

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3881971号
(P3881971)

(45) 発行日 平成19年2月14日(2007.2.14)

(24) 登録日 平成18年11月17日(2006.11.17)

(51) Int. Cl. F I
G 1 O L 13/06 (2006.01) G 1 O L 13/06 2 4 O C
G 1 O L 13/04 (2006.01) G 1 O L 13/04 Z

請求項の数 8 (全 16 頁)

(21) 出願番号	特願2003-297306 (P2003-297306)	(73) 特許権者	393031586 株式会社国際電気通信基礎技術研究所
(22) 出願日	平成15年8月21日(2003.8.21)		京都府相楽郡精華町光台二丁目2番地2
(65) 公開番号	特開2005-70214 (P2005-70214A)	(74) 代理人	100099933 弁理士 清水 敏
(43) 公開日	平成17年3月17日(2005.3.17)	(72) 発明者	河井 恒 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
審査請求日	平成16年6月24日(2004.6.24)	(72) 発明者	津崎 実 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
		(72) 発明者	戸田 智基 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

最終頁に続く

(54) 【発明の名称】 声質差評価テーブル作成装置、音声コーパスの声質差評価テーブル作成システム、及び音声合成システム

(57) 【特許請求の範囲】

【請求項1】

複数種類の発話音声データの間の声質の差異を表す声質差評価値テーブルを作成する声質差評価テーブル作成装置であって、前記複数種類の発話音声データの各々には、声質の差異の評価用の音声波形データが含まれており、

前記複数種類の発話音声データから、第1及び第2の発話音声データの任意の組合せを決定するための組合せ決定手段と、

前記組合せ決定手段により決定された組合せの前記第1及び第2の発話音声データの各々から、前記声質の差異の評価用の音声波形データを抽出するための音声波形データ抽出手段と、

前記第1及び第2の発話音声データの声質の差異の評価用の音声波形データをそれぞれ再生した音声からなる第1及び第2の音声刺激を、対比して被験者に提示し、前記第1及び第2の音声刺激に対する聴感上の声質の差異の大きさを評定する入力を被験者より受ける知覚試験を行ない、入力された評定の値から、当該組合せの発話音声データの間の声質の差異の大きさに関する評価を数値で表す評価値を導出するための知覚試験手段とを含み、

前記評定する入力は、声質の差異の大きさが声質の差異の大きさを表す複数段階のカテゴリのいずれに属するかを示すものであり、

前記知覚試験手段により導出された評価値を、当該発話音声データの組合せと対応付けて格納することにより、前記声質差評価値テーブルを作成するためのテーブル作成手段と

を含む、声質差評価テーブル作成装置。

【請求項 2】

前記知覚試験手段は、

前記第 1 及び第 2 の発話音声データの声質の差異の評価用の音声波形データを再生し、再生された音声からなる前記第 1 及び第 2 の音声刺激を知覚試験の対象となる被験者に対して対比して呈示するための刺激呈示手段と、

前記被験者による、前記第 1 及び第 2 の音声刺激の間の聴感上の声質の差異の大きさに関する評定の入力を受けるための入力手段とを含む、請求項 1 に記載の声質差評価テーブル作成装置。

【請求項 3】

前記入力手段は、

前記第 1 及び第 2 の音声刺激の間の聴感上の差異の大きさが、前記複数段階のカテゴリのいずれに属するかを示す評定を前記被験者に入力させるための手段と、

前記被験者により、前記第 1 及び第 2 の音声刺激の間の聴感上の差異の大きさが、前記複数段階のカテゴリのいずれに属するかを示す評定が入力されたことに回答し、当該入力により示されるカテゴリに対して予め割当てられている数値を、前記第 1 及び第 2 の音声刺激の差異の評価値として出力するための手段とを含む、請求項 2 に記載の声質差評価テーブル作成装置。

【請求項 4】

複数種類の発話音声データの間の声質の差異を表す声質差評価値テーブルを作成する声質差評価テーブル作成装置であって、前記複数種類の発話音声データの各々には、声質の差異の評価用の音声波形データが含まれており、

前記複数種類の発話音声データに含まれる前記声質の差異の評価用の音声波形データの各々を所定の手順により定められる順番で抽出するための抽出手段と、

前記抽出手段により抽出された音声波形データの各々に対して、再生時の聴感上の声質を、声質に関する複数個の互いに異なる評定尺度の各々について被験者により評定させる知覚試験を行ない、被験者による評定から、当該発話音声データの声質に関する、前記複数個の評定尺度に対応する複数の評価値を導出することにより、当該複数の評価値を要素とする声質ベクトルを前記複数の発話音声データの各々に対して算出するための知覚試験手段と、

前記知覚試験により、前記複数種類の発話音声データの各々に対して算出される前記声質ベクトルの間に定義される距離として、前記複数種類の発話音声データの任意の組合せの間の声質差に関する評価値を算出し、当該発話音声データの組合せと対応付けて格納することにより前記声質差評価値テーブルを作成するためのテーブル作成手段とを含む、声質差評価テーブル作成装置。

【請求項 5】

前記知覚試験手段は、

前記抽出手段により抽出された、前記声質の差異の評価用の音声波形データの各々を再生して得られる音声を被験者に提示するための再生手段と、

前記再生手段により前記被験者に提示された音声に対する、前記複数個の評定尺度の各々に関する被験者の評定値を取得するための評定値取得手段とを含み、

前記複数個の評定尺度の各々は、複数個のカテゴリにより表され、

前記被験者の評定値は、前記複数個の評定尺度の各々について、当該評定尺度について前記再生手段により前記被験者に提示された音声の声質が、当該評定尺度を表す複数個のカテゴリのうち、どのカテゴリに属するかを示す値であり、

各評定値には、対応の評価値が割り当てられており、

前記知覚試験手段はさらに、前記複数個の評定尺度の各々に関する被験者の評定値に対応するカテゴリに割り当てられている評価値を要素として、前記声質ベクトルを算出するためのベクトル作成手段を含む、請求項 4 に記載の声質差評価テーブル作成装置。

【請求項 6】

10

20

30

40

50

前記テーブル作成手段は、

前記複数個の発話音声データの任意の組合せに対して、当該組合せに属する発話音声データに対して前記ベクトル生成手段により得られるベクトルの間のユークリッド距離を算出するための距離算出手段と、

前記距離算出手段により算出されたユークリッド距離を、対応する発話音声データの組合せと対応付け、当該発話音声データ間の声質の差異の大きさとして前記声質差評価テーブルに格納するための手段とを含む、請求項 5 に記載の声質差評価テーブル作成装置。

【請求項 7】

異なる複数の収録条件下で収録された複数種類の発話音声データからなる音声コーパスと、

前記音声コーパスに含まれる前記複数種類の発話音声データを入力とする、請求項 1 から請求項 6 のいずれかに記載の声質差評価テーブル作成装置とを含み、

前記複数種類の発話音声データの各々には、声質の差異の評価用の音声波形データが含まれている、音声コーパスの声質差評価テーブル作成システム。

【請求項 8】

複数種類の発話音声データからなり、それぞれ音声素片に分離可能な複数の発話音声データを含む音声コーパスと、

予め定める入力情報を取得し、当該入力情報から音声を作成するための音声素片を選択するための評価関数により、一連の音声素片を前記音声コーパスより選択し、抽出するための手段と、

前記抽出するための手段が抽出した一連の音声素片を接続して、発話音声を作成するための手段と、

前記音声コーパスに含まれる前記複数種類の発話音声データを入力として、請求項 1 から請求項 7 のいずれかに記載の声質差評価テーブル作成装置により作成された声質差評価テーブルとを含む音声合成システムであって、

前記複数種類の発話音声データの各々には、声質の差異の評価用の音声波形データが含まれており、

前記抽出するための手段は、前記声質差評価テーブルに格納された、前記複数種類の発話音声データ間の声質差の評価値を前記評価関数の入力として用いて、前記入力情報から音声を作成するための音声素片を選択する、音声合成システム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、音声合成技術に関し、特に、音声コーパスから音声素片を選択し、接続することにより自然な発話に近い音声を合成する音声素片接続型音声合成技術に関する。

【背景技術】

【0002】

コンピュータ技術及びデータコミュニケーション技術の発達に伴い、人間と機械との間のインターフェイスが重要となっている。人間にとっては、人と話をするのと同様に機械とのコミュニケーションを行なうことが望ましく、そのための技術開発が進められている。

【0003】

人間から機械への情報の伝達としては、音声認識、画像認識等の認知技術が主として用いられる。また機械から人間への情報の伝達方法は種々あるが、中でも音声合成技術が用いられる機会が増加しており、音声応答システム、音声翻訳システム等が代表的な応用例である。さらに、近年のロボット等の開発の進展に伴い、音声認識及び画像認識と音声合成とを組合せることで、人間とロボットとのコミュニケーションを人間同士のコミュニケーションと同様に実現することが期待される。

【0004】

図 8 に、音声コーパスを用いる音声素片接続型音声合成システムのブロック図を示す。

10

20

30

40

50

図 8 を参照して、音声コーパスを用いる音声素片接続型音声合成システムでは、人間による自然な発話の音声を収録し、発話の音声素片を音声コーパス 4 0 として予めコーパス化しておく。

【 0 0 0 5 】

このシステムに対して、入力テキスト 4 2 が与えられると、音声素片選択部 4 4 は音声コーパス 4 0 の中から、入力テキストを構成する音声に対応する音声素片を、音声の合成に用いる音声素片の候補として選択する。音声素片選択部 4 4 は、選択した音声素片を評価関数 5 0 によって評価し、その結果に従って、音声の合成に用いるのに最適な音声素片を決定する。このようにして、入力テキスト 4 2 を構成する音声にそれぞれ対応する音声素片を抽出する。音声素片接続部 4 6 が、これら一連の音声素片を接続することにより、

10

【 0 0 0 6 】

評価関数 5 0 には、音声素片選択部 4 4 より、これ以前に選択された音声素片及び候補となっている音声素片について観測可能な物理量に変数として与えられる。評価関数 5 0 は、与えられた物理量に関する評価値を従属変数として出力する。音声素片選択部 4 4 は、評価関数 5 0 により出力された評価値に基づいて、選択した複数の音声素片から、直前の音声素片に接続するのに好適な音声素片を決定する。

【 0 0 0 7 】

このようにして合成された音声は、人間が実際に発声した音声を用いて合成されたものである。そのため、いわゆる「機械音らしさ」を感じさせない比較的自然的な音声を合成することができ

20

【 0 0 0 8 】

一方、音声素片接続型音声合成システムでは、それぞれ別個に収録された音声からそれぞれ抽出した音声素片を接続して、連続的な音声情報を合成する。そのため、接続される音声素片の声質が均質であることが求められる。さもなければ合成した音声に不連続感が生じ、合成された音声の音質は劣化する。よって多くの場合、単一の話者の音声からなる

30

【 0 0 0 9 】

この場合、話者が一人であるため、大規模な音声コーパスを構築するには、長期間かけて音声を収録する必要がある。そのため、1 人の話者の発声を、複数回の収録期間（以下、この収録期間を「セッション」と呼ぶ。）に分けて収録する。場合によってはその収録に数ヶ月から数年の期間を必要とする。

【 0 0 1 0 】

このように、単一の話者による音声を収録した場合であっても、上記の通り、大規模な音声コーパスを作成するには、数ヶ月に及ぶ収録期間が必要となることがある。このように長い収録期間中、音声の収録条件を毎日一定に保つことは極めて困難である。とりわけ話者が体調を一定に保つことは極めて難しい。そのため、これら収録条件の変化に起因して、収録された音声データの声質が、セッション毎に異なるものとなることは、避けられない。よって合成された音声の声質にばらつきが生じ、不連続感が生じるという問題がある。従って、音声波形素片を接続する際の不自然さを解消する技術が望まれている。

40

【 0 0 1 1 】

音声素片を接続する際の不自然さを解消するために、合成に用いる音声波形素片をどのようにして評価し、選択するかが問題となる。通常、各音声波形素片に関連する何らかの音響特徴量を算出し、所定の条件に合致する音声波形素片が選択される。不自然さを小さくするためには、知覚特性にできるだけ一致した尺度を用いて素片選択を行なうことが重要である。

50

【 0 0 1 2 】

後掲の非特許文献 1 では、知覚特性を反映した「コスト関数」と呼ばれる評価関数を用いて、候補の音声素片についてコストを算出し、その算出されたコストが最小となる波形素片を選択している。このようなコスト関数を用いて波形素片を選択することで、より自然な音声を合成できると期待される。

【 0 0 1 3 】

しかし、コスト関数として、どのような物理尺度を用いれば、波形接続時に知覚される不自然さが解消されるかについては明らかではない。即ち、物理尺度と合成音声の自然さとの間の対応関係は明らかでない。そのため非特許文献 1 では、コスト関数を様々な要因に対応する複数のサブコスト関数に分けている。

10

【 0 0 1 4 】

サブコスト関数は、それぞれ対応の物理量（観測可能なもの）が与えられると、その関数としてサブコストを出力する。これらサブコストに重みを乗算し、加算することにより評価値となるコストが算出される。

【 0 0 1 5 】

非特許文献 1 では、韻律に関するサブコスト関数、F0（基本周波数）の不連続に関するサブコスト関数、音素環境代替におけるサブコスト関数、スペクトルの不連続に関するサブコスト関数、音素の適合性に関するサブコスト関数を用いている。そして、これらサブコスト関数のうち、特に知覚評価との関係が比較的分かりやすい要因である音素環境代替に関しては、知覚評価と物理量との間のマッピングを行なっている。しかしその他の要因については知覚評価を用いていない。

20

【 0 0 1 6 】

【非特許文献 1】戸田 智基、河井 恒、津崎 実、鹿野 清宏、「素片接続型日本語テキスト音声合成における音素単位とダイフオン単位に基づく素片選択」、電子情報通信学会論文誌、Vol.J85-D-11., No.12, pp.1760-1770, Dec.2002.

【発明の開示】

【発明が解決しようとする課題】

【 0 0 1 7 】

非特許文献 1 に記載の技術では、音素環境代替による自然性劣化を知覚評価により評価し、その結果をサブコスト関数に反映している。しかし、セッション毎に異なる話者の声質による合成音声の自然性劣化については非特許文献 1 では考慮されていない。これは、セッションの違いに起因する声質の差異と、音声の物理的特徴との間の対応関係が不明か、それを特定するのが極めて困難であるためである。

30

【 0 0 1 8 】

それゆえに、本発明の目的は、収録時期の異なる音声データより抽出された音声素片同士を接続して音声を合成する際に、話者の体調変化等に起因する声質の変化によって生じる、合成された音声の音質劣化を軽減するための声質差評価テーブル作成装置及び音声合成システムを提供することである。

【 0 0 1 9 】

本発明の別の目的は、収録時期の異なる音声データの、話者の体調変化等に起因する声質の差異を良好な感度で評価するための声質差評価テーブル作成装置及び音声合成システムを提供することである。

40

【課題を解決するための手段】

【 0 0 2 0 】

本発明の第 1 の局面に係る声質差評価テーブル作成装置は、複数種類の発話音声データの間の声質の差異を表す声質差評価値テーブルを作成する装置である。この声質差評価テーブル作成装置は、複数種類の発話音声データから、第 1 及び第 2 の発話音声データの任意の組合せを抽出するための抽出手段と、抽出手段により抽出された任意の組合せの発話音声データの各々に基づいて、第 1 及び第 2 の音声刺激を生成するための刺激生成手段と

50

、任意の組合せの発話音声データの各々に対して、第1及び第2の音声刺激に対する聴覚上の声質の差異に関する知覚試験を行ない、当該組合せの発話音声データの間の声質の差異に関する評価値を導出するための知覚試験手段と、知覚試験の結果に基づいて、任意の組合せの発話音声データの間の声質の差異を表わす評価値を、当該発話音声データの組合せと対応付けて格納することにより、声質差評価値テーブルを作成するためのテーブル作成手段とを含む。

【0021】

異なる収録条件により収録された発話音声データの間の声質の差異を、知覚試験によって評価することができる。よって、物理的尺度によって評価することが困難な発話音声の声質の差異を、良好な感度で評価することが可能となる。

10

【0022】

好ましくは、知覚試験手段は、知覚試験の対象となる被験者に対して、第1及び第2の音声刺激を対比して呈示するための刺激呈示手段と、被験者による第1及び第2の音声刺激の間の聴感上の声質の差異の大きさに関する評定結果を、知覚試験の結果として取得するための取得手段とを含む。

【0023】

知覚試験を行なう際に、声質の差異を評価する対象となる発話音声の声質を、被験者に比較させることにより、比較対象となる両者の発話音声データ間の声質の差異を、被験者が明確に評定することができる。

【0024】

より好ましくは、取得手段は、予め定める評定尺度を用いて被験者により評定された聴感上の差異の大きさを得るための手段と、被験者による評定により得られた評定尺度の値に基づいて評価値を導出するための評価値導出手段とを含む。

20

【0025】

評定尺度を用いて、知覚試験を行なうことにより、より明確な知覚試験の結果を得ることができるようになる。

【0026】

評価値導出手段は、評定尺度による評定の結果を予め定められた換算基準により換算することにより評価値を導出するための手段を含んでもよい。

【0027】

さらに、換算基準を定め、これを用いて評価値を導出することにより、聴感上の声質の差異の大きさを数値化することができる。

30

【0028】

本発明の第2の局面に係る声質差評価テーブル作成装置は、複数種類の発話音声データの間の声質の差異を表す声質差評価値テーブルを作成する装置である。この声質差評価テーブル作成装置は、複数種類の発話音声データに含まれる発話音声データの各々を所定の手順により定められる順番で抽出するための抽出手段と、抽出手段により抽出された発話音声データの各々に対して、聴覚上の声質に関する知覚試験を行ない、当該発話音声データの声質に関する評価値を導出するための知覚試験手段と、知覚試験により得られる、複数種類の発話音声データの声質に関する評価値に基づいて、発話音声データの任意の組合せの間の声質差に関する評価値を、当該発話音声データの組合せと対応付けて格納することにより、声質評価値テーブルを作成するためのテーブル作成手段とを含む。

40

【0029】

被験者に、音声刺激に関する聴感上の印象をそれぞれ評定させるため、知覚試験によって、より多角的な評定を得ることができる。また、少ない試験回数で、声質差に関する評価値を導出することができる。

【0030】

好ましくは、知覚試験手段は、抽出手段により抽出された発話音声データの各々に基づいて音声刺激を生成するための刺激生成手段と、刺激生成手段により生成された音声刺激に対する、予め定められる複数個の評定尺度の各々に関する被験者の評価値を取得するた

50

めの評価値取得手段とを含む。

【0031】

複数個の評定尺度を用いて、被験者に評定を行なわせることにより、より多様な聴感上の印象をもとに、声質の差異に関する評価値を導出することが可能になる。

【0032】

より好ましくは、評価値取得手段は、刺激生成手段により生成された音声刺激に対する、複数個の評定尺度の各々に関する被験者の評価値を、それぞれ予め定められた複数個の段階を表わす離散的な値として被験者から取得するための手段と、取得するための手段が取得した、複数個の評定尺度の各々に関する被験者の離散的な評価値を要素とするベクトルを生成するためのベクトル生成手段とを含む。

10

【0033】

被験者による評価をもとにベクトルを生成することにより、被験者による聴感上の印象に関する評価を、数値化することができる。

【0034】

さらに好ましくは、テーブル作成手段は、複数個の発話音声データの任意の組合せに対して、当該組合せに属する発話音声データに対して、ベクトル生成手段により得られるベクトルの間の距離を所定の算出方法に従って算出するための距離算出手段と、距離算出手段により算出された距離を、対応する発話音声データの組合せと対応付けて声質差評価テーブルに格納するための手段とを含む。

【0035】

20

ベクトル生成手段より得られるベクトル間の距離を、声質の差異に関する評価値として用いることにより、被験者による多面的な評定をもとに、声質の差異を数値化することが可能になる。

【0036】

本発明の第3の局面に係る音声コーパスの声質差評価テーブル作成システムは、異なる複数の収録条件下で収録された複数種類の発話音声データからなる音声コーパスと、音声コーパスに含まれる複数種類の発話音声データを入力とする、本発明の第1の局面又は第2の局面に係る声質差評価テーブル作成装置とを含む。

【0037】

この声質差評価テーブル作成システムにより、音声コーパスを用いて行なう様々な音声処理技術において、声質の差異に関する評価に基づく処理を行なうことが可能となる。

30

【0038】

本発明の第4の局面に係る音声合成システムは、複数種類の発話音声データからなり、それぞれ音声素片に分離可能な複数の発話音声データを含む音声コーパスと、予め定める入力情報を取得し、当該入力情報と所定の関係にある音声素片を、音声コーパスより選択し、抽出するための手段と、抽出するための手段が抽出した一連の音声素片を接続して、発話音声を合成するための手段と、音声コーパスに含まれる、これら複数種類の発話音声データを入力として、本発明の第1の局面から第3の局面のいずれかに係る声質差評価テーブル作成装置により作成された声質差評価テーブルとを含む音声合成システムである。この音声合成システムの、抽出するための手段は、声質差評価テーブルに格納された複数種類の発話音声データの間の声質差の評価値に基づいて、所定の関係にある音声素片を選択する。

40

【0039】

この音声合成システムは、複数種類の発話音声データからなる音声コーパスをもとに音声を合成する際に、声質の差異が大きな音声素片を接続することを防止できる。よって、声質の差異が大きな音声素片を接続することにより生じる、合成音声の音質劣化を軽減することができる。

【発明を実施するための最良の形態】

【0040】

以下、図面を参照しつつ、本発明の実施の形態に係る音声素片接続型音声合成システム

50

について説明する。

【 0 0 4 1 】

[第 1 の実施の形態]

図 1 に本発明の一実施の形態に係るシステムの機能的構成をブロック図形式で示す。図 1 を参照して、この音声素片選択システム 7 0 は、図 8 に示す従来技術の音声合成システムに用いられるものと同様の、評価関数 7 2 を用いた評価により音声素片を選択する音声素片選択部 4 4 と、評価関数 7 2 が参照する、セッション間の声質の差異を表わす情報を格納した声質差評価テーブル 1 0 8 を作成する声質差評価テーブル作成装置 1 0 0 とを含む。

【 0 0 4 2 】

声質差評価テーブル作成装置 1 0 0 は、音声素片選択システム 7 0 が音声合成に用いる一連の音声素片 6 0 を選択する際の、選択候補となる音声のデータを格納する音声コーパス 1 0 2 と、音声コーパス 1 0 2 に接続され、音声コーパス 1 0 2 内の音声データを用いて、被験者 1 0 6 に対して知覚試験を行なうことにより、セッション間での声質の差異に関する評価値を導出し、導出した声質差に関する評価値をまとめて声質差評価テーブル 1 0 8 を作成する声質差評価装置 1 0 4 とを含む。声質差評価テーブル 1 0 8 は、評価関数 7 2 による評価の際に参照される。

【 0 0 4 3 】

音声コーパス 1 0 2 には、単一の話者により発話され、収録された 1 セッション分の音声波形信号からなる音声データ 1 1 0 A , 1 1 0 B , ... , 1 1 0 N が、複数セッション分記憶されている。これらの 1 セッション分の音声データには、それぞれセッションを識別するための識別番号 1 1 2 A , 1 1 2 B , ... , 1 1 2 N が付与されている。

【 0 0 4 4 】

図 2 に、声質差評価装置 1 0 4 の構成をブロック図形式で示す。図 2 を参照して、声質差評価装置 1 0 4 は、音声コーパス 1 0 2 に接続され、音声コーパス 1 0 2 に記憶されている各セッションの音声データの識別番号をもとに、被験者 1 0 6 に対して呈示する刺激対の組合せを決定する処理を行なう試験処理部 1 2 2 と、音声コーパス 1 0 2 及び試験処理部 1 2 2 に接続され、試験処理部 1 2 2 からの命令に従い、被験者 1 0 6 に対して提示する刺激対を、音声コーパス 1 0 2 内の音声波形のデータより抽出する刺激抽出部 1 2 0 とを含む。

【 0 0 4 5 】

声質差評価装置 1 0 4 はさらに、刺激抽出部 1 2 0 が抽出した音声波形のデータを再生し、被験者 1 0 6 に対して音声刺激を呈示する刺激呈示部 1 2 4 と、被験者 1 0 6 に対して呈示した音声刺激についての被験者 1 0 6 の評定を取得する評定取得部 1 2 6 と、試験処理部 1 2 2 及び評定取得部 1 2 6 に接続され、試験処理部 1 2 2 が決定する刺激対の組合せと、評定取得部 1 2 6 が取得する被験者 1 0 6 の評定とをもとに、声質差評価テーブル 1 0 8 を作成するテーブル作成部 1 2 8 とを含む。

【 0 0 4 6 】

なお、テーブル作成部 1 2 8 は、一つの刺激対に対する処理が完了すると、試験処理部 1 2 2 に対して完了信号を送る機能を有する。また、試験処理部 1 2 2 は、完了信号を受けると、次の刺激対の処理を開始する機能を有する。

【 0 0 4 7 】

図 3 に、声質差評価テーブル 1 0 8 の構成の一例を示す。図 3 を参照して、声質差評価テーブル 1 0 8 は、セッションごとのエントリ 1 3 0 A , 1 3 0 B , ... を含む。各エントリには、そのセッションの識別番号と、そのセッションで収録された音声と他の各セッションで収録された音声との間の声質の差異を示す評価値がそれぞれ格納される。なお、これらの評価値は、複数の被験者 1 0 6 から得られた評定の結果をセッションの組合せごとに統合したものである。

【 0 0 4 8 】

本実施の形態に係るシステムは、以下のように動作する。なお、本実施の形態に係る声

10

20

30

40

50

質差評価装置 104 が知覚試験を行なうにあたり、被験者 106 は予め適切な方法で選ばれているものとする。また被験者には、予め十分かつ適切な教示が与えられているものとする。

【0049】

図1を参照して、話者が自然に発話した音声は、音声コーパス102に音声データ110A、110B、...、110Nとしてセッション毎に格納されている。これらの音声データ110A、110B、...、110Nにはそれぞれ識別番号が付与されている。また、できるだけ声質差の評価を正確にするため、各セッションの最初には同じ文を読むこととし、これを声質又は声質の差異の評価に用いるものとする。

【0050】

図2を参照して、声質差評価装置104が起動すると、試験処理部122は、音声コーパス102に格納されている音声データの識別番号112A、112B、...、112Nを、音声コーパス102より読出す。試験処理部122は、読出した識別番号をテーブル作成部128に与える。テーブル作成部128は、与えられた識別番号をもとに、声質差評価テーブル108の作成準備を行なう。即ち、図3に示す声質差評価テーブル108で、データの入っていないものを、図示しない記憶装置上に作成する。声質差評価テーブル108の作成準備が完了すると、テーブル作成部128は、試験処理部122に完了信号を与える。

【0051】

試験処理部122は、テーブル作成部128より完了信号を受けたことに応答して、音声コーパス102より取得した識別番号112A、112B、...、112Nの中から、声質の差異を比較させる対象となる2セッションの音声データの識別番号を選び、刺激抽出部120及びテーブル作成部128に与える。この際、どの2つのセッションを選ぶかには、様々な方法がある。比較対象となるべきセッションの対が全て抽出できるものであれば、どのような方法であってもよい。

【0052】

刺激抽出部120は、2セッションの音声データの識別番号が与えられると、音声コーパス102に格納されている音声データの中から、与えられた識別番号にそれぞれ対応する音声データをそれぞれ特定し、特定した音声データから、上記した声質評価用の1発声分の音声波形のデータをそれぞれ抽出する。刺激抽出部120は、抽出した音声波形のデータを対にして刺激呈示部124に与える。

【0053】

刺激呈示部124は、与えられた一对の音声波形のデータを対比する形で再生し、被験者106に呈示する。刺激の呈示が終了すると、刺激呈示部124は、呈示が終了したことを示す信号を評定取得部126に与える。

【0054】

被験者106には、刺激対を構成する2つの刺激間の聴感上の差異を評定するための評定尺度が予め示されている。図4に、評定尺度の一例を示す。図4を参照して、評定尺度160は、5段階のカテゴリによって構成された尺度である。これらのカテゴリには、それぞれカテゴリに対応する、数値で表わされている評定値が予め付与されている。なお、この評定尺度160は、印刷された評定用紙によって与えられてもよい。また、ディスプレイ装置によって被験者に対して表示してもよい。

【0055】

刺激呈示部124より一对の刺激が呈示されると、被験者106は、それら刺激を比較する。被験者106はさらに、それらの間の差異を表わすのに、評定尺度160のカテゴリの内どれが最も適当であるかを判断し、そのカテゴリを選択する。本実施の形態では、評定取得部126は一般的な入力装置(例えばキーボード)を含んでおり、被験者106は、選択したカテゴリに対応するキーを押す。

【0056】

図2を参照して、評定取得部126は、被験者106がカテゴリを選択すると、そのカ

10

20

30

40

50

テゴリに付与されている評定値をテーブル作成部 1 2 8 に与える。被験者 1 0 6 が複数である場合、評定取得部 1 2 6 は、複数の被験者 1 0 6 による選択に対応する評定値をそれぞれ特定し、特定された評定値の平均値を、統合された評定値としてテーブル作成部 1 2 8 に与える。

【 0 0 5 7 】

テーブル作成部 1 2 8 は、テーブル作成部 1 2 8 が予め準備した声質差評価テーブル 1 0 8 の、試験処理部 1 2 2 より与えられた識別番号の組合せに該当する項目に、評定取得部 1 2 6 より与えられた評定値を格納する。テーブル作成部 1 2 8 は、声質評価テーブル 1 0 8 への評定値の格納が完了したことに応答して、試験処理部 1 2 2 に完了信号を与える。

10

【 0 0 5 8 】

試験処理部 1 2 2 は、この完了信号に応答して、音声コーパス 1 0 2 より取得した識別番号の中から、別の 2 セッションの音声データの識別番号を前述した選択方法によって選り、刺激抽出部 1 2 0、及びテーブル作成部 1 2 8 に与える。

【 0 0 5 9 】

声質差評価装置 1 0 4 は、上記した動作を繰返し、全てのセッションの組合せについて、被験者 1 0 6 に対する知覚試験を行なう。テーブル作成部 1 2 8 は、評定取得部 1 2 6 より与えられる評定値を声質差評価テーブル 1 0 8 に格納する。全ての組合せについて評定値の格納が終了すると、声質差評価装置 1 0 4 は一連の動作を終了する。

【 0 0 6 0 】

20

このようにして作成された声質差評価テーブル 1 0 8 は、図 1 に示す音声素片選択部 4 4 による評価関数 7 2 の値の算出時に、音声素片選択部 4 4 により参照される。即ち、図 1 を参照して、音声素片選択部 4 4 は、評価関数 7 2 に対して、観測可能な物理量に加えて、候補となる音声素片が属するセッションと、その直前に音声合成に使用された音声素片が属するセッションとの 2 つのセッションの識別番号の組を変数として与える。評価関数 7 2 は、与えられた物理量と、与えられた識別番号の組に対応する声質差評価テーブル 1 0 8 の項目に格納されている評定値とに基づく評価値を出力する。この評価値には、セッションごとの声質の差異に関する評定値が反映されることとなる。

【 0 0 6 1 】

音声素片選択部 4 4 はこの評価値をもとに一連の音声素片を決定する。その結果、図 8 30 に示す音声素片接続部 4 6 が一連の音声素片を接続することによって合成される音声には、接続される音声素片同士の聴感上の声質のばらつきが少なくなる。そのため、このようにして合成された音声は、不連続感が軽減し、自然に聞こえるものとなる。

【 0 0 6 2 】**[第 2 の実施の形態]**

第 1 の実施の形態に係るシステムでは、被験者 1 0 6 は、異なる 2 つのセッションで収録された音声波形のデータを比較し、それらの音声波形のデータにおける声質の差異を評定した。しかし、本発明は、このような実施の形態には限定されない。

【 0 0 6 3 】

第 2 の実施の形態に係る声質差評価装置は、2 つのセッションの音声データを対比するのではなく、セッションごとの音声データをもとに、知覚検査により各セッションの声質を表わす特徴ベクトルを作成する。セッションのベクトル間の距離によりセッション間の声質差が表現される。本実施の形態に係る知覚試験では、声質を評定するための複数の評価語対からなる評価語セットを予め準備し、被験者に与える。被験者はこの評価語セットに基づいて、呈示された刺激の声質に関する聴感上の印象を評定する。

40

【 0 0 6 4 】

図 5 に本実施の形態に係る声質差評価装置 2 0 4 の構成をブロック図形式で示す。図 5 を参照して、声質差評価装置 2 0 4 は、被験者 1 0 6 に呈示する音声刺激の抽出元となる音声データのセッションを 1 つずつ、所定の順序で決定する処理を行なう試験処理部 2 2 2 と、試験処理部 2 2 2 により決定されたセッションの音声データから、特定の声質評価

50

用の音声データを抽出する刺激抽出部 220 と、刺激抽出部 220 が抽出した音声刺激を被験者 106 に対して呈示する刺激呈示部 224 と、刺激呈示部 224 により呈示された音声刺激に対して被験者 106 が、前述した評価語セットを用いて行なう聴感上の印象の評定を取得する評定取得部 226 とを含む。

【0065】

声質差評価装置 204 はさらに、試験処理部 222 及び評定取得部 226 に接続され、試験処理部 222 が決定するセッションの識別番号と、評定取得部 226 が取得する被験者の評定とをもとに、各セッションで収録された音声の声質を表わす声質ベクトルを作成する声質ベクトル作成部 228 と、声質ベクトルを格納する声質ベクトルテーブル 230 と、声質ベクトルテーブル 230 に格納された声質ベクトルに基づき、任意のセッション間の声質差に関する評価値を算出する声質差算出部 232 と、声質差算出部 232 が算出する評価値をもとに声質差評価テーブル 236 を作成するテーブル作成部 234 とを含む。声質差評価テーブル 236 は、図 3 に示す第 1 の実施の形態に係る声質差評価テーブル 108 と同様の構成である。

10

【0066】

図 6 に、本実施の形態に係る知覚検査において、被験者に与えられる評価語セットの一例を示す。図 6 を参照して、評価語セット 240 は、対義語となる一对の形容詞からなる評価語対（「張りがある 張りが無い」、「濁った 澄んだ」、「明るい 暗い」、など）を複数含む。これらの評価語対について左側の評価語から右側の評価語に向かって、7 段階のカテゴリ（「非常に」、「かなり」、…、「かなり」、「非常に」）が与えられている。これらのカテゴリには、第 1 の実施の形態に係る評価尺度のカテゴリと同様に、それぞれ数値からなる評定値が評価語対ごとに付与されている。

20

【0067】

図 7 に、本実施の形態に係る声質ベクトルテーブルの構成の一例を示す。図 7 を参照して、声質ベクトルテーブル 230 は、各セッションに対応するエン트리 280A, 280B, … によって構成されたテーブルである。声質ベクトルテーブル 230 の各エント리는、セッションの識別番号の項目 272 と、当該セッションの音声データの声質を示す声質ベクトル 274 とを含む。声質ベクトル 274 は、図 6 に示す評価語セットの各評価語対に関する被験者の評定をそれぞれ数値化したものを成分とするベクトルである。

【0068】

本実施の形態に係るシステムは以下のように動作する。

30

【0069】

図 5 を参照して、最初に試験処理部 222 は、音声コーパス 102 より、全てのセッションの識別番号を讀出し、テーブル作成部 234 及び声質ベクトル作成部 228 に与える。声質ベクトル作成部 228 は与えられた識別番号を声質ベクトルテーブル 230（図 7 参照）に格納し、声質ベクトルテーブル 230 の準備を行なう。声質ベクトルテーブル 230 の準備が完了すると、声質ベクトル作成部 228 は、試験処理部 222 に完了信号を与える。

【0070】

試験処理部 222 は、声質ベクトル作成部 228 より完了信号を受けたことに応答して、音声コーパス 102 より讀出した識別番号の中から、1セッションの音声データの識別番号を選び、刺激抽出部 220 及び声質ベクトル作成部 228 に与える。

40

【0071】

この識別番号が与えられると、刺激抽出部 220 は、音声コーパス 102 に格納されている、与えられた識別番号に対応するセッションの音声データから、前述した声質評価用の 1 発声分の音声波形データを抽出し、刺激呈示部 224 に与える。

【0072】

刺激呈示部 224 は、音声波形データを再生し、再生が完了すると、刺激呈示部 224 は、評定取得部 226 に対し再生完了を示す信号を与える。

【0073】

50

被験者106には、予め図6に示す評価語セットが与えられている。被験者106は、刺激呈示部224より音声刺激を受けると、与えられた評価語セット内のある各評価語対について、刺激によって受けた聴感上の印象が評定尺度のカテゴリの内のどのカテゴリに属するかを判断する。被験者106は、最も適当であると判断したカテゴリを当該評価語対に対応する評定尺度より選択する。被験者106はこの選択を全ての評価語対について行なう。

【0074】

図5を参照して、評定取得部226は、刺激呈示部224よりの再生完了を示す信号に
10 応答して、被験者106による評定を取得する。評定取得部226は被験者106により
選択されたカテゴリを評価語対ごとに特定し、対応する評定値をそれぞれ声質ベクトル作
成部228に与える。被験者106が複数である場合、評定取得部226は、第1の実施
の形態と同様に、複数の被験者106による選択からそれぞれ特定し、特定された複数の
評定値の平均値を、統合された評定値として声質ベクトル作成部228に与える。

【0075】

声質ベクトル作成部228は、評価語対に対応する評定値をもとに、これらの評定値を
成分とする声質ベクトルを作成する。声質ベクトル作成部228は、声質ベクトルテー
ブル230のうち、試験処理部222から与えられた識別番号のエントリの声質ベクトルの
項目に、作成した声質ベクトルを格納する。声質ベクトル作成部228は、声質ベクトル
の格納が完了すると、試験処理部222に完了信号を与える。試験処理部222は、この
完了信号に
20 応答して、次のセッションを決定する。

【0076】

以上の動作を繰返すことにより、全てのセッションについての声質ベクトルが声質ベ
クトルテーブル230に格納される。

【0077】

声質ベクトルの作成及び格納が完了すると、声質差算出部232は、セッション間の声
質の差異に関する評価値の算出を次のようにして行なう。即ち、声質差算出部232は、
声質ベクトルテーブル230から、任意の2つのセッションの識別番号と声質ベクトルと
を
30 読出す。声質差算出部232は、読出した2組分の声質ベクトルをもとに、評価語セッ
トの各評価語対における評定尺度をそれぞれ軸とする多次元空間における、読出した声質
ベクトル間の距離を算出する。ここで算出するベクトル間の距離は、例えばベクトル間の
ユークリッド距離であってもよい。声質差算出部232は、このようにして算出したベク
トル間の距離を
読出した2セッション分の識別番号と共にテーブル作成部234に与える
。

【0078】

テーブル作成部234は、声質差評価テーブル236の、声質差算出部232より与え
られた識別番号に対応する項目に、声質差算出部232より与えられるベクトル間の距離
を格納する。なお、声質差評価テーブル236は、図3に示す声質差評価テーブル108
と同様の構成である。

【0079】

声質差算出部232とテーブル作成部234とは、全てのセッションの組合せについて
40 これら一連の動作を実行し、声質差評価テーブル236を作成する。

【0080】

このようにして作成された声質差評価テーブル236は、第1の実施の形態と同様に、
図1に示す評価関数72の計算の中に、1つの評価尺度として組込まれる。

【0081】

このように第2の実施の形態に係る声質差評価装置を含むシステムでは、セッションご
との音声データの間の声質の差異の大きさを被験者に直接評価させるのではなく、声質に
関する評価語セットによって表現される被験者の聴覚上の印象をもとに声質差評価テー
ブルを作成する。そのため、より多角的な評定結果を被験者より得ることが可能となる。

【0082】

10

20

30

40

50

この結果、音声素片選択部 44 がこの評価値をもとに選択した音声素片を接続することにより合成される音声では、接続される音声素片同士の声質が類似するものとなるので、不連続感が生じることが少なくなる。従って、このようにして合成された音声は、自然なものとなる。また、セッションごとに 1 回だけ評価すればよいので、2 つのセッションの組合せごとに評価する第 1 の実施の形態と比較して、評価に要する時間を短縮することができる。

【0083】

以上のブロック図形式で説明した各機能部は、いずれもコンピュータハードウェア及び当該コンピュータ上で実行されるプログラムにより実現することができる。このコンピュータとしては、音声を扱う設備を持ったものであれば、汎用のハードウェアを有するものを用いることができる。また、上で説明した装置の各機能ブロックは、この明細書の記載に基づき、当業者であればプログラムで実現することができる。そうしたプログラムもまた 1 つのデータであり、記憶媒体に記憶させて流通させることができる。

10

【0084】

なお、上記した実施の形態における被験者の人数は問わない。十分な教示及び訓練を受けた被験者であれば、単一又は小人数の被験者であっても、十分な精度の試験結果を得ることができる。また、多人数の被験者に対して知覚試験を行なうことにより、合成された音声を聞く一般的なユーザの評定に近い評定値に基づく声質差評価テーブルを作成することができる。

【0085】

また、上記した実施の形態では、複数の被験者を用いて知覚試験を行なう場合、各被験者から得られた評定値の平均値によって、評定値を統合した。しかし本発明は、このような実施の形態には限定されない。平均値以外にも中央値、最大値、及び最小値などによって評定値を統合してもよい。

20

【0086】

なお、上記した実施の形態では、刺激抽出部 120、220 は、特定の 1 発声分の音声波形のデータを音声データよりそれぞれ抽出した。しかし、本発明はこのような実施の形態には限定されない。例えば複数発話分を用いてもよい。また不特定の発声を用いて評価することもできるが、その場合には、評価の精度が十分に保証される知覚試験の方法を用いることが望ましい。

30

【0087】

なお、上記した発明の実施の形態では、被験者に与える評定尺度の各カテゴリに、予め評定値が付与されていた。しかし、本発明は、このような実施の形態には限定されない。例えば、事前に別の知覚試験を行ない、各カテゴリに対応する評定値を決定してもよい。また、知覚試験によって得られる試験結果をもとに、各カテゴリの評定値を統計的に求めることも可能である。

【0088】

また、上記した発明の実施の形態では、声質差の評価を行なうための知覚試験の方法及び評定値の決定方法として、被験者に音声刺激を呈示し、音声刺激に対する評定をカテゴリ尺度によって行なわせるものであった。しかし、本発明はこのような実施の形態には限定されない。抽出した音声波形のデータを呈示刺激とし、知覚試験によって聴感上の印象に基づく尺度を構成する方法であれば、どのような方法を用いてもよい。

40

【0089】

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味及び範囲内のすべての変更を含む。

【図面の簡単な説明】

【0090】

【図 1】本発明の第 1 の実施の形態に係るシステムの構成を示すブロック図である。

50

【図2】本発明の第1の実施の形態に係る声質差評価装置の構成を示すブロック図である。

【図3】本発明の実施の形態に係る声質差評価テーブル108の一例を示す図である。

【図4】本発明の第1の実施の形態に係る評定尺度の一例を示す図である。

【図5】本発明の第2の実施の形態に係る声質差評価装置の構成を示すブロック図である。

【図6】本発明の第2の実施の形態に係る評定尺度の一例を示す図である。

【図7】本発明の第2の実施の形態に係る声質ベクトルテーブルの一例を示す図である。

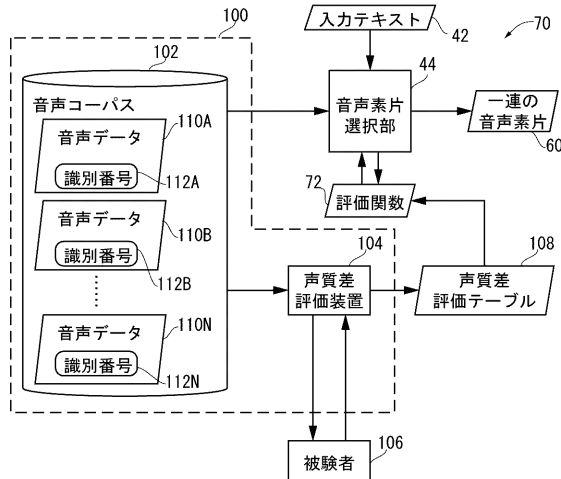
【図8】音声素片接続型音声合成システムの基本構成を示すブロック図である。

【符号の説明】

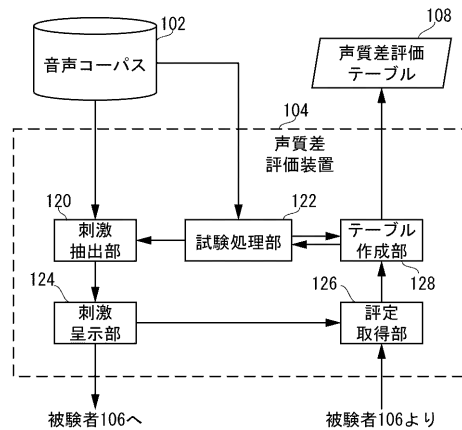
【0091】

40, 102 音声コーパス、42 入力テキスト、44 音声素片選択部、46 音声素片接続部、48 合成音声データ、50, 72 評価関数、70 音声素片選択システム、100 声質差評価テーブル作成装置、104, 204 声質差評価装置、106 被験者、108, 236 声質差評価テーブル、110A, 110B, ..., 110N 音声データ、112A, 112B, ..., 112N 識別番号、120, 220 刺激抽出部、122, 222 試験処理部、124, 224 刺激呈示部、126, 226 評定取得部、128, 234 テーブル作成部、228 声質ベクトル作成部、230 声質ベクトルテーブル、232 声質差算出部

【図1】



【図2】



【図3】

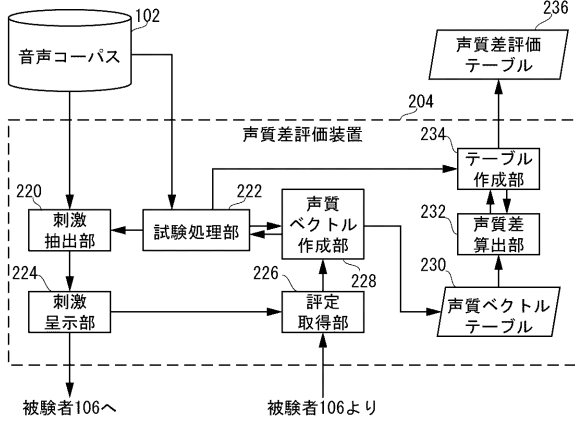
		108				
		001	002	003	004	005
130A	001		2.3	3.7	0.9	3.8
	002	2.3		0.2	4.2	1.5
130B	003	3.7	0.2		4.6	2.1
	004	0.9	4.2	4.6		3.3
	005	3.8	1.5	2.1	3.3	

【 図 4 】

同じ	似ている	どちらでもない	異なる	極端に異なる

160

【 図 5 】



【 図 6 】

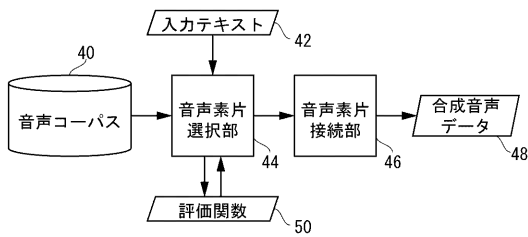
	非常に	かなり	やや	でも	やや	かなり	非常に
				どちら			
				もない			
張りがある							張りがない
濁った							澄んだ
明るい							暗い
⋮							⋮
⋮							⋮

240

【 図 7 】

識別番号	声質ベクトル
001	(2, -1, 3, ...)
002	(0, 2, -1, ...)
003	(-2, 0, 1, ...)
004	(-2, 1, 1, ...)
005	(3, 2, 0, ...)

【 図 8 】



フロントページの続き

審査官 荏原 雄一

(56)参考文献 特開平10 - 254471 (JP, A)
特開平05 - 108002 (JP, A)
特開2003 - 108824 (JP, A)

(58)調査した分野(Int.Cl., DB名)
G10L 13/00 - 13/08