

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号  
特許第4965371号  
(P4965371)

(45) 発行日 平成24年7月4日(2012.7.4)

(24) 登録日 平成24年4月6日(2012.4.6)

(51) Int.Cl.

F I

G 1 O L 21/04 (2006.01)

G 1 O L 19/00 (2006.01)

G 1 O L 21/04 1 1 O Z

G 1 O L 19/00 3 1 2 E

請求項の数 3 (全 45 頁)

(21) 出願番号	特願2007-195708 (P2007-195708)	(73) 特許権者	000005821
(22) 出願日	平成19年7月27日 (2007.7.27)		パナソニック株式会社
(65) 公開番号	特開2008-58956 (P2008-58956A)		大阪府門真市大字門真1006番地
(43) 公開日	平成20年3月13日 (2008.3.13)	(74) 代理人	110001276
審査請求日	平成22年7月27日 (2010.7.27)		特許業務法人 小笠原特許事務所
(31) 優先権主張番号	特願2006-208689 (P2006-208689)	(72) 発明者	前田 芽衣子
(32) 優先日	平成18年7月31日 (2006.7.31)		大阪府門真市大字門真1006番地 松下
(33) 優先権主張国	日本国 (JP)		電器産業株式会社内
		(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	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

(57) 【特許請求の範囲】

【請求項 1】

入力されるオーディオ信号を所定の再生時間で再生するために前記オーディオ信号に速度変換処理を適用して圧縮伸長するための目標圧伸比が設定される音声再生装置であって、

前記オーディオ信号に対して、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手段と、

前記判別手段において判別された前記音声区間に基づいて、前記オーディオ信号中に含まれる音声区間の時間比率を示す音声含有率を算出する音声含有率算出手段と、

(1) 前記音声含有率と、

(2) 前記目標圧伸比と、

(3) 前記音声区間の平均速度比と前記非音声区間の平均速度比とが満たすべき算出条件を示す複数の音声速度比算出パターンを有する音声速度比算出分布と

を用いて、前記音声区間の平均速度比と、前記非音声区間の平均速度比とを算出し、算出したそれぞれの平均速度比を速度比条件として設定する速度比条件設定手段と、

前記速度比条件設定手段により設定された前記速度比条件に基づいて、前記音声区間及び前記非音声区間のそれぞれの平均圧伸比に対する圧伸比変化量の和が、前記音声区間又は前記非音声区間内においてゼロとなるように、前記音声区間の平均速度比及び前記非音声区間の平均速度比を決定する速度比決定手段と、

前記オーディオ信号に含まれる音声区間及び非音声区間の再生速度を前記速度比決定手

段で決定された前記音声区間の速度比および前記非音声区間の速度比に基づいてそれぞれ変換する速度変換手段とを備え、

前記複数の音声速度比算出パターンは、互いに異なる前記算出条件を有し、

前記速度比条件設定手段は、前記音声含有率と前記目標圧伸比に応じて、前記複数の音声速度比算出パターンの中から一つの音声速度比算出パターンが定まる前記音声速度比算出分布に基づいて、前記速度比条件を設定することを特徴とし、

前記判別手段は、予め定められた特定音を含む特定イベント区間と、前記特定音を含まない非特定イベント区間とを判別し、前記非特定イベント区間に対して前記音声区間と前記非音声区間とを判別し、

前記判別手段により判別された前記特定イベント区間に基づいて、前記オーディオ信号中に含まれる特定イベント区間の時間比率を示す特定イベント含有率を算出する特定イベント含有率算出手段をさらに備え、

前記速度比条件設定手段は、

(4) 前記特定イベント含有率と、

(5) 前記目標圧伸比と、

(6) 前記特定イベント区間の平均速度比と前記非特定イベント区間の平均速度比が満たすべき算出条件を示す複数の特定イベント速度比算出パターンを有する特定イベント速度比算出分布とを用いて、

前記特定イベント区間の平均速度比と前記非特定イベント区間の平均速度比を算出し、算出したそれぞれの平均速度比を前記速度比条件としてさらに設定し、

前記複数の特定イベント速度比算出パターンは、互いに異なる前記算出条件を有し、

前記速度比条件設定手段は、前記特定イベント含有率と前記目標圧伸比に応じて、前記複数の特定イベント速度比算出パターンの中から一つの特定イベント速度比算出パターンが定まる前記特定イベント速度比算出分布に基づいて、前記速度比条件を設定することを特徴とする、音声再生装置。

#### 【請求項2】

入力されるオーディオ信号を所定の再生時間で再生するために前記オーディオ信号に速度変換処理を適用して圧縮伸長するための目標圧伸比が設定される音声再生装置であって、

前記オーディオ信号に対して、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手段と、

前記判別手段において判別された前記音声区間に基づいて、前記オーディオ信号中に含まれる音声区間の時間比率を示す音声含有率を算出する音声含有率算出手段と、

(1) 前記音声含有率と、

(2) 前記目標圧伸比と、

(3) 前記音声区間の平均速度比と前記非音声区間の平均速度比とが満たすべき算出条件を示す複数の音声速度比算出パターンを有する音声速度比算出分布と

を用いて、前記音声区間の平均速度比と、前記非音声区間の平均速度比とを算出し、算出したそれぞれの平均速度比を速度比条件として設定する速度比条件設定手段と、

前記速度比条件設定手段により設定された前記速度比条件に基づいて、前記音声区間及び前記非音声区間のそれぞれの平均圧伸比に対する圧伸比変化量の和が、前記音声区間又は前記非音声区間内においてゼロとなるように、前記音声区間の平均速度比及び前記非音声区間の平均速度比を決定する速度比決定手段と、

前記オーディオ信号に含まれる音声区間及び非音声区間の再生速度を前記速度比決定手段で決定された前記音声区間の速度比および前記非音声区間の速度比に基づいてそれぞれ変換する速度変換手段とを備え、

前記複数の音声速度比算出パターンは、互いに異なる前記算出条件を有し、

前記速度比条件設定手段は、前記音声含有率と前記目標圧伸比に応じて、前記複数の音声速度比算出パターンの中から一つの音声速度比算出パターンが定まる前記音声速度比算出分布に基づいて、前記速度比条件を設定することを特徴とし、

10

20

30

40

50

前記オーディオ信号の内容を示すコンテンツ情報を取得するコンテンツ取得手段を、さらに備え、

前記速度比条件設定手段は、

前記コンテンツ情報が示す内容に応じた複数の前記音声速度比算出分布を有し、

前記コンテンツ取得手段により取得された前記コンテンツ情報に基づいて前記複数の音声速度比算出分布の中から一つを選択し、

選択した前記音声速度比算出分布に基づいて、前記速度比条件を設定する

ことを特徴とし、

前記判別手段は、予め定められた特定音を含む特定イベント区間と、前記特定音を含まない非特定イベント区間とを判別し、前記非特定イベント区間に対して前記音声区間と前記非音声区間とを判別し、

前記判別手段により判別された前記特定イベント区間に基づいて、前記オーディオ信号中に含まれる特定イベント区間の時間比率を示す特定イベント含有率を算出する特定イベント含有率算出手段をさらに備え、

前記速度比条件設定手段は、

(4) 前記特定イベント含有率と、

(5) 前記目標圧伸比と、

(6) 前記特定イベント区間の平均速度比と前記非特定イベント区間の平均速度比が満たすべき算出条件を示す複数の特定イベント速度比算出パターンを有する特定イベント速度比算出分布とを用いて、

前記特定イベント区間の平均速度比と前記非特定イベント区間の平均速度比を算出し、算出したそれぞれの平均速度比を前記速度比条件としてさらに設定し、

前記複数の特定イベント速度比算出パターンは、互いに異なる前記算出条件を有し、

前記速度比条件設定手段は、前記特定イベント含有率と前記目標圧伸比に応じて、前記複数の特定イベント速度比算出パターンの中から一つの特定イベント速度比算出パターンが定まる前記特定イベント速度比算出分布に基づいて、前記速度比条件を設定する

ことを特徴とする、音声再生装置。

#### 【請求項3】

前記判別手段は、複数種類の前記特定音をそれぞれ含む複数の前記特定イベント区間と、前記複数種類の特定音を全て含まない非特定イベント区間とを判別し、

前記速度比条件設定手段は、

前記複数の特定イベント区間に対応する複数の前記特定イベント含有率にそれぞれ対応する複数の前記特定イベント速度比算出分布を有し、

判別された特定イベント区間に対応する前記複数の特定イベント速度比算出分布に基づいて、前記速度比条件を設定する

ことを特徴とする、請求項1または2に記載の音声再生装置。

#### 【発明の詳細な説明】

#### 【技術分野】

#### 【0001】

本発明は、音声再生装置に関し、より特定的には、オーディオ信号の再生速度を変えて再生する音声再生装置に関する。

#### 【背景技術】

#### 【0002】

従来、オーディオ信号の再生速度を変えて再生する音声再生装置として、音声を含む音声区間の再生速度と、音声を含まない非音声区間の再生速度とを別々に変える音声再生装置が提案されている（例えば、特許文献1参照）。以下、図33を参照して従来の音声再生装置について説明する。図33は、従来の音声再生装置の構成を示したブロック図である。

#### 【0003】

図33に示す従来の音声再生装置において、ユーザがオーディオ信号全体の再生時間に

10

20

30

40

50

対して目標時間を設定する。この目標時間は、オーディオ信号全体を等倍の再生速度で再生したときの再生時間よりも短い時間とする。音響分析部 9 1 は、入力されるオーディオ信号を音声区間及び非音声区間に分離する。速度変換部 9 2 は、一定時間長以上の非音声区間に挟まれた音声区間のオーディオ信号に対して、その冒頭部分が所定の再生速度よりも遅くなり、かつ末尾に向けて次第に所定の再生速度に戻るよう速度変換を行っている。ここで、話速変換部 9 2 における上記速度変換処理によって、音声区間の再生時間が長くなり、結果的にオーディオ信号全体の再生時間が目標時間に対して遅延してしまうという問題があった。そこで、非音声区間長制御部 9 3 は、速度変換部 9 2 から出力される遅延時間情報を参照して、非音声区間に対して当該遅延時間を短くするための処理を行う。具体的には、非音声区間長制御部 9 3 は、非音声区間を削除したり、圧縮したりする処理を行って、遅延時間を短くしている。速度変換部 9 2 で速度変換された音声区間のオーディオ信号と、非音声区間長制御部 9 3 で処理された非音声区間のオーディオ信号は、合成部 9 4 で合成され、合成部 9 4 から出力される。

10

【特許文献 1】特開 2 0 0 1 - 2 2 2 3 0 0 号公報（第 1 - 2 頁、図 1）

【発明の開示】

【発明が解決しようとする課題】

【0 0 0 4】

ここで、入力されるオーディオ信号に含まれる音声区間の比率は、入力されるオーディオ信号に応じて異なっている。しかしながら、従来の音声再生装置では、音声区間が含まれる比率に関わらず、音声区間に対しては上記速度変換処理を一律に行い、非音声区間に対しては目標時間を達成するための削除や圧縮を行っている。したがって、例えば入力されるオーディオ信号が音声区間を多く含む信号である場合、音声区間に対しては上記速度変換処理が一律に行われるので、話速変換部 9 2 において生じる遅延時間が長くなってしまふ。そして遅延時間が長くなれば、非音声区間に対する区間の削除量や圧縮量も大きくなってしまい、情報の欠落が大きくなったり、再生が聞き取り難くなったりする。このように従来の音声再生装置では、目標時間を達成しつつ、入力されるオーディオ信号に応じた適切な速度変換を行うことができなかった。

20

【0 0 0 5】

それ故、本発明の目的は、目標時間を達成しつつ、入力されるオーディオ信号に応じた適切な速度変換を行うことが可能な音声再生装置を提供することを目的とする。

30

【課題を解決するための手段】

【0 0 0 6】

第 1 の発明は、入力されるコンテンツのオーディオ信号の再生速度を変えて所定の再生時間で再生する音声再生装置であって、オーディオ信号に対して、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手段と、判別手段において判別された判別結果に基づいて、オーディオ信号に含まれる音声区間の比率を示す音声含有率を算出する音声含有率算出手段と、オーディオ信号の再生時間が所定の再生時間となるように、オーディオ信号に予め設定された再生速度に対する音声区間及び非音声区間の速度比を音声含有率に基づいてそれぞれ算出する速度比算出手段と、オーディオ信号を入力とし、当該オーディオ信号に含まれる音声区間及び非音声区間の再生速度を速度比に基づいてそれぞれ変換する速度変換手段とを備える。

40

【0 0 0 7】

第 2 の発明は、上記第 1 の発明において、速度比算出手段は、オーディオ信号に予め設定された再生速度で再生される再生時間を所定の再生時間に圧伸する比率を示す圧伸比と、音声含有率と、音声区間の平均速度比を示す音声平均速度比及び非音声区間の平均速度比を示す非音声平均速度比の算出方法との対応を示す対応情報を用いて、音声平均速度比及び非音声平均速度比をそれぞれ算出して速度比条件として設定する速度比条件設定手段と、音声区間が細分化された各区間における速度比を音声平均速度比に基づく速度比に決定するとともに、非音声区間が細分化された各区間における速度比を非音声平均速度比に基づく速度比に決定して音声区間及び非音声区間の速度比をそれぞれ算出する速度比決定

50

手段とを有する。

【 0 0 0 8 】

第 3 の発明は、上記第 2 の発明において、音声再生装置は、判別手段において判別された各音声区間の開始時刻から終了時刻までの時間を音声区間長としてそれぞれ算出する音声区間長算出手段をさらに備え、速度比条件設定手段は、音声平均速度比に応じた音声区間の終了時刻における終了速度比を速度比条件としてさらに設定し、速度比決定手段は、判別手段において判別された各音声区間に対して、音声平均速度比、音声区間長、及び終了速度比に基づいて音声区間の各区間における速度比を決定する。

【 0 0 0 9 】

第 4 の発明は、上記第 2 の発明において、速度比決定手段は、判別手段において判別された各音声区間に対して、音声区間の開始時刻から経過した時間を音声区間長で除算して得られる経過割合に応じて音声区間の各区間における速度比をそれぞれ決定する。

10

【 0 0 1 0 】

第 5 の発明は、上記第 2 の発明において、速度比決定手段は、音声区間の開始時刻から時間が経過するにつれて再生速度が速くなるように、音声区間の各区間における速度比を決定する。

【 0 0 1 1 】

第 6 の発明は、上記第 2 の発明において、速度比条件設定手段は、音声平均速度比及び非音声平均速度比の算出方法を少なくとも 1 種類含む対応情報であって、オーディオ信号によって構成されるコンテンツの種類に応じて異なる対応を示す対応情報を用いて、音声平均速度比及び非音声平均速度比を算出する。

20

【 0 0 1 2 】

第 7 の発明は、上記第 2 の発明において、速度比条件設定手段は、音声平均速度比及び非音声平均速度比がユーザによって指定された範囲内となるように対応情報を作成し、当該対応情報を用いて音声平均速度比及び非音声平均速度比を算出する。

【 0 0 1 3 】

第 8 の発明は、上記第 2 の発明において、対応情報は、音声含有率の大きさに応じて当該音声含有率と音声平均速度比及び非音声平均速度比の算出方法とが異なる対応を示す情報である。

【 0 0 1 4 】

30

第 9 の発明は、上記第 2 の発明において、速度比条件設定手段は、音声平均速度比及び非音声平均速度比の算出方法を少なくとも 1 種類含む対応情報であって、ユーザの使用目的に応じて異なる対応を示す対応情報を用いて、音声平均速度比及び非音声平均速度比を算出する。

【 0 0 1 5 】

第 10 の発明は、上記第 2 の発明において、対応情報は、圧伸比の大きさに応じて当該圧伸比と音声平均速度比及び非音声平均速度比の算出方法とが異なる対応を示す情報である。

【 0 0 1 6 】

第 11 の発明は、上記第 1 の発明において、コンテンツ全体を構成するオーディオ信号と、当該コンテンツ全体を構成するオーディオ信号に対して判別手段において判別された判別結果とを予め蓄積する蓄積手段をさらに備え、音声含有率算出手段は、蓄積手段に予め蓄積された判別結果に基づいて、コンテンツ全体を構成するオーディオ信号に含まれる音声区間の比率を示す音声含有率を算出する。

40

【 0 0 1 7 】

第 12 の発明は、上記第 1 の発明において、音声含有率算出手段は、判別手段において過去に判別された判別結果に基づいて、速度比算出手段が算出するときに用いる音声含有率を逐次算出する。

【 0 0 1 8 】

第 13 の発明は、上記第 12 の発明において、音声含有率算出手段は、速度比算出手段

50

が算出するときから第１の所定時間分だけ過去に判別された判別結果に基づいて、速度比算出手段が算出するときに用いる音声含有率を第１の所定時間以下の時間である第２の所定時間毎に逐次算出し、音声再生装置は、速度変換手段に入力されるデータ量及び速度変換手段から出力されるデータ量と、オーディオ信号に予め設定された再生速度で再生される再生時間を所定の再生時間に圧伸する比率を示す圧伸比とに基づいて、第２の所定時間毎の圧伸比を逐次算出する圧伸比算出手段をさらに備え、速度比算出手段は、第２の所定時間毎の圧伸比と、第２の所定時間毎の音声含有率と、音声区間の平均速度比を示す音声平均速度比及び非音声区間の平均速度比を示す非音声平均速度比の算出方法との対応を示す対応情報を用いて、音声平均速度比及び非音声平均速度比をそれぞれ算出し、算出した音声平均速度比及び非音声平均速度比を速度比条件として第２の所定時間毎に逐次設定する速度比条件設定手段と、第２の所定時間内に含まれる音声区間及び非音声区間に対して、音声区間が細分化された各区間における速度比を音声平均速度比に基づく速度比に決定するとともに、非音声区間が細分化された各区間における速度比を非音声平均速度比に基づく速度比に決定して、第２の所定時間毎に逐次設定される速度比条件に基づいて音声区間及び非音声区間の速度比をそれぞれ算出する速度比決定手段とを有し、速度変換手段は、オーディオ信号に含まれる音声区間及び非音声区間の再生速度を速度比決定手段において算出された音声区間及び非音声区間の速度比に基づいてそれぞれ変換する。

10

**【００１９】**

第１４の発明は、上記第１３の発明において、音声再生装置は、速度比決定手段において決定される音声区間の各区間における速度比が示す第２の所定時間毎の変化を抑制するための統計量を算出する統計量算出手段をさらに備え、速度比条件設定手段は、統計量、第２の所定時間毎の音声含有率、及び第２の所定時間毎の圧伸比に基づいて音声平均速度と当該音声平均速度に基づく音声区間の終了速度比とを算出し、算出した音声平均速度及び終了速度比を速度比条件として第２の所定時間毎に逐次設定する。

20

**【００２０】**

第１５の発明は、上記第１４の発明において、統計量算出手段は、コンテンツの開始時刻から速度比算出手段が算出するときまでの判別結果に基づいて、コンテンツの開始時刻から速度比算出手段が算出するときまでに含まれる音声区間の比率を示す音声含有率を統計量として算出する。

**【００２１】**

30

第１６の発明は、上記第１３の発明において、音声再生装置は、判別手段によって過去に判別された判別結果に基づいて、速度比算出手段が算出するときに用いる音声区間の開始時刻から終了時刻までの時間である音声区間長を逐次算出する音声区間長算出手段をさらに備え、速度比条件設定手段は、音声平均速度比に応じた音声区間の終了時刻における終了速度比を速度比条件としてさらに第２の所定時間毎に逐次設定し、速度比決定手段は、判別手段において判別された各音声区間に対して、音声平均速度比、音声区間長、及び終了速度比に基づいて音声区間の各区間における速度比を決定する。

**【００２２】**

第１７の発明は、上記第１６の発明において、速度比算出手段が算出するときに用いる音声区間長は、判別手段において過去に判別された音声区間の開始及び終了時刻から算出される音声区間長のうち、所定区間長以上の音声区間長のみに基づいて算出される。

40

**【００２３】**

第１８の発明は、上記第１６の発明において、速度比算出手段が算出するときに用いる音声区間長は、判別手段において過去に判別された音声区間の開始及び終了時刻から算出される音声区間長の最大値に基づいて算出される。

**【００２４】**

第１９の発明は、上記第１の発明において、所定時間分のオーディオ信号と、当該所定時間分のオーディオ信号に対して判別手段において判別された判別結果とを予め蓄積する蓄積手段をさらに備え、音声含有率算出手段は、蓄積手段に予め蓄積された判別結果に基づいて、所定時間分のオーディオ信号に含まれる音声区間の比率を示す音声含有率を算出

50

する。

【 0 0 2 5 】

第 2 0 の発明は、上記第 1 の発明において、判別手段は、オーディオ信号に対して、特定イベント音を含む特定イベント区間と、当該特定イベント区間以外の音声区間及び非音声区間とを判別し、音声再生装置は、判別手段において判別された判別結果に基づいて、オーディオ信号に含まれる特定イベント区間の比率を示す特定イベント含有率を算出する特定イベント含有率算出手段をさらに備え、速度比算出手段は、オーディオ信号に予め設定された再生速度に対する特定イベント区間の速度比を特定イベント含有率に基づいて算出するとともに、オーディオ信号の再生時間が所定の再生時間となるように、オーディオ信号に予め設定された再生速度に対する特定イベント区間以外の音声区間及び非音声区間の速度比を音声含有率に基づいてそれぞれ算出し、速度変換手段は、オーディオ信号に含まれる特定イベント区間と当該特定イベント区間以外の音声区間及び非音声区間との再生速度を速度比に基づいて変換する。

10

【 0 0 2 6 】

第 2 1 の発明は、入力されるコンテンツのオーディオ信号の再生速度を変えて所定の再生時間で再生する音声再生方法であって、オーディオ信号に対して、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、判別ステップにおいて判別された判別結果に基づいて、オーディオ信号に含まれる音声区間の比率を示す音声含有率を算出する音声含有率算出ステップと、オーディオ信号の再生時間が所定の再生時間となるように、オーディオ信号に予め設定された再生速度に対する音声区間及び非音声区間の速度比を音声含有率に基づいてそれぞれ算出する速度比算出ステップと、オーディオ信号を入力とし、当該オーディオ信号に含まれる音声区間及び非音声区間の再生速度を速度比に基づいてそれぞれ変換する速度変換ステップとを含む。

20

【 0 0 2 7 】

第 2 2 の発明は、入力されるコンテンツのオーディオ信号の再生速度を変えて所定の再生時間で再生する音声再生装置のコンピュータに実行させるためのプログラムであって、オーディオ信号に対して、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、判別ステップにおいて判別された判別結果に基づいて、オーディオ信号に含まれる音声区間の比率を示す音声含有率を算出する音声含有率算出ステップと、オーディオ信号の再生時間が所定の再生時間となるように、オーディオ信号に予め設定された再生速度に対する音声区間及び非音声区間の速度比を音声含有率に基づいてそれぞれ算出する速度比算出ステップと、オーディオ信号を入力とし、当該オーディオ信号に含まれる音声区間及び非音声区間の再生速度を速度比に基づいてそれぞれ変換する速度変換ステップとを、コンピュータに実行させるプログラムである。

30

【 0 0 2 8 】

第 2 3 の発明は、上記第 2 2 の発明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。

【 0 0 2 9 】

第 2 4 の発明は、入力されるコンテンツのオーディオ信号の再生速度を変えて所定の再生時間で再生する集積回路であって、オーディオ信号に対して、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手段と、判別手段において判別された判別結果に基づいて、オーディオ信号に含まれる音声区間の比率を示す音声含有率を算出する音声含有率算出手段と、オーディオ信号の再生時間が所定の再生時間となるように、オーディオ信号に予め設定された再生速度に対する音声区間及び非音声区間の速度比を音声含有率に基づいてそれぞれ算出する速度比算出手段と、オーディオ信号を入力とし、当該オーディオ信号に含まれる音声区間及び非音声区間の再生速度を速度比に基づいてそれぞれ変換する速度変換手段とを備える。

40

【発明の効果】

【 0 0 3 0 】

上記第 1 の発明によれば、入力されるオーディオ信号の音声含有率を算出することによ

50

り、所定の再生時間を達成しつつ、入力されるオーディオ信号に対して当該オーディオ信号の音声含有率に基づいた最適な音声区間及び非音声区間の速度比をそれぞれ算出することができる。つまり、音声含有率を算出することで、入力されるオーディオ信号に含まれる音声区間の比率を知ることができ、音声区間及び非音声区間の両方について所定の再生時間を達成するための最適な速度比を算出することができる。これにより、どのようなオーディオ信号が入力されても、再生内容の不連続性や情報の欠落による不快感などを低減させた、聞き取り易い再生を実現することができる。

【 0 0 3 1 】

上記第2の発明によれば、所定の再生時間に圧伸する比率を示す圧伸比と音声含有率とに基づく音声平均速度比及び非音声平均速度比をそれぞれ算出して、音声区間及び非音声区間の各区間における速度比が決定されることで、所定の再生時間を達成しつつ、入力されるオーディオ信号の音声含有率に基づいた最適な音声区間及び非音声区間の速度比をそれぞれ算出することができる。

10

【 0 0 3 2 】

上記第3の発明によれば、音声平均速度比、音声区間長、及び終了速度比に基づいて音声区間の各区間における速度比が決定されることで、音声区間の終了時刻における速度比を終了速度比で一定にしつつ、文頭から文末まで音声区間長に適した速度比を決定することができ、音声区間末の再生の高速化による聞き取り難さや不自然さを低減することができる。

【 0 0 3 3 】

20

上記第4の発明によれば、音声区間の各区間における速度比を経過割合に応じて決定することで、音声区間長の長短に関わらず、簡単な関数を用いて音声区間の各区間における速度比を決定することができる。

【 0 0 3 4 】

上記第5の発明によれば、音声区間の冒頭部分の再生速度が他の部分と比べて相対的に遅くなるので、冒頭部分の聞き逃しによって再生内容の理解度が低下することを防ぐことができる。

【 0 0 3 5 】

上記第6の発明によれば、コンテンツの種類に応じて異なる音声平均速度比及び非音声平均速度比を算出することができ、音声区間及び非音声区間の各区間における速度比をコンテンツに応じたより精度の高いものにすることができる。

30

【 0 0 3 6 】

上記第7の発明によれば、音声区間及び非音声区間の各区間における速度比がユーザによって指定された範囲内の音声平均速度比及び非音声平均速度比に基づく速度比となり、ユーザの聞き取り能力や好みに応じた速度変換処理を行うことができる。

【 0 0 3 7 】

上記第8の発明によれば、所定の再生時間を示す圧伸比を達成しつつ、入力されるオーディオ信号の音声含有率に基づいた最適な音声区間及び非音声区間の速度比をそれぞれ算出することができる。

【 0 0 3 8 】

40

上記第9の発明によれば、ユーザの使用目的に応じて異なる音声平均速度比及び非音声平均速度比の算出方法を変更することが可能になり、早聞き再生や遅聞き再生だけではなく、挿入や一時停止等を含んだオーディオ信号の出力時間に関する様々な制御について取り扱うことができる。これにより、コンテンツを視聴用、概要把握用、語学学習用、書き起こし用など用途に分けて個別に作成する必要がなく、同一のコンテンツを様々な目的で利用可能となる。

【 0 0 3 9 】

上記第10の発明によれば、圧伸比の大きさに応じて音声平均速度比及び非音声平均速度比の算出方法との対応が異なることで、例えば圧伸比が低い場合は概要把握用途に、圧伸比が高い場合は学習用途等に、ユーザの目的に応じた速度比を決定することができる。

50



これにより、ユーザは目的に応じて機器の使い分けを意識せずに使用でき、またユーザの視聴要望に即した速度変換処理を行うことができる。

【 0 0 4 0 】

上記第 1 1 の発明によれば、コンテンツ全体についての音声含有率を算出することにより、精度の高い音声区間及び非音声区間の速度比を算出することができる。

【 0 0 4 1 】

上記第 1 2 の発明によれば、蓄積手段を設けることなく処理が可能なため、オーディオ信号が蓄積手段に蓄積されるのを待つ必要がなく、リアルタイムで速度変換処理を行うことができる。

【 0 0 4 2 】

上記第 1 3 の発明によれば、第 1 の所定時間分の音声含有率が第 2 の所定時間に反映されることとなり、音声含有率の変動をすぐに反映した音声区間及び非音声区間の速度比の算出が可能になる。また、音声含有率を第 1 の所定時間分から算出することで再生時間に誤差が生じた場合であっても、第 2 の所定時間毎の圧伸比を用いて速度変換処理を行うので、当該誤差をこれから先の速度変換処理において解消させることができる。

【 0 0 4 3 】

上記第 1 4 の発明によれば、音声区間の各区間における速度比が示す前記第 2 の所定時間毎の変化を抑制するための統計量を用いることで、第 1 の所定時間分の音声含有率が局所的に高くなった場合でも、音声区間の各区間における速度比が上がりすぎることを防ぐことができる。その結果、音声含有率が異なる様々なコンテンツに対応した速度変換処理が可能となる。

【 0 0 4 4 】

上記第 1 5 の発明によれば、第 1 の所定時間分の音声含有率と時間変化の傾向が異なる音声含有率を統計量として利用することで、第 1 の所定時間分の音声含有率が局所的に変動した場合であっても、その変動は抑制され、聞き取り易い速度変換処理が可能となる。

【 0 0 4 5 】

上記第 1 6 の発明によれば、音声区間長が逐次算出されることで、音声区間の終了時刻が分からなくても、音声区間に対して適切な速度変換処理を行うことができ、より精度の高いリアルタイム処理を実現することができる。

【 0 0 4 6 】

上記第 1 7 の発明によれば、所定区間長以上の音声区間長のみに基づいて算出されることで、「はい」や「うん」など相槌や、「えー」などのフィラーなどを除いた平均的な音声区間長を算出することができる。

【 0 0 4 7 】

上記第 1 8 の発明によれば、速度比算出手段が算出するときに用いる音声区間長が過去の音声区間長の最大値に基づき算出されることで、音声区間の終了時刻では終了速度比で変換される割合が低下し、音声区間の平均速度比が更に下がる効果があり、聞き易い速度変換処理を実現することができる。

【 0 0 4 8 】

上記第 1 9 の発明によれば、所定時間分のオーディオ信号単位で速度変換を行うことができる。これにより、コンテンツの録画中であっても、全体の録画終了を待たずに速度変換処理を行うことができる。また、音声区間及び非音声区間の判別結果が蓄積手段に蓄積されることにより、音声含有率の実測値を算出することができ、より最適な速度比で速度変換を行うことができる。

【 0 0 4 9 】

上記第 2 0 の発明によれば、特定イベント含有率を算出して特定イベント区間の速度比を算出することで、例えば特定イベント音が音楽である場合、音楽番組などのオーディオ信号に対して音楽区間を音声区間及び非音声区間よりも遅い再生速度で再生を行うことができる。その結果、音楽を重視した速度変換処理を行うことができる。

【発明を実施するための最良の形態】

10

20

30

40

50

## 【 0 0 5 0 】

以下、本発明の実施形態について、図面を参照しながら説明する。

## 【 0 0 5 1 】

(第1の実施形態)

まず、図1を参照して本発明の第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、第1の実施形態に係る音声再生装置の構成例を示すブロック図である。図1において、本音声再生装置は、音声非音声判別部11、蓄積部12、音声含有率算出部13、速度比条件設定部14、音声区間長算出部15、速度比決定部16、及び速度変換部17で構成される。なお、本実施形態では、速度変換対象となるオーディオ信号をコンテンツ単位で予め蓄積部12に蓄積し、この蓄積したオーディオ信号を用いて再生速度を変えた再生処理を行う音声再生装置について説明する。また以下の説明において、音声が含まれる区間を音声区間とする。また、音声区間以外の区間、つまり音声を含まない区間を非音声区間とする。

10

## 【 0 0 5 2 】

音声非音声判別部11は、オーディオ信号を入力として、音声区間と非音声区間とを判別する。また音声非音声判別部11は、この判別結果と共に音声区間の始末端時刻(開始時刻及び終了時刻)を出力する。入力されるオーディオ信号は、CD、DVD、メモリ、又はハードディスクなどに記録されたオーディオ信号である。なお、オーディオ信号は、インターネットなどの通信回線を介して配信されたオーディオ信号や放送により受信したオーディオ信号などであってもよい。また、オーディオ信号は音声合成などその場で生成したものや、マイクで収録したもの、電話などの通信機器を通じて出力されるものでもよい。ここで、音声区間及び非音声区間を判別する方法としては、例えばオーディオ信号のパワーを算出し、閾値により判別を行う方法が挙げられる。また例えば「C e p s t r u m F l u xを用いた音声と音楽のセグメンテーション」<SP2000-1、内田貴之、山下昌毅、杉山雅英による、信学技報、SP2000-17>に記載されるように、ケプストラムの変化度合いを計測して判別を行う方法もある。ケプストラムの変化度合いを計測する方法では、BGMが重畳した音声であっても判別が可能である。

20

## 【 0 0 5 3 】

蓄積部12は、ハードディスク、DVD、又はメモリ媒体(例えばSDカード)などの読み書き可能な記録媒体で構成される。蓄積部12には、音声非音声判別部11に入力されるのと同じオーディオ信号がコンテンツ単位で蓄積される。また蓄積部12には、音声非音声判別部11から出力された、判別結果と音声区間の始末端時刻とが蓄積される。ここで例えば、TV放送を録画する場合を考える。この場合、TV放送を構成するオーディオ信号及びビデオ信号は蓄積部12に蓄積される。またこの蓄積と共に音声非音声判別部11において判別処理が行われ、判別結果や音声区間の始末端時刻が蓄積部12に蓄積される。蓄積部12には、コンテンツ1つに対して、オーディオ信号、ビデオ信号、判別結果、及び音声区間の始末端時刻が対応付けされて蓄積される。なお、オーディオ信号及びビデオ信号のフォーマットは、どのようなフォーマットであってもかまわない。

30

## 【 0 0 5 4 】

音声含有率算出部13は、コンテンツのオーディオ信号に含まれる音声区間の比率を示す音声含有率を算出する。具体的には、音声含有率算出部13は、蓄積部12に蓄積された各コンテンツに対して、それぞれに対応する判別結果や音声区間の始末端時刻を用いて音声含有率を算出する。音声含有率は、所定時間のオーディオ信号に含まれる音声区間長の和を当該所定時間で除算したものである。本実施形態では、コンテンツのオーディオ信号に含まれる音声区間長の和をコンテンツ長で除算したものを音声含有率とする。ここでコンテンツとは、速度変換を行う番組全体を意味する。したがって、コンテンツ長は通常、番組長に等しく、30分や1時間といったものが多い。なお、ユーザが番組の一部を速度変換対象として指定した場合、その一部をコンテンツとしてもよい。

40

## 【 0 0 5 5 】

音声含有率は、コンテンツによって異なる。例えば図2に示すようにコンテンツをジャ

50

ジャンル別に見た場合、音声含有率はジャンルによって異なることがわかる。図2は、ジャンル別の音声含有率を示した図である。図2において横軸はジャンルを示し、縦軸は音声含有率を示している。また図2に示す音声含有率は、同一週に放送された番組のうち、ジャンル別の視聴率の上位6位までの番組を抽出して、抽出した番組ごとの音声含有率をジャンル別に集計して平均化したものである。ニュースの音声含有率は、約60%であり、5ジャンルの中で最も高い値となっている。スポーツや音楽の音声含有率は、約40%となり、ニュースと比べて20%近い開きがある。また同じジャンルにおいても、音声含有率には図3に示すような多少のばらつきが存在する。図3は、各ジャンルの音声含有率の平均と標準偏差とを示した図である。ドラマやアニメでは標準偏差が16.2となり、他のジャンルに比べて高くなっている。このように音声含有率は、コンテンツによって異なる。

10

#### 【0056】

したがって、音声含有率を考慮しない従来技術では、上述したように、音声含有率の高いニュース番組などで目標時間からの遅延時間が長くなる。その結果、遅延時間を解消するために部分的に音声区間の高速再生や削除を行い、再生されるオーディオ信号が聞き取り難くなるという問題があった。これに対し、本実施形態では、コンテンツの音声含有率を算出する。これにより、目標時間から遅れることなく、コンテンツに応じた最適な音声及び非音声区間の速度比の算出が可能となり、部分的に偏ることなく聞き取り易い再生を実現することができる。なお、音声含有率を用いた速度比の算出方法については、後述にて詳述する。

20

#### 【0057】

速度比条件設定部14は、音声含有率及び目標圧伸比を入力とし、音声区間の平均速度比、非音声区間の平均速度比、及び音声区間の終端速度比を算出し、これらを速度比条件として設定する。

#### 【0058】

圧伸比とは、速度変換処理後の再生時間長を速度変換処理前の再生時間長で除算したものである。等倍速の再生では、圧伸比は1となる。2倍速の再生では、圧伸比は0.5となる。圧伸比が0から1までの値をとるとき、再生時間長は圧縮され、等倍速よりも速い速度で再生される。圧伸比が1より大きな値をとるとき、再生時間長は伸張され、等倍速よりも遅い速度で再生される。また目標圧伸比とは、速度変換を行いたいコンテンツの再生時間長をどれくらい圧縮もしくは伸張するかを示したものである。目標圧伸比は、圧縮の場合は0から1までの値をとり、伸張の場合は1以上の値をとる。目標圧伸比は、ユーザによって入力されてもよいし、予め装置に設定されていてもよい。また、ユーザが目標圧伸比を直接入力しなくてもよい。この場合、コンテンツ再生の目標時間を入力する。ユーザが目標時間を入力した場合、目標時間を速度変換処理前の再生時間長で除算することで、目標圧伸比を得ることができる。また速度比とは、等倍速に対する速度の比率を示したものである。速度比は、圧伸比の逆数で表される。また音声区間の終端速度比とは、音声区間の終端時刻における速度比を意味する。

30

#### 【0059】

次に、音声及び非音声区間の平均速度比を算出する方法について説明する。速度比条件設定部14は、予め設定された速度比算出分布を用いて平均速度比を算出する。速度比算出分布とは、音声含有率及び目標圧伸比に応じて、どの算出パターンで平均速度比を算出するかを示した分布である。換言すれば、速度比算出分布は、音声含有率と、目標圧伸比と、音声及び非音声区間の平均速度比を算出する方法との対応を示した対応情報である。

40

#### 【0060】

以下、算出パターンについて説明する。音声及び非音声区間の平均速度比は、目標圧伸比を達成するように算出される。具体的には、式(1)を満たすように算出される。

【数 1】

$$\frac{S}{V_{m1}} + \frac{1-S}{V_{m2}} = E \quad \cdots (1)$$

なお、 $S$  は音声含有率、 $V_{m1}$  は音声区間の平均速度比、 $V_{m2}$  は非音声区間の平均速度比、 $E$  は目標圧伸比を示す。算出パターンとしては、図 4 に示すように 5 種類の算出パターン a ~ e が考えられる。図 4 は、5 種類の算出パターンを示した図である。算出パターン a ~ e の条件は、以下ようになる。

a : 非音声区間の平均速度比  $V_{m2} = A_n$  (固定値) として、与えられる音声含有率  $S$  と目標圧伸比  $E$  から式 (1) を満たすように、音声区間の平均速度比  $V_{m1}$  を算出する。但し、 $V_{m1} < A_n$  を算出条件とする。

10

b : 音声区間の平均速度比  $V_{m1} = B_s$  (固定値) として、与えられる音声含有率  $S$  と目標圧伸比  $E$  から式 (1) を満たすように、非音声区間の平均速度比  $V_{m2}$  を算出する。但し、 $V_{m2} < B_s$  を算出条件とする。

c : 音声及び非音声区間の平均速度比を  $V_{m1} = V_{m2}$  として、与えられる音声含有率  $S$  と目標圧伸比  $E$  から式 (1) を満たすように、音声及び非音声区間の平均速度比  $V_{m1}$  及び  $V_{m2}$  を算出する。

d : 非音声区間の平均速度比  $V_{m2} = D_n$  (固定値) として、与えられる音声含有率  $S$  と目標圧伸比  $E$  から式 (1) を満たすように、音声区間の平均速度比  $V_{m1}$  を算出する。但し、 $V_{m1} < D_n$  を算出条件とする。

20

e : 音声区間の平均速度比  $V_{m1} = E_s$  (固定値) として、与えられる音声含有率  $S$  と目標圧伸比  $E$  から式 (1) を満たすように、非音声区間の平均速度比  $V_{m2}$  を算出する。但し、 $V_{m2} < E_s$  を算出条件とする。

【0061】

このように、平均速度比の算出パターンが異なれば、同じ音声含有率及び目標圧伸比であっても、音声区間の平均速度比と非音声区間の平均速度比の組み合わせは異なることとなる。そこで、音声含有率及び目標圧伸比に応じて、どの算出パターンを選択するか速度比算出分布を用いて決定する。以下、速度比算出分布について説明する。

【0062】

30

図 5 に速度比算出分布の一例を示す。図 5 において、縦軸は音声含有率、横軸は目標圧伸比を示しており、算出パターン a ~ c の領域の分布が示されている。ここで、速度比算出分布は、上述した算出パターンから所定の算出パターンを選択し、選択した算出パターンに対して上述した条件を満足しつつ、音声及び非音声区間の平均速度比の取り得る値を設定することで得られる。図 5 に示す速度比算出分布では、上述した算出パターンのうち算出パターン a ~ c が選択されている。算出パターン a では、非音声区間の平均速度比が最大値である  $A_n = 4$ 、音声区間の平均速度比  $V_{m1}$  の取り得る値が  $1.3 < V_{m1} < 2$  と設定されている。この取り得る値は、算出パターン a の算出条件 ( $V_{m1} < A_n$ ) を満足している。算出パターン b では、音声区間の平均速度比が  $B_s = 1.3$ 、非音声区間の平均速度比  $V_{m2}$  の取り得る値が  $1.3 < V_{m2} < 4$  と設定されている。この取り得る値は、算出パターン b の算出条件 ( $V_{m2} < B_s$ ) を満足している。算出パターン c では、音声及び非音声区間の平均速度比を  $V_{m1} = V_{m2}$  ( $1.3 < V_{m1} < 1.3$ ) と設定されている。このように算出パターンを選択し、音声及び非音声区間の平均速度比の取り得る値を設定することで、図 5 の速度比算出分布を得ることができる。音声含有率と目標圧伸比が算出パターン a の領域内にある場合、音声及び非音声区間の平均速度比は、算出パターン a で算出される。算出パターン b、c についても同様である。このように、音声含有率及び目標圧伸比に応じてどの算出パターンで算出するかが、速度比算出分布によって決まることとなる。なお、図 5 の一番左側にある処理不可の領域は、音声含有率に対して目標圧伸比が極端に小さく、音声及び非音声区間の平均速度比をユーザが聞き取り可能な範囲で最大にしても、目標圧伸比を達成できない領域である。

40

50

## 【 0 0 6 3 】

なお、図 5 の速度比算出分布では、音声含有率が高いほど、目標圧伸比に対して算出パターン a によって算出される割合が高くなる。図 5 の算出パターン a では、非音声区間の平均速度比が最大値 ( $A_n = 4$ ) に設定されている。これにより、目標圧伸比を達成する上で音声区間の平均速度比  $V_{m1}$  を可能な限り遅くすることができる。

## 【 0 0 6 4 】

また、図 5 の速度比算出分布では、音声含有率が低いほど、目標圧伸比に対して算出パターン b によって算出される割合が高くなる。図 5 の算出パターン b では、音声区間の平均速度比が  $B_s = 1.3$  (固定値) に設定されている。これにより、目標圧伸比を達成する上で非音声区間の平均速度比  $V_{m2}$  を可能な限り遅くすることができる。このように、図 5 の速度比算出分布は、音声含有率の大きさに応じて音声含有率と算出方法との対応が異なるものとなる。

## 【 0 0 6 5 】

また、目標圧伸比が大きいとき、算出パターン c が選択される。つまり、音声と非音声の平均速度比を等しくしている。目標圧伸比が大きいときは、音声と非音声が同じ平均速度比である方が、より自然に再生することができる。このように、図 5 の速度比算出分布は、目標圧伸比の大きさに応じて目標圧伸比と算出方法とが異なるものとなる。

## 【 0 0 6 6 】

また、図 5 に示す速度比算出分布では、音声含有率と目標圧伸比で領域が一意に定まるように、音声及び非音声区間の平均速度比の取り得る値が設定されている。つまり、音声及び非音声区間の平均速度比の取り得る値は、隣り合う算出パターン間の境界で平均速度比の値が連続するように設定されている。具体的には、算出パターン a では音声区間の平均速度比  $V_{m1}$  の最下限が  $1.3$  であり、隣り合う算出パターン b の音声区間の平均速度比  $B_s$  が  $1.3$  である。これにより、算出パターン a 及び b の境界において平均速度比の値が連続することとなる。また、算出パターン b では音声区間の平均速度比  $B_s$  が  $1.3$  であり、隣り合う算出パターン c では音声区間の平均速度比  $V_{m1}$  の最上限が  $1.3$  である。これにより、算出パターン b 及び c の境界において平均速度比の値が連続することとなる。

## 【 0 0 6 7 】

なお、目標圧伸比が大きくなるにつれて音声及び非音声区間の平均速度比がどのように連続して変化するかという観点から説明すると、次のようになる。目標圧伸比が処理不可の領域内の値をとるとき、音声及び非音声区間の平均速度比は算出されない。算出パターン a の領域内の値をとるとき、目標圧伸比が大きくなるにつれて、音声区間の平均速度比は 2 から  $1.3$  まで小さくなる。このとき、非音声区間の平均速度比は、4 で一定である。算出パターン b の領域内の値をとるとき、音声区間の平均速度比は  $1.3$  で一定となり、非音声区間の平均速度比は目標圧伸比が大きくなるにつれて 4 から  $1.3$  まで小さくなる。算出パターン c の領域内の値をとるとき、音声及び非音声区間の平均速度比は、共に同じ値となりながら、 $1.3$  から 1 まで小さくなる。

## 【 0 0 6 8 】

このように、速度比算出分布が、算出パターンの切り替わる境界の平均速度比が連続値となるように設定されることで、平均速度比が不連続な値をとる際に急激な速度変換が起こり、違和感が生じるという問題を回避することができる。

## 【 0 0 6 9 】

図 6 は、音声含有率が  $0.5$  のときの目標圧伸比、音声区間の平均速度比、非音声区間の平均速度比を示している。上述した図 5 に示されるように、目標圧伸比が  $0.375$  から  $0.510$  までの値をとるとき、算出パターン a が選択される。目標圧伸比が  $0.501$  から  $0.769$  までの値をとるとき、算出パターン b が選択される。目標圧伸比が  $0.769$  から 1 までの値をとるとき、算出パターン c が選択される。ここで、図 5 に示した各算出パターンには、上述したように、平均速度比の取り得る値が設定されている。従って、上述した算出パターン及び式 (1) により、図 6 に示すような平均速度比が算出され

10

20

30

40

50

る。

【0070】

目標圧伸比が0.1及び0.3の値をとるとき、図5に示す速度比算出分布からも分かるように、処理不可の領域となるので、音声及び非音声区間の速度比は算出されない。目標圧伸比が0.4及び0.5の値をとるとき、共に算出パターンaによって算出される。なお、算出パターンaによって算出される場合、目標圧伸比が増加するにつれて音声区間の平均速度比が小さくなっていることが分かる。目標圧伸比が0.6及び0.7の値をとるとき、共に算出パターンbによって算出される。なお、算出パターンbによって算出される場合、目標圧伸比が増加するにつれて非音声区間の平均速度比が小さくなっていることが分かる。目標圧伸比が0.9及び1.0の値をとるとき、音声及び非音声区間の平均速度比が等しくなり、目標圧伸比が増加するにつれて音声及び非音声区間の平均速度比が小さくなっていることが分かる。

10

【0071】

次に、音声区間の終端速度比を算出方法について説明する。速度比条件設定部14は、算出した音声区間の平均速度比から、音声区間の終端速度比を算出する。図7は、音声及び非音声区間の速度比変化を示した模式図である。音声区間1の区間長は、音声区間2の区間長よりも短くなっている。縦軸は変換速度比であり、横軸は経過時間である。変換速度比とは、速度変換部17の速度変換処理に用いられる速度比を示しており、音声及び非音声区間をそれぞれ細分化した各区間における速度比によって示される。変換速度比の決定方法については、後述にて説明する。図7に示すように、音声区間長が異なっても、音声区間の終端時刻での速度比は等しくなっている。速度比条件設定部14は、この終端時刻の速度比を音声区間の終端速度比として算出している。図7からも明らかとなり、音声区間は終端速度比よりも遅い速度比が設定されている。音声区間の終端速度比 $V_{end}$ は音声区間の平均速度比 $V_{m1}$ に を加算したものとす。つまり、終端速度比 $V_{end} = V_{m1} +$  とする。なお、聴取実験により、 を0.2とし、 $V_{end}$ は2.0を超えないものが好ましいことが分かった。なお、図7では、非音声区間の速度比は平均速度比 $V_{m2\ end}$ で一定である。

20

【0072】

以上のように、速度比条件設定部14は、音声含有率及び目標圧伸比を入力とし、音声区間の平均速度比、非音声区間の平均速度比、及び音声区間の終端速度比を算出し、これらを速度比条件として設定する。

30

【0073】

音声区間長算出部15は、音声区間の始終時刻を入力とし、音声区間長を算出する。速度比決定部16は、速度比条件設定部14で設定された速度比条件と音声区間長とに基づき、音声及び非音声区間の変換速度比を決定する。ここで、変換速度比とは、上述したように、速度変換部17の速度変換処理に用いられる速度比を示しており、音声及び非音声区間をそれぞれ細分化した各区間における速度比によって示される。ただし、音声区間中の速度比を一定にする場合や、一定時間ごとに速度比を切り替える場合は音声区間長を必ずしも算出する必要はないため、音声区間長算出部15を設けなくてもよい。たとえ、音声区間長算出部15を設けなかったとしても、音声含有率によって音声区間の平均速度比を設定しているため、従来の方法よりも聞き易くなる。これに対し、音声区間長算出部15を設けることによって、音声区間長を算出し、音声区間の細分化された各区間に対して速度比を設定することで、更に聞き易くなる効果がある。

40

【0074】

速度変換部17は、オーディオ信号を入力とし、速度比決定部16で決定された変換速度比に従って速度変換を行う。速度変換の方法としては、例えば「高品質音声速度変換方式のDSPによる実現」<鈴木，三崎，電子情報通信学会 音声研究会資料 SP90-34、(1990.8.23)>、特許第3189587号などに記載された公知の方法を用いるとする。このような方法により、1倍速以下の遅い速度比での再生や、1倍速以上の速い速度比での再生が可能となる。また、速度変換の方法はこれに限らず、音を合成

50

したり、区間の削除や挿入などを行ったり、速度比決定部 16 で決定された変換速度比を満たすような処理を行っているものであれば方法は問わない。例えば、変換速度比が 0.5 である場合を仮定すると、ある入力区間に対して出力再生時間が 2 倍となればよく、音を引き伸ばしたり、無音区間を追加したり、新たに音を合成してもよい。このように、速度変換部 17 はある区間に対する入力と出力との関係が対応付けられており、変換速度比を満たすような処理を行っているものであれば、速度変換の方法として含まれる。

【0075】

以下、図 8 を参照して、第 1 の実施形態に係る音声再生装置の処理について説明する。図 8 は、第 1 の実施形態に係る音声再生装置の処理の流れを示すフローチャートである。

【0076】

まず、ユーザが入力装置（図示なし）においてコンテンツを録画する指示をしたとき、当該コンテンツのオーディオ信号及びビデオ信号が蓄積部 12 に蓄積される。このとき、音声非音声判別部 11 はそのコンテンツのオーディオ信号について音声区間と非音声区間とを判別する（ステップ S101）。なお、ステップ S101 において判別された判別結果と音声区間の始終端時刻についても、蓄積部 12 に蓄積される。

【0077】

ステップ S101 の次に、入力装置において、ユーザが所望のコンテンツを再生する指示をしたか否かが判断される（ステップ S102）。ユーザの指示があった場合（ステップ S102 で Yes）、音声含有率算出部 13 は、指示されたコンテンツの音声含有率を算出する（ステップ S103）。

【0078】

ステップ S103 の次に、ユーザが入力装置（図示なし）において目標圧伸比を設定する（ステップ S104）。速度比条件設定部 14 は、ステップ S103 で算出された音声含有率と、予め設定された速度比算出分布とから、ステップ S104 で設定された目標圧伸比が処理不可の領域内にあるか否かを判断する（ステップ S105）。処理不可の領域内に目標圧伸比が設定された場合（ステップ S105 で No）、処理はステップ S104 に戻る。ステップ S104 において、速度比条件設定部 14 は、目標圧伸比を処理可能な値に再設定する。図 5 の場合、音声含有率が 0.5 のとき処理可能な最小の目標圧伸比は、0.375 となる。したがって速度比条件設定部 14 は、最も近い領域境界の値 0.375 を目標圧伸比として再設定する。なお、速度比条件設定部 14 が自動で再設定するのではなく、目標圧伸比の入力を再度ユーザに求めるようにしてもよい。

【0079】

ステップ S105 の次に、速度比条件設定部 14 は、ステップ S103 で算出された音声含有率、ステップ S104 及び S105 で設定された目標圧伸比に基づいて、音声区間の平均速度比、非音声区間の平均速度比、及び音声区間の終端速度比を算出し、速度比条件を設定する（ステップ S106）。なお、速度比条件の算出方法については、上述したとおりである。

【0080】

ステップ S106 の次に、音声区間長算出部 15 は、音声区間の始終端時刻を入力とし、音声区間長を算出する（ステップ S107）。音声区間長は、図 9 及び図 10 に示すように、同じコンテンツ内でも長短様々なものが含まれているが、ジャンルによっても大きく異なる。図 9 は、ニュース番組に含まれる音声区間長とその頻度を示した図である。図 10 は、野球番組に含まれる音声区間長とその頻度を示した図である。図 9 及び図 10 において、横軸は音声区間長であり、縦軸は番組中に発生した頻度である。

【0081】

ここで、上述した従来技術では、音声区間長を考慮せず始端からの経過時間のみで速度比を設定しており、音声区間長が長いものでは経過時間に伴って速度比が段々速くなり、聞きにくくなるという課題があった。これに対し、本実施形態では、音声区間長を考慮することで、図 7 に示したように、音声区間長の長短に関わらず、音声区間の終端での速度比が等しくなるように速度比を決定することができる。これにより、音声区間の速度比が

10

20

30

40

50

始端から徐々に速くなるが、終端が聞き取り可能な速度比までしか速くならないため、従来技術に比べ、聞き易さが大きく改善した。

【 0 0 8 2 】

ステップ S 1 0 7 の次に、速度比決定部 1 6 は、蓄積部 1 2 に蓄積された音声区間の始終端時刻を参照して、コンテンツの始端から順に所定の単位時間毎に音声区間であるか否かを判断する（ステップ S 1 0 8）。音声区間と判断した場合、速度比決定部 1 6 は、音声区間の始終端時刻と、ステップ S 1 0 7 で算出された音声区間長とに基づき、音声区間における経過割合を算出する（ステップ S 1 0 9）。音声区間の経過割合とは、音声区間の始端を 0、終端を 1 として、始端からの経過時間を音声区間長で除算したものである。

【 0 0 8 3 】

ステップ S 1 0 9 の次に、速度比決定部 1 6 は、音声区間の経過割合から、音声区間の変換速度比を決定する（ステップ S 1 1 0）。以下、ステップ S 1 1 0 の具体的な処理例について説明する。変換速度比の算出処理の一例としては、音声区間の平均圧伸比に対する圧伸比変化量の和が 0 になるように、変換速度比を算出する方法が挙げられる。図 1 1 は、音声区間の圧伸比の変化を示した図である。図 1 1 において、 $x$  は経過割合、 $v_x$  は経過割合が  $x$  のときの変換圧伸比、 $v_s$  は始端圧伸比、 $v_e$  は終端圧伸比、 $v_{m1}$  は平均圧伸比とする。ここで、始端圧伸比  $v_s$  と終端圧伸比  $v_e$  とを結ぶ圧伸比の変化カーブは、式 ( 2 ) で表現される。

【数 2】

$$v_x = (v_e - v_s)x + v_s \quad \cdots (2)$$

平均圧伸比  $v_{m1}$  は、音声区間の平均速度比  $V_{m1}$  の逆数である。終端圧伸比  $v_e$  は、終端速度比  $V_{end}$  の逆数である。ここで、圧伸比変化量は、音声区間の平均圧伸比  $v_{m1}$  を 0 と想定したとき、 $v_{m1}$  に対して増減する量（図 1 1 の網掛け部分の面積）を意味する。したがって、この量の和が 0 となるためには、図 1 1 に示したように、 $x = 0.5$  のときに変換圧伸比  $v_x$  が平均圧伸比  $v_{m1}$  となるようにすればよい。平均圧伸比  $v_{m1}$  は、音声区間の平均速度比  $V_{m1}$  から求まる値であり、終端圧伸比  $v_e$  は、終端速度比  $V_{end}$  から求まる値である。したがって、式 ( 3 ) を満たすように始端圧伸比  $v_s$  を設定すればよいこととなる。

【数 3】

$$v_{m1} = \frac{1}{2}(v_e - v_s) + v_s \quad \cdots (3)$$

なお、経過割合が  $x$  のときの変換圧伸比  $v_x$ 、始端圧伸比  $v_s$ 、終端圧伸比  $v_e$ 、平均圧伸比  $v_{m1}$  を速度比で表すと、式 ( 4 ) のようになる。ここで、 $V_x$  は経過割合が  $x$  のときの変換速度比、 $V_s$  は始端速度比、 $V_{end}$  は上述した終端速度比、 $V_{m1}$  は上述した平均速度比を示す。

【数 4】

$$v_{m1} = \frac{1}{V_{m1}}, \quad v_e = \frac{1}{V_{end}}, \quad v_s = \frac{1}{V_s}, \quad v_x = \frac{1}{V_x} \quad \cdots (4)$$

そして、式 ( 2 ) に式 ( 3 ) および ( 4 ) を代入すると、式 ( 5 ) が得られる。



【数 5】

$$V_x = \frac{V_{end} V_s}{(V_s - V_{end})x + V_{end}} \quad \cdots (5)$$

なお、速度変換後の音声区間長は、平均速度比  $V_{m1}$  で一様に変換した時間長と等しくなることから、式 (6) が成り立つ。

【数 6】

$$V_s = 2V_{m1} - V_{end} \quad \cdots (6)$$

10

このようにステップ S 1 1 0 において、速度比決定部 1 6 は、式 (5) に音声区間の経過割合  $x$  を代入することで、音声区間の変換速度比  $V_x$  を決定することができる。このステップ S 1 1 0 で算出した音声区間の変換速度比は、上述した図 7 のような変化となる。つまり、音声区間の冒頭部分を遅くし、終端に向かって徐々に速めていくように、変換速度比を音声区間長に応じて変化させることができる。

【0084】

ステップ S 1 0 8 において非音声区間と判断した場合、速度比決定部 1 6 は、非音声区間の始端から終端まで、速度比条件設定部 1 4 で設定された非音声区間の平均速度比を変換速度比として決定する。つまり、図 7 に示したように、速度比決定部 1 6 は、平均速度比で一定となるように、非音声区間の始端から終端までの変換速度比を決定する。

20

【0085】

ステップ S 1 1 0 及び S 1 1 1 の次に、速度比決定部 1 6 は、コンテンツの終端まで変換速度比を算出したか否かを判断する (ステップ S 1 1 2)。終端ではないとき、処理はステップ S 1 0 8 へ戻る。このように、コンテンツの終端までの変換速度比が算出されるまで、速度比決定部 1 6 においてステップ S 1 0 8 ~ S 1 1 2 までの処理が繰り返される。ステップ S 1 1 2 においてコンテンツの終端まで変換速度比が算出されたと判断された場合、速度変換部 1 7 において変換速度比に従ってオーディオ信号の速度変換が行われ、速度変換後のオーディオ信号の再生が開始される (ステップ S 1 1 3)。入力装置 (図示なし) が本装置の処理を終了するか否かの指示を受け付ける (ステップ S 1 1 4)。ユーザが他のコンテンツについて速度変換処理を行う場合 (ステップ S 1 1 4 で No)、処理はステップ S 1 0 2 へ戻る。

30

【0086】

以上のように、本実施形態に係る音声再生装置によれば、コンテンツの音声含有率を算出することにより、コンテンツに応じた速度比条件を設定することができる。これにより、目標圧伸比、つまり目標時間を達成しつつも、音声区間及び非音声区間の速度比をコンテンツに応じた最適な速度比にそれぞれ設定することができる。その結果、どのようなコンテンツのオーディオ信号が入力されても、聞き取り易い速度で再生することが可能となり、再生内容の不連続性や情報の欠落による不快感などを低減させた再生を行うことができる。

40

【0087】

また本実施形態に係る音声再生装置によれば、速度比決定部 1 6 において図 1 1 に示すように圧伸比の変化を示す関数として、1 次関数を用いるとした。つまり、本装置に入力されたオーディオ信号に対して、一次直線で速度比を設定している。ここで、日本語はモーラリズムの言語と言われており、個々のモーラが同じ長さになるように話す傾向がある。モーラは言葉を話すときの長さの単位であり、日本語では俳句や短歌で数える拍に相当する。「かな」でいえば、一文字に相当している。このモーラ毎に速度比を変化させることが望ましいが、入力されるオーディオ信号に対して一次直線で速度比の算出を行っているため、モーラ毎に速度比を変化させなくても、十分に自然な再生を実現することができ

50

る。さらに、音声区間の始端から終端までの速度比は、一次関数によって、細かく切り替えられている。これにより、知覚される時間よりも短い間隔で速度を変化させることとなり、違和感が少ない自然な再生を提供することができる。

#### 【0088】

また本実施形態に係る音声再生装置によれば、音声区間の始端から終端まで音声区間長に応じて速度比を設定している。これにより、音声区間の終端時刻において予め設定した終端速度比よりも速い速度比になることなく、音声区間の終端時刻付近において速度比が速くなりすぎて聞き取り難くなることを防ぐことができる。

#### 【0089】

なお、上述では図11に示したように、圧伸比の変化を一次関数によって表すようにしたが、他の関数によって表されても構わない。例えば、上に凸または下に凸の指数関数であってもよい。また例えば、予め用いることができる速度比が限られている場合は、2段階や数段階の速度比で変換をおこなっても、音声区間長の経過割合に応じた速度変換を行うことで、不自然さを低減させた再生を提供することができる。図12は、2段階の変換速度比を算出した場合を示す図である。図12において、より好ましくは、音声区間全体に対して、最初の変換速度比が始端から2～3割の範囲を占めるようにする。これにより、より自然な再生を実現することが聴取実験で明らかとなった。また例えば、音声区間の速度比を非音声区間と同様に一定の速度比としてもよい。この場合であっても、音声含有率により、適切な速度比が設定されるため、従来技術のような非音声区間の削除や極端な高速化が行われずに済み、聞き易い再生を提供することができる。

#### 【0090】

なお、上述では、速度比条件設定部14は、図5に示した速度比算出分布を用いるとしたが、これに限定されない。ユーザ自身が、所望の算出パターンを選択して音声及び非音声区間の平均速度比の取り得る値を所望の値に設定し、速度比算出分布を作成するようにしてもよい。つまり、速度比条件設定部14が用いる速度比算出分布は、予め設定されているものに限らず、ユーザによって設定されるものであってもよい。例えば、図5では音声区間の平均速度比が2.0までとり得る。しかし、高齢者や語学学習者では2.0よりももっと遅い平均速度比での聴き取りを希望する場合もある。その際に、ユーザが望む平均速度比を超えないように、音声及び非音声区間の平均速度比の取り得る値を設定することで、ユーザの聴き取り能力に応じた再生処理が可能となる。また、高齢者や語学学習者では通常の再生速度よりさらに遅くして聞きたい場合が存在する。音声区間は通常の平均速度比1.0よりも遅い速度にし、非音声区間を通常の平均速度比1.0より高速化して通常の再生時間と同じ時間で収めたい、あるいはもっと短い時間で視聴したいといった要望に答えるためにも、速度比条件設定部14の速度比算出分布は用途に応じて切り替えることを可能にしている。

#### 【0091】

また、図5に示した速度比算出分布は、ジャンル毎に予め用意されていてもよい。この場合、EPG等のジャンル情報やユーザの指示によって、いずれの速度比算出分布を用いるかを選択する。ここで音声含有率以外にも、画像の動きの激しさ等はジャンルによって異なる。例えば、ドキュメンタリーなどの静止画像が多いジャンルでは、非音声区間を高速で再生しても、画像の高速化による情報の欠落は少ない。また、非音声区間を高速で再生できるので、音声区間を1倍速に近づけることができる。その結果、番組内容を理解しやすい再生を行うことができる。ここで、ドキュメンタリーなどの静止画像が多いジャンルについての速度比算出分布の例を図13に示す。図13に示すように、速度比算出分布は、算出パターンa及びbの領域で構成される。このうち、算出パターンaでは、非音声区間の平均速度比が最大値である $A_n = 4$ 、音声区間の平均速度比 $V_{m1}$ の取り得る値が $1 \leq V_{m1} \leq 2$ と設定されている。算出パターンbでは、音声区間の平均速度比が $B_s = 1$ 、非音声区間の平均速度比 $V_{m2}$ の取り得る値が $1 \leq V_{m2} \leq 4$ と設定されている。

#### 【0092】

図14は、図13に示す速度比算出分布において、音声含有率が0.5のときの目標圧

10

20

30

40

50

伸比、音声区間の平均速度比、非音声区間の平均速度比を示している。目標圧伸比が 0.1 及び 0.3 の値をとるとき、図 13 に示す速度比算出分布からも分かるように、処理不可の領域となるので、音声及び非音声区間の速度比は算出されない。目標圧伸比が 0.4、0.5、及び 0.6 の値をとるとき、共に算出パターン a によって算出される。目標圧伸比が 0.7、0.9、及び 1 の値をとるとき、共に算出パターン b によって算出される。図 14 から分かるように、図 13 に示す速度比算出分布を用いた場合、算出パターンが 2 つのパターンで構成されるので、音声及び非音声の平均速度比の差が大きくなる。換言すれば、非音声区間を高速化し、音声区間を 1 倍速にすることができ、番組内容をより理解し易い再生を実現することができることを意味する。

#### 【0093】

また、例えばスポーツなど動きの激しいシーンが多いジャンルでは、音声と非音声の平均速度比に大きな差をつけないほうがよい。なぜならば、動きの激しいシーンが多いジャンルは、動きの少ないジャンルに比べ、番組の内容理解に対して音声以外の部分が与える影響が大きいので、非音声区間の聞き取り易さや見易さを向上させる必要があるからである。ここで、スポーツなど動きの激しいシーンが多いジャンルについての速度比算出分布の例を図 15 に示す。図 15 に示すように、速度比算出分布は、2 つの算出パターン a、2 つの算出パターン b、及び算出パターン c の領域で構成される。このうち、一番左側の算出パターン a では、非音声区間の平均速度比が  $A_n = 3$ 、音声区間の平均速度比  $V_{m1}$  の取り得る値が  $1.8 \leq V_{m1} \leq 2.5$  と設定されている。この算出パターン a と隣り合う算出パターン b では、音声区間の平均速度比が  $B_s = 1.8$ 、非音声区間の平均速度比  $V_{m2}$  の取り得る値が  $2.5 \leq V_{m2} \leq 3$  と設定されている。この算出パターン b と隣り合う算出パターン a では、非音声区間の平均速度比が  $A_n = 2.5$ 、音声区間の平均速度比  $V_{m1}$  の取り得る値が  $1.5 \leq V_{m1} \leq 1.8$  と設定されている。この算出パターン a と隣り合う算出パターン b では、音声区間の平均速度比が  $B_s = 1.5$ 、非音声区間の平均速度比  $V_{m2}$  の取り得る値が  $1.5 \leq V_{m2} \leq 2.5$  と設定されている。この算出パターン b と隣り合う算出パターン c では、音声及び非音声区間の平均速度比を  $V_{m1} = V_{m2} (1 \leq V_{m1} \leq 1.5)$  と設定されている。

#### 【0094】

図 16 は、図 15 に示す速度比算出分布において、音声含有率が 0.5 のときの目標圧伸比、音声区間の平均速度比、非音声区間の平均速度比を示している。目標圧伸比が 0.1 及び 0.3 の値をとるとき、図 15 に示す速度比算出分布からも分かるように、処理不可の領域となるので、音声及び非音声区間の速度比は算出されない。目標圧伸比が 0.4 及び 0.5 の値をとるとき、共に算出パターン a によって算出される。目標圧伸比が 0.6 の値をとるとき、算出パターン b によって算出される。目標圧伸比が 0.7、0.9、及び 1 の値をとるとき、算出パターン c によって算出される。図 15 から分かるように、目標圧伸比に対して算出パターンは多数切り替わっている。これにより、図 15 に示す速度比算出分布を用いた場合、音声と非音声の平均速度比に大きな差が生じない。その結果、非音声区間の聞き取り易さ及び見易さが向上する。また、図 6 で示した速度比算出分布よりも、非音声区間の速度比を若干遅めに設定している。これにより、非音声区間において生じる動きが激しいシーンが多いジャンルに対して、非音声区間の聞き取り易さ及び見易さをさらに向上させることができる。このように、速度比算出分布をジャンル毎に準備しておくことで、よりの確な速度変換処理が可能になる。

#### 【0095】

なお、ジャンル毎だけではなく、動きの激しさなどを示す画像情報や、音響的な特徴に応じた速度比算出分布が予め用意されていてもよい。このような速度比算出分布は、例えば、音楽やある特定の人物の音声などユーザが着目したい音に対して個別に速度を制御したい場合に、有効である。

#### 【0096】

また、上述では、速度比算出分布を構成する領域として、算出パターン a ~ c を用いて場合について説明したが、目的に合わせて算出パターン d 及び e を用いてもよい。音楽番

10

20

30

40

50

組において音楽を重視して再生したい場合、音楽は非音声区間であるため、出来るだけ非音声区間の速度比を遅くするとよい。その代わり、音声区間の速度比を速くする必要がある。したがって、このような音楽番組などのジャンルに対しては、上述した算出パターン d を用いるのが好適である。これは音楽に限らず、ユーザが着目したい音に対して再生を行う場合に有効である。また、ユーザがコンテンツに含まれる非音声区間をサーチする場合、音声区間を出来るだけ高速で再生することが望まれる。したがって、この場合、上述した算出パターン e を用いるのが好適である。

【 0 0 9 7 】

また、速度比条件設定部 1 4 は、音声区間長算出部 1 5 で得られた音声区間長から、図 9 及び図 1 0 に示したような統計的な分布を求め、その統計的な分布に基づく速度比算出分布を用いてもよい。音声区間長とその生起頻度はコンテンツの属性を示している。このため、統計的な分布に基づく速度比算出分布を用いることで、コンテンツの属性に応じた速度変換処理が可能になる。例えば音声含有率が同じであったとしても、音声区間長が短いものが多く音声区間の生起頻度が高いコンテンツや、音声区間長が長いものが多く音声区間の生起頻度が低いコンテンツが存在する。後者のコンテンツでは、一つの音声区間あたりに含まれる情報量が多く、理解にかかるユーザの負荷が高いことが想定される。したがって、このようなコンテンツに対しては、音声区間に対してより重点的に遅い速度比を配分するような速度比算出分布を用いる。このように速度比条件設定部 1 4 は、音声区間長の統計的な分布に基づく速度比算出分布を用いて、速度比条件を設定してもよい。このことは、プライベートコンテンツについて特に有効である。プライベートコンテンツは放送コンテンツと異なり編集等の処理を施していないものが多いため、音声含有率や音声区間長もコンテンツごとにばらつきが大きい。そのため、様々な速度比算出分布を用意することで、プライベートコンテンツなどコンテンツ間で音声区間長や音声含有率のばらつきが大きいものにおいても適切な速度比条件を設定することが可能になる。

【 0 0 9 8 】

また、上述では、目標圧伸比が 0 から 1 となる場合についてのみ説明を行ったが、速度制御後の再生時間が通常再生時間と同じかそれよりも長い時間で視聴を行う遅聞きや遅見再生など目標圧伸比 1 以上の場合についても、同様に速度比算出分布を予め用意しておくことで、速度比条件を設定することが可能である。また、一つの音声区間ごとに音声区間長と同じ長さの非音声区間を設け発音練習を促すような発音練習モードの出力制御も可能となる。例えば、発音練習モードでは、非音声区間の平均速度比は直前の音声区間長によって定めるとすると、音声区間や非音声区間の速度比は以下の式から算出できる。

【 数 7 】

$$V_{m1} = \frac{2 * S}{E} \quad \dots (7)$$

【 数 8 】

$$V_{m2}(n) = \frac{N(n) * V_{m1}}{M(n)} \quad \dots (8)$$

なお、S は音声含有率、V m 1 は音声区間の平均速度比、V m 2 は非音声区間の平均速度比、E は目標圧伸比を示す。n 番目の音声区間の音声区間長を M ( n ) とし、n 番目の音声区間に後続する非音声区間の非音声区間長を N ( n ) とする。音声区間長と同じ長さの非音声区間長を設けるため、音声区間の平均比 V m 1 は式 ( 7 ) のように表せる。また発音練習を行うには、音声区間と同じ長さの非音声区間を必要とする。このため、後続の非音声区間の速度比は、音声区間長に応じて算出する必要があるため、式 ( 8 ) のようになる。このように、音声区間長と同じ長さの非音声区間を設け、一定の時刻内で再生するという発音練習モードにおいても、音声含有率や音声区間長を利用することで、適切な速度

比を設定することが可能になる。なお、今回は音声区間と同じ長さの非音声区間を設けたが、非音声区間の長さは学習の進み具合などに応じて変更してもよい。このように、語学学習用に新たにコンテンツを作成しなくても、音声含有率と音声区間長の利用によって、発音練習に適した速度で音声を提示することが可能になる。また、学習に費やしたい時間を始めに指定することで、コンテンツ長から圧伸比を算出し、学習時間内におさまるように速度を制御することが可能になる。また、学習のレベルに応じて速度比を変えることも可能となる。

#### 【 0 0 9 9 】

他の使用目的としては、音声の書き起こしを行うときに用いる書き起こしモードが考えられる。この場合も同様に速度比条件設定部 1 4 の速度比算出分布を変えることで対応可能である。音声を書き起こすには書き起こす人の書き込み能力、例えば、紙に記入する場合、一定時間で何文字記入可能かということや、キーボードで打ち込む場合、一定時間に何タイプ可能かなど各ユーザの書き込み能力に応じて、再生速度を変える必要がある。書き込み能力より再生速度が速ければ、すぐに書き込み部分が追い越され、書き込み部分より先の部分が再生されてしまう。そのため、一時停止や巻き戻しなどの操作が必要となる。また、そのような再生制御の操作は書き込み処理を中断させるため、二重に時間を無駄に消費させることになる。書き込み能力より再生速度が遅ければ再生部分に追い越されることはないが、書き込み後も次の音声が始まるまで待ち時間が発生し無駄な時間を消費することになりは無い。そこで、音声の書き起こしを行うときはユーザの書き込み能力に応じた速度で再生する必要がある。従来の方法では音声区間長や音声含有率が不明なため、非音声区間も音声区間と同じ速度で再生されたり、速度を遅くした場合どの程度の時間がかかるかは事前にわからなかったりした。今回速度条件設定部 1 4 で以下のような処理をおこなえば、音声区間は聞きやすく非音声区間を省くような再生が可能となる。音声区間の平均速度比は式 ( 9 ) のように表せる。

#### 【 数 9 】

$$V_{m1} = \frac{Q * S}{Q * E - U * P} \quad \dots (9)$$

#### 【 数 1 0 】

$$V_{m2} = \infty \quad \dots (10)$$

なお、S は音声含有率、V m 1 は音声区間の平均速度比、V m 2 は非音声区間の平均速度比、E は目標圧伸比を示す。コンテンツに含まれる音声区間の総数を U とし、コンテンツの全長を Q としている。P は音声区間を再生後に、書き込み行の変更など書き起こし作業に必要な時間として設けた一定時間である。従って作業によっては P の値は 0 でも構わない。非音声区間は再生されないため、非音声区間の速度比 V m 2 は式 ( 1 0 ) のようになる。ここで、音声区間の平均速度比 V m 1 をユーザの能力に応じた書き込み速度に設定すれば、書き込み作業中に一時停止や早送り、巻き戻しなどの機器操作をすることなく、書き起こし作業が可能となる。このように速度条件設定部 1 4 で音声区間と非音声区間の平均速度比を設定することで、ユーザの書き込み能力に応じた速度比で音声区間のみを再生するような書き起こし作業モードが可能となる。

#### 【 0 1 0 0 】

その他にも、学習開始時には音声区間を遅い速度で再生を行い、学習経過時間に応じて徐々に速度比を早くしながら学習時間の合計は所定の時間になるように、音声や非音声区間の速度を制御するような聴き取り練習用の制御も可能になる。例えば、一回の学習の中で、学習開始時には学習終了時に比べ、遅い速度で音声区間の再生を行い徐々に音声区間の再生速度を上げていくような制御を行うこともできる。また長期的な学習の中で、最初

に学習を開始した時点からの経過時間や前回の操作履歴などから、今回の学習における音声区間の再生速度を制御してもよい。このように音声区間や非音声区間の再生速度を変更したり、音声区間の直前や直後に一時停止や無音やオーディオ区間を付与したり、画面を挿入したりしながらコンテンツ全体の時間調整をすることにも対応可能である。

#### 【0101】

(第2の実施形態)

図17を参照して、本発明の第2の実施形態に係る音声再生装置について説明する。図17は、第2の実施形態に係る音声再生装置の構成例を示すブロック図である。図17において、本音声再生装置は、音声非音声判別部11、音声含有率予測部18、速度比条件設定部14、音声区間長予測部19、速度比決定部16、速度変換部17、及び圧伸比算出部20で構成される。本実施形態は、第1の実施形態に係る音声再生装置に対し、オーディオ信号の再生速度を変えた再生処理をリアルタイムで行う点で異なる。以下、異なる点を中心に説明する。

10

#### 【0102】

音声含有率予測部18は、音声非音声判別部11から出力されたフレーム毎の判別結果から、現時点より過去数分の音声含有率を算出する。そして、算出した音声含有率を用いて、現時点より1つ進んだセクションにおける音声含有率を予測する。現時点より過去数分の音声含有率を $X(z-1)$ とすると、音声含有率の予測値 $Y(z)$ は式(11)で表現される。 $Y(z-1)$ は一つ前のセクションでの音声含有率の予測値である。以下、音声含有率の予測値を予測音声含有率と称す。は0から1までの値で、シミュレーションにより最適な値を求めた。

20

#### 【数11】

$$Y_{(z)} = \alpha Y_{(z-1)} + (1-\alpha) X_{(z-1)} \quad \cdots (11)$$

式(11)において、初期値 $X(0) = Y(0)$ は、一般的なコンテンツの音声含有率の平均値とする。また音声含有率 $X(z)$ が算出されるまでの間予測音声含有率 $Y(z)$ は $Y(0)$ を維持するとする。図18に、式(11)で示される予測音声含有率 $Y(z)$ を求める方法を模式的に示す。図18に示すように、現時点より過去数分の音声含有率は、フレーム毎に移動しながら算出される。また、予測音声含有率 $Y(z)$ や $Y(z-1)$ がセクション毎に予測される。図18に示すように、音声含有率予測部18は、一つ前のセクションでの予測音声含有率 $Y(z-1)$ と、現時点より過去数分の音声含有率 $X(z-1)$ を用いて予測している。

30

#### 【0103】

ここで音声含有率 $X(z)$ の算出時間が短すぎると、音声区間長が長いものであれば、音声含有率が1となってしまう。また、音声と音声の間にある短いポーズ区間のみを抽出してしまい、音声含有率が0に近くなるなど、極端な値をとる可能性がある。また、算出時間が長すぎると、平滑化してしまい、音声含有率の予測に利用できなくなる。そのため、音声含有率 $X(z)$ の算出時間は、コンテンツの音声区間の集まり具合を適度に表す必要があり、実験の結果、1分以上が望ましいことがわかった。したがって、上述では過去数分としている。図19に音声含有率 $X(z)$ と予測音声含有率 $Y(z)$ の算出結果の一例を示す。このグラフは、縦軸に音声含有率を横軸に番組の経過時間を示したものである。また図19においては、30分番組のニュース番組について、音声含有率 $X(z)$ と式(11)によって算出した予測音声含有率 $Y(z)$ を図示したものである。図19に示すように、予測音声含有率 $Y(z)$ は、実際の音声含有率 $X(z)$ とほぼ同じ値に推移することが分かる。

40

#### 【0104】

圧伸比算出部20は、セクション単位で圧伸比を算出する。具体的には、まず現時刻tにおける圧伸比を算出する。現時刻tにおける圧伸比は、現時刻tにおいて速度変換部1

50

7 から出力された出力データ量を速度変換部 17 に入力された入力データ量で除算することと求められる。次に圧伸比算出部 20 は、現時刻  $t$  がセクション境界に達したどうかを判断する。セクション境界に達したとき、次のセクションでの速度比条件を設定するため、次のセクションにおける圧伸比（以下、セクション圧伸比と称す）を算出する。セクション圧伸比とは、次のセクションをどのくらいの圧伸比で変換するかを定めたものである。圧伸比算出部 20 は、ユーザ又は機器が予め設定した目標圧伸比と、現時刻  $t$  における圧伸比とから、式 (12) を用いて算出する。式 (12) において、 $R_t$  は目標圧伸比、 $R(t)$  は現時刻  $t$  における圧伸比、次のセクションの時間長を  $T(z)$ 、 $R_s(z)$  は次のセクションの圧伸比とする。

【数 12】

10

$$R_{s(z)} = R_t + (R_t - R(t)) \frac{t}{T(z)} \quad \cdots (12)$$

【0105】

速度比条件設定部 14 は、音声含有率予測部 18 で算出された予測音声含有率と、圧伸比算出部 20 で算出されたセクション圧伸比を入力とする。速度比条件設定部 14 は、第 1 の実施形態で説明した速度比条件設定部 14 と同様の方法で、セクション毎に速度比条件を設定する。

【0106】

20

ここで、上述した逐次的な処理を行う場合、従来技術では局所的な音声区間の偏りが生じた場合に、その偏りが生じた箇所において再生時間を達成しようと処理するので、非音声区間の極端な削除や高速化が局所的に生じていた。これに対し、本実施形態では、セクション圧伸比がセクション単位で算出される。また速度比条件設定部 14 は、セクション単位で、速度比条件を設定する。つまり、速度比条件は、セクション圧伸比によってセクション単位に更新される。これにより、偏りが生じた箇所において再生時間を達成しなくても、次以降のセクションへ持ち越すことができるので、目標時間を達成しつつ、聞き取り易い再生を実現した逐次的な処理を行うことができる。

【0107】

なお、予測音声含有率をそのまま利用すると、予測音声含有率の増減と平均速度比の増減が直結してしまう。予測音声含有率が高い部分で平均速度比が速くなることは聞こえに影響を与える可能性がある。なぜならば、例えば発話速度が同じと仮定すると、音声含有率が高いほど、情報量が多いため、文章を理解するのが難しいからである。そこで、予測音声含有率が増加すると平均速度比を下げ、予測音声含有率が減少すると平均速度比が上がるように、予測音声含有率を下記のように調整してもよい。

【数 13】

30

$$W_{(z)} = W_{(z-1)} + \gamma(Y_{(z-1)} - Y_{(z)}) \quad \cdots (13)$$

40

$W(z)$  は、セクション  $z$  における、調整後の予測音声含有率である。以下、調整音声含有率と呼ぶ。 $W(z-1)$  は、セクション  $z$  の 1 つ前のセクションにおける調整音声含有率である。 $Y(z)$  は、セクション  $z$  における予測音声含有率である。 $Y(z-1)$  は、1 つ前のセクションにおける予測音声含有率である。  $\gamma$  は、加算する際の係数である。

【0108】

この調整音声含有率を用いて速度比条件を設定することで、予測音声含有率が 1 つ前のセクションにおける予測音声含有率よりも上げれば、音声区間の平均速度比が下がる。また予測音声含有率が 1 つ前のセクションにおける予測音声含有率と同じ場合、音声区間の平均速度比は変化しない。

【0109】

50

このような調整音声含有率を利用することで、音声含有率の高いセクションではより遅くすることが可能となり、情報量に応じた速度条件の設定が可能となる。なお、このような調整を行うことで、セクション圧伸比を達成できずに誤差が生じる恐れがあるが、式(12)に示すように、次以降のセクションにおいて誤差を解消することができる。これにより、目標圧伸比は達成することができる。例えば、図19に示したように、予測音声含有率は一定ではなく、高いところもあれば低いところもある。そのため、予測音声含有率が高いセクションで再生速度が遅くなり、目標圧伸比と差が広がったとしても、次以降の予測音声含有率が低いセクションでこの差を解消することができる。

#### 【0110】

音声区間長予測部19は、音声非音声判別部11の過去の判別結果から、音声区間長を算出し、音声区間長の予測を行う。音声区間長の予測値は、一つ一つの音声区間長の予測値ではなく、文を話す際の平均的な音声区間長を代表値として予測する。一つ一つの音声区間長は、話者交替や会話内容などの様々な要因によって、予測することは難しい。そこで、平均的な音声区間長を予測値として利用して、コンテンツに適した音声区間の速度比制御を行う。図9及び図10に示したように、音声区間長の分布はジャンル毎に大きく異なる。図9は、ニュース番組での各音声区間長の頻度を示したものであるが、500msをピークとし、4000ms辺りまでゆるやかに頻度が減っている。一方、図10は、野球番組での各音声区間長の頻度を示したものであるが、同じく500msをピークとしているものの、急激に頻度が減少していく様子がみられる。この減少の仕方の違いにより、これらの番組を視聴した際の印象としては、ニュース番組では長めの音声区間が続き、野球番組では短い音声区間が続くように聞こえる。そのため、音声区間長を予測せずに、固定長により音声区間の速度比制御を行うと、音声区間の速度比が必要以上に遅くなりすぎたり、速すぎたりする恐れがある。これは、音声非音声判別部11が逐次的に音声か非音声かを判別しているために、音声区間の終端時刻が把握できないからである。また、音声区間長の違いは、話者や会話内容の違いによるものであり、コンテンツ毎に異なるものとなる。このような理由から、音声区間長の予測が必要となる。

#### 【0111】

そこで、処理を開始してから $n$  ( $n$ は自然数)番目の音声区間がもつ音声区間長の予測値 $L(n)$ は、一つ前の $n-1$ 番目の音声区間がもつ音声区間長の予測値 $L(n-1)$ と実測値 $M(n-1)$ とから、式(14)で表現される。

#### 【数14】

$$L_{(n)} = \frac{n-1}{n} L_{(n-1)} + \frac{1}{n} M_{(n-1)} \quad \cdots (14)$$

なお、音声区間には、「はい」や「うん」など相槌や、「えー」などのフィラーなどが含まれる。これらは言語的な理解が容易なため、速度比に関わらず聞き取り易い。そこで、式(14)の音声区間長の予測値 $L(n)$ の算出には、所定の閾値以上の音声区間長をもつ音声区間を利用するとした。ここでは、所定の閾値の一例として、1000msを採用する。

#### 【0112】

式(14)に基づいて予測した音声区間長を図20に示す。図20において、縦軸は音声区間長であり、横軸は経過時間を示している。図20では、X軸において音声区間長の始端時刻の位置に、その音声区間長を示している。また図20では、音声区間長の実測値 $M(n)$ 、直前の音声区間長の実測値 $M(n-1)$ 、及び予測音声区間長 $L(n)$ と図示している。このように、相槌やフィラーなどを除いた音声区間長から予測音声区間長を算出することで、速い速度比では聞き取りにくくなる音声区間長が長い音声区間に適した速度比を設定することができる。また、コンテンツに応じて音声区間長が逐次的に算出されることで、音声区間長の分布の違いに応じた速度比の設定が可能となる。速度比決定部16は、逐次的な速度比条件と予測音声区間長とに基づき、音声及び非音声区間の変換速度



比を決定する。

【0113】

以下、図21を参照して、第2の実施形態に係る音声再生装置の処理について説明する。図21は、第2の実施形態に係る音声再生装置の処理の流れを示すフローチャートである。

【0114】

まず、入力装置（図示なし）においてユーザによるコンテンツを再生する指示を受け付けたか否かが判断される（ステップS201）。ユーザがコンテンツを再生する指示をしたとき、音声非音声判別部11にオーディオ信号が入力される。音声非音声判別部11は、入力されたコンテンツのオーディオ信号について音声区間と非音声区間とをフレーム毎に判別する（ステップS202）。 10

【0115】

ステップS202の次に、音声含有率予測部18は、音声非音声判別部11から出力されたフレーム毎の判別結果から、現時点より過去数分の音声含有率を算出し、算出した音声含有率を用いて、現時点より1つ進んだセクションにおける音声含有率を予測する（ステップS203）。

【0116】

ステップS203の次に、速度比条件設定部14は、音声含有率予測部18で算出された予測音声含有率と、圧伸比算出部20で算出されたセクション圧伸比を入力とし、セクション毎に速度比条件を設定する（ステップS204）。また、音声区間長予測部19は、音声非音声判別部11の過去の判別結果から、音声区間長を算出し、音声区間長の予測を行う（ステップS205）。 20

【0117】

ステップS205の次に、速度比決定部16は、音声非音声判別部11から出力される音声区間の始末端時刻を参照して、所定の単位時間毎に音声区間であるか否かを判断する（ステップS206）。音声区間と判断した場合、速度比決定部16は、音声区間における経過割合を算出する（ステップS207）。音声区間の経過割合とは、音声区間の始端を0、終端を1として、始端からの経過時間を音声区間長で除算したものである。本実施形態では、音声非音声判別が逐次的に行われているため、音声区間の始端時刻は把握できるが、音声区間の終端時刻は現時点では分からない。そこで、音声区間長予測部19で予測された予測音声区間長を音声区間長として用いる。これにより、速度比決定部16は、音声区間における経過割合を算出することができる。なお、音声区間長として実際の値ではなく、予測音声区間長を用いるので、実際の音声区間の経過割合とは必ずしも一致しない。従って、音声区間の経過割合が1以下であっても音声区間の終端時刻となる可能性がある。そこで、速度比決定部16は、音声区間の経過割合が1を越えていないかどうかを判断する（ステップS208）。1を越えていない場合、速度比決定部16は、音声区間の経過割合から音声区間の変換速度比を決定する（ステップS209）。ステップS209の処理は、第1の実施形態の速度比決定部16で説明した処理と同様であるので、説明を省略する。ステップS208において経過割合が1を超えた場合、処理はステップS210へ進み、速度比条件である終端速度比を変換速度比に算出する。この場合、音声区間長予測部19で予測された音声区間長を超過した状態であるため、音声区間の終端速度比を変換速度比として算出する必要がある。 30 40

【0118】

ステップS206において非音声区間と判断した場合、速度比決定部16は、非音声区間の変換速度比を決定する（ステップS211）。ステップS209、S210、及びS211の次に、速度比決定部16は、速度変換対象となるコンテンツの終端時刻まで変換速度比を決定したか否かを判断する（ステップS212）。終端時刻ではないとき、処理はステップS202へ戻る。このように、速度変換対象となるコンテンツの終端時刻までの変換速度比が算出されるまで、速度比決定部16においてステップS202～S212までの処理がセクション単位で繰り返される。ステップS212においてコンテンツの終 50

端時刻まで変換速度比が算出されたと判断された場合、入力装置（図示なし）が本装置の処理を終了するか否かの指示を受け付ける（ステップS213）。ユーザが他のコンテンツについて速度変換処理を行う場合（ステップS213でNo）、処理はステップS202へ戻る。

#### 【0119】

以上のように、本実施形態に係る音声再生装置によれば、リアルタイムで処理を行いながら速度変換を行うことができる。また第1の実施形態に比べ、音声含有率や音声区間長に予測値を用いている。このため、実測値との誤差が生じるが、この誤差は圧伸比設定部28で設定されるセクション圧伸比によって解消される。これにより、本実施形態に係る音声再生装置によれば、リアルタイムで処理を行いながら、目標圧伸比を達成しつつ、区間削除や音声区間の極端な高速化をせずに、速度変換を行うことができる。

10

#### 【0120】

##### （第3の実施形態）

図22を参照して、本発明の第3の実施形態に係る音声再生装置について説明する。図22は、第3の実施形態に係る音声再生装置の構成例を示すブロック図である。図22において、本音声再生装置は、音声非音声判別部11、音声含有率予測部18、速度比条件設定部14、音声区間長予測部19、速度比決定部16、速度変換部17、圧伸比算出部20、及び統計量算出部21で構成される。本実施形態は、第2の実施形態に係る音声再生装置に対し、統計量算出部21を新たに備え、速度比条件設定部14の処理が異なる。以下、統計量算出部21と、速度比条件設定部14の処理を中心に説明する。

20

#### 【0121】

統計量算出部21は、音声区間の上限速度比を修正するための統計量を算出している。例えば、コンテンツの始端から現時点までの音声含有率を利用する。このような音声含有率を以下、長期音声含有率と称す。コンテンツ毎の長期音声含有率の時間変化を図23に示す。図23において、縦軸は長期音声含有率を示し、横軸は始端（0分）からの経過時間を示している。また、コンテンツ毎の予測音声含有率を図24に示す。予測音声含有率は、第2の実施形態において説明した予測音声含有率と同じである。ここでは、算出間隔を1分としている。予測音声含有率のグラフは、音声区間が密集している部分や疎の部分で反映され、山谷がはっきりしたグラフとなっている。長期音声含有率のグラフは、始端付近で多少変動があるものの、概ね平坦であり、第1の実施形態で用いたコンテンツ全体に対する音声含有率に近いグラフとなる。

30

#### 【0122】

そこで、この長期音声含有率を用いて音声区間の上限速度比の修正を行うことを考える。音声区間の上限速度比の修正を逐次行っていくことで、予測音声含有率が局所的に高くなった場合でも、音声区間の速度比が上がりすぎることを防ぐことができる。

#### 【0123】

速度比条件設定部14は、音声含有率予測部18で算出された予測音声含有率、統計量算出部21で算出された長期音声含有率、及び圧伸比算出部20で算出されたセクション圧伸比を入力とする。上述した図21のステップS204において、速度比条件設定部14は図25に示す処理を行う。図25は、第3の実施形態に係る速度比条件設定部14の処理を示すフローチャートである。

40

#### 【0124】

図25において、速度比条件設定部14は、入力される音声含有率が予測音声含有率であるか否かを判断する（ステップS301）。予測音声含有率が入力された場合、処理はステップS302へ進み、速度比条件設定部14は、予測音声含有率を用いて音声及び非音声の平均速度比を算出する。また速度比条件設定部14は、算出した音声区間の平均速度比から終端速度比を算出する（ステップS303）。なお、予測音声含有率に基づく終端速度比をVend1とする。ステップS302及びS303の処理は、上述した第1の実施形態と同様の処理である。

#### 【0125】

50

一方、ステップ S 3 0 1 において予測音声含有率が入力されない場合、つまり長期音声含有率が入力された場合、速度比条件設定部 1 4 は、長期音声含有率を用いて音声及び非音声の平均速度比を算出する（ステップ S 3 0 4）。また速度比条件設定部 1 4 は、算出した音声区間の平均速度比から終端速度比を算出する（ステップ S 3 0 5）。なお、長期音声含有率に基づく終端速度比を  $V_{end2}$  とする。ステップ S 3 0 4 及び S 3 0 5 の処理は、上述した第 1 の実施形態と同様の処理である。

#### 【 0 1 2 6 】

ステップ S 3 0 5 の次に、速度比条件設定部 1 4 は、長期音声含有率に基づいて算出した音声区間の終端速度比  $V_{end2}$  を上限速度比として設定する（ステップ S 3 0 6）。ステップ S 3 0 3 及び S 3 0 6 の次に、速度比条件設定部 1 4 は、予測音声含有率に基づき終端速度比  $V_{end1}$  と、長期音声含有率に基づく上限速度比  $V_{end2}$  とを比較する（ステップ S 3 0 7）。終端速度比  $V_{end1}$  が上限速度比  $V_{end2}$  を超える場合（ $V_{end1} > V_{end2}$ ）、速度比条件設定部 1 4 は終端速度比  $V_{end1}$  を上限速度比  $V_{end2}$  に修正する（ステップ S 3 0 8）。またこの修正に併せて、速度比条件設定部 1 4 は音声区間の平均速度比も長期音声含有率によって算出された値に修正する。つまり、音声及び非音声区間の平均速度比、音声区間の終端速度比の 3 つの速度比を表す速度比条件のうち、非音声区間の平均速度比についてのみ予測音声含有率で算出された値を用いる。それ以外の音声区間の平均速度比及び終端速度比は、長期音声含有率によって求められた値を用いる。

#### 【 0 1 2 7 】

このように、音声区間の平均速度比及び終端速度比の修正を逐次行っていくことで、予測音声含有率が局所的に高くなって音声区間の平均速度比が高めに設定され得る場合でも、聞き易い速度比での再生を行うことができる。

#### 【 0 1 2 8 】

音声区間長予測部 1 9 は、音声非音声判別部 1 1 の過去の判別結果から、予測音声区間長を算出する。本実施形態では、予測音声区間長として音声区間長の最大値を利用する。これは、聞き取り易さ重視の観点から、どのような音声区間であっても漏らさずに終端まで速度比制御を行えるようにするためである。

#### 【 0 1 2 9 】

図 2 6 に示すように、音声区間長の分布は経過時間によって大きく異なる。このため、音声区間長の予測が必要となる。図 2 6 は、音声区間長の実測値と、直前の音声区間長の実測値と、予測音声区間長の分布を示した図である。図 2 6 において、縦軸は音声区間長を示し、横軸はコンテンツの始端からの経過時間を示している。また音声区間長は音声区間の始端時刻に表示している。ここで、 $n$  番目の音声区間長の実測値を  $M(n)$  とする。予測音声区間長  $Lm(n)$  は式 (15) ~ 式 (17) のように表現される。

#### 【 数 1 5 】

$$Lm_{(1)} = \max \quad \cdots (15)$$

#### 【 数 1 6 】

$$Lm_{(n)} = Lm_{(n-1)} + \beta \quad , \quad M_{(n-1)} < Lm_{(n-1)} \quad \cdots (16)$$

#### 【 数 1 7 】

$$Lm_{(n)} = M_{(n-1)} \quad , \quad M_{(n-1)} \geq Lm_{(n-1)} \quad \cdots (17)$$

10

20

30

40

50

式(15)において、maxはコンテンツに含まれる音声区間長のうち最大の音声区間長を複数のコンテンツについて平均した値である。事前にジャンル情報が得られる場合は、ジャンル情報毎に上記平均値を算出し、テーブルを用意しておく。

【0130】

は、予測音声区間長 $L_m$ が次の音声区間までの経過時間とともに減少するように設定された値である。 $n-1$ 番目の音声区間の始端時刻を0とし、 $n$ 番目の音声区間の始端時刻を $t$ とすると、式(18)のように表せる。 $k$ は正の値をとるものとする。

【数18】

$$\beta = -kt \quad \dots (18)$$

10

なお、 $\beta$ は指数関数でもよく、経過時間 $t$ の減少関数であればよい。

【0131】

式(15)～式(18)により予測された予測音声区間長は、図26に示すようになる。図26に示すように、予測音声区間長が音声区間長の実測値よりも長いものが多い。この値を速度比算出時に用いることで、音声区間の終端時刻では終端速度比で変換される割合が低下し、音声区間の平均速度比が更に下がる効果を有する。その結果、聞き取り易い再生を提供することができる。

【0132】

20

(第4の実施形態)

図27を参照して、本発明の第4の実施形態に係る音声再生装置について説明する。図27は、第4の実施形態に係る音声再生装置の構成例を示すブロック図である。図27において、本音声再生装置は、音声非音声判別部11、一時蓄積部22、音声含有率算出部13、速度比条件設定部14、音声区間長算出部15、速度比決定部16、速度変換部17、及び圧伸比算出部20で構成される。本実施形態は、第1の実施形態に係る音声再生装置に対し、蓄積部12よりも蓄積量が少ない一時蓄積部22と、第2の実施形態で説明した圧伸比算出部20を備える点で異なる。以下、異なる点を中心に説明する。

【0133】

一時蓄積部22は、ハードディスク、DVD、又はメモリ媒体(例えばSDカード)などの読み書き可能な記録媒体で構成される。一時蓄積部22には、音声非音声判別部11に入力されるのと同じオーディオ信号がセクション1個分もしくは数個分蓄積される。そして、一時蓄積部22において蓄積されたオーディオ信号が速度変換処理された後、その速度変換処理されたセクションのオーディオ信号は消去され、新しいセクションのオーディオ信号が蓄積される。ここで、本実施形態に係るセクションとは、所定間隔で区切られた区間だけではなく、所定のイベントで区切られた区間でもよい。例えば、イベントをCMとすると、CM区間と、CMとCMに挟まれた番組区間の2種類のセクションができる。イベントが音楽であれば、音楽区間と、音楽と音楽に挟まれた区間の2種類のセクションができる。また、セクションは、ユーザによって指示された区間であってもよい。

30

【0134】

40

なお、第1の実施形態と同様、1つのセクションを構成するオーディオ信号及びビデオ信号が一時蓄積部22に蓄積されるとき、音声非音声判別部11において判別処理が行われ、当該セクションの判別結果や音声区間の始末端時刻も一時蓄積部22に蓄積される。また、一時蓄積部22には、オーディオ信号やビデオ信号と、判別結果及び音声区間の始末端時刻とが対応付けられて蓄積される。なお、オーディオ信号及びビデオ信号のフォーマットは、どのようなフォーマットであってもかまわない。また、本実施形態に係る音声再生装置が、上述した蓄積部12をさらに備えていてもよい。この場合、蓄積部12においてコンテンツ単位で蓄積されたオーディオ信号や音声非音声判別結果が、セクション単位で読み出され、一時蓄積部22に蓄積されるようにする。

【0135】

50

このような一時蓄積部 22 を設けることで、蓄積されたセクションで実際の音声区間長（実測値）を算出することができる。これにより、実際の音声区間にあわせた速度制御が可能になる。また蓄積されたセクションの実際の音声含有率を求めることができるため、第 2 の実施形態で説明した予測音声含有率を用いる場合に比べて、局所的な変動が少なく、コンテンツ全体の音声含有率と近い値となる。

#### 【0136】

音声含有率算出部 13 は、一時蓄積部 22 で蓄積された判別結果や音声区間の始末端時刻から、セクションの音声含有率を算出する。このセクション内に含まれる音声区間長の和を求め、セクション全体の時間長（以下、セクション長と称す）で除算したものが本実施形態の音声含有率となる。

10

#### 【0137】

以下、図 28 を参照して、第 4 の実施形態に係る音声再生装置の処理について説明する。図 28 は、第 4 の実施形態に係る音声再生装置の処理の流れを示すフローチャートである。

#### 【0138】

まず、入力装置（図示なし）において、ユーザが所望のコンテンツを再生する指示をしたか否かが判断される（ステップ S 401）。ユーザの指示があった場合、コンテンツのオーディオ信号及びビデオ信号がセクション分だけ一時蓄積部 22 に蓄積され、音声非音声判別部 11 は、セクション内のオーディオ信号について音声区間と非音声区間とを判別する（ステップ S 402）。なお、ステップ S 402 において判別された判別結果と音声区間の始末端時刻についても、一時蓄積部 22 に蓄積される。

20

#### 【0139】

ステップ S 402 の次に、音声含有率算出部 13 は、セクションの音声含有率を算出する（ステップ S 403）。速度比条件設定部 14 は、ステップ S 403 で算出された音声含有率、圧伸比算出部 20 で算出されたセクション圧伸比に基づいて、音声区間の平均速度比、非音声区間の平均速度比、及び音声区間の終端速度比を算出する（ステップ S 404 及び S 405）。この処理は、第 1 の実施形態と同様の処理である。次に、ステップ S 406 において速度比条件設定部 14 は、ステップ S 405 で算出した終端速度比  $V_{end1}$  と、ユーザによって指定された又は予め装置に設定された終端速度比の上限速度比  $V_{end2}$  とを比較する。ステップ S 406 において終端速度比  $V_{end1}$  が上限速度比  $V_{end2}$  よりも大きいと判断された場合、速度比条件設定部 14 は、終端速度比  $V_{end1}$  を上限速度比  $V_{end2}$  に修正する（ステップ S 407）。ステップ S 406 において終端速度比  $V_{end1}$  が上限速度比  $V_{end2}$  よりも小さいと判断された場合、処理はステップ S 408 へ進む。

30

#### 【0140】

ここで、セクション音声含有率は、コンテンツの一部を構成するセクション内での値である。したがって、音声区間が局所的に集中するセクションなどが存在すれば、局所的にセクション音声含有率の値が大きくなる場合がある。以上のステップ S 406 及び S 407 の処理を行うことで、セクション音声含有率の値が大きくなり、音声区間の終端速度比が大きくなり過ぎることを防ぐことができる。なお、第 3 の実施形態で説明した統計量算出部 21 で長期音声含有率を算出し、上限速度比を修正するようにしてもよい。

40

#### 【0141】

ステップ S 407 の次に、音声区間長算出部 15 は、音声区間の始末端時刻を入力とし、音声区間長を算出する（ステップ S 408）。速度比決定部 16 は、一時蓄積部 22 に蓄積された音声区間の始末端時刻を参照して、セクションの始端から順に所定の単位時間毎に音声区間であるか否かを判断する（ステップ S 409）。音声区間と判断した場合、速度比決定部 16 は、音声区間の始末端時刻と、ステップ S 408 で算出された音声区間長とに基づき、音声区間における経過割合を算出する（ステップ S 410）。

#### 【0142】

ステップ S 410 の次に、速度比決定部 16 は、音声区間の経過割合から、音声区間の

50

変換速度比を決定する（ステップS 4 1 1）。ステップS 4 1 1の処理は、第1の実施形態と同様である。ステップS 4 0 9において非音声区間と判断した場合、速度比決定部16は、非音声区間の始端から終端まで、速度比条件設定部14で設定された非音声区間の平均速度比を変換速度比として決定する（ステップS 4 1 2）。

#### 【0143】

ステップS 4 1 1及びS 4 1 2の次に、速度比決定部16は、セクションの終端まで変換速度比を算出したか否かを判断する（ステップS 4 1 3）。終端ではないとき、処理はステップS 4 0 9へ戻る。このように、セクションの終端までの変換速度比が算出されるまで、速度比決定部16においてステップS 4 0 9～S 4 1 3までの処理が繰り返される。ステップS 4 1 3においてセクションの終端まで変換速度比が算出されたと判断された場合、速度変換部17において変換速度比に従ってオーディオ信号の速度変換が行われ、速度変換後のオーディオ信号の再生が開始される（ステップS 4 1 4）。速度比決定部16は、速度変換対象となるコンテンツの終端時刻まで再生されたか否かを判断する（ステップS 4 1 5）。終端時刻ではないとき、次のセクション分のオーディオ信号が一時蓄積部22に蓄積され、処理はステップS 4 0 2へ戻る。ステップS 4 1 5においてコンテンツの終端時刻まで再生されたと判断された場合、入力装置（図示なし）が本装置の処理を終了するか否かの指示を受け付ける（ステップS 4 1 6）。ユーザが他のコンテンツについて速度変換処理を行う場合（ステップS 4 1 6でNo）、処理はステップS 4 0 1へ戻る。

#### 【0144】

以上のように、本実施形態に係る音声再生装置によれば、セクション単位で速度変換を行うことができる。ここでユーザが、例えば放送中の番組を最初から録画していたが、その番組放送の途中から視聴可能になったとする。このときユーザは、その番組の冒頭を見逃したのでその冒頭を速度変換して視聴しようとするとき、ユーザは速度変換処理の開始時点を冒頭の時点に指定する。なお、速度変換処理の終了時点は、最新の録画がなされた時点である。つまり、冒頭から最新の録画がなされた時点までの区間が1つのセクションとなる。これにより、ユーザは、冒頭から最新の録画がなされた時点まで速度変換した視聴をすることができ、その後においては通常再生によって視聴を続けることができる。このように、本実施形態に係る音声再生装置によれば、セクション単位で速度変換処理を行うので、コンテンツの録画中であっても、全体の録画終了を待たずに速度変換処理を行うことができる。

#### 【0145】

また本実施形態に係る音声再生装置によれば、一時蓄積部22を備えることにより、音声含有率の実測値を算出することができ、音声含有率として予測値を用いる第2及び第3の実施形態に比べて、より最適な速度比で速度変換を行うことができる。また、一時蓄積部22を備えることにより、音声区間長の実測値を算出することができる。音声区間長に実測値を用いる限り、音声区間長が分からないことによる終端速度比の上がり過ぎは生じず、音声区間長として予測値を用いる第2及び第3の実施形態に比べて、より最適な速度比で速度変換を行うことができる。

#### 【0146】

（第5の実施形態）

図29を参照して、本発明の第5の実施形態に係る音声再生装置について説明する。図29は、第5の実施形態に係る音声再生装置の構成例を示すブロック図である。図29において、本音声再生装置は、音判別部23、蓄積部12、音声含有率算出部13、速度比条件設定部14、音声区間長算出部15、速度比決定部16、速度変換部17、及び特定イベント含有率算出部24で構成される。

#### 【0147】

なお、上述した第1の実施形態では、音声区間及び非音声区間の速度比を算出したが、本実施形態では、コンテンツに含まれる特定イベント区間についてさらに個別の速度比を算出することが可能な音声再生装置について説明する。また本実施形態に係る音声再生装

置は、第1の実施形態に係る音声再生装置に対し、音声非音声判別部11の代わりに音判別部23を備える点と、特定イベント含有率算出部24をさらに備える点で大きく異なる。

#### 【0148】

音判別部23は、オーディオ信号を入力として、特定イベント音を含む特定イベント区間、当該特定イベント区間以外の音声区間及び非音声区間を判別する。特定イベント音とは、個別の音源からの音であってもよいし、複数の音源からの音を一まとめにしたものであってもよい。個別の音源からの音としては、例えば、話者Aからの音声、楽器Bからの音、機器Cからの特定音などが挙げられる。複数の音源からの音を一まとめにしたものとしては、例えば、複数の話者からの音声を一まとめにしたものや、音楽、雑音などが挙げられる。また、特定イベント音は1つとは限らず、複数であってもよい。特定イベント音が複数ある場合、音判別部23は、オーディオ信号を入力として、複数の特定イベント区間、当該各特定イベント区間以外の音声区間及び非音声区間を判別することになる。また、特定イベント音が話者Aや話者Bなどの音声である場合、音判別部23が判別する音声区間は特定イベント区間以外の音声区間を意味することになる。以下では、特定イベント音を音楽と仮定して説明する。

10

#### 【0149】

音判別部23は、オーディオ信号を入力として、特定イベント区間である音楽区間、当該音楽区間以外の音声区間及び非音声区間を判別する。これらの区間を判別する方法としては、例えば「MPEG符号化データからのオーディオインデキシング」<中島康之，陸洋，菅野勝，柳原広昌，米山暁夫（KDDI研究所）2000、信学論D-II，Vol. J83-D-II，No. 5，pp. 1361-1371>に記載された公知の方法がある。この方法では、まずオーディオ信号を有音部と無音部に分類する。そして有音部についてさらに、ベイズ推定を用いて音声・音楽・歓声の3つのカテゴリに分類する。このような方法で、音判別部23は、音楽区間、当該音楽区間以外の音声区間及び非音声区間を判別する。なお、上記歓声は、音楽区間以外の非音声区間に含まれるとする。音判別部23で判別された判別結果や、音声区間の始末端時刻、音楽区間の始末端時刻は、蓄積部12に蓄積される。

20

#### 【0150】

特定イベント含有率算出部24は、蓄積部12に蓄積された特定イベント区間の始末端時刻から特定イベントの含有率を算出する。特定イベント含有率は、コンテンツのオーディオ信号に含まれる特定イベント区間（ここでは音楽区間）の比率を示したものである。以下の説明では、特定イベント含有率を音楽含有率と言い換えて説明する。音楽含有率は、具体的には、所定時間のオーディオ信号に含まれる音楽区間長の和を当該所定時間で除算したものである。ここでは、コンテンツ全体に含まれる音楽区間長の和をコンテンツ長で除算したものとする。

30

#### 【0151】

速度比条件設定部14は、まず目標圧伸比と特定イベント含有率算出部24で算出された音楽含有率から、音楽区間の平均速度比（以下、音楽速度比と称す）を算出する。なお、目標圧伸比は、ユーザによって設定されたものでもよいし、予め装置に設定されたものでもよい。具体的には、速度比条件設定部14は、目標圧伸比に応じて異なる音楽含有率と音楽速度比との対応を示したテーブルや対応関数に基づいて、音楽速度比を算出する。このテーブルや対応関数は、予め用意されているとする。図30は、音楽含有率と音楽速度比との対応を示したテーブルの例を示す図である。図30に示すテーブルは、目標圧伸比が0.5のときの対応関係を示したものである。

40

#### 【0152】

速度比条件設定部14は、算出した音楽速度比に基づいて、音楽区間以外の音声区間及び非音声区間の圧伸比を算出する。ここで音楽含有率を $S_m$ 、音楽速度比を $F$ 、目標圧伸比を $E$ 、音楽区間以外の音声区間及び非音声区間の平均速度比を $G$ とすると、平均速度比 $G$ は式(19)となる。

50

【数 19】

$$G = \frac{F - FSm}{FE - Sm} \quad \dots (19)$$

例えば目標圧伸比が 0.5、音楽含有率が 10% の場合、図 30 により、音楽速度比は 1 倍速となる。したがってこの場合、式 (19) に  $E = 0.5$ 、 $Sm = 0.1$ 、 $F = 1$  を代入すると、 $G = 2.25$  となる。

【0153】

圧伸比は速度比の逆数で表せる。平均速度比  $G$  が 2.25 であるため、音楽区間以外の音声区間及び非音声区間の圧伸比はその逆数 0.44 となる。そこで、この圧伸比 (0.44) を音声区間及び非音声区間についての目標圧伸比とすれば、第 1 の実施形態と同様の方法で、音声区間の平均速度比、非音声区間の平均速度比、音声区間の終端速度比を算出することができる。なお、速度比条件設定部 14 が用いる速度比算出分布は、音楽含有率に応じて設定されるようにしてもよい。

10

【0154】

以下、図 31 を参照して、第 5 の実施形態に係る音声再生装置の処理について説明する。図 31 は、第 5 の実施形態に係る音声再生装置の処理の流れを示すフローチャートである。

【0155】

20

まず、ユーザが入力装置 (図示なし) においてコンテンツを録画する指示をしたとき、当該コンテンツのオーディオ信号及びビデオ信号が蓄積部 12 に蓄積される。このとき、音判別部 23 は、音楽区間、当該音楽区間以外の音声区間及び非音声区間を判別する (ステップ S501)。なお、ステップ S501 において判別された判別結果、音声区間の始終端時刻、及び音楽区間の始終端時刻についても、蓄積部 12 に蓄積される。

【0156】

ステップ S501 の次に、入力装置において、ユーザが所望のコンテンツを再生する指示をしたか否かが判断される (ステップ S502)。ユーザの指示があった場合 (ステップ S502 で Yes)、音声含有率算出部 13 は、指示されたコンテンツの音声含有率を算出する (ステップ S503)。また、特定イベント含有率算出部 24 は、指示されたコ

30

【0157】

ステップ S504 の次に、速度比条件設定部 14 は、目標圧伸比と特定イベント含有率算出部 24 で算出された音楽含有率から、音楽速度比を算出する (ステップ S505)。速度比条件設定部 14 は、算出した音楽速度比に基づいて、式 (19) を用いて音楽区間以外の音声区間及び非音声区間の圧伸比を算出する。そして、算出した圧伸比を用いて、音声区間の平均速度比、非音声区間の平均速度比、及び音声区間の終端速度比を速度比条件として設定する (ステップ S506)。音声区間長算出部 15 は、音声区間の始終端時刻を入力とし、音声区間長を算出する (ステップ S507)。

【0158】

40

ステップ S507 の次に、速度比決定部 16 は、蓄積部 12 に蓄積された音楽区間の始終端時刻を参照して、コンテンツの始端から順に所定の単位時間毎に音楽区間であるか否かを判断する (ステップ S508)。音楽区間と判断した場合、速度比決定部 16 は、ステップ S505 で算出した音楽速度比を変換速度比として決定する (ステップ S509)。つまり、音楽区間の始端から終端までの変換速度比は、音楽速度比で一定となる。

【0159】

ステップ S508 において音楽区間でないと判断した場合、速度比決定部 16 は、音声区間の始終端時刻を参照して、音声区間であるか否かを判断する (ステップ S510)。音声区間と判断した場合、速度比決定部 16 は、ステップ S506 で設定された音声区間の平均速度比と、音声区間の始終端時刻と、ステップ S507 で算出された音声区間長と

50



に基づき、音声区間における経過割合を算出する（ステップS511）。速度比決定部16は、音声区間の経過割合から、音声区間の変換速度比を決定する（ステップS512）。音声区間でないと判断した場合、速度比決定部16は、ステップS506で設定された非音声区間の平均速度比を変換速度比として決定する（ステップS513）。つまり、非音声区間の始端から終端までの変換速度比は、当該非音声区間の平均速度比で一定となる。なお、ステップS511～S513の処理は、第1の実施形態と同様である。

#### 【0160】

ステップS509、S512、及びS513の次に、速度比決定部16は、コンテンツの終端まで変換速度比を算出したか否かを判断する（ステップS514）。終端ではないとき、処理はステップS508へ戻る。このように、コンテンツの終端までの変換速度比が算出されるまで、速度比決定部16においてステップS508～S514までの処理が繰り返される。ステップS514においてコンテンツの終端まで変換速度比が算出されたと判断された場合、速度変換部17において変換速度比に従ってオーディオ信号の速度変換が行われ、速度変換後のオーディオ信号の再生が開始される（ステップS515）。入力装置（図示なし）が本装置の処理を終了するか否かの指示を受け付ける（ステップS516）。ユーザが他のコンテンツについて速度変換処理を行う場合（ステップS516でNo）、処理はステップS502へ戻る。

#### 【0161】

以上のように、本実施形態に係る音声再生装置によれば、特定イベント含有率を算出して特定イベント区間の速度比を設定することで、音楽番組などを視聴するに際し、特定イベント区間である音楽区間をそれ以外の音声区間及び非音声区間よりも遅い速度で再生することができる。これにより、速度変換処理において、音楽を重視した再生を行うことができる。また、特定イベント音を音楽ではなく、コンテンツ中に登場するある話者Aの音声とした場合、話者Aの音声に重点がおかれ、特定イベント区間以外の音声区間よりも遅い速度で話者Aの音声速度変換処理される。例えば、何度も視聴しているコンテンツに対して、話者Aの発言内容を確認したいときなどに、話者Aの発言内容だけ遅い速度で再生を行うことは有用である。また、セキュリティカメラのように長時間記録し続けている場合、雑音部分を特定イベント音として識別し、その部分を高速再生することで、冗長なシーンを見る時間を減らすような使い方も可能となる。このように、ある特定イベント音に対して、個別の速度を設定することにより、その部分の重視を促すように遅い速度で再生したり、冗長な部分を低減するために速い速度で再生を行ったり、用途に応じた速度設定が可能になる。

#### 【0162】

なお、上述した第1～第5の実施形態に係る音声再生装置は、一般的なコンピュータシステム50に音声再生プログラムを実行させることによって実現されてもよい。図32は、音声再生装置がコンピュータシステム50によって実現される構成例を示すブロック図である。

#### 【0163】

図32において、コンピュータシステム50は、CPU51、メモリ52、ハードディスク53、ディスクドライブ装置54、モニタ55、スピーカ56、及び入力装置57で構成される。CPU51は、音声再生プログラムを実行させることによって、上述した蓄積部12及び一時蓄積部22以外の第1～第5の実施形態に係る音声再生装置の各構成部と同一の機能を実現する。メモリ52やハードディスク53は、音声再生プログラムを実行させることによって、蓄積部12及び一時蓄積部22と同一の機能を実現する。

#### 【0164】

ディスクドライブ装置54は、コンピュータシステム50を音声再生装置として機能させるための音声再生プログラムが記憶された記録媒体58から、当該音声再生プログラムを読み出す。音声再生プログラムが任意のコンピュータシステム50にインストールされることにより、コンピュータシステム50を上述した音声再生装置として機能させることができる。

## 【0165】

なお、記録媒体58は、例えばフレキシブルディスクや光ディスクなどのディスクドライブ装置54によって読み取り可能な形式の記録媒体である。また音声再生プログラムは、コンピュータシステム50に予めインストールされていてもかまわない。また音声再生プログラムは、インターネットなどの電気通信回線によって提供されてもよい。また音声再生処理は、全部または一部をハードウェアによって処理される形態であってもよい。

## 【0166】

モニタ55は、ディスクドライブ装置54を介して読み込んだ記録媒体58に記録されたビデオ信号や、ハードディスク53に記録されたビデオ信号などを表示する。スピーカ56は、ディスクドライブ装置54を介して読み込んだ記録媒体58に記録されたオーディオ信号、ハードディスク53に記録されたオーディオ信号、速度変換処理後のオーディオ信号を音に変換して再生する。入力装置57は、例えばキーボードやマウスなどで構成され、目標圧伸比の入力などを受け付ける。

## 【0167】

このように、上述した第1～第5の実施形態に係る音声再生装置は、一般的なコンピュータシステム50に音声再生プログラムを実行させることによって実現される。

## 【0168】

また、上述した第1～第5の実施形態に係る音声再生装置は、LSIなどの集積回路や、専用の信号処理回路を用いて1チップ化したものによって実現されてもよい。また上述した第1～第5の実施形態に係る音声再生装置は、音声再生装置を構成する各構成部の機能に相当するものをそれぞれチップ化したものによって実現されてもよい。なお、ここでは、LSIとしたが、集積度の違いにより、IC、システムLSI、スーパーLSI、ウルトラLSIと呼称されることもある。また集積回路化の手法は、LSIに限るものではなく、専用回路又は汎用プロセッサで実現してもよい。LSI製造後に、プログラムすることが可能なFPGA(Field Programmable Gate Array)や、LSI内部の回路セルの接続や設定を再構成可能なりコンフィギュラブル・プロセッサを利用してもよい。さらには、半導体技術の進歩又は派生する別技術によりLSIに置き換わる集積回路化の技術が登場すれば、当然、その技術を用いて機能ブロックの集積化を行ってもよい。

## 【産業上の利用可能性】

## 【0169】

本発明に係る音声再生装置は、目標時間を達成しつつ、入力されるオーディオ信号に応じた適切な速度変換を行うことが可能なハードディスクレコーダーやDVDレコーダー等のAVコンテンツ視聴用機器、パソコンや携帯電話等のモバイル機器上で動作するアプリケーション等に有用である。また、視聴用途だけではなく、学習コンテンツ再生システム等において内容の理解を容易にするための用途や、セキュリティカメラで撮影された映像等の長時間のコンテンツについて概要の把握を容易にするための用途等にも有用である。

## 【図面の簡単な説明】

## 【0170】

【図1】第1の実施形態に係る音声再生装置の構成例を示すブロック図

【図2】ジャンル別の音声含有率を示した図

【図3】各ジャンルの音声含有率の平均と標準偏差とを示した図

【図4】5種類の算出パターンを示した図

【図5】速度比算出分布の一例を示す図

【図6】音声含有率が0.5のときの目標圧伸比、音声区間の平均速度比、非音声区間の平均速度比を示す図

【図7】音声及び非音声区間の速度比変化を示した模式図

【図8】第1の実施形態に係る音声再生装置の処理の流れを示すフローチャート

【図9】ニュース番組に含まれる音声区間長とその頻度を示した図

【図10】野球番組に含まれる音声区間長とその頻度を示した図

【図 1 1】音声区間の圧伸比の変化を示した図

【図 1 2】2 段階の変換速度比を算出した場合を示す図

【図 1 3】ドキュメンタリーなどの静止画像が多いジャンルについての速度比算出分布の例を示す図

【図 1 4】図 1 3 に示す速度比算出分布において、音声含有率が 0.5 のときの目標圧伸比、音声区間の平均速度比、非音声区間の平均速度比を示す図

【図 1 5】スポーツなど動きの激しいシーンが多いジャンルについての速度比算出分布の例を示す図

【図 1 6】図 1 5 に示す速度比算出分布において、音声含有率が 0.5 のときの目標圧伸比、音声区間の平均速度比、非音声区間の平均速度比を示す図

10

【図 1 7】第 2 の実施形態に係る音声再生装置の構成例を示すブロック図

【図 1 8】式 (11) で示される予測音声含有率  $Y(z)$  を求める方法を模式的に示す図

【図 1 9】音声含有率  $X(z)$  と予測音声含有率  $Y(z)$  の算出結果の一例を示す図

【図 2 0】式 (14) に基づいて予測した音声区間長を示す図

【図 2 1】第 2 の実施形態に係る音声再生装置の処理の流れを示すフローチャート

【図 2 2】第 3 の実施形態に係る音声再生装置の構成例を示すブロック図

【図 2 3】コンテンツ毎の長期音声含有率の時間変化を示す図

【図 2 4】コンテンツ毎の予測音声含有率を示す図

【図 2 5】第 3 の実施形態に係る速度比条件設定部 14 の処理を示すフローチャート

【図 2 6】音声区間長の実測値と、直前の音声区間長の実測値と、予測音声区間長の分布を示した図

20

【図 2 7】第 4 の実施形態に係る音声再生装置の構成例を示すブロック図

【図 2 8】第 4 の実施形態に係る音声再生装置の処理の流れを示すフローチャート

【図 2 9】第 5 の実施形態に係る音声再生装置の構成例を示すブロック図

【図 3 0】音楽含有率と音楽速度比との対応を示したテーブルの例を示す図

【図 3 1】第 5 の実施形態に係る音声再生装置の処理の流れを示すフローチャート

【図 3 2】音声再生装置がコンピュータシステム 50 によって実現される構成例を示すブロック図

【図 3 3】従来の音声再生装置の構成を示したブロック図

【符号の説明】

30

【0171】

11 音声非音声判別部

12 蓄積部

13 音声含有率算出部

14 速度比条件設定部

15 音声区間長算出部

16 速度比決定部

17 速度変換部

18 音声含有率予測部

19 音声区間長予測部

40

20 圧伸比算出部

21 統計量算出部

22 一時蓄積部

23 音判別部

24 特定イベント含有率算出部

50 コンピュータシステム

51 CPU

52 メモリ

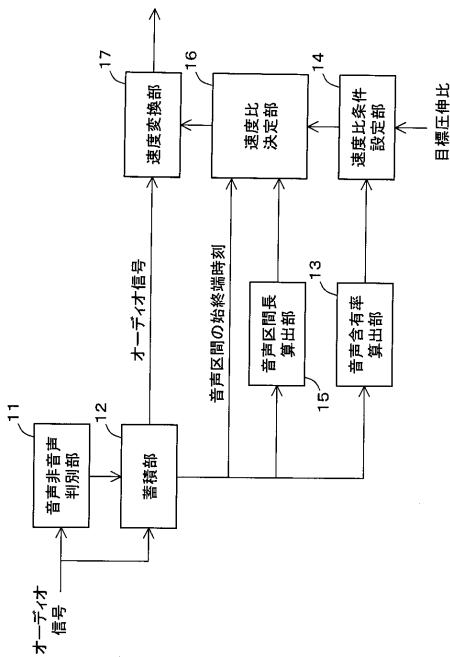
53 ハードディスク

54 ディスクドライブ装置

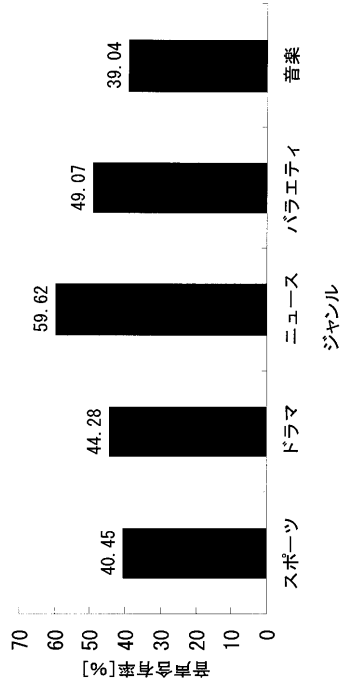
50

- 5 5      モニタ
- 5 6      スピーカ
- 5 7      入力装置

【 図 1 】



【 図 2 】



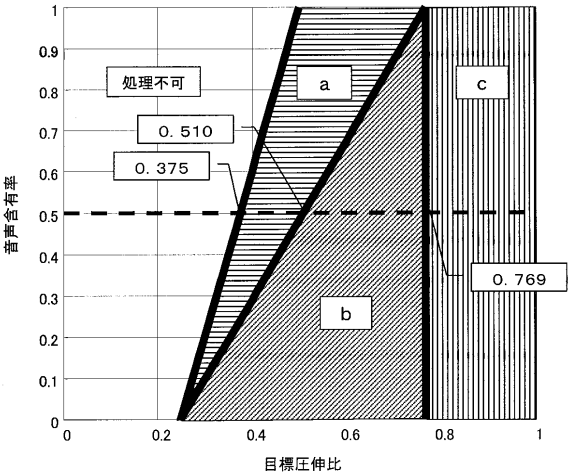
【図 3】

スポーツ	音声含有率	ドラマ・映画・アニメ	音声含有率	ニュース	音声含有率
平均	40.4	平均	44.3	平均	59.6
標準偏差	5.3	標準偏差	16.2	標準偏差	8.2

【図 4】

平均速度比算出 パターン	a	b	c	d	e
音声区間	$Ym1 = \frac{Ar \cdot S}{E - 1 + S}$ $Ym2 = An$	$Ym1 = Bs$ $Ym2 = (1 - S) \cdot \frac{Bs}{Bs \cdot E - S}$	$Ym1 = \frac{1}{E}$ $Ym2 = \frac{1}{E}$	$\frac{Ym1 = Dn \cdot S / (Dn \cdot E - 1 + S)}{Ym2 = Dn}$	$Ym1 = Es$ $Ym2 = (1 - S) \cdot \frac{Es}{Es \cdot E - S}$
非音声区間	$Ym1 \leq Ym2$ $Ym2 = An(一定)$	$Ym1 \leq Ym2$ $Ym1 = Bs(一定)$	$Ym1 = Ym2$	$Ym1 \geq Ym2$ $Ym2 = Dn(一定)$	$Ym1 \geq Ym2$ $Ym2 = Es(一定)$

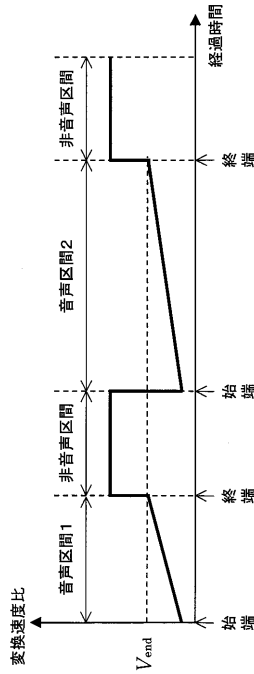
【図 5】



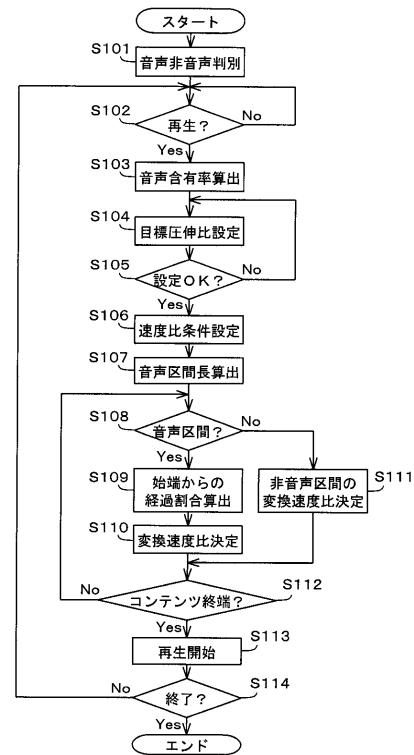
【図 6】

音声含有率0.50のときの平均速度比									
目標圧伸比		Ar=4				Bs=1.3			
		01	03	04	05	06	07	09	1
算出/音声		不可	不可	a	a	b	b	c	c
音声区間		不可	不可	1.818	1.333	1.300	1.300	1.111	1.000
非音声区間		不可	不可	4.000	4.000	2.321	1.585	1.111	1.000
平均速度比									

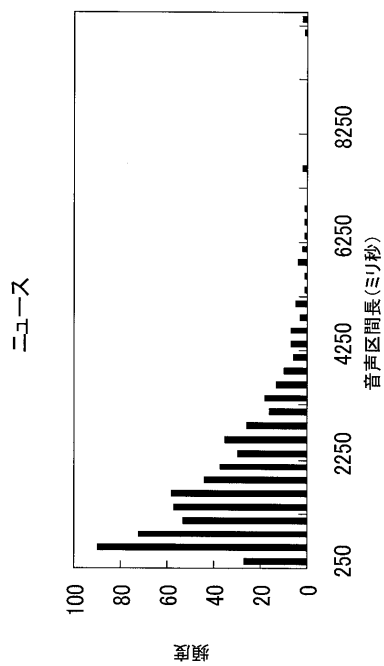
【図 7】



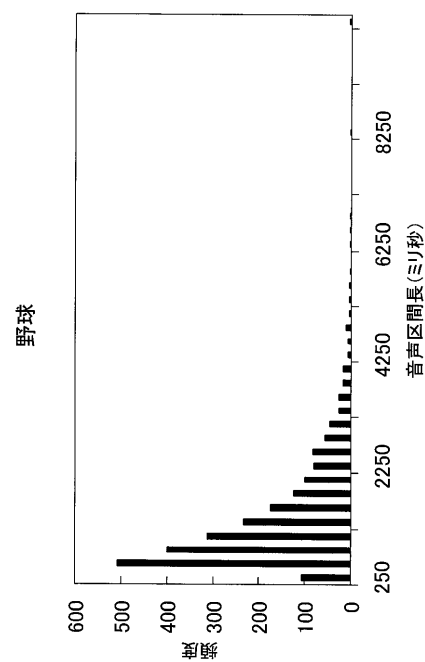
【図 8】



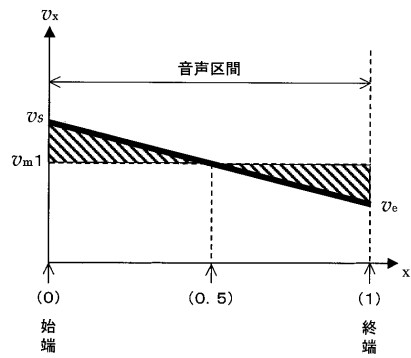
【図 9】



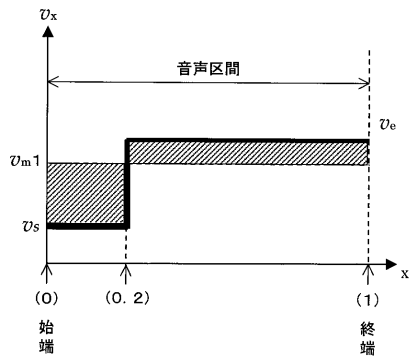
【図 10】



【図 1 1】



【図 1 2】

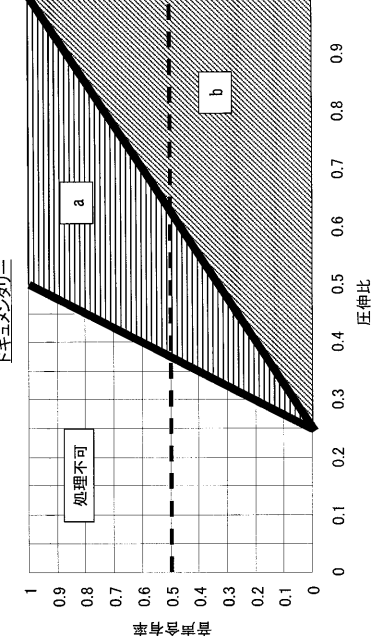


【図 1 4】

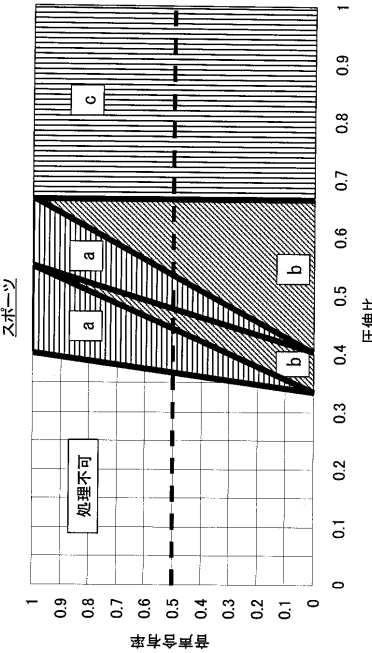
ドキュメンタリー

音声含有率0.5のときの平均速度比									
A <sub>0</sub> = 4					B <sub>0</sub> = 1				
目標伸比	0.1	0.3	0.4	0.5	0.6	0.7	0.9	1	
算出/与え	不可	不可	a	a	a	b	b	b	
平均速度比	不可	不可	1818	1333	1053	1000	1000	1000	
音声区間	不可	不可	4000	4000	4000	2500	1250	1000	
非音声区間	不可	不可							

【図 1 3】



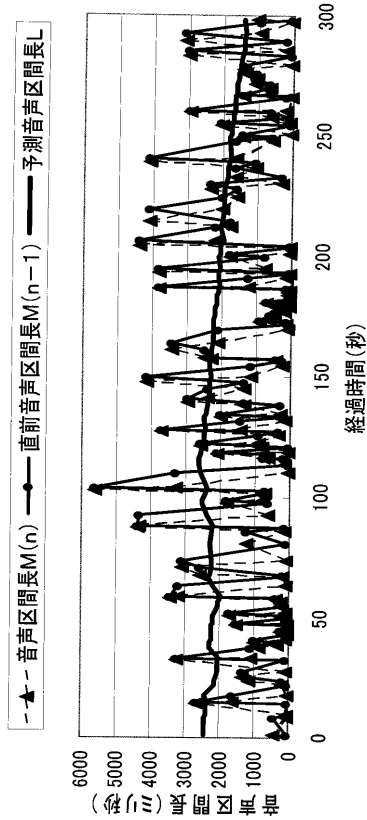
【図 1 5】



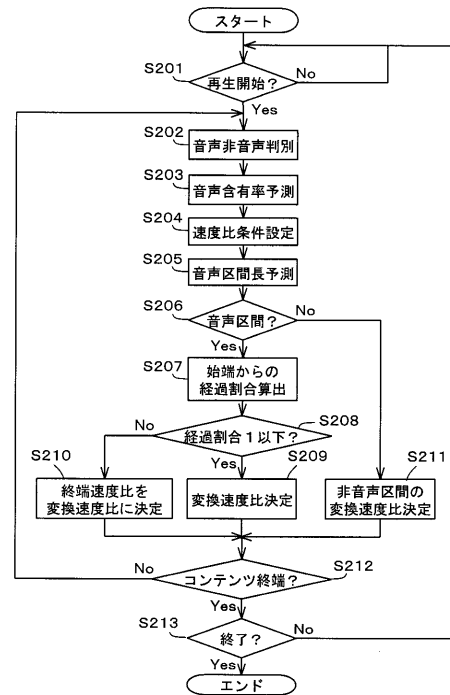




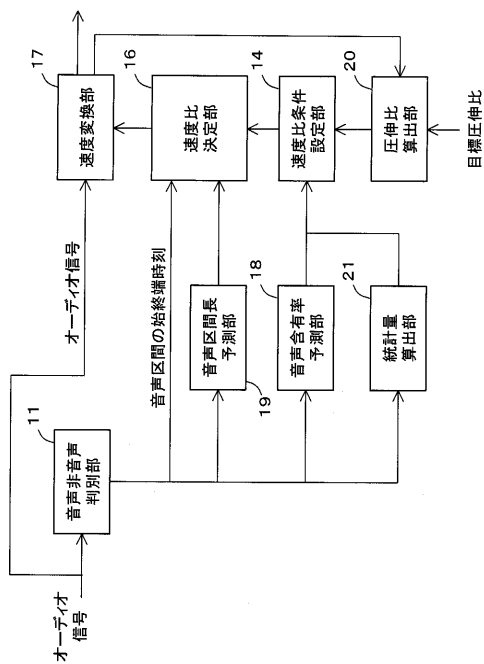
【図 20】



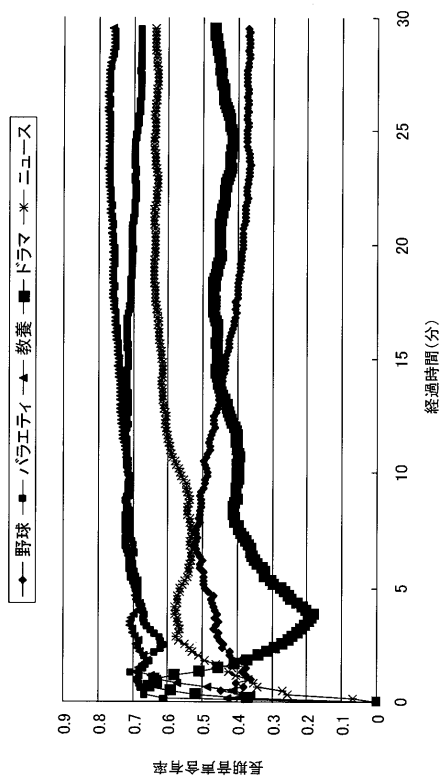
【図 21】



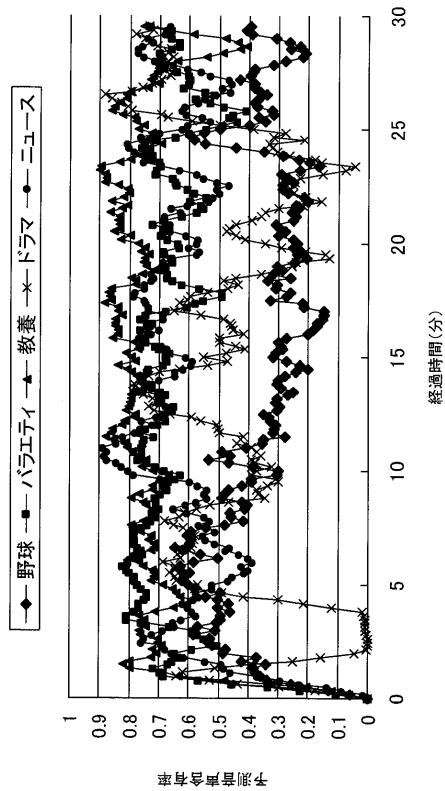
【図 22】



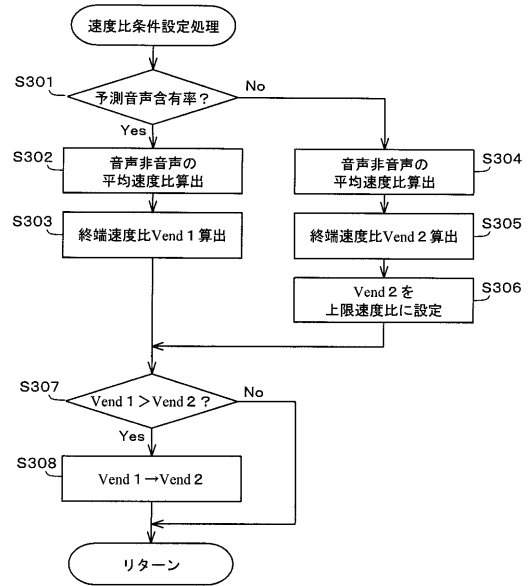
【図 23】



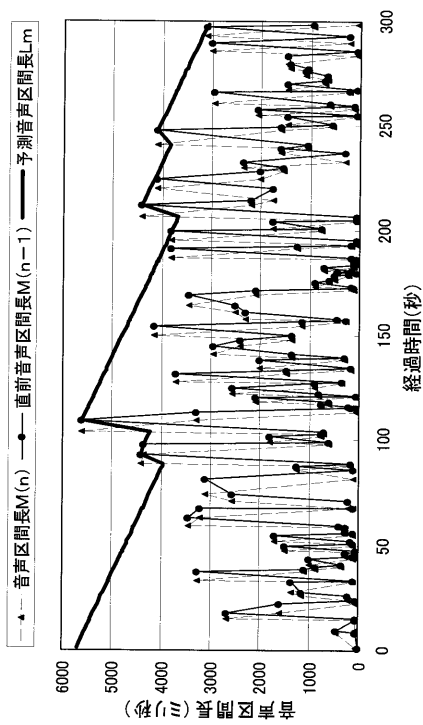
【図 24】



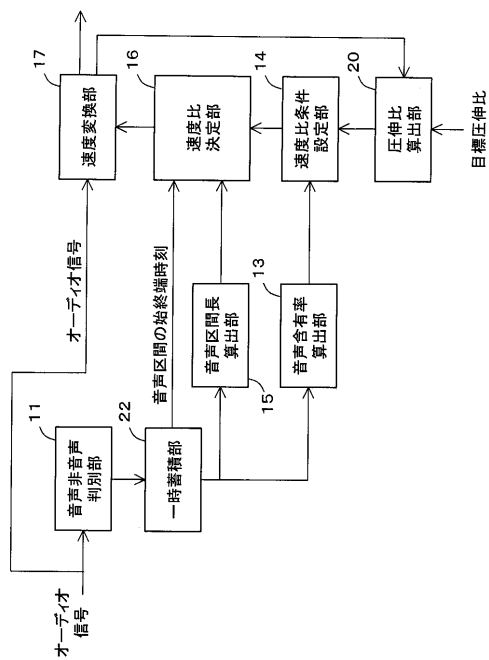
【図 25】



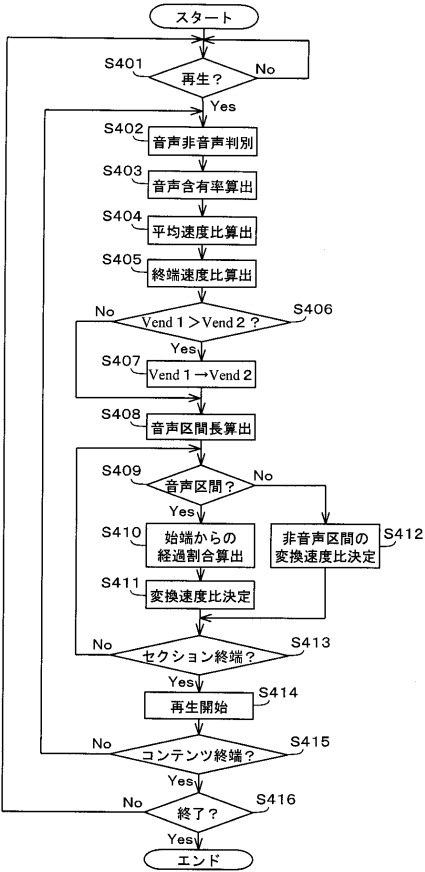
【図 26】



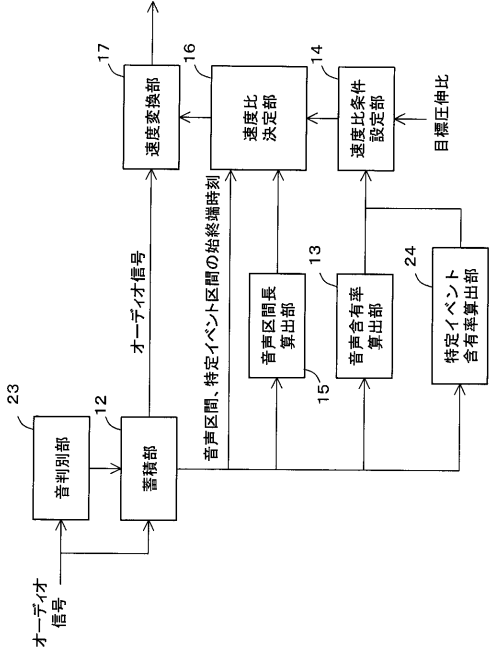
【図 27】



【図 28】



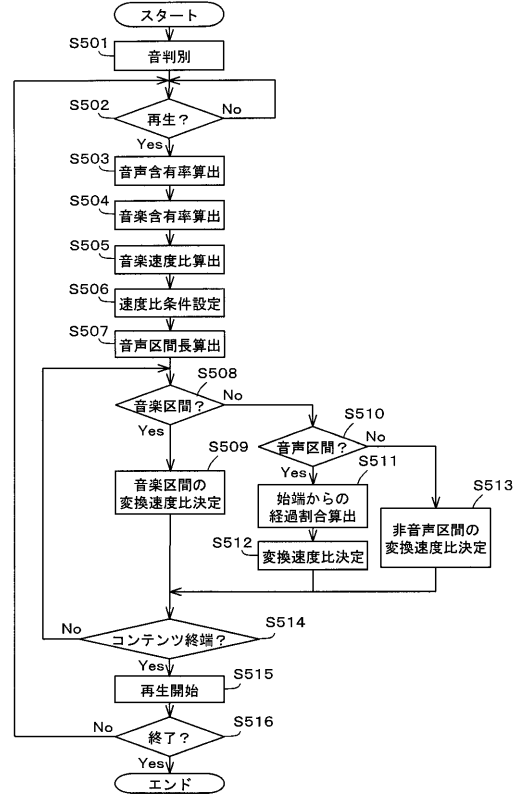
【図 29】



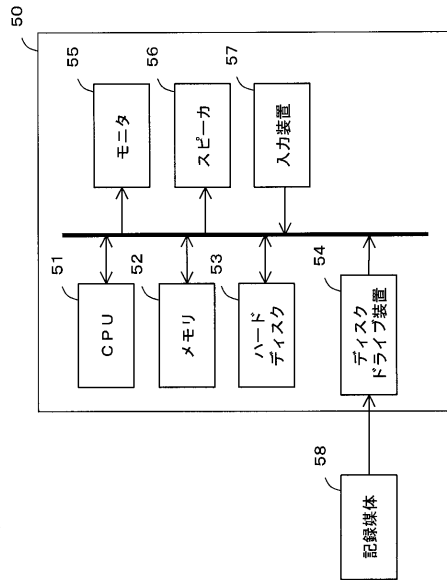
【図 30】

音楽含有率	0%~30%	~40%	~50%	~60%	~70%	~100%
音楽速度比	1	1.1	1.2	1.3	1.5	2

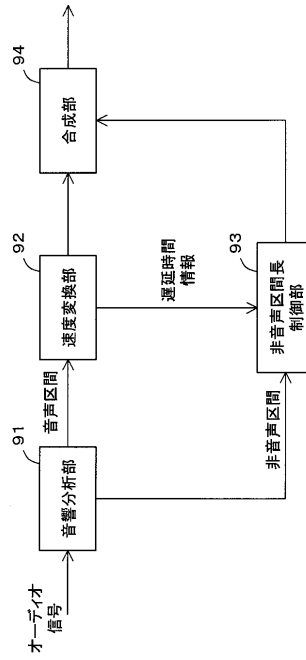
【図 31】



【図 3 2】



【図 3 3】



---

フロントページの続き

審査官 菊池 智紀

- (56)参考文献 特開平08-254992(JP,A)  
特開平04-367898(JP,A)  
特開平06-337696(JP,A)  
特開平08-328586(JP,A)  
特開2001-184100(JP,A)  
特開平09-146587(JP,A)  
特開2005-148307(JP,A)  
特開2005-266571(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G10L 19/00 - 21/06  
G11B 20/10  
JSTPlus(JDreamII)