値および前記標準偏差によって前記算出用フレーム長毎に得られる値を用いて算出される係数を前記速度比の基準値に乗ずることで、前記算出用フレーム長の音声含有率が相対的に高いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より小さな値に変更し、前記算出用フレーム長の音声含有率が相対的に低いときに当該算出用フレーム長における前記音声区間の速度比を当該基準値より大きな値に変更し、音声含有率の時間的な変動に対して適応的に前記音声区間の速度比を調整する速度制御を行う音声録音再生装置。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声再生装置、音声録音再生装置、およびそれらの方法、記録媒体、集積回路に関し、より特定のには、再生速度を変換して再生する音声再生装置、音声録音再生装置、およびそれらの方法、記録媒体、集積回路に関する。

【背景技術】

【0002】

従来、予め記録された音声を再生する音声再生装置において、声の高さを変えることなく、より高速に再生する方法が知られている（例えば、特許文献1参照）。特許文献1に開示された音声再生装置では、音声信号全体を指定速度で再生するとき、音声区間については部分的に再生速度比を低速化している。これにより、特許文献1に開示された従来の音声再生装置は、情報の欠落が少なく、聴き取りやすい再生音声を提供することができる。

20

【特許文献1】特開2001-222300号公報

【0003】

以下、図11を参照して、上記特許文献1に開示された従来の音声再生装置9について、具体的に説明する。図11は、従来の音声再生装置9の構成を示すブロック図である。図11において、従来の音声再生装置9は、音響分析部91、話速変換部92、非音声区間長制御部93、および合成部94を備える。

30

【0004】

音響分析部91は、入力される音声データに対して、予め設定されているパワー閾値に基づき音声区間および非音声区間を判別する。そして、音響分析部91は、音声区間および非音声区間の時間情報をそれぞれ求める。図11に示す従来の音声再生装置9では、音響分析部91において判別された音声区間および非音声区間に対して、異なる再生処理を適用する。音響分析部91で判別された音声区間の音声データおよび上記各時間情報は、話速変換部92に出力される。音響分析部91で判別された非音声区間の音声データは、非音声区間長制御部93に出力される。

【0005】

40

話速変換部92は、まず音声区間の音声データと上記各時間情報とに基づいて、一定時間長以上の非音声区間に挟まれた音声区間を特定する。そして、話速変換部92は、当該音声区間の冒頭部分の速度比を所定速度比より遅く、末尾に向けて次第に所定速度比に戻すような速度比制御を行う。速度比が制御された音声区間の音声データは、合成部94に出力される。また、話速変換部92は、波形の伸長処理によって生じる音声区間の遅延時間情報を非音声区間長制御部93に出力する。

【0006】

一方、非音声区間長制御部93では、話速変換部92から出力された上記遅延時間情報に基づいて、非音声区間の音声データに対して削除および圧縮する処理を適宜行う。つまり、非音声区間長制御部93では、目標の指定速度比に合うように、かつ、話速変換部9

50

2で生じた音声区間の遅延を解消するような処理が行われる。非音声区間長制御部93において処理された非音声区間の音声データは、合成部94に出力される。

【0007】

合成部94は、話速変換部92から出力された音声区間の音声データと、非音声区間長制御部93から出力された非音声区間の音声データとを合成する。そして、合成部94は、速度比が変換された音声区間と非音声区間とが合成された音声データを変換音声データとして、最終的な再生音声を出力する。

【0008】

上記従来の音声再生装置9では、例えば指定速度として m 倍速(m は1以上の正数)が与えられたとき、音声区間の冒頭部分では m 倍速より遅い速度比で再生する。そして、従来の音声再生装置9は、音声区間の末尾に向かって次第に再生速度比を速くする。ここで、一般的に音声区間の冒頭部分には、重要な情報が含まれている場合が多い。したがって、従来の音声再生装置9によれば、音声区間の冒頭部分にある重要な情報を欠落させることなく、聴きとりやすい再生を実現することができる。このように従来の音声再生装置9では、音声区間については聴き取りやすい処理が、非音声区間については指定速度比に適應するような処理がそれぞれ行われている。

【発明の開示】

【発明が解決しようとする課題】

【0009】

ここで、高速再生時には、音声の発話速度が速くなり、ユーザにとって内容を理解するための負荷が大きくなる。さらに、番組全体の中で音声区間が偏って集中すると(音声が続続的に発声されると)、ユーザにとってさらに理解が困難になる。しかしながら、上記従来の音声再生装置9では、一つの音声区間の中で再生速度比を変更することのみを想定している。つまり、上記従来の音声再生装置9では、例えばテレビ番組などの全体を通して、同一の速度比制御処理が適用される。したがって、従来の音声再生装置9においては、音声区間が偏って集中する部分で相対的に音声の内容の聴き取りが困難になるという本質的課題があった。

【0010】

それ故、本発明の目的は、テレビなどの番組全体を考慮した最適な速度比制御を行って、より聴き取りやすい再生を実現する音声再生装置、音声録音再生装置、およびそれらの方法、記録媒体、および集積回路を提供することを目的とする。

【課題を解決するための手段】

【0011】

第1の発明は、入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための音声再生装置であって、音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別部と、音声区間および非音声区間に基づき、所定時間長に対して当該音声区間が含まれる比率を示す音声含有率を少なくとも算出する音声情報算出部と、音声速度変換処理の比率が1以上の速度比を基準値として予め設定し、所定時間長の音声含有率が相対的に高いときに当該所定時間長における音声区間の速度比を当該基準値より小さな値に変更し、所定時間長の音声含有率が相対的に低いときに当該所定時間長における音声区間の速度比を当該基準値より大きな値に変更する速度比算出部とを備え、音声含有率の時間的な変動に対して適応的に音声区間の速度比を調整する速度制御を行うものである。

【0012】

第2の発明は、上記第1の発明において、速度比算出部は、ユーザの操作に応じて短縮された再生時間を設定し、音声区間における変更した速度比に基づいて、音声信号の再生時間が設定された再生時間となるように非音声区間の速度比を算出することを特徴とする。

【0013】

第3の発明は、上記第2の発明において、速度比算出部は、設定された再生時間内にお

10

20

30

40

50

いて非音声区間の速度比を一定に算出することを特徴とする。

【0014】

第4の発明は、上記第1の発明において、速度比算出部は、所定時間長の設定値を選択することによって音声含有率の時間的な変動に対する音声区間の速度比の適応度合いを可変することを特徴とする。

【0015】

第5の発明は、上記第1の発明において、音声再生装置は、入力される音声信号のうち、少なくとも所定時間長の音声信号を含むように当該音声信号を順次更新しながら記録するバッファと、バッファに記録された音声信号に対して音声速度変換処理を行って出力する速度変換部とを、さらに備え、判別部は、バッファに記録された所定時間長の音声信号に対して音声区間と非音声区間とを判別し、音声情報算出部は、さらに、所定時間長の音声含有率に関する統計値を算出して、これまでに算出した統計値を単位時間毎に順次更新し、速度比算出部は、単位時間ごとに更新される統計値および当該更新時の所定時間長に設定された音声含有率に応じて音声区間の速度比を算出し、速度変換部は、バッファで順次更新される音声信号に対して、単位時間ごとに算出された音声区間の速度比を用いて順次音声速度変換処理を行うことを特徴とする。

10

【0016】

第6の発明は、上記第1の発明において、音声情報算出部は、所定時間長の音声含有率に関する統計値をさらに算出し、速度比算出部は、統計値および所定時間長の音声含有率に応じて音声区間の速度比を算出することを特徴とする。

20

【0017】

第7の発明は、上記第5または6の発明において、統計値は、所定時間長の音声含有率の平均値および標準偏差であることを特徴とする。

【0018】

第8の発明は、上記第7の発明において、速度比算出部は、所定時間長における音声含有率の平均値に対する差分値および標準偏差によって得られる値を用いて算出される係数を速度比の基準値に乗じて、音声区間の速度比を算出することを特徴とする。

【0019】

第9の発明は、上記第7の発明において、音声情報算出部は、それぞれ時間長が異なる所定時間長を複数設定してそれぞれ音声含有率および統計値を算出し、速度比算出部は、所定時間長における音声含有率の平均値に対する差分値および標準偏差によって所定時間長毎に得られる値を用いて算出される係数を速度比の基準値に乗じて算出することを特徴とする。

30

【0020】

第10の発明は、入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための音声再生方法であって、音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別ステップと、音声区間および非音声区間に基づき、所定時間長に対して当該音声区間が含まれる比率を示す音声含有率を少なくとも算出する音声情報算出ステップと、音声速度変換処理の比率が1以上の速度比を基準値として予め設定し、所定時間長の音声含有率が相対的に高いときに当該所定時間長における音声区間の速度比を当該基準値より小さな値に変更し、所定時間長の音声含有率が相対的に低いときに当該所定時間長における音声区間の速度比を当該基準値より大きな値に変更する速度比算出ステップとを含む。

40

【0021】

第11の発明は、入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するためのコンピュータで実行される音声再生プログラムを記録した当該コンピュータで読み取り可能な記録媒体であって、コンピュータに、音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別ステップと、音声区間および非音声区間に基づき、所定時間長に対する当該音声区間が含まれる比率を示す音声含有率を少なくとも算出する音声情報算出ステップと、音声速度変換処理の

50

比率が1以上の速度比を基準値として予め設定し、所定時間長の音声含有率が相対的に高いときに当該所定時間長における音声区間の速度比を当該基準値より小さな値に変更し、所定時間長の音声含有率が相対的に低いときに当該所定時間長における音声区間の速度比を当該基準値より大きな値に変更する速度比算出ステップとを実行させるためのプログラムを記録した、コンピュータに読み取り可能な記録媒体である。

【0022】

第12の発明は、入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための集積回路であって、音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別部と、音声区間および非音声区間に基づき、所定時間長に対して当該音声区間が含まれる比率を示す音声含有率を少なくとも算出する音声情報算出部と、音声速度変換処理の比率が1以上の速度比を基準値として、所定時間長の音声含有率が相対的に高いときに当該所定時間長における音声区間の速度比を当該基準値より小さな値に変更し、所定時間長の音声含有率が相対的に低いときに当該所定時間長における音声区間の速度比を当該基準値より大きな値に変更する速度比算出部とを備える。

10

【0023】

第13の発明は、入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための音声録音再生装置であって、入力される音声信号を記録する情報記録部と、情報記録部に記録される前の音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別部と、音声区間および非音声区間に基づき、所定時間長に対して当該音声区間が含まれる比率を示す音声含有率を少なくとも算出する音声情報算出部と、音声速度変換処理の比率が1以上の速度比を基準値として予め設定し、所定時間長の音声含有率が相対的に高いときに当該所定時間長における音声区間の速度比を当該基準値より小さな値に変更し、所定時間長の音声含有率が相対的に低いときに当該所定時間長における音声区間の速度比を当該基準値より大きな値に変更する速度比算出部とを備える。

20

【0024】

第14の発明は、上記第13の発明において、情報記録部には、音声信号が記録される際に判別部が判別した結果が記録され、音声情報算出部は、情報記録部に記録された結果に基づいて、所定時間長の音声含有率を算出することを特徴とする。

30

【0025】

第15の発明は、上記第13の発明において、情報記録部には、音声信号が記録される際に、判別部が判別した結果および所定時間長の音声含有率が記録され、速度比算出部は、情報記録部に記録された所定時間長の音声含有率を用いて、音声区間の速度比を算出することを特徴とする。

【0026】

第16の発明は、入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するための音声録音再生方法であって、入力される音声信号を記録する情報記録ステップと、情報記録ステップに記録される前の音声信号に対して音声を含む音声区間と音声を含まない非音声区間とを判別する判別ステップと、音声区間および非音声区間に基づき、所定時間長に対して当該音声区間が含まれる比率を示す音声含有率を少なくとも算出する音声情報算出ステップと、音声速度変換処理の比率が1以上の速度比を基準値として予め設定し、所定時間長の音声含有率が相対的に高いときに当該所定時間長における音声区間の速度比を当該基準値より小さな値に変更し、所定時間長の音声含有率が相対的に低いときに当該所定時間長における音声区間の速度比を当該基準値より大きな値に変更する速度比算出ステップとを含む。

40

【0027】

第17の発明は、入力される音声信号に音声速度変換処理を適用して通常より再生時間を短縮して当該音声信号を聴取するためのコンピュータで実行される音声録音再生プログラムを記録した記録媒体であって、コンピュータに、入力される音声信号を記録部に記録

50

する情報記録ステップと、記録部に記録される前の音声信号に対して音声を含む音声区間と音声を含まない非音声区間と判別する判別ステップと、音声区間および非音声区間に基づき、所定時間長に対して当該音声区間が含まれる比率を示す音声含有率を少なくとも算出する音声情報算出ステップと、音声速度変換処理の比率が1以上の速度比を基準値として、所定時間長の音声含有率が相対的に高いときに当該所定時間長における音声区間の速度比を当該基準値より小さな値に変更し、所定時間長の音声含有率が相対的に低いときに当該所定時間長における音声区間の速度比を当該基準値より大きな値に変更する速度比算出ステップとを実行させるためのプログラムを記録した、コンピュータに読み取り可能な記録媒体である。

【発明の効果】

10

【0028】

第1の発明によれば、音声含有率の変動に応じた音声区間の速度比を算出することで、入力された音声信号の速度変換後の再生音質を音声含有率の変動に応じた了解性の優れたものにすることができる。

【0029】

第2の発明によれば、設定された再生時間となるように、重要な音声情報が含まれていない非音声区間の速度比を音声区間の速度比とは別に算出することで、音声区間の速度比をユーザが聴取可能な範囲内の速度比に調整することができる。

【0030】

第3の発明によれば、重要な音声情報が含まれていない非音声区間の速度比を一定の速度比とすることで、能率のよい速度変換をした再生が可能となる。

20

【0031】

第4の発明によれば、例えば所定時間長が長い場合には、調整される音声区間の速度比が音声含有率の変動に対して大局的でより正確性の高い値となる。また例えば、所定時間長が短い場合には、調整される音声区間の速度比が音声含有率の変動に対して敏感でより追従性のよい値となる。つまり、所定時間長の設定値を選択して音声含有率の変動に対する音声区間の速度比の適応度合いを可変することによって、正確性または追従性を自由に選択することができる。

【0032】

第5の発明によれば、統計値を単位時間毎に更新することで、音声信号の入力に応じて即時に速度変換処理をして再生することができる。

30

【0033】

第6の発明によれば、音声区間の速度比の算出に対して、統計値を用いることで、より実際の音声含有率の変動に即した音声区間の速度比を算出することができ、結果的に速度変換後の再生音質をより了解性のある自然なものにすることができる。

【0034】

第7の発明によれば、音声区間の存在の偏り度合いを考慮した音声区間の速度比を算出することができる。

【0035】

第8の発明によれば、音声区間の存在の偏り度合いに即した音声区間の速度比を算出することができる。

40

【0036】

第9の発明によれば、音声含有率の敏感な変動および大局的な変動の双方に対応した最適な音声区間の速度比を算出することができる。

【0037】

第13の発明によれば、音声含有率の変動に応じた音声区間の速度比を算出することで、記録した音声信号の速度変換後の再生音質を音声含有率の変動に応じた了解性の優れたものにすることができる。

【0038】

第14の発明によれば、音声信号を記録後、速度変換した再生が行われる前までの処理

50

時間を判別部における処理時間分だけ短縮することができる。

【 0 0 3 9 】

第 1 5 の発明によれば、音声信号を記録後、速度変換した再生が行われる前までの処理時間を判別部および音声情報算出部における処理時間分だけ短縮ことができ、音声信号を記録後、即時に速度変換をした再生を行うことができる。

【発明を実施するための最良の形態】

【 0 0 4 0 】

(第 1 の実施形態)

図 1 を参照して、本発明における第 1 の実施形態に係る音声再生装置について説明する。図 1 は、本発明における第 1 の実施形態に係る音声再生装置 1 の構成を示すブロック図である。図 1 において、音声再生装置 1 は、音声 / 非音声判別部 1 1、音声情報算出部 1 2、音声情報記録部 1 3、速度比算出部 1 4、および音声速度変換部 1 5 を有する。なお、本実施形態に係る音声再生装置 1 は、記録メディアなどに録音された音声信号を速度変換して再生する前に一旦、録音された音声信号全体について読み出し可能であることを想定した装置である。ここで、録音対象としては、例えばテレビやラジオ番組が挙げられる。また記録メディアは、例えば映画などが予め収録された D V D 等の記録メディアであってもよい。以下の説明では、一例として、第 1 の実施形態に係る音声再生装置 1 が、録音されたテレビ番組の音声信号に対して速度変換処理を行うとする。

【 0 0 4 1 】

記録メディアなどに録音された音声信号が読み出され、音声 / 非音声判別部 1 1 に入力される。音声 / 非音声判別部 1 1 は、入力された音声信号のパワーの包絡値や周期性などの分析を行う。そして、音声 / 非音声判別部 1 1 は入力された音声信号に対して音声区間および非音声区間を時間軸上で判別する。音声信号の時間軸上で判別された音声区間および非音声区間の情報（以下、判別情報という）は、速度変換した再生を行う前に音声情報算出部 1 2 に出力される。

【 0 0 4 2 】

音声情報算出部 1 2 は、音声 / 非音声区間の判別情報に基づいて、音声区間および非音声区間の速度比を算出するために必要な音声情報を算出する。音声情報としては、音声含有率、音声含有率の平均値、および標準偏差などを算出する。具体的には、音声情報算出部 1 2 は、録音された番組全体を通して音声含有率を算出した後に、音声含有率の平均値と標準偏差とを算出する。音声情報算出部 1 2 で算出された音声含有率、音声含有率の平均値、および標準偏差は、音声情報記録部 1 3 にそれぞれ記録される。以下、音声含有率、音声含有率の平均値、および標準偏差について説明する。

【 0 0 4 3 】

音声含有率は、所定数（少なくとも 1 つ以上）のフレームに対して音声区間が含まれる時間比率を示すものである。音声含有率はフレーム毎に算出される。ここでフレームとは、入力される音声信号を分析するための処理単位であり、当該フレームの時間長をフレーム長とする。当該フレームには、音声区間および / または非音声区間が含まれる。また、音声含有率の算出に用いられる少なくとも 1 つ以上のフレームを算出用フレームとし、その時間長を算出用フレーム長とする。以下の説明では、一例として、1 フレームの時間長（1 フレーム長）を 1 分とする。また、音声含有率を算出するための算出用フレーム長を n （ n は正数）分とする。つまり、1 フレーム長を 1 分としたので、算出用フレームは n 個のフレームから構成されることとなる。また、録音された番組全体のフレーム数が N （ N は正数）個あるとする。そして、フレームナンバーを k （ $k = 1 \sim N$ ）として、フレームナンバーが k のときのフレームを「第 k フレーム」とする。このとき、第 k フレームの音声含有率 $R_{i s _ n}(k)$ は、数式（1）で表現される。

【数 1】

$$Ris_n(k) = \frac{(k \text{ フレームを含むそれから先の } n \text{ 分間における音声区間の時間長})}{(k \text{ フレームを含むそれから先の } n \text{ 分間におけるフレーム長})} \quad \cdot \cdot \quad (1)$$

つまり、数式(1)によって算出される第kフレームの音声含有率 $Ris_n(k)$ は、算出用フレーム長に対して音声区間が含まれる時間比率を示す。

【0044】

ここで、図2～図4を参照して、上記音声含有率 $Ris_n(k)$ の算出例を挙げる。図2～図4では、一例として、テレビ放送のドキュメンタリ番組(30分間)の音声含有率を算出するとし、1分、5分、および10分の3種類の算出用フレーム長で算出している。図2は、算出用フレーム長が1分のときの音声含有率 $Ris_1(k)$ の算出例を示す図である。図3は、算出用フレーム長が5分のときの音声含有率 $Ris_5(k)$ の算出例を示す図である。図4は、算出用フレーム長が10分のときの音声含有率 $Ris_10(k)$ の算出例を示す図である。なお、図2～図4において、横軸はフレームナンバー(k)を示し、縦軸は音声含有率(%)を示す。

10

【0045】

図2において、第1フレーム($k=1$)の音声含有率 $Ris_1(1)$ は、算出用フレーム長を1分としたので、数式(1)より第1フレームの音声含有率そのものとなる。図3においては、数式(1)より算出される第1フレームの音声含有率 $Ris_5(1)$ は、図2の第1～第5フレームの音声含有率を平均したものである。図4においては、数式(1)より算出される第1フレームの音声含有率 $Ris_10(1)$ は、図2の第1～第10フレームの音声含有率を平均したものである。

20

【0046】

図2～図4に示すように、各算出用フレーム長で音声含有率の変動の様子が異なることが分かる。具体的には、算出用フレーム長が短い場合(図2)には、音声含有率のフレーム間の変動が比較的大きくなる。つまり、算出用フレーム長が短い場合には、音声含有率の実際の変動が敏感に反映されたものとなる。これに対し、図3および図4に示すように、算出用フレーム長が長くなるにつれて、音声含有率のフレーム間の変動が比較的小さくなる。これは、上述したように、算出用フレーム長が長くなるにつれて各フレームの音声含有率が平均化されるためである。つまり、算出用フレーム長が長い場合には、平均化によって小さい変動が吸収され、音声含有率の変動が大局的に反映される。また、各算出用フレーム長の分散および標準偏差も、音声含有率の変動の違いにより、異なる値となる。

30

【0047】

次に音声含有率の平均値および標準偏差について説明する。音声含有率の平均値は、音声含有率 $Ris_n(k)$ を番組全体において平均した値である。上述した図2でいえば、 $Ris_1(1)$ から $Ris_1(30)$ の音声含有率を平均した値である。つまり、算出用フレーム長 n (n は正数) で表現すれば、音声含有率の平均値は、 $Ris_n(1)$ から $Ris_n(N)$ までの音声含有率の平均である。また、標準偏差は、音声含有率 $Ris_n(k)$ と音声含有率の平均値とを用いて算出される値である。ここで、上記図2～図4に示した音声含有率 $Ris_n(k)$ の値をもとに、各算出用フレーム長について、それぞれ音声含有率の平均値と標準偏差とを求めると図5に示すような値となる。図5は、各算出用フレーム長の音声含有率の平均値および標準偏差の算出結果を示す図である。図5において、算出用フレーム長が1分である音声含有率の平均値 A_1 は0.506と、算出用フレーム長が5分である音声含有率の平均値 A_5 は0.498と、算出用フレーム長が10分である音声含有率の平均値 A_{10} は0.488となる。また、図5において、平均値 A_1 に対する標準偏差 S_1 は0.161と、平均値 A_5 に対する標準偏差 S_5 は0.073と、平均値 A_{10} に対する標準偏差 S_{10} は0.028となる。

40

【0048】

このように、図5に示すように、標準偏差においては、算出用フレーム長が短い場合に

50

は、変動が大きく（ばらつきが大きく）なるために標準偏差の値が大きくなる。算出用フレーム長が長い場合には、変動が小さく（ばらつきが小さく）なるために標準偏差の値が小さくなる。つまり、標準偏差は、算出用フレーム長の長さによって大きな影響を受ける値であり、一般的には番組全体における音声区間の存在の偏りを示す値と考えることができる。

【 0 0 4 9 】

次に、入力される音声信号を速度変換して再生する段階において、速度比算出部 1 4 は、音声情報記録部 1 3 に記録された音声情報（音声含有率、音声含有率の平均値、および標準偏差）を用いて、音声区間の存在の偏りに応じた音声区間の速度比をフレーム毎に算出する。そして、速度比算出部 1 4 は、上記音声区間の速度比とユーザなどが入力する所望再生時間とに基づいて、非音声区間の速度比を算出する。そして、速度比算出部 1 4 は、音声 / 非音声判別部 1 1 において判別された判別情報に対して、フレーム毎の速度比を設定して音声速度変換部 1 5 へ出力する。なお、ここでは算出された各フレームの音声区間の速度比は、当該フレーム内に存在する音声区間に一律に適用されたとする。また、非音声区間の速度比は、後述するように例えば一定の速度比でフレーム内の非音声区間に適用されたとする。

【 0 0 5 0 】

ここで、速度比の算出方法を説明する前に、音声の再生速度と聴き取りやすさの関係について説明する。通常の再生時間より短い時間で音声信号を聴取するために、通常の再生時間に対する再生時間長の設定値である目標再生時間比 R_t ($0 < R_t < 1$) が与えられたとする。例えばユーザが通常の再生時間に対して半分の再生時間で聴取しようとする、目標再生時間比 R_t は $R_t = 0.5$ となる。このような目標再生時間比 R_t は、数式 (2) で表現される。数式 (2) において、音声含有率の平均値を A_0 と、音声含有率が一定であるときの音声区間の速度比を SR_{s0} と、および音声含有率が一定であるときの非音声区間の速度比を SR_{ns0} とする。

【 数 2 】

$$R_t = \frac{A_0}{SR_{s0}} + \frac{1-A_0}{SR_{ns0}} \quad \cdots (2)$$

数式 (2) より、目標再生時間比 R_t および音声含有率の平均値 A_0 が与えられ、音声区間の速度比 SR_{s0} および非音声区間の速度比 SR_{ns0} のうち、いずれか一方が決まれば残りの他方が算出されることが分かる。

【 0 0 5 1 】

数式 (2) に示す音声区間の速度比 SR_{s0} は、一般的に通常速（等倍速）である 1.0 に近い値ほど聴き取りやすい。音声区間の速度比 SR_{s0} の値が大きくなるほど、単位時間当たりの情報量が増大するので、ユーザにとって聴取が難しくなる。また、音声区間の速度比 SR_{s0} の値が 2.0 程度になると、ユーザが聴き取りに集中しなければ内容を理解することが困難となる。このように、音声区間の速度比 SR_{s0} が大きい場合、長時間の聴取にかなりの困難さが生じてくる。したがって、音声区間の速度比 SR_{s0} は、目標再生時間比 R_t にある程度左右されることなく、ユーザの聴取可能な範囲内で設定されるのが最適である。これに基づき、通常は音声区間の速度比 SR_{s0} が 1 ~ 1.8 程度となる範囲を利用する。また、音声 / 非音声判別を利用しない一定速度比での再生であれば、実用上は速度比を 1.3 ~ 1.5 とすることが多い。

【 0 0 5 2 】

本実施形態においては、上記音声区間の速度比 SR_{s0} の最適な設定範囲を考慮しつつ、上述したように標準偏差が番組全体における音声区間の存在の偏りの度合いを示すと考え、現在のフレームにおける音声含有率と音声含有率の平均値との差および標準偏差とを用いて音声区間の速度比 SR_{s0} を可変する。すなわち、速度比 SR_{s0} を基準値として、音声区間が集中して音声含有率が上記音声含有率の平均値より高い部分に関しては当該基準値より音声区間の速度比を小さな値に変更し、逆に音声含有率が上記音声含有率の平

均値より低い部分に関しては当該基準値より音声区間の速度比を大きな値に変更する。

【 0 0 5 3 】

ここで、番組全体のフレーム数を N と、算出用フレーム長が n 分のときの標準偏差を S_n と、算出用フレーム長が n 分のときの第 k フレームにおける音声含有率を $Ris_n(k)$ と、第 k フレームにおける音声区間の速度比を $SRs(k)$ と、算出用フレーム長が n 分のときの音声含有率の平均値を A_n と、算出用フレーム長ごとに異なる重み係数を C_n と、非音声区間の速度比を $SRns$ と、および音声含有率が一定と仮定したときの基準値の速度比を $SRs0$ とする。なお、非音声区間の速度比 $SRns$ は、ここではフレームの音声含有率に依存せず一定値とする。このとき、音声区間の存在の偏りに応じた音声区間の速度比 $SRs(k)$ は、例えば数式(3)と表現される。

10

【数3】

$$SRs(k) = SRs0 * (1 + Sn * Cn * (An - Ris_n(k))) \quad \cdot \cdot \quad (3)$$

【 0 0 5 4 】

さらに、音声区間の速度比 $SRs(k)$ を音声含有率の大局的な変動および短期的な変動の双方が反映した値として算出する場合には、それぞれ時間長が異なる複数種類の算出用フレーム長の音声情報を用いて算出する。つまり、複数種類の算出用フレーム長の音声情報を多重に用いて音声区間の速度比を算出する。ここで、 M 種類の算出用フレーム長の音声情報を用いるとすると、第 k フレームの音声区間の速度比 $SRs(k)$ は、数式(4)となる。

20

【数4】

$$SRs(k) = SRs0 * (1 + \sum_{n=1}^M Sn * Cn * (An - Ris_n(k))) \quad \cdot \cdot \quad (4)$$

数式(4)において、 C_n は、算出用フレーム長ごとに異なる重み係数であり、各算出用フレーム長の音声含有率の偏差を音声区間の速度比 $SRs0$ に反映させる度合いを示すものである。

【 0 0 5 5 】

ここで、多重の音声情報として、算出用フレーム長が1分、5分、10分のときの各音声情報を用いたとき、音声区間の速度比 $SRs(k)$ は、数式(5)となる。

30

【数5】

$$SRs(k) = SRs0 * \{1 + S1 * C1 * (A1 - Ris_1(k)) + S5 * C5 * (A5 - Ris_5(k)) + S10 * C10 * (A10 - Ris_10(k))\} \quad \cdot \cdot \quad (5)$$

ここで、数式(5)により音声情報を多重に用いた速度比の算出結果の一例を図6に示す。図6は、音声情報を多重に用いた速度比の算出結果の一例を示す図である。なお、図6に示す算出例は、数式(5)において $SRs0 = 1.5$ 、 $C1 = 1$ 、 $C5 = 10$ 、 $C10 = 20$ として算出し、短期的変動よりも長期的な変動に重点を置いた速度比を算出することを意図した例である。また、 $A1$ 、 $A5$ 、 $A10$ 、 $S1$ 、 $S5$ 、 $S10$ 、 $Ris_1(k)$ 、 $Ris_5(k)$ 、および $Ris_10(k)$ は、それぞれ図2～図5に示した値である。また、図6では、数式(5)により音声情報を多重に用いた速度比の他に、数式(3)を用いて算出フレーム長(1分、5分、および10分)に基づく音声情報から算出された各速度比を比較のために示している。

40

【 0 0 5 6 】

図6において、菱形のプロットで描かれたグラフは、音声情報を多重に用いて算出された音声区間の速度比を示す。また、丸のプロットで描かれたグラフは、算出用フレーム長が1分のときの音声情報のみを用いて算出された音声区間の速度比を示す。四角のプロットで描かれたグラフは、算出用フレーム長が5分のときの音声情報のみを用いて算出された音声区間の速度比を示す。三角のプロットで描かれたグラフは、算出用フレーム長が10分のときの音声情報のみを用いて算出された音声区間の速度比を示す。

50

【 0 0 5 7 】

図 6 に示すように、音声情報を多重に用いて算出された音声区間の速度比は、それぞれ単独の算出用フレーム長の音声情報のみを用いて算出された速度比と比べて、音声含有率の短期的な変動および長期的な変動の双方が反映された値であることが分かる。つまり、多重の音声情報を用いて算出された音声区間の速度比は、番組全体を通して音声区間の存在の偏りに応じた速度比であり、聴き取りやすい速度となるよう考慮された速度比である。

【 0 0 5 8 】

速度比算出部 1 4 は、上述した方法で音声区間の速度比 SRs を算出後、入力される再生時間から設定される目標再生時間比 Rt を達成するように非音声区間の速度比 $SRns$ を算出する。なお、非音声区間の速度比 $SRns$ は、上述したように例えば可変とせず一定の速度比とする。これは、有益な情報の大部分が音声区間に含まれているため、音声区間の速度比の調整を重視したことに基づくものである。これにより、本実施形態に係る音声再生装置は、能率良い再生を実現できる。以下、非音声区間の速度比 $SRns$ の算出方法について説明する。

【 0 0 5 9 】

目標再生時間比 Rt は、数式 (4) に基づいて算出されたフレーム毎の音声区間の速度比 $SRs(k)$ を用いて、数式 (6) と表現される。なお、 $Ris(k)$ は、音声含有率を求める算出用フレーム長の最も短いものとする。上述の例で考えると、3 種類の算出用フレーム長のうち最も短いのは、1 分の算出用フレーム長である。

【 数 6 】

$$Rt = \frac{1}{N} \sum_{k=1}^N \left\{ \frac{Ris(k)}{SRs(k)} + \frac{(1-Ris(k))}{SRns} \right\} \quad \cdot \cdot \quad (6)$$

【 0 0 6 0 】

したがって、非音声区間の速度比 $SRns$ は、数式 (6) を整理して数式 (7) となる。

【 数 7 】

$$SRns = \frac{\sum_{k=1}^N (1-Ris(k))}{N * Rt - \sum_{k=1}^N \frac{Ris(k)}{SRs(k)}} \quad \cdot \cdot \quad (7)$$

なお、数式 (7) から分かるように、音声区間の速度比 $SRs(k)$ がフレーム毎に算出されるのに対して、非音声区間の速度比 $SRns$ は、フレームには依存せず (k には依存せず) 一定速度比として算出される。ここで、非音声区間の速度比 $SRns$ の算出例を挙げる。例えば音声区間の速度比が 1 分、5 分、10 分の多重な音声情報を用いて算出されるとする。また、数式 (4) において、 $SRs0$ を 1.5 と、重み係数を $C1 = 1$ 、 $C5 = 10$ 、 $C10 = 20$ とする。このとき、図 6 に示したように、音声情報を多重に用いて算出された音声区間の速度比 $SRs(k)$ は 1.23 ~ 1.68 の範囲の値となる。ここで、目標再生時間比 Rt を例えば 0.5 とする。このとき、非音声区間の速度比 $SRns$ は、数式 (7) より、3.177 となる。つまり、非音声区間の速度比 $SRns$ は、音声区間の速度比 (例えば図 6 に示す 1.23 ~ 1.68) より高速の速度比に設定される。このように、速度比算出部 1 4 は、音声情報記録部 1 3 に記録された音声情報を用いて、音声含有率の変動に応じた音声区間の速度比をフレーム毎に算出し、非音声区間の速度比をフレームに関係なく一定の速度比で算出する。そして、算出された音声区間および非音声区間の速度比の情報は、音声速度変換部 1 5 に出力される。

【 0 0 6 1 】

音声速度変換部 15 は、速度比算出部 14 において算出された音声区間および非音声区間の速度比の情報に基づいて、入力される記録メディアなどに録音された音声信号に対して、速度変換処理を行う。速度変換処理の方法としては、例えば入力される音声信号を時間軸上に圧縮伸長して速度変換を行う方法などがある。しかし、この方法に限定されず、その他の公知方法を用いて速度変換処理が行われてもよい。このように、本実施形態の音声速度変換部 15 において速度変換された音声信号は、音声 / 非音声判別部 11 の判別結果と音声含有率に応じて動的に可変する速度比で変換された音声信号である。

【0062】

次に、図 7 を参照して、本実施形態に係る音声再生装置 1 の処理の流れについて説明する。図 7 は、本実施形態に係る音声再生装置 1 の処理の流れを示すフローチャートである。図 7 において、まず、ユーザが例えば記録メディアに記録された番組全体の記録時間に対して目標とする再生時間を設定する（ステップ S1）。これにより、目標再生時間比 R_t ($0 < R_t < 1$) が設定される。次に、記録メディアなどに録音された番組全体が読み出され、音声 / 非音声判別部 11 において、再生前に番組全体を通して音声区間および非音声区間を判別する（ステップ S2）。そして、音声情報算出部 12 において、ステップ S2 で判別された音声 / 非音声区間の情報に基づいて、複数種類の算出用フレーム長について音声含有率がそれぞれ算出される（ステップ S3）。次に、音声情報算出部 12 において、ステップ S3 で算出された各算出用フレーム長の音声含有率を用いて、音声含有率の平均値および標準偏差がそれぞれ算出される（ステップ S4）。そして、ステップ S3 および S4 で算出された音声情報（音声含有率、音声含有率の平均値および標準偏差）が音声情報記録部 13 に記録される（ステップ S5）。ここまですが再生前に行われる処理である。番組全体を通して音声情報が算出された後、速度変換をする再生が開始される。再生される段階で、速度比算出部 14 は、音声情報記録部 13 に記録された音声情報に基づいて、音声区間の存在の偏りに応じた音声区間の速度比をフレーム毎に算出する（ステップ S6）。次に、速度比算出部 14 において、ステップ S6 で算出された音声区間の速度比と、ステップ S1 で設定された目標再生時間比 R_t とに基づいて、非音声区間の速度比が算出される（ステップ S7）。そして、音声 / 非音声判別部 11 において判別された音声 / 非音声区間の判別情報に対して、フレーム毎の速度比を設定して音声速度変換部 15 へ出力する。ステップ S7 の次に、ステップ S6 および S7 で算出された音声区間および非音声区間の速度比の情報に基づいて、入力される記録メディアなどに録音された音声信号に対して、速度変換処理を行う（ステップ S8）。以上で本実施形態に係る音声再生装置 1 の処理の流れについての説明を終了する。

【0063】

以上のように、本実施形態に係る音声再生装置によれば、音声含有率を音声信号全体に対して算出後、統計値として音声含有率の平均値と標準偏差とを算出して番組中の音声区間の存在の偏り度合いを予め求め、これらの音声情報を用いて音声区間の速度比を算出することで、音声含有率の変動に応じて動的に可変する音声区間の速度比を算出することができる。つまり、本実施形態に係る音声再生装置は、音声が集まる部分には速度比を低減し、音声が集まっていない部分には速度比を増加させる処理を行う。これにより、本実施形態に係る音声再生装置によれば、テレビ番組や映画など全体を通して音声の了解性を保つことができる。また、非音声区間の速度比は、所定の再生時間となるように音声区間の速度比に基づいて一定速度比として算出される。これにより、能率のよい再生速度での再生が可能となる。また、各算出用フレーム長の音声情報を多重して平均値などの統計値を求めることで、音声含有率の長期的な変動や短期的な変動に対して、追従性の高い、より滑らかな速度比の制御を実現することが可能となる。

【0064】

なお、上述した速度比算出部 14 では、各算出用フレーム長の音声情報を多重して音声区間の速度比 $S_{Rs}(k)$ を算出したが、これに限定されない。例えば、音声区間の速度比 $S_{Rs}(k)$ が単独の算出用フレーム長のみ用いて算出されたものでもよい。時間長が長い算出用フレーム長を用いて算出した場合には、算出された音声区間の速度比は、変化

する音声含有率に対して大局的な値であり、より正確性のある値となる。時間長が短い算出用フレーム長を用いて算出した場合には、算出された音声区間の速度比は、変動する音声含有率に対してより追従性のよい値となる。

【0065】

また、上述した速度比算出部14では、音声区間の速度比を算出するための音声情報として、音声含有率 $R_{is_n}(k)$ 、音声含有率の平均値 A_n 、標準偏差 S_n を用いるとしたが、これに限定されない。例えば、上記標準偏差の代わりに、分散や偏差平均など、標準偏差と同等の統計値が用いられてもよい。つまり、音声区間の速度比を算出するための音声情報としては、音声含有率 $R_{is_n}(k)$ 以外に、音声含有率の平均値 A_n および標準偏差と同等の統計値が含まれる。

10

【0066】

また、上述した速度比算出部14では、音声区間の速度比をフレーム毎に算出するとしたが、フレーム内の音声区間1つ1つに対して、さらに文頭、文中、文末などの区分に分け、各区分で速度比を可変してもよい。例えば、ある音声区間の文頭では、速度比算出部14で算出された音声区間の速度比に対してやや速度比を小さくする。そして、文末になるにつれて速度比が大きくなるように設定する。これにより、重要な情報を多く含む文頭部分がユーザにとってより聴き取りやすいものとなる。このように、速度比算出部14は、1つの音声区間中の各区分について速度比を可変するものであってもよい。

【0067】

なお、上述した第1の実施形態で説明した音声/非音声判別部11、音声情報算出部12、速度比算出部14、および音声速度変換部15は、例えば音声信号を入力とし、音声速度変換部15で速度変換された音声信号を出力とする一般的なコンピュータシステム等の情報処理装置で実現可能である。この場合、上述した動作をコンピュータに実行させるプログラムを所定の情報記録媒体に格納し、当該情報記録媒体に格納されたプログラムをコンピュータが読み出して実行することによって、本発明の実現が可能となる。この場合、上記情報処理装置に接続されたキーボードなどの入力部を用いて、ユーザが所望する再生時間を入力する。また、音声情報算出部12で算出される音声情報は、例えば情報処理装置内のハードディスクなどに記録される。また、上記プログラムを格納する情報記録媒体は、例えば、ROMまたはフラッシュメモリのような不揮発性半導体メモリやCD-ROM、DVD、あるいはそれらに類する光学式ディスク状記録媒体である。また、プログラムを他の媒体や通信回線を通じて上記情報処理装置に供給してもかまわない。また、音声情報算出部12で算出される音声情報は情報処理装置内のハードディスクに記録されるとしたが、情報処理装置内のメモリや情報処理装置外の他の記録媒体に記録されてもよい。

20

30

【0068】

(第2の実施形態)

図8を参照して、本発明における第2の実施形態に係る音声再生装置について説明する。図8は、本発明における第2の実施形態に係る音声再生装置2の構成を示すブロック図である。図8において、音声再生装置2は、入力バッファ21、音声/非音声判別部11、音声情報逐次更新部22、速度比算出部14、および音声速度変換部15を有する。

40

【0069】

なお、本実施形態に係る音声再生装置2は、例えばテレビ番組や映画などの音声信号全体が既に記録メディアなどに録音済みであり、録音された音声信号全体のうち一部(所定時間分)の音声信号を一時的に保存しながら逐次的に音声情報を算出して、音声信号の入力に応じて即座に速度変換した再生を行うことを想定した装置である。そのため、本実施形態に係る音声再生装置2は、上述した第1の実施形態に係る音声再生装置1に対して、入力バッファ21を新たに有し、音声情報逐次更新部22において音声情報を逐次更新する点で大きく異なる。以下、異なる点を中心に説明する。また、音声/非音声判別部11、速度比算出部14、および音声速度変換部15は、上述した第1の実施形態と同様であるので、同一の符号を付して、詳細な説明を省略する。

50

【 0 0 7 0 】

記録メディアなどに録音された音声信号が入力バッファ 2 1 に入力される。入力バッファ 2 1 は、入力された音声信号を適宜バッファする。つまり、入力バッファ 2 1 では、音声情報逐次更新部 2 2 で音声情報を逐次更新するために必要な所定時間分の音声信号のデータが一時的に記録される。一時的に保存された所定時間分の音声信号は、音声 / 非音声判別部 1 1 および音声速度変換部 1 5 にそれぞれ出力される。音声 / 非音声判別部 1 1 は、入力された所定時間分の音声信号に対して音声区間および非音声区間を判別する。音声 / 非音声判別部 1 1 において判別された音声 / 非音声区間の情報は、音声情報逐次更新部 2 2 および速度比算出部 1 4 にそれぞれ出力される。

【 0 0 7 1 】

音声情報逐次更新部 2 2 は、音声 / 非音声区間の判別情報に基づいて音声情報を逐次更新する。なお、第 1 の実施形態では数式 (3) および数式 (4) において、音声含有率 $R_{is_n}(k)$ を音声信号全体について一旦算出した後に、統計値である音声含有率の平均値 A_n および標準偏差 S_n を算出していた。これに対し、本実施形態では、音声信号の入力に応じて即座に速度変換した再生を行うために、統計値である上記音声含有率の平均値 A_n および標準偏差 S_n の初期値を予め記録部 (図示しない) などにそれぞれ記録設定して、当該統計値を記録部などに逐次記録しながら更新していく。以下、音声情報である音声含有率の平均値および標準偏差の更新方法について説明する。

【 0 0 7 2 】

音声含有率の平均値 A_n は、起動に際して初期値が設定される。そして、音声含有率の平均値 A_n は、音声信号が入力されるフレーム毎に逐次更新される。上記初期値は、例えば再生する番組のジャンルなどによって異なり、当該ジャンルに合わせて適宜設定される。例えば、頻繁にアナウンサが話す機会の多いテレビのニュース番組などの場合は、音声含有率の平均値が 8 5 % 程度となる。また、話者の話す機会が少ない様々な映像シーンを多用するドキュメンタリ番組などの場合は、音声含有率の平均値が 5 0 % 程度になる。

【 0 0 7 3 】

ここで、入力バッファ 2 1 に記録される音声信号の所定時間分を例えば上述した算出用フレーム長 (n 分) とする。そして、入力バッファ 2 1 は、算出用フレーム長 (n 分) 分の音声信号を確保しながら、例えば 1 フレーム分の音声信号を順次記録更新していくとする。また、音声情報逐次更新部 2 2 は、例えば音声 / 非音声判別部 1 1 で 1 フレーム分の音声 / 非音声区間が判別される毎に、音声含有率の平均値 A_n の逐次更新を行うとする。この場合、音声含有率の平均値 A_n はフレーム毎に更新され、 k フレーム目の逐次更新される音声含有率の平均値の更新値 (以下、音声含有率の更新平均値とする) を $A_n(k)$ とする。このとき、音声含有率の更新平均値 $A_n(k)$ は、数式 (8) で表現される。

【 数 8 】

$$A_n(k) = \alpha_1 * A_n(k-1) + \beta_1 * R_{is_n}(k) \quad \cdots (8)$$

なお、数式 (8) において、 α_1 および β_1 は音声含有率の更新平均値 $A_n(k)$ の更新速度を規定するパラメータである。すなわち、 α_1 の値が大きいほど k フレームの 1 つ前のフレームの更新平均値 $A_n(k-1)$ の占める割合が高くなり、更新平均値 $A_n(k)$ の更新速度が緩やかになる。また、 β_1 の値が大きいほど k フレームの音声含有率 $R_{is_n}(k)$ の占める割合が高くなり、更新平均値 $A_n(k)$ の更新速度が速くなる。数値例としては、例えば $\alpha_1 = 0.98$ 、 $\beta_1 = 0.02$ としてもよい。

【 0 0 7 4 】

また、標準偏差 S_n も上記音声含有率の平均値と同様に、起動に際して初期値が設定される。そして、標準偏差 S_n は、フレーム毎に逐次更新される。上記初期値は、音声含有率の平均値 A_n と同様に、例えば再生する番組のジャンルなどによって異なり、当該ジャンルに合わせて適宜設定される。具体的には標準偏差 S_n は、上記初期値と、更新平均値 $A_n(k)$ と、 k フレームの音声含有率 $R_{is_n}(k)$ とを用いて更新される。ここで、 k フレーム目の標準偏差の更新値を $S_n(k)$ とすると、標準偏差の更新値 $S_n(k)$

は、数式(9)で表現される。

【数9】

$$S_n(k) = \sqrt{\alpha^2 * (S_n(k-1))^2 + \beta^2 * (R_{is_n}(k) - A_n(k))^2} \quad \cdot \cdot \quad (9)$$

なお、数式(9)において、 α^2 および β^2 は標準偏差の更新値 $S_n(k)$ の更新速度を規定するパラメータである。数値例としては、例えば $\alpha^2 = 0.98$ 、 $\beta^2 = 0.02$ としてもよい。

【0075】

次に、速度比算出部14は、音声含有率 $R_{is_n}(k)$ と、フレーム毎に更新された音声含有率の更新平均値 $A_n(k)$ および標準偏差の更新値 $S_n(k)$ とに基づいて、上述した第1の実施形態と同様に、数式(3)～数式(5)に基づいて音声区間の速度比 $S_{Rs}(k)$ を算出する。また、速度比算出部14は、算出した音声区間の速度比 $S_{Rs}(k)$ と目標再生時間比 R_t とに基づいて非音声区間の速度比 S_{Rns} を算出する。そして、速度比算出部14は、音声/非音声判別部11から入力される音声/非音声区間の判別情報に対して、フレーム毎の速度比を設定して音声速度変換部15へ出力する。音声速度変換部15は、速度比算出部14において算出された音声区間および非音声区間の速度比の情報に基づいて、入力バッファ21から入力される音声信号に対してフレーム毎に逐次速度変換処理を行う。

【0076】

以上のように、本実施形態に係る音声再生装置2は、統計値である音声含有率の平均値および標準偏差を逐次更新する。これにより、本実施形態に係る音声再生装置2は、音声情報を番組全体に対して事前に算出することなく、音声信号の入力に応じて即時に速度変換処理を行うことができる。

【0077】

なお、上述した第2の実施形態で説明した音声再生装置2は、音声/非音声判別部11、音声情報逐次更新部22、速度比算出部14、および音声速度変換部15は、例えば音声信号を入力とし、音声速度変換部15で速度変換された音声信号を出力とする一般的なコンピュータシステム等の情報処理装置で実現可能である。この場合、上述した動作をコンピュータに実行させるプログラムを所定の情報記録媒体に格納し、当該情報記録媒体に格納されたプログラムをコンピュータが読み出して実行することによって、本発明の実現が可能となる。また、上記情報処理装置に接続されるキーボードなどの入力部において、ユーザが所望する再生時間や上述した初期値を入力する。また、入力バッファ21は、例えば情報処理装置内のハードディスク内で構成される。また、上記プログラムを格納する情報記録媒体は、例えば、ROMまたはフラッシュメモリのような不揮発性半導体メモリやCD-ROM、DVD、あるいはそれらに類する光学式ディスク状記録媒体である。また、プログラムを他の媒体や通信回線を通じて上記情報処理装置に供給してもかまわない。また、入力バッファ21を例えば情報処理装置内のハードディスク内で構成されたとしたが、情報処理装置内のメモリや情報処理装置外の他の記録媒体で構成されてもよい。

【0078】

(第3の実施形態)

図9を参照して、本発明における第3の実施形態に係る音声録音再生装置について説明する。図9は、本発明における第3の実施形態に係る音声録音再生装置3の構成を示すブロック図である。図9において、音声録音再生装置3は、音声/非音声判別部11、情報記録部31、音声情報算出部12、音声情報記録部13、速度比算出部14、および音声速度変換部15を有する。

【0079】

なお、本実施形態に係る音声録音再生装置3は、情報記録部31に音声記録して再生する音声録音再生装置であって、入力される音声信号を情報記録部31に記録すると同時に、音声/非音声判別部11で判別された音声区間や非音声区間の情報も情報記録部31

10

20

30

40

50

に記録することを特徴とする装置である。以下、この特徴を中心に説明する。また、音声／非音声判別部 1 1、音声情報算出部 1 2、音声情報記録部 1 3、速度比算出部 1 4、および音声速度変換部 1 5 は、上述した第 1 の実施形態と同様であるので、同一の符号を付して、詳細な説明を省略する。

【0080】

録音対象となる音声信号が音声／非音声判別部 1 1 および情報記録部 3 1 にそれぞれ入力される。音声／非音声判別部 1 1 は、入力された音声信号に対して音声区間および非音声区間を判別する。音声／非音声判別部 1 1 において判別された音声／非音声区間の判別情報は、情報記録部 3 1 に出力される。情報記録部 3 1 において、入力された録音対象である音声信号と音声／非音声区間の判別情報とがそれぞれ記録される。

10

【0081】

音声情報算出部 1 2 は、情報記録部 3 1 に記録された音声信号全体についての音声／非音声区間の情報を読み出して、音声情報を算出する。具体的には、音声情報算出部 1 2 は、記録された音声信号全体を通して音声含有率を算出した後に、音声含有率の平均値および標準偏差を算出する。そして、音声情報算出部 1 2 で算出された音声含有率、音声含有率の平均値、および標準偏差は、音声情報記録部 1 3 にそれぞれ記録される。

【0082】

そして、再生される段階において、速度比算出部 1 4 は、音声情報記録部 1 3 に記録された音声情報を用いて、音声含有率の変動に応じた音声区間の速度比をフレーム毎に算出する。また、速度比算出部 1 4 は、音声区間の速度比と目標再生時間比 R_t とに基づいて非音声区間の速度比を算出する。そして、記録された音声／非音声区間の判別情報に対して、フレーム毎の速度比を設定して音声速度変換部 1 5 へ出力する。音声速度変換部 1 5 は、速度比算出部 1 4 において算出された音声区間および非音声区間の速度比の情報に基づいて、情報記録部 3 1 に記録された音声信号に対して速度変換処理を行う。

20

【0083】

以上のように、本実施形態に係る音声録音再生装置 3 は、入力される音声信号を情報記録部 3 1 に記録するとともに、音声／非音声判別部 1 1 で判別された音声区間や非音声区間の情報も情報記録部 3 1 に記録している。これにより、本実施形態に係る音声録音再生装置 3 によれば、音声信号全体を記録した段階で音声信号全体についての音声区間や非音声区間の判別が終了しているため、再生前に行われる音声情報の算出時間を短縮することができる。

30

【0084】

なお、上述した情報記録部 3 1 において、音声／非音声判別部 1 1 で判別された音声区間や非音声区間の判別情報に加え、さらに音声情報算出部 1 2 で算出された音声情報が記録されてもよい。この場合、図 1 0 に示すように、音声情報記録部 1 3 は省略される。図 1 0 は、情報記録部 3 1 に音声区間や非音声区間の情報と音声情報とを記録する音声録音再生装置 4 の構成を示すブロック図である。図 1 0 において、音声録音再生装置 4 は、音声／非音声判別部 1 1、情報記録部 3 1、音声情報算出部 1 2、速度比算出部 1 4、および音声速度変換部 1 5 を有する。

【0085】

図 1 0 において、情報記録部 3 1 では、入力された録音対象である音声信号と、音声／非音声判別部 1 1 において判別された音声／非音声区間の情報と、音声情報算出部 1 2 で算出された音声情報とがそれぞれ記録される。つまり、音声録音再生装置 4 は、記録とともに音声／非音声区間の判別情報および音声情報が情報記録部 3 1 に記録される。これにより、音声録音再生装置 4 によれば、記録後において再生時間が入力されれば、即時に速度比を算出することができる。その結果、音声録音再生装置 4 は、速度変換した再生音声を短時間で出力することができる。

40

【0086】

なお、上述した第 3 の実施形態で説明した音声／非音声判別部 1 1、音声情報算出部 1 2、音声情報記録部 1 3、速度比算出部 1 4、および音声速度変換部 1 5 は、例えば音声

50

信号を入力とし、音声速度変換部 15 で速度変換された音声信号を出力とする一般的なコンピュータシステム等の情報処理装置で実現可能である。この場合、上述した動作をコンピュータに実行させるプログラムを所定の情報記録媒体に格納し、当該情報記録媒体に格納されたプログラムをコンピュータが読み出して実行することによって、本発明の実現が可能となる。また、上記情報処理装置に接続されるキーボードなどの入力部において、ユーザが所望する再生時間が入力される。また、情報記録部 31 および音声情報記録部 13 は、例えば情報処理装置内のハードディスク内で構成される。また、上記プログラムを格納する情報記録媒体は、例えば、ROM またはフラッシュメモリのような不揮発性半導体メモリや CD-ROM、DVD、あるいはそれらに類する光学式ディスク状記録媒体である。また、プログラムを他の媒体や通信回線を通じて上記情報処理装置に供給してもかま

10

【0087】

また、上述した第 1 ～ 第 3 の実施形態で説明した音声 / 非音声判別部 11、音声情報算出部 12、音声情報記録部 13、速度比算出部 14、音声情報逐次更新部 22 および音声速度変換部 15 は、例えば音声信号、再生時間情報、および上述した初期値などを入力とし、音声速度変換部 15 で速度変換された音声信号を出力とする集積回路でも実現可能である。この場合、第 1 の実施形態における音声情報記録部 13、第 2 の実施形態における入力バッファ 21、第 3 の実施形態における音声情報記録部 13 および情報記録部 31 は、例えば集積回路内のメモリで構成される。そして、上述した機能を果たす電気回路を 1 つの小型パッケージに集積して、音声信号の処理等を行う音声信号処理回路 DSP (Digital Signal Processor) 等を構成することによって、本発明の実現が可能となる。なお、第 1 の実施形態における音声情報記録部 13、第 2 の実施形態における入力バッファ 21、第 3 の実施形態における音声情報記録部 13 および情報記録部 31 は、上記集積回路とは別の他の記録媒体で構成されてもよい。

20

【産業上の利用可能性】

【0088】

本発明に係る音声再生装置、音声録音再生装置、およびそれらの方法、記録媒体、および集積回路は、音声含有率の変動に応じた最適な速度比制御を行って、より聴き取りやすい再生を実現する DVD プレーヤ、HDD プレーヤ、CD プレーヤ等にも有用である。

30

【図面の簡単な説明】

【0089】

【図 1】本発明における第 1 の実施形態に係る音声再生装置 1 の構成を示すブロック図

【図 2】算出用フレーム長が 1 分のときの音声含有率 $R_{is_1}(k)$ の算出例を示す図

【図 3】算出用フレーム長が 5 分のときの音声含有率 $R_{is_5}(k)$ の算出例を示す図

【図 4】算出用フレーム長が 10 分のときの音声含有率 $R_{is_10}(k)$ の算出例を示す図

【図 5】各算出用フレーム長の音声含有率の平均値および標準偏差の算出結果を示す図

【図 6】多重の音声情報を用いた速度比の算出結果の一例を示す図

40

【図 7】本実施形態に係る音声再生装置 1 の処理の流れを示すフローチャート

【図 8】本発明における第 2 の実施形態に係る音声再生装置 2 の構成を示すブロック図

【図 9】本発明における第 3 の実施形態に係る音声録音再生装置 3 の構成を示すブロック図

【図 10】情報記録部 31 に音声区間や非音声区間の情報と音声情報とを記録する音声録音再生装置 4 の構成を示すブロック図

【図 11】従来の音声再生装置 9 の構成を示すブロック図

【符号の説明】

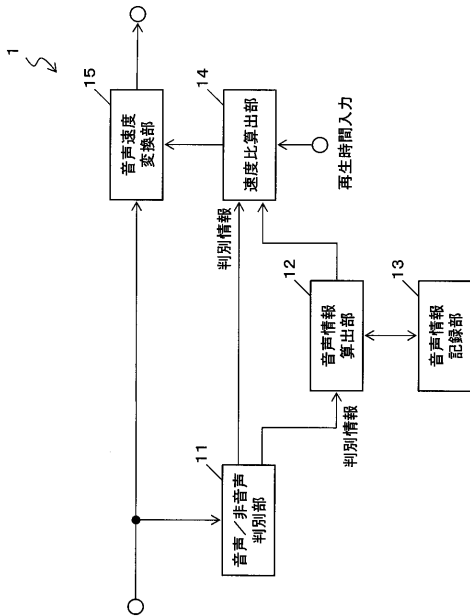
【0090】

1、2 音声再生装置

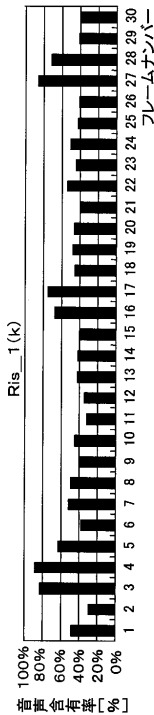
50

- 3、4 音声録音再生装置
- 1 1 音声 / 非音声判別部
- 1 2 音声情報算出部
- 1 3 音声情報記録部
- 1 4 速度比算出部
- 1 5 音声速度変換部
- 2 1 入力バッファ
- 2 2 音声情報逐次更新部
- 3 1 情報記録部

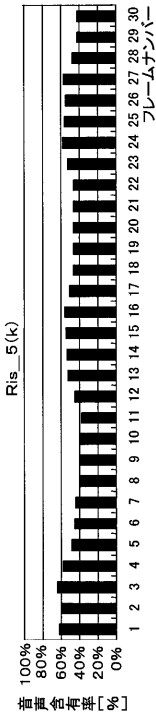
【図 1】



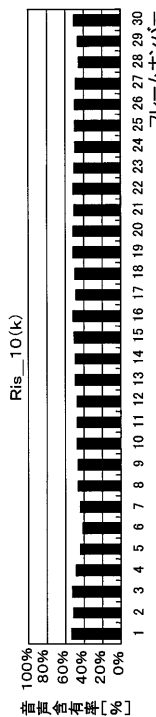
【図 2】



【図 3】



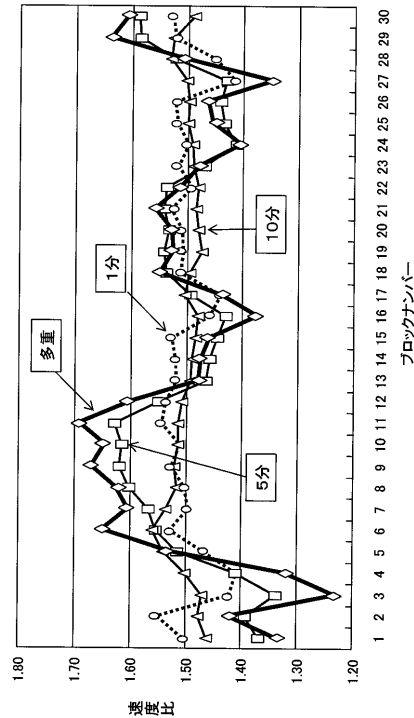
【図 4】



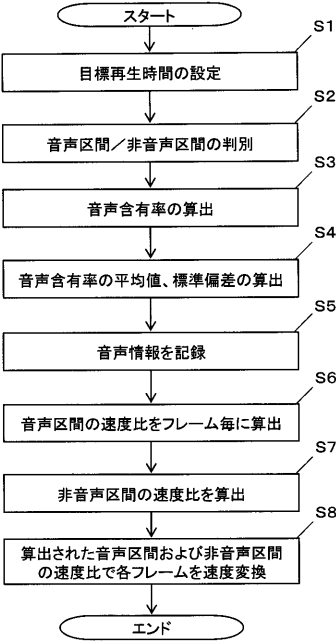
【図 5】

| 算出用フレーム長[分] | 平均値 | 標準偏差 |
|-------------|------------|------------|
| 1 | A1: 0.506 | S1: 0.161 |
| 5 | A5: 0.498 | S5: 0.073 |
| 10 | A10: 0.488 | S10: 0.028 |

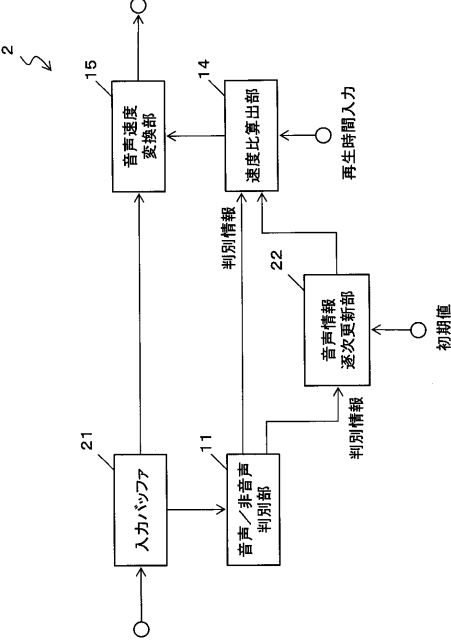
【図 6】



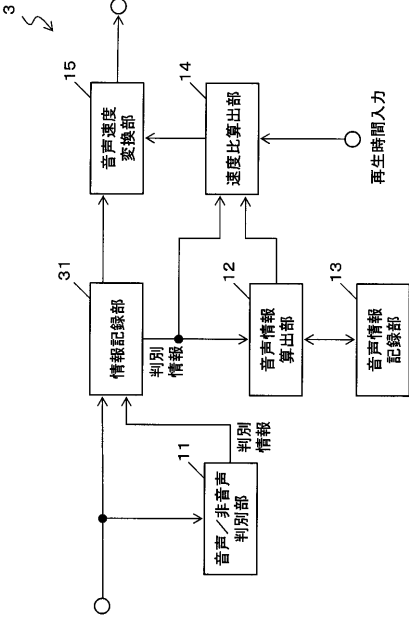
【図 7】



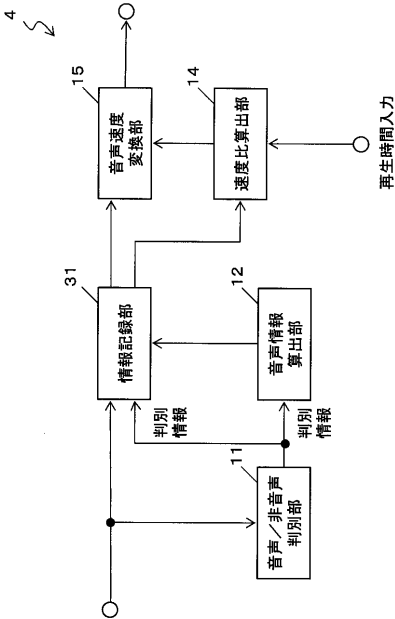
【図 8】



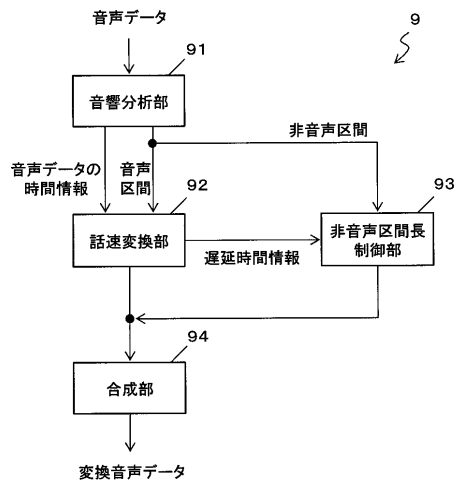
【図 9】



【図 10】



【図 11】



フロントページの続き

審査官 安田 勇太

(56)参考文献 特開平 0 4 - 3 6 7 8 9 8 (J P , A)
特開 2 0 0 1 - 2 2 2 3 0 0 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)
G 1 0 L 2 1 / 0 4
G 1 0 L 1 9 / 0 0
G 1 0 L 1 1 / 0 0