

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6436806号  
(P6436806)

(45) 発行日 平成30年12月12日(2018.12.12)

(24) 登録日 平成30年11月22日(2018.11.22)

(51) Int.Cl.	F I
<b>G 1 0 L 13/10 (2013.01)</b>	<b>G 1 0 L 13/10 1 1 3 B</b>
<b>G 1 0 L 13/06 (2013.01)</b>	<b>G 1 0 L 13/06 2 1 0</b>

請求項の数 6 (全 21 頁)

(21) 出願番号	特願2015-19009 (P2015-19009)	(73) 特許権者	000233169
(22) 出願日	平成27年2月3日(2015.2.3)		株式会社日立超エル・エス・アイ・システムズ
(65) 公開番号	特開2016-142936 (P2016-142936A)		東京都立川市緑町7番地1
(43) 公開日	平成28年8月8日(2016.8.8)	(74) 代理人	100091096
審査請求日	平成29年3月22日(2017.3.22)		弁理士 平木 祐輔
		(74) 代理人	100105463
			弁理士 関谷 三男
		(74) 代理人	100102576
			弁理士 渡辺 敏章
		(72) 発明者	孫 慶華
			東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		審査官	上田 雄

最終頁に続く

(54) 【発明の名称】 音声合成用データ作成方法、及び音声合成用データ作成装置

(57) 【特許請求の範囲】

【請求項 1】

メモリから各種プログラムを読みだして実行し、音声合成処理で用いられる音声合成用データを作成するプロセッサが、第1種音声のコーパスデータの入力を受け付け、当該コーパスデータに含まれる音声データから第1韻律を抽出するステップと、

前記プロセッサが、予め用意されている第2種音声の韻律モデルを前記コーパスデータに含まれるテキストに対して適用し、前記テキストに対応する第2韻律を予測するステップと、

前記プロセッサが、前記第1韻律と前記第2韻律の差分値を算出するステップと、

前記プロセッサが、前記テキストに含まれる文字のうち、前記差分値が所定の閾値よりも大きい文字を判定するステップと、

前記プロセッサが、前記判定するステップの結果に基づいて、前記第1韻律と前記第2韻律との違いに起因する特徴テキスト部分に対応するデータを抽出するステップと、  
を含み、

前記第1種音声は口語調音声であり、前記第2種音声は読み上げ調音声であり、

前記第2種音声の韻律モデルは、読み上げ韻律・音韻予測モデルであり、

前記抽出された特徴テキスト部分に対応するデータは、前記特徴テキスト部分の音声波形データ、韻律・音韻情報、及びテキストデータを含み、

さらに、前記プロセッサが、前記特徴テキスト部分のテキストデータを用いて、与えられるテキストデータにおける口語調表現を抽出するためのルールを生成するステップを含

10

20

むことを特徴とする音声合成用データ作成方法。

【請求項 2】

請求項 1 において、

前記プロセッサは、前記特徴テキスト部分のテキストデータに加えて、当該テキストデータが含まれる口語調テキストの前後のコンテキスト情報を用いて前記ルールを生成することを特徴とする音声合成用データ作成方法。

【請求項 3】

メモリから各種プログラムを読みだして実行し、音声合成処理で用いられる音声合成用データを作成するプロセッサが、口語調音声のコーパスデータの入力を受け付け、当該コーパスデータに含まれる音声データから口語調韻律データを抽出するステップと、

前記プロセッサが、予め用意されている読み上げ調の韻律モデルを前記コーパスデータに含まれるテキストに対して適用し、前記テキストに対応する読み上げ調韻律データを予測するステップと、

前記プロセッサが、前記口語調韻律データと前記読み上げ調韻律データの差分値を算出するステップと、

前記プロセッサが、前記差分値に基づいて、前記テキストのセグメントに対して、当該セグメントの口語調の程度を示す口語調度を算出し、前記口語調韻律データに付与するステップと、

前記プロセッサが、前記口語調度が付与された前記口語調韻律データを用いて、前記音声合成用データを生成するステップと、

を含むことを特徴とする音声合成用データ作成方法。

【請求項 4】

請求項 3 において、

前記音声合成用データを生成するステップは、前記プロセッサが、前記口語調度が付与された前記口語調韻律データを用いて、入力テキストの口語調度を予測するための統計モデルである口語調度予測モデルを生成することを含むことを特徴とする音声合成用データ作成方法。

【請求項 5】

各種プログラムを格納するメモリと、

前記メモリから前記各種プログラムを読みだして実行し、音声合成処理で用いられる音声合成用データを作成するプロセッサと、を有し、

前記プロセッサは、

第 1 種音声のコーパスデータの入力を受け付け、当該コーパスデータに含まれる音声データから第 1 韻律を抽出する処理と、

予め用意されている第 2 種音声の韻律モデルを前記コーパスデータに含まれるテキストに対して適用し、前記テキストに対応する第 2 韻律を予測する処理と、

前記第 1 韻律と前記第 2 韻律の差分を算出する処理と、

前記テキストに含まれる文字のうち、前記差分が所定の閾値よりも大きい文字を判定する処理と、

前記第 1 韻律と前記第 2 韻律との違いに起因する特徴テキスト部分を抽出する処理と

、  
を実行し、

前記第 1 種音声は口語調音声であり、前記第 2 種音声は読み上げ調音声であり、

前記第 2 種音声の韻律モデルは、読み上げ韻律・音韻予測モデルであり、

前記抽出された特徴テキスト部分に対応するデータは、前記特徴テキスト部分の音声波形データ、韻律・音韻情報、及びテキストデータを含み、

前記プロセッサは、さらに、前記特徴テキスト部分のテキストデータを用いて、与えられるテキストデータにおける口語調表現を抽出するためのルールを生成する処理を実行する音声合成用データ作成装置。

【請求項 6】

請求項 5 において、

前記プロセッサは、前記特徴テキスト部分のテキストデータに加えて、当該テキストデータが含まれる口語調テキストの前後のコンテキスト情報を用いて前記ルールを生成することを特徴とする音声合成用データ作成装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声合成用データ作成方法、及び音声合成用データ作成装置に関し、例えば、収録した口語調音声から、読み上げ調との差異（口語調度）を自動的に測定する技術に関するものである。

10

【背景技術】

【0002】

テキストを音声に変換して読み上げるテキスト音声合成技術、およびそれを利用したテキスト音声合成システムがある。このような技術、システムの応用先として、例えば、カーナビゲーションでのガイド音声、携帯電話・スマートフォンでのメール読み上げや音声対話インタフェース、視覚障害者向けのスクリーンリーダー、電子書籍の読み上げ機能などが存在する。

【0003】

近年、音声合成技術はその重要性が増大している。例えば、カーナビゲーションや携帯電話・スマートフォンの普及に加えて障害者や高齢者を対象としたユニバーサルデザインの志向によって音声を使ったユーザインターフェースが今まで以上に使われるようになってきている。また、近年では、電子書籍端末の普及が始まり、音声読上げに必須な音声合成技術のニーズも拡大している。

20

【0004】

最近、音声合成技術は、カーナビや携帯電話などの音声対話処理で多く用いられるようになってきている。加えて、そのような応用例では、ユーザが会話を自然に進められるように、話し言葉（以降、口語調）での音声を合成したいというニーズが増えてきている。

【0005】

従来の音声合成技術であっても、合成音声の韻律（抑揚やリズム、強さなど）を工夫することで口語調の音声を生成できるものは存在している。例えば、標準的な発話特性を有する韻律特徴（読み上げ調韻律）に、感情や個性、発話スタイルなどの発話特性の違いに応じて補正するため修正韻律（口語調韻律）を重畳することで、口語調を含む多様な韻律を実現する手法が提案されている（特許文献 1 参照）。

30

【0006】

また、入力されたテキストについて、口語調音声の特徴を担う部分（口語表現部分）とそれ以外の部分を分けて、前者に関しては口語調音声コーパスから得られた韻律や音声を用いて合成し、後者に関しては従来読み上げ調音声合成手法で合成するという技術も考案されている（特許文献 2 参照）。このような手法では、従来培ってきた読み上げ音声合成技術を最大限に生かしたうえ、自然性が高く、安定した口語調音声合成が実現できると考えられる。

40

【先行技術文献】

【特許文献】

【0007】

【特許文献 1】特開 2003 - 337592 号公報

【特許文献 2】特開 2014 - 062970 号公報

【発明の概要】

【発明が解決しようとする課題】

【0008】

しかし、特許文献 1 の技術では、音声合成時に大きく韻律を制御する（変形させる）ため、合成音声の品質が劣化するという課題が指摘されている。

50

## 【 0 0 0 9 】

また、特許文献 1 及び 2 に開示の何れの手法においても、口語調での音声合成を実現するために、収録した口語調音声コーパスから、口語調音声データの作成が必須となる。例えば、特許文献 2 では、収録した口語音声から、熟練したラベリング作業者の経験に基づいて、手作業で口語調表現を抽出していた。しかし、この作業は、非常に時間がかかるうえ、作業者ごとに判定基準が異なり、抽出した口語調表現の一貫性を保つことが難しい（例えば、同じ音声に対しても、作業者ごとに、抽出された口語調表現が異なる。）。さらに、特許文献 1 の音声合成方法を用いる場合、音声のセグメントごとに、口語調表現らしさをより正確に定量的に評価することが望まれる。しかし、作業者の手作業ではこのような定量的な評価はほぼ不可能である。

10

## 【 0 0 1 0 】

本発明はこのような状況に鑑みてなされたものであり、口語調音声のように、読み上げ音声（平静音声）の特徴とは異なる韻律や声質の特徴を持つ音声に対して、その特徴部分のデータを自動的に抽出する技術を提供するものである。

## 【課題を解決するための手段】

## 【 0 0 1 1 】

上記課題を解決するために、本発明では、まず、予め用意されている第 2 種音声の韻律モデルをコーパスデータに含まれるテキストに対して適用し、テキストに対応する第 2 韻律を予測し、当該コーパスに含まれる音声データから抽出された第 1 韻律と第 2 韻律の差分値を算出する。次に、コーパスデータのテキストに含まれる文字のうち、差分値が所定の閾値よりも大きい文字を判定する。そして、この判定結果に基づいて、第 1 韻律と第 2 韻律との違いに起因する特徴テキスト部分に対応するデータを抽出する。

20

## 【 0 0 1 2 】

本発明に関連する更なる特徴は、本明細書の記述、添付図面から明らかになるものである。また、本発明の態様は、要素及び多様な要素の組み合わせ及び以降の詳細な記述と添付される特許請求の範囲の様態により達成され実現される。

## 【 0 0 1 3 】

本明細書の記述は典型的な例示に過ぎず、本発明の特許請求の範囲又は適用例を如何なる意味に於いても限定するものではないことを理解する必要がある。

## 【発明の効果】

30

## 【 0 0 1 4 】

本発明によれば、口語調音声のように、読み上げ音声（平静音声）の特徴とは異なる韻律や声質の特徴を持つ音声に対して、その特徴部分のデータを自動的に抽出することができるようになる。

## 【図面の簡単な説明】

## 【 0 0 1 5 】

【図 1】本発明の実施形態による音声合成システムのハードウェア構成例を示す図である。

【図 2】本発明の第 1 の実施形態による音声合成システムの機能ブロックを示す図である。

40

【図 3】本発明の実施形態による読み上げ調韻律・音韻予測部によって予測された韻律特徴量の例を示す図である。

【図 4】本発明の実施形態による韻律・音韻特徴抽出部によって抽出された韻律特徴量の例を示す図である。

【図 5】本発明の実施形態による口語調度算出部によって算出された韻律差分値の例を示す図である。

【図 6】本発明の第 1 の実施形態による口語調表現部分の自動抽出部によって計算された音節（モーラ）レベルでの口語調度の例を示す図である。

【図 7】本発明の第 1 の実施形態による口語調表現部分の自動抽出部によって計算された形態素レベルでの口語調度の例を示す図である。

50

【図 8】本発明の第 1 の実施形態による口語調表現部分の自動抽出部によって計算されたアクセント句レベルでの口語調度の例を示す図である。

【図 9】本発明の第 1 の実施形態による口語調表現部分の自動抽出部によって抽出された口語調音声データの例を示す図である。

【図 10】本発明の第 1 の実施形態による口語調表現部分の自動抽出部によって抽出された口語調韻律データの例を示す図である。

【図 11】本発明の第 1 の実施形態による口語調表現部分の自動抽出部によって抽出された口語調テキストデータの例を示す図である。

【図 12】本発明の第 1 の実施形態による口語調表現抽出ルール作成部で用いられるデータの例を示す図である。

【図 13】本発明の第 1 の実施形態による口語調表現抽出ルール作成部で生成されたルールの例を示す図である。

【図 14】本発明の第 1 の実施形態による、英語 F0 パターンによる口語調度の計算を説明する図である。

【図 15】本発明の第 1 の実施形態による、英語単語レベルで計算した口語調度の例を示す図である。

【図 16】本発明の第 1 お実施形態による口語調データ作成処理を説明するためのフローチャートである。

【図 17】本発明の第 2 の実施形態による音声合成システムの機能ブロックを示す図である。

【図 18】本発明の第 1 の実施形態による口語調表現抽出ルール作成部で用いられるデータ（口語調度が連続値）の例を示す図である。

【図 19】本発明の第 2 の実施形態による口語調度予測 & 韻律混合比決定部で生成した口語調度および口語調混合比の例を示す図である。

【図 20】本発明の第 2 の実施形態による韻律混合部で生成した韻律の例を示す図である。

【図 21】本発明の第 2 の実施形態による韻律混合処理の概念を示す図である。

【図 22】本発明の第 2 の実施形態による口語調データ作成処理を説明するためのフローチャートである。

【発明を実施するための形態】

【0016】

本発明は、従来手作業で行っていた口語調データ（口語調表現抽出ルール又は口語調度モデル、口語調韻律・音韻モデル、及び口語調音声 D B）の作成を、収録口語調音声コーパスと読み上げ韻律・音韻予測モデルを用いて自動化するものである。具体的には、本発明は、収録した口語調音声から抽出した特徴量と、収録した音声のテキストデータを読み上げモデルで読み上げた場合の特徴量とを比較して口語調表現部分を抽出する。口語の言い回しで特徴ある部分（例えば、「今日は雨かなあ」の「かなあ」の部分）以外は、収録音声と読み上げ音声とでは差がないが、特徴部分に関しては特徴量に差異が生じるという性質を利用したものである。この特徴ある部分の情報から口語調表現抽出ルール等、口語調韻律・音韻モデル、口語調音声データが作成され、口語調データとして登録される。音声合成装置では、この口語調データを用いて、例えば特許文献 2 で示された方法で口語調の音声合成データを出力する。

【0017】

以下、添付図面を参照して本発明の実施形態について説明する。添付図面では、機能的に同じ要素は同じ番号で表示される場合もある。なお、添付図面は本発明の原理に則った具体的な実施形態と実装例を示しているが、これらは本発明の理解のためのものであり、決して本発明を限定的に解釈するために用いられるものではない。

【0018】

本実施形態では、当業者が本発明を実施するのに十分詳細にその説明がなされているが、他の実装・形態も可能で、本発明の技術的思想の範囲と精神を逸脱することなく構成・

10

20

30

40

50

構造の変更や多様な要素の置き換えが可能であることを理解する必要がある。従って、以降の記述をこれに限定して解釈してはならない。

【 0 0 1 9 】

更に、本発明の実施形態は、後述されるように、汎用コンピュータ上で稼動するソフトウェアで実装しても良いし専用ハードウェア又はソフトウェアとハードウェアの組み合わせで実装しても良い。

【 0 0 2 0 】

なお、以後の説明では「テーブル」形式によって本発明の各情報について説明するが、これらの情報は必ずしもテーブルによるデータ構造で表現されていなくても良く、リスト、DB、キュー等のデータ構造やそれ以外で表現されていても良い。そのため、データ構造に依存しないことを示すために「テーブル」、「リスト」、「DB」、「キュー」等について単に「情報」と呼ぶことがある。

【 0 0 2 1 】

また、各情報の内容を説明する際に、「識別情報」、「識別子」、「名」、「名前」、「ID」という表現を用いることが可能であり、これらについてはお互いに置換が可能である。

【 0 0 2 2 】

以下では「プロセッサ」を主語（動作主体）として本発明の実施形態における各処理について説明を行うが、プロセッサはプログラムを実行することで定められた処理をメモリ及び通信ポート（通信制御装置）を用いながら行うため、「プログラム」を主語とした説明としてもよい。プログラムの一部または全ては専用ハードウェアで実現してもよく、また、モジュール化されていても良い。

【 0 0 2 3 】

（ 1 ）第 1 の実施形態

< ハードウェア構成 >

図 1 は、本発明の実施形態による音声合成システムのハードウェア構成例を示す図である。音声合成システム 1 は、各種プログラムを実行するプロセッサ（CPU：Central Processing Unit）101と、各種プログラムを格納するメモリ102と、各種データを格納する記憶装置103と、出力装置104と、入力装置105と、必要に応じて通信デバイス106と、を有している。

【 0 0 2 4 】

メモリ102は、少なくとも口語調データ（口語調表現抽出ルール、口語調韻律・音韻モデル、及び口語調音声DBを含む）を作成する口語調データ作成プログラム1021と、読み上げ調データ（読み上げ韻律・音韻予測モデル、及び読み上げ調音声DBを含む）を作成する読み上げ調データ作成プログラム1022と、音声合成処理を実行する音声合成プログラム1023と、を格納する。

【 0 0 2 5 】

記憶装置103は、読み上げ調データを作成するために用いられる、収録した読み上げ調音声データ1031と、収録した読み上げ調音声データから抽出された特徴量を学習して得られる読み上げ韻律・音韻予測モデル1032と、収録した読み上げ調音声データ1031から生成された読み上げ調音声DB1033と、を格納する。また、記憶装置103は、口語調データを作成する元データである収録した口語調音声コーパスデータ1034と、口語調音声コーパスデータ1034から抽出された口語調抽出データ1035と、口語調抽出データ1035に基づいて生成された口語調表現抽出ルール1036と、口語調抽出データ1035の特徴を学習して得られる口語調韻律・音韻モデル1037と、口語調抽出データ1035から生成される口語調音声DB1038と、を格納している。

【 0 0 2 6 】

出力装置104は、音声合成処理して得られる音声出力するデバイスである。例えば、スピーカが該当する。また、出力装置104は、口語調データ作成処理によって生成されるデータを画面上に表示したりしても良い。

## 【0027】

入力装置105は、テキストや音声を入力するためのデバイスである。例えば、テキスト入力であればキーボード、マウス、入力すべきテキストデータを取り込んで自動的に入力するソフトウェア等が該当し、音声入力であれば、マイクや入力すべき音声データを取り込んで自動的に入力するソフトウェア等が該当する。

## 【0028】

通信デバイス106は、必要に応じて設けられるデバイスであり、処理に必要なデータを受信したり、処理済のデータを他の装置に送信したりする場合に用いられる。

## 【0029】

なお、図1に示される音声合成システム1は、読み上げ調データ作成処理10と、口語調データ作成処理11と、音声合成処理12と、を実行するシステムとして構成されているが、本発明の特徴である口語調データ作成処理11のみを実行する口語調データ作成装置として構成しても良い。この場合、メモリ102に格納されるプログラムは口語調データ作成プログラムのみということになる。また、各処理を別のコンピュータで実行させるようなシステムを音声合成システムとして構成しても良い。この場合、読み上げ調データ作成処理10を実行するコンピュータ（読み上げデータ作成装置）、口語調データ作成処理11を実行するコンピュータ（口語調データ作成装置）、音声合成処理12を実行するコンピュータ（音声合成装置）がそれぞれ独立に存在していても良い（それぞれがネットワークに接続され、遠隔的に設置されていても良い）。

## 【0030】

コンピュータが口語調データ作成装置を構成する場合には、メモリ102は、口語調データ作成プログラム1021を格納する。また、この場合、記憶装置103は、収録した口語調音声コーパスデータ1034、口語調抽出データ1035、口語調表現抽出ルール1036、口語調韻律・音韻モデル1037、及び口語調音声DB1038を格納することとなる。

## 【0031】

## &lt;機能ブロックと処理内容&gt;

図2は、本発明の第1の実施形態による音声合成システム1の機能ブロック図である。音声合成システム1は、読み上げ調データ作成処理10と、口語調データ作成処理（口語調の音声合成用データを作成する処理）11と、口語調データを用いた音声合成処理12

## 【0032】

読み上げ調データ作成処理10は、韻律・音韻抽出部、音声DB作成部、韻律・音韻特徴自動学習部などで構成されるが、これらによる処理は一般的な音声合成データ作成に用いられる処理であり、本発明の特徴となるものではないので、以降、その説明は原則省略する。

## 【0033】

## (i) 口語調データ作成処理

口語調データ作成処理11は、収録した口語調音声と、読み上げ韻律・音韻予測モデルを入力すると、口語調音声から口語調表現部分を自動抽出し、口語調音声合成に必要な口語調表現抽出ルール、口語調韻律・音韻モデル、口語調音声DB（データベースの略）を生成する処理である。この口語調データ作成処理11を実現するために、韻律・音韻特徴抽出部202、韻律・音韻自動学習部206、音声DB作成部207という通常の音声データ作成装置が持つ処理単位を備える。これに加えて、本発明に特徴的な、読み上げ調韻律・音韻予測部201、口語調度算出部203、口語調表現部分の自動抽出部204が設けられ、さらに、入力テキストから口語調テキスト表現を自動検出できる口語調表現抽出ルールを生成する口語調表現抽出ルール作成部205が設けられる。

## 【0034】

音声合成処理12は、口語調データ作成処理によって生成された口語調表現抽出ルール1036、口語調韻律・音韻モデル1037、及び口語調音声DB1038を用いて、入

力されたテキストデータを処理し、合成音声を生成する処理である。当該音声合成処理は、特許文献 2 に開示された内容と同様であるので、以降、詳細な説明は原則省略する。

#### 【0035】

以上のように、読み上げ調データ作成処理 10 と音声合成処理 12 は、通常の音声合成システムにも存在する処理単位であり、口語調データ作成処理 11 が本発明の特徴となる処理に相当する。従って、以下では口語調データ作成処理 11 を中心に説明することとする。

#### 【0036】

なお、以降の説明において、収録した音声として「お願いします」や「今日は雨かなあ」などの単文を用いているが、複数の文を結合した長い文書の収録も可能である。また、図 2 では入力テキストを漢字かな文としているが、もちろん、英語や中国語などの外国語でも構わない。その場合は、内部処理もその外国語に対応したプログラム・データ（たとえば、収録した口語調音声、収録した読み上げ調音声）を用いなければならないことは言うまでもない。

#### 【0037】

読み上げ調韻律・音韻予測部 201 は、収録した口語調音声コーパスデータ 1034 から、収録した口語調音声の発話テキストを読み出し、それに対して読み上げ韻律・音韻予測モデル 1032 を適用し、読み上げ調の韻律特徴量および音韻特徴量を予測する。つまり、ここでは、発話者がこのテキストに対して、読み上げ調スタイルで発話した場合は、韻律・音韻特徴がどのようなものであるかが分かる。ただし、韻律特徴量は、発話速度を表す特徴量（例えば、音素継続長、音節継続長など）、声の高さを表す特徴量（例えば、基本周波数の時間変化パターン（F0 パターン）など）、音の大きさを表す特徴量（例えば、短時間平均パワーなど）等である。音韻特徴量は、声道形状を表す特徴量（例えば、ケプストラム、LPC 係数など）が考えられる。また、これらの情報をすべて用いる必要がないが、口語調の特徴に最も寄与する基本周波数を用いることが好ましい。ただし、以降、本明細書では、音韻特徴量についての説明を省略し、単に韻律特徴量と記載した場合でも、韻律特徴量と音韻特徴量と両方を意味するものとする。また、韻律特徴量についても、理解しやすい F0 パターンを中心に説明を行うこととする。例えば、収録した口語調音声「今日は雨かなあ」のテキストに対して予測された韻律特徴は図 3 に示されるようなものとなる。なお、読み上げ韻律・音韻予測モデル 1032 は、口語調音声と同じ話者の読み上げ調音声から学習したものをを用いることが望ましいが、別の話者から学習したモデルを口語調音声話者に適用したもので良い。

#### 【0038】

韻律・音韻特徴抽出部 202 は、収録した収録した口語調音声コーパスデータ 1034 から収録音声データを読み出し、その音声の韻律・音韻特徴量を抽出する。つまり、発話者がこのテキストに対して、実際に口語調スタイルで発話した場合の韻律・音韻特徴がどのようなものであるかが分かる。ただし、収録した口語調音声には、事前に音素セグメンテーション情報が、自動および手動で付与されているものとする。なお、抽出する特徴量は、読み上げ調韻律・音韻予測部で予測された特徴量と同じである。例えば、収録した口語調音声「今日は雨かなあ」の音声波形に対して、抽出した韻律特徴量は、図 4 に示されるようなものとなる。

#### 【0039】

口語調度算出部 203 は、読み上げ調韻律・音韻予測部 201 で予測された特徴量と、韻律・音韻特徴抽出部 202 で抽出した収録口語調音声の特徴量とを比較し、口語調への寄与度（口語調度）を計算する。例えば、口語調度算出部 203 は、単純に音素ごとに韻律特徴量の差分を取り、下記式 1 を用いて音素ごとの口語調度を計算する。図 5 は、継続長係数 = 0.3；高さ係数 = 0.5；強さ係数 = 0.2 の場合、計算された口語調度を示している。

#### 【0040】

口語調度 = | 継続長係数 \* 継続長差分 | + | 高さ係数 \* 高さ差分 |

10

20

30

40

50



$$+ | \text{強さ係数} * \text{強さ差分} | \dots \dots \dots \quad (\text{式} 1)$$

ここで、“ $|A|$ ”は“ $A$ ”の絶対値を示すものとする。

#### 【0041】

口語調表現部分の自動抽出部204は、口語調度算出部203で計算された口語調度を用いて、音声を構成する各セグメントについて、セグメントの口語調度を計算し、口語調度が所定の閾値以上を示すセグメントを口語調表現部分として自動抽出する。抽出された口語調表現部分は、口語調抽出データ1035に格納される。口語調表現のセグメント単位は、合成時に用いる韻律モデルにも依存するが、日本語であれば音節単位、形態素単位、アクセント句単位などが適切だと考えられる。例えば、セグメント単位が音節（モーラ）の場合、口語調特徴が母音のみに現れることを仮定すると、音節口語調度は式2のようになる。つまり、音節に含まれる母音の口語調度が音節口語調度として与えられる。式2に従うと、音節口語調度は、図6のようになる。閾値が“20”の場合は、7番目の“ナ”と8番目の“ァ”が抽出される。閾値が“10”の場合には、6番目の“カ”も口語調データとして、抽出されることになる。なお、例えば、閾値は経験値で定められる値であり、予め決めておく。

10

#### 【0042】

$$\text{音節口語調度} = \text{母音口語調度} \dots \dots \dots \quad (\text{式} 2)$$

#### 【0043】

また、例えば、セグメント単位が形態素であるとする場合、口語調度は式3のように表される。つまり、音節に含まれる音節の口語調度の平均値が形態素口語調度として与えられる。式3に従うと、形態素口語調度は、図7のようになる。閾値が“25”の場合は、4番目の“かなあ”が口語調データとして、抽出される。

20

#### 【0044】

$$\text{形態素口語調度} = \text{音節平均口語調度} \dots \dots \dots \quad (\text{式} 3)$$

#### 【0045】

さらに、セグメント単位がアクセント句の場合、口語調度は式4のように表される。つまり、アクセント句に含まれる形態素の形態素口語調度のうち、最大値が口語調度として与えられる。式4に従うと、アクセント句口語調度は図8のようになる。閾値が“20”の場合は、2番目の“雨かなあ”が口語調データとして、抽出される。

30

#### 【0046】

$$\text{アクセント句口語調度} = \text{形態素最大口語調度} \dots \dots \dots \quad (\text{式} 4)$$

#### 【0047】

音声DB作成部207は、口語調表現部分の自動抽出部204によって抽出された口語調抽出データの音声波形を蓄積し、音声合成に用いる口語調音声DB1038を作成する。口語調音声DB1038は、音声合成装置による音声合成処理に適合する所定のフォーマットで作成される。例えば、「今日は雨かなあ」から抽出した口語調表現部分の音声波形は、図9のようになる。

#### 【0048】

韻律音韻特徴自動学習部206は、口語調表現部分の自動抽出部204から抽出された口語調抽出データの韻律・音韻情報（図10参照）を用いて、音声合成に用いる口語調韻律・音韻モデル1037を作成する。口語調韻律・音韻モデル1037は、音声合成装置による音声合成処理に適合する所定のフォーマットで作成される。口語調韻律・音韻モデル1037は、コンテキストから韻律・音韻情報を推定する統計モデルでも良いし、口語調のデータとして抽出され肉声の韻律・音韻情報をそのまま蓄積したモデルでも良い。例えば、「今日は雨かなあ」から抽出した口語調表現部分の韻律（F0パターン）は、図10のようになる。

40

#### 【0049】

口語調表現抽出ルール作成部205は、口語調表現部分の自動抽出部204によって抽出された口語調抽出データのテキスト（図11）を用いて、口語調表現抽出ルールを作成する。最も簡単な口語調表現抽出ルールは、「“かなあ”という文字列がマッチした場合

50

、その部分を口語調表現とする。」のように、文字列表現のみを用いた文字列マッチングルールである。ただし、このようなルール作成手法では、例えば「お願いします」の口語調音声に対して、“します”の部分を口語調表現として抽出されたとすると、「します」という文字列がマッチした場合、その部分を口語調表現とする。」というルールを作成されてしまう。このルールは明らかに不適切である。従って、作成したルールには、前後のコンテキスト情報を考慮した方が良くと考えられる。例えば、「お願いします」「今日は雨かなあ」の文に対して、口語調テキストを形態素単位（アクセント句単位など、形態素より大きい言語単位でも良い）に分解し、それぞれコンテキストと口語調度（“Yes”と“No”の2値）を付与すると、図12のようになる。このデータに対して、機械学習手法を用いて、口語調表現抽出ルールを自動作成することができる。例えば、2分岐決定木を自動構築した場合、図13のようになる。もちろん、ニューラルネットワーク、スーパーベクトルマシンなどのカテゴリを推測する手法を用いても良い。図13は、「お願いします」「今日は雨かなあ」の文に対して、口語調テキストを形態素単位（アクセント句単位など、形態素より大きい言語単位でも良い）に分解し、それぞれコンテキストと口語調度（口語調表現部分抽出部で口語調抽出に用いる口語調度の連続値）を付与した場合の2分岐決定木（図12を基に学習したツリー）を示している。このデータに対して、機械学習手法を用いて、口語調表現度予測モデルを自動作成することができる。例えば、重回帰解析などの連続値を推測する統計手法を用いることができる。そして、合成時に文を構成する各形態素について、口語調度を予測し、ある閾値を超えた形態素を「口語調表現」とし、一方、予測した口語調度がその閾値より小さい形態素を「口語調表現でない」とする。

#### 【0050】

以上のように、読み上げ調韻律・音韻予測部201、韻律・音韻特徴抽出部202、口語調度算出部203、及び口語調表現部分の自動抽出部204については、日本語「今日は雨かなあ」を適用した場合を例に説明したが、英語や中国語などの外国語でも構わない。例えば、口語調音声は英語「Oh, It's raining.」である場合、図14で示すように、F0観測値（収録した口語調音声からのF0値）とF0予測値（読み上げ調音声からのF0値）が得られたとする。英語の場合は、口語調表現のセグメント単位は、音素や音節より、単語や韻律語を用いたほうが良い。例えば、セグメント単位が単語の場合、口語調度は式5のようになる。ただし、この式は一例であり、上記式1を用いても構わない。式5に従うと、各単語の口語調度が図15のようになる。閾値を50と設定した場合、一番最初の“Oh”が、口語調表現として抽出される。

#### 【0051】

単語口語調度 = | 予測した単語最大F0値 - 観測した単語最大F0値 | ……  
(式5)

ここで、“|A|”は“A”の絶対値を示している。

#### 【0052】

##### (ii) 音声合成処理

音声合成処理では、まずテキスト入力部に音声合成すべきテキスト（例えば、かな漢字文）がユーザによって入力され、テキスト解析部で解析される。

#### 【0053】

口語調表現自動抽出部は、テキスト解析部で解析されたコンテキスト情報と口語調データ作成処理11の口語調表現抽出ルール作成部205で作成された口語調表現抽出ルール1036を用いて、入力テキストを「口語調表現」部分と「口語調表現でない」部分に分割する。「口語調表現」部分は、口語調部分の韻律・音韻作成部に出し、「口語調表現でない」部分は、読み上げ部分の韻律・音韻予測部に出し、出力する。ただし、入力テキストに必ず「口語調表現」部分と「口語調表現でない」部分と両方存在すると限らないので、入力テキストが必ず分割されると限らない。

#### 【0054】

例えば、テキスト「今日は晴れかなあ」が入力された場合、図13で示した口語調表現

抽出ルールに従い、口語調表現が抽出される。この例では、形態素「かなあ」が「口語調表現」として抽出され、残りの形態素が「口語調表現でない」と判断される。そのため、入力テキスト「今日は晴れかなあ」は、口語調表現でない部分の「今日は晴れ」と口語調表現部分の「かなあ」と分割される。また、例えば、テキスト「掃除します」が入力された場合、図 13 で示した口語調表現抽出ルールには「掃除」「します」の両方とも口語調表現として登録されていないため、口語調表現は抽出されず、文分割は行われない。

【0055】

韻律生成部は、口語調部分の韻律・音韻作成部で生成された韻律・音韻特徴量と読み上げ部分の韻律・音韻予測部で生成された韻律・音韻特徴量を合併し、文全体の韻律・音韻特徴量ターゲットを生成する。

10

【0056】

そして、波形生成部は、読み上げ調音声 DB 1033 を参照して、口語調ではない部分のテキストについて声質を考慮した処理を実行し、読み上げ調部分のテキストについて音声波形を生成する。また、口語調音声生成部は、口語調音声 DB 1038 を参照して、口語調部分のテキストについて音声波形を生成する。

【0057】

波形接続部は、口語調部分の音声波形と口語調ではない部分（読み上げ調部分）の音声波形を接続し、音声出力部は、最終的な合成音声を出力する。

【0058】

<口語調データ作成処理のフローチャート>

20

図 16 は、本発明の第 1 の実施形態による口語調データ作成処理を説明するためのフローチャートである。

【0059】

(i) ステップ 1601

プロセッサ 101 は、収録した口語調音声コーパスデータ 1034 の入力を受け付ける。当該データには、収録音声データとそれに対応するテキストデータ（発話テキスト）がセットとなっている。

【0060】

(ii) ステップ 1602

プロセッサ 101 は、収録した収録した口語調音声コーパスデータ 1034 の収録音声データから、その音声の韻律・音韻特徴量を抽出する。詳細については上述した通りである。

30

【0061】

(iii) ステップ 1603

プロセッサ 101 は、収録した口語調音声コーパスデータ 1034 の発話テキストに対して読み上げ韻律・音韻予測モデル 1032 を適用し、読み上げ調の韻律特徴量および音韻特徴量を予測する。つまり、ここでは、発話者がこのテキストに対して、読み上げ調スタイルで発話した場合は、韻律・音韻特徴がどのようなものであるかが分かる。詳細は上述した通りである。

【0062】

40

(iv) ステップ 1604

プロセッサ 101 は、ステップ 1602 で抽出した収録口語調音声の特徴量と、ステップ 1603 で予測した韻律・音韻特徴量とを比較し、口語調への寄与度（口語調度）を計算する。

【0063】

(v) ステップ 1605

プロセッサ 101 は、ステップ 1604 で得られた口語調度を用いて、音声を構成する各セグメントについて、セグメントの口語調度を計算し、口語調度が所定の閾値以上を示すセグメントを口語調表現部分として自動抽出する。抽出された口語調表現部分は、口語調抽出データ 1035 に格納される。詳細は上述した通りである。

50

## 【 0 0 6 4 】

## (vi) ステップ 1 6 0 6

プロセッサ 1 0 1 は、ステップ 1 6 0 5 で得られた口語調抽出データの音声波形を蓄積し、音声合成に用いる口語調音声 D B 1 0 3 8 を作成する

## 【 0 0 6 5 】

## (vii) ステップ 1 6 0 7

プロセッサ 1 0 1 は、口語調抽出データの韻律・音韻情報（図 1 0 参照）を用いて、音声合成に用いる口語調韻律・音韻モデル 1 0 3 7 を作成する。詳細は上述した通りである。

## 【 0 0 6 6 】

## (viii) ステップ 1 6 0 8

プロセッサ 1 0 1 は、ステップ 1 6 0 5 で得られた口語調抽出データのテキスト（図 1 1）を用いて、口語調表現抽出ルール 1 0 3 6 を作成する。詳細は上述した通りである。

## 【 0 0 6 7 】

## ( 2 ) 第 2 の実施形態

第 2 の実施形態は、特許文献 1 のような音声合成装置に用いる口語調音声合成用データを作成することを想定したものである。ハードウェア構成は第 1 の実施形態と同様であるので、説明は省略する。ただし、記憶装置 1 0 3 は、口語調抽出データ 1 0 3 5 の代わりに口語調度付き口語調音声データ 1 7 0 2、口語調表現抽出ルール 1 0 3 6 の代わりに口語調度予測モデル 1 7 0 4 を格納する。

## 【 0 0 6 8 】

## &lt; 機能ブロックと処理内容 &gt;

図 1 7 は、本発明の第 2 の実施形態による音声合成システムの機能ブロックを示す図である。第 2 の実施形態では、従来手作業による音声の口語調度ラベリングに代わって、収録した口語調音声にセグメントごとに、口語調度の定量的な評価を実現し、入力テキストの各セグメントに対する口語調度を予測する。この予測した口語調度によって、口語調音声から学習した韻律・音韻モデルと読み上げ調音声から学習した読み上げ調韻律・音韻モデルと、セグメント毎の混合割合を計算し、文全体の韻律・音韻特徴の予測を行う。第 1 の実施形態とは異なり、入力テキストを分割することがないので、分割された口語調表現部分と口語調表現でない部分と接続するときの不連続感を低減できると考えられる。

## 【 0 0 6 9 】

以下では、第 1 の実施形態とは異なる部分のみ説明することとする。

## (i) 口語調データ作成処理

口語調度付与部 1 7 0 1 は、口語調度算出部 2 0 3 で算出された韻律特徴の差分情報を用いて、収録した口語調音声の各セグメントに口語調度を付与し、口語調度付き口語調音声データ 1 7 0 2 を生成する。ここで、セグメントの単位は、音素、音節、形態素、アクセント句、フレーズ、文などが考えられるが、口語調音声の特徴を担う最小単位として、形態素を用いたことが好ましい。各セグメントの口語調度の計算については、口語調度算出部 2 0 3 で算出された韻律特徴の差分情報から求められるが、その具体例については、第 1 の実施形態で説明したので、ここでは詳細については省略する。「お願いします」「今日は雨かなあ」の文に対して、口語調テキストを形態素単位に分解し、それぞれコンテキストと口語調度を付与すると、図 1 8 のようになる。

## 【 0 0 7 0 】

口語調度予測モデル学習部 1 7 0 3 は、口語調度付与部 1 7 0 1 が生成した口語調度付き口語調音声データ 1 7 0 2 を用いて、口語調度を予測する統計モデル（口語調度予測モデル）1 7 0 4 を生成する。第 1 の実施形態では、入力文（テキスト）を「口語調」の部分と「口語調でない」の部分と分割するためのルールを作成しているが、第 2 の実施形態では、入力文を構成するすべてのセグメントについて、口語調度を予測するための統計モデルを作成することになる。

## 【 0 0 7 1 】

## (ii) 音声合成処理

第2の実施形態では、テキスト解析部が入力テキストを解析した後、口語調度予測&韻律混合比決定部が、口語調度予測モデル1704を用いて、テキスト文を構成する各セグメントについて、口語調度を予測する。さらに、口語調度予測&韻律混合比決定部は、この予測した口語調度に基づいて、口語調韻律と読み上げ調韻律の混合比率を計算する。例えば、「今日は晴れかなあ」というテキスト文が入力された場合、すべての形態素について口語調を予測した結果は、図19のようになる。ここで、口語調混合比を式6のように定義した場合（口語調下限値 = 0，口語調上限値 = 50とする）、口語調混合比は、図19に示される値となる。

【0072】

口語調混合比 =  $\text{MIN}(100\%, (\text{口語調度} - \text{口語調下限値}) / (\text{口語調上限値} - \text{口語調下限値}))$

・・・・・・ (式6)

ここで、 $\text{MIN}(A, B)$ は、AとBとの間で小さい方の値を選ぶことを意味するものとする。

【0073】

読み上げ調韻律・音韻予測部201は、読み上げ韻律・音韻予測モデル1032を参照し、入力テキストの読み上げ調の音声データを予測する。口語調韻律・音韻作成部は、口語調韻律・音韻モデル1037を参照し、入力テキストの口語調の音声データを作成する。

【0074】

韻律混合部は、口語調度予測&韻律混合比決定部で生成された口語調混合比を用いて、口語調の音声データと読み上げ調の音声データの韻律混合処理を実施する。例えば、入力テキスト「今日は晴れかなあ」に対して、図20に示されるように口語調韻律と読み上げ韻律が予測された場合、式7を用いて、韻律生成を行う。

【0075】

韻律 =  $\text{口語調韻律} * \text{口語調混合比} + \text{読み上げ調韻律} * (1 - \text{口語調混合比})$ ・・・・・・ (式7)

【0076】

図21は、当該韻律混合処理の概念を示す図である。図21に示されるように、「今日」については読み上げ調音声データが90%、口語調音声データが10%用いられる。「は」についてはそれぞれ84%、16%用いられ、「晴れ」についてはそれぞれ66%、34%用いられる。そして、「かな」については読み上げ調音声データが4%、口語調音声データが96%用いられて、混合韻律が生成される。

【0077】

最後に、音声生成部は、韻律混合部で生成した韻律をターゲットとして音声を生成し、音声出力部がこれを出力する。

【0078】

<口語調データ作成処理のフローチャート>

図22は、本発明の第2の実施形態による口語調データ作成処理を説明するためのフローチャートである。

【0079】

(i) ステップ2201

プロセッサ101は、収録した口語調音声コーパスデータ1034の入力を受け付ける。当該データには、収録音声データとそれに対応するテキストデータ（発話テキスト）がセットとなっている。

【0080】

(ii) ステップ2202

プロセッサ101は、収録した収録した口語調音声コーパスデータ1034の収録音声データから、その音声の韻律・音韻特徴量を抽出する。詳細については上述した通りであ

10

20

30

40

50

る。

【 0 0 8 1 】

( iii ) ステップ 2 2 0 3

プロセッサ 1 0 1 は、収録した口語調音声コーパスデータ 1 0 3 4 の発話テキストに対して読み上げ韻律・音韻予測モデル 1 0 3 2 を適用し、読み上げ調の韻律特徴量および音韻特徴量を予測する。つまり、ここでは、発話者がこのテキストに対して、読み上げ調スタイルで発話した場合は、韻律・音韻特徴がどのようなものであるかが分かる。詳細は上述した通りである。

【 0 0 8 2 】

( iv ) ステップ 2 2 0 4

プロセッサ 1 0 1 は、ステップ 2 2 0 2 で抽出した収録口語調音声の特徴量と、ステップ 2 2 0 3 で予測した韻律・音韻特徴量とを比較し、口語調への寄与度（口語調度）を計算する。

【 0 0 8 3 】

( v ) ステップ 2 2 0 5

プロセッサ 1 0 1 は、ステップ 2 2 0 4 で算出された口語調度（韻律特徴の差分情報）を用いて、収録した口語調音声の各セグメントに口語調度を付与し、口語調度付き口語調音声データ 1 7 0 2 を生成する。詳細は上述した通りである。

【 0 0 8 4 】

( vi ) ステップ 2 2 0 6

プロセッサ 1 0 1 は、ステップ 2 2 0 5 で得られた口語調度付き口語調音声データ 1 7 0 2 の音声波形を蓄積し、音声合成に用いる口語調音声 D B 1 0 3 8 を作成する。

【 0 0 8 5 】

( vii ) ステップ 2 2 0 7

プロセッサ 1 0 1 は、口語調度付き口語調音声データ 1 7 0 2 の韻律・音韻情報（図 1 0 参照）を用いて、音声合成に用いる口語調韻律・音韻モデル 1 0 3 7 を作成する。詳細は上述した通りである。

【 0 0 8 6 】

( viii ) ステップ 2 2 0 8

プロセッサ 1 0 1 は、ステップ 2 2 0 4 で得られた口語調度付き口語調音声データ 1 7 0 2 を用いて、口語調予測モデル（口語調度予測モデル） 1 7 0 4 を生成する。詳細は上述した通りである。

【 0 0 8 7 】

( 3 ) まとめ

( i ) 第 1 の実施形態では、口語調音声データから韻律特徴量を抽出し、一方、当該口語調音声データに対応するテキストデータに対して読み上げ韻律・音韻予測モデルを適用して読み上げ調の韻律特徴量を予測する。次に、これらの韻律特徴量の差分を取り、差分値が所定の閾値（経験から設定される値）よりも大きい箇所を口語調の特徴部分（音声合成に用いる口語調データ）として抽出する。これらの処理は、収録した口語調音声コーパスと読み上げ韻律・音韻予測モデルを与えれば自動的に実行される。このように、口語調音声を始めとする韻律や声質の変化が大きい発話スタイルの合成音声から、その特徴を担う部分（口語調音声の場合は、口語調表現部分）を自動的に抽出するので、作業コストを抑えることができるうえ、異なった作業による基準の不統一を改善できる。

【 0 0 8 8 】

第 1 の実施形態では、口語調表現抽出ルールが生成される。このルールは、口語調の特徴部分のテキストデータを用いて、与えられるテキストデータにおける口語調表現を抽出するためのルールとして生成される。この場合、特徴部分のテキストデータに加えて、当該テキストデータが含まれる口語調テキストの前後のコンテキスト情報を用いて当該ルールを生成するようにしても良い。このようなルールを作成することにより、このルールに従って生成された合成音声を、より自然で安定的な口語調音声とすることができるように

10

20

30

40

50

なる。

【 0 0 8 9 】

第2の実施形態では、口語調韻律データと読み上げ調韻律データの差分値に基づいて、テキストのセグメントに対して、当該セグメントの口語調の程度を示す口語調度を算出し、これを口語調音声データに付与する。そして、この口語調度が付与された口語調音声データを用いて、音声合成用データが生成される。第2の実施形態による音声合成用データは、音声合成すべき入力テキストの口語調度を予測するための統計モデル（口語調度予測モデル）となっている。第2の実施形態によっても上述の第1の実施形態と同様の技術的効果を期待することができる。

【 0 0 9 0 】

(ii) 本発明は、実施形態の機能を実現するソフトウェアのプログラムコードによっても実現できる。この場合、プログラムコードを記録した記憶媒体をシステム或は装置に提供し、そのシステム或は装置のコンピュータ（又はCPUやMPU）が記憶媒体に格納されたプログラムコードを読み出す。この場合、記憶媒体から読み出されたプログラムコード自体が前述した実施形態の機能を実現することになり、そのプログラムコード自体、及びそれを記憶した記憶媒体は本発明を構成することになる。このようなプログラムコードを供給するための記憶媒体としては、例えば、フレキシブルディスク、CD-ROM、DVD-ROM、ハードディスク、光ディスク、光磁気ディスク、CD-R、磁気テープ、不揮発性のメモリカード、ROMなどが用いられる。

【 0 0 9 1 】

また、プログラムコードの指示に基づき、コンピュータ上で稼動しているOS（オペレーティングシステム）などが実際の処理の一部又は全部を行い、その処理によって前述した実施の形態の機能が実現されるようにしてもよい。さらに、記憶媒体から読み出されたプログラムコードが、コンピュータ上のメモリに書きこまれた後、そのプログラムコードの指示に基づき、コンピュータのCPUなどが実際の処理の一部又は全部を行い、その処理によって前述した実施の形態の機能が実現されるようにしてもよい。

【 0 0 9 2 】

さらに、実施の形態の機能を実現するソフトウェアのプログラムコードを、ネットワークを介して配信することにより、それをシステム又は装置のハードディスクやメモリ等の記憶手段又はCD-RW、CD-R等の記憶媒体に格納し、使用時にそのシステム又は装置のコンピュータ（又はCPUやMPU）が当該記憶手段や当該記憶媒体に格納されたプログラムコードを読み出して実行するようにしてもよい。

【 0 0 9 3 】

ここで述べたプロセス及び技術は本質的に如何なる特定の装置に関連することはなく、コンポーネントの如何なる相応しい組み合わせによっても実装できる。更に、汎用目的の多様なタイプのデバイスがここで記述内容に従って使用可能である。ここで述べた方法のステップを実行するのに、専用の装置を構築するのも有益である。また、実施形態に開示されている複数の構成要素の適宜な組み合わせにより、種々の発明を形成することもできる。例えば、実施形態に示される全構成要素から幾つかの構成要素を削除してもよい。さらに、異なる実施形態にわたる構成要素を適宜組み合わせてもよい。本発明は、具体例に関連して記述したが、これらは、すべての観点において限定の為ではなく説明のためである。本分野にスキルのある者であれば、本発明を実施するのに相応しいハードウェア、ソフトウェア、及びファームウェアの多数の組み合わせがあることを理解できるものと考えられる。例えば、記述したソフトウェアは、アセンブラ、C/C++、perl、Shell、PHP、Java（登録商標）等の広範囲のプログラム又はスクリプト言語で実装できる。

【 0 0 9 4 】

さらに、上述の実施形態において、制御線や情報線は説明上必要と考えられるものを示しており、製品上必ずしも全ての制御線や情報線を示しているとは限らない。全ての構成が相互に接続されていてもよい。

【符号の説明】

**【 0 0 9 5 】**

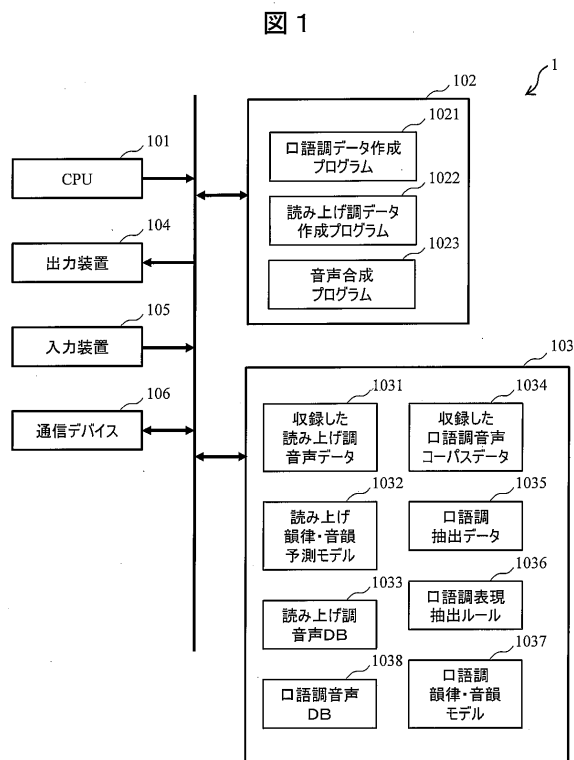
- ```

1  音声合成システム
1 0  読み上げ調データ作成処理
1 1  口語調データ作成処理
1 2  音声合成処理
1 0 1  CPU
1 0 2  メモリ
1 0 3  記憶装置
1 0 4  出力装置
1 0 5  入力装置
1 0 6  通信デバイス

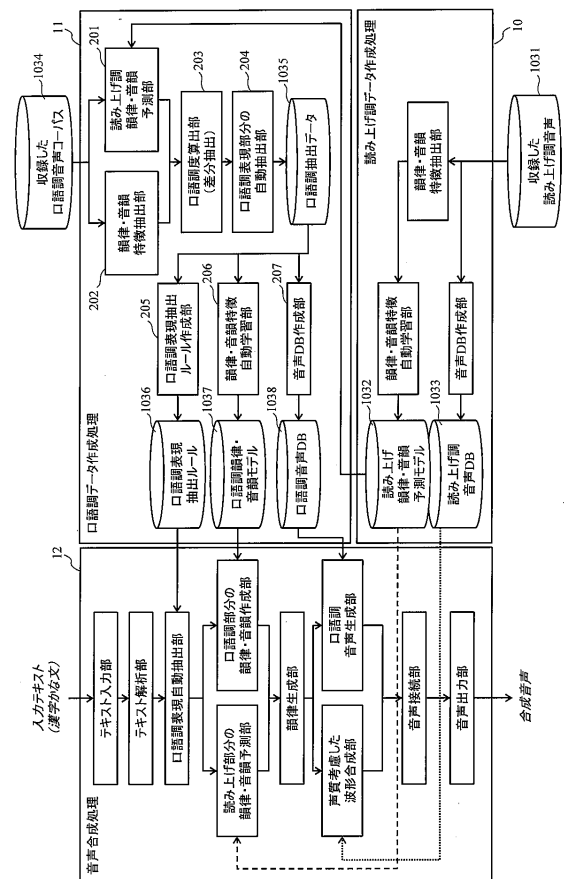
```

10

【 図 1 】



【圖 2】





【図 3】

図 3

| ID | 音素 | 長さ(ms) | 高さ(Hz) | 強さ(dB) |
|----|----|--------|--------|--------|
| 1  | KY | 25     | 232    | -40    |
| 2  | O  | 83     | 248    | -20    |
| 3  | O  | 75     | 274    | -20    |
| 4  | W  | 32     | 262    | -30    |
| 5  | A  | 75     | 239    | -25    |
| 6  | A  | 62     | 230    | -25    |
| 7  | M  | 21     | 239    | -30    |
| 8  | E  | 55     | 263    | -20    |
| 9  | K  | 27     | 252    | -35    |
| 10 | A  | 68     | 237    | -30    |
| 11 | N  | 33     | 225    | -35    |
| 12 | A  | 50     | 200    | -40    |
| 13 | A  | 30     | 180    | -40    |

【図 4】

図 4

| ID | 音素 | 長さ(ms) | 高さ(Hz) | 強さ(dB) |
|----|----|--------|--------|--------|
| 1  | KY | 25     | 232    | -40    |
| 2  | O  | 83     | 248    | -20    |
| 3  | O  | 75     | 274    | -20    |
| 4  | W  | 32     | 262    | -30    |
| 5  | A  | 75     | 239    | -25    |
| 6  | A  | 62     | 230    | -25    |
| 7  | M  | 21     | 239    | -30    |
| 8  | E  | 55     | 263    | -20    |
| 9  | K  | 27     | 252    | -35    |
| 10 | A  | 68     | 207    | -25    |
| 11 | N  | 33     | 215    | -25    |
| 12 | A  | 74     | 248    | -20    |
| 13 | A  | 89     | 259    | -25    |

【図 5】

図 5

| ID | 音素 | 長さ差分<br>(ms) | 高さ差分<br>(Hz) | 強さ差分<br>(dB) | 口語調度 |
|----|----|--------------|--------------|--------------|------|
| 1  | KY | 0            | 0            | 0            | 0    |
| 2  | O  | 0            | 0            | 0            | 0    |
| 3  | O  | 0            | 0            | 0            | 0    |
| 4  | W  | 0            | 0            | 0            | 0    |
| 5  | A  | 0            | 0            | 0            | 0    |
| 6  | A  | 0            | 0            | 0            | 0    |
| 7  | M  | 0            | 0            | 0            | 0    |
| 8  | E  | 0            | 0            | 0            | 0    |
| 9  | K  | 0            | 0            | 0            | 0    |
| 10 | A  | 0            | -30          | 5            | 14   |
| 11 | N  | 0            | -10          | 10           | 3    |
| 12 | A  | 24           | 48           | 20           | 35.2 |
| 13 | A  | 59           | 79           | 15           | 60.2 |

【図 7】

図 7

| ID | 形態素 | 口語調度 |
|----|-----|------|
| 1  | 今日  | 0    |
| 2  | は   | 0    |
| 3  | 雨   | 0    |
| 4  | かなあ | 36.4 |

【図 8】

図 8

| ID | アクセント句 | 口語調度 |
|----|--------|------|
| 1  | 今日 は   | 0    |
| 2  | 雨かなあ   | 36.4 |

【図 6】


図 6

| ID | 音節<br>(モーラ) | 口語調度 |
|----|-------------|------|
| 1  | キョ          | 0    |
| 2  | ウ           | 0    |
| 3  | ワ           | 0    |
| 4  | ア           | 0    |
| 5  | メ           | 0    |
| 6  | カ           | 14   |
| 7  | ナ           | 35.2 |
| 8  | ア           | 60.2 |

【図 9】

図 9


抽出した口語調音声データ

| ID | 文字列 | 口語調音声DB                                                                               |
|----|-----|---------------------------------------------------------------------------------------|
| 1  | かなあ |  |

【図 10】

図 10

抽出した口語調韻律データ

| ID | 文字列 | 口語調韻律データ(F0)                                                                      |
|----|-----|-----------------------------------------------------------------------------------|
| 1  | かなあ |  |

【図 11】

図 11

抽出した口語調テキスト

| ID | 文字列 | 所属文書    |
|----|-----|---------|
| 1  | かなあ | 今日は雨かなあ |

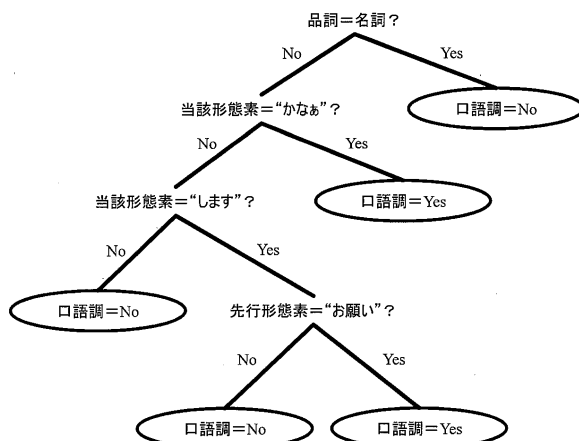
【図 12】

図 12

| ID | 表記  | 口語調 | コンテキスト                                                     |
|----|-----|-----|------------------------------------------------------------|
| 1  | 今日  | No  | 先行形態素=なし;<br>当該形態素="今日";<br>後続形態素="は";<br>品詞=名詞;<br>...    |
| 2  | は   | No  | 先行形態素="今日";<br>当該形態素="は";<br>後続形態素="雨";<br>品詞=助詞;<br>...   |
| 3  | 雨   | No  | 先行形態素="は";<br>当該形態素="雨";<br>後続形態素="かなあ";<br>品詞=名詞;<br>...  |
| 4  | かなあ | Yes | 先行形態素="雨";<br>当該形態素="かなあ";<br>後続形態素=なし;<br>品詞=助詞;<br>...   |
| 5  | お願い | No  | 先行形態素=なし;<br>当該形態素="お願い";<br>後続形態素="します";<br>品詞=名詞;<br>... |
| 6  | します | Yes | 先行形態素="お願い";<br>当該形態素="します";<br>後続形態素=なし;<br>品詞=動詞;<br>... |

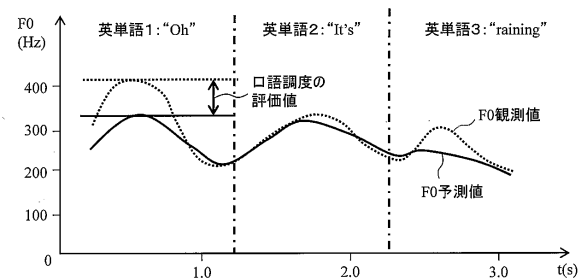
【図 13】

図 13



【図 14】

図 14

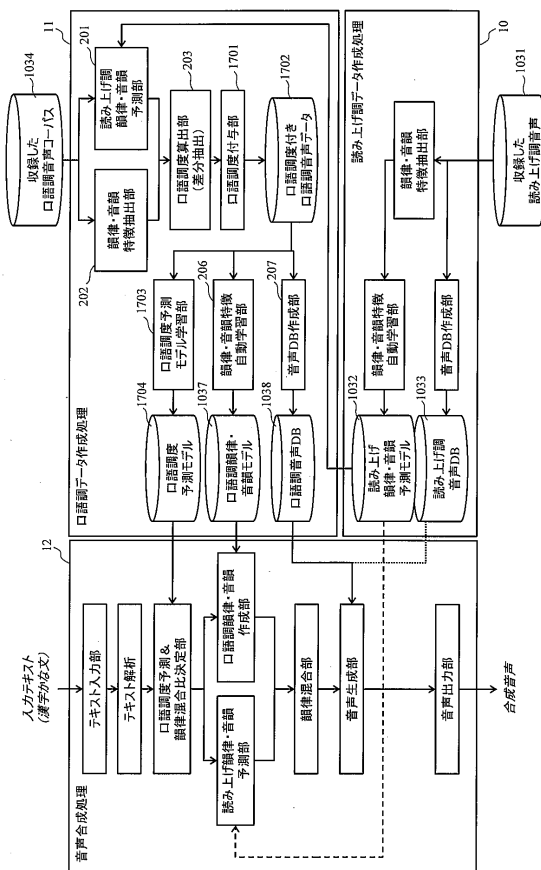


【図 15】

図 15

| ID | 単語      | 観測した最大F0値 | 予測したd最大F0値 | 口語調度 |
|----|---------|-----------|------------|------|
| 1  | Oh      | 405       | 325        | 80   |
| 2  | It's    | 335       | 330        | 5    |
| 3  | Raining | 325       | 285        | 40   |

【 ㄨ 1 7 】



【 図 1 9 】



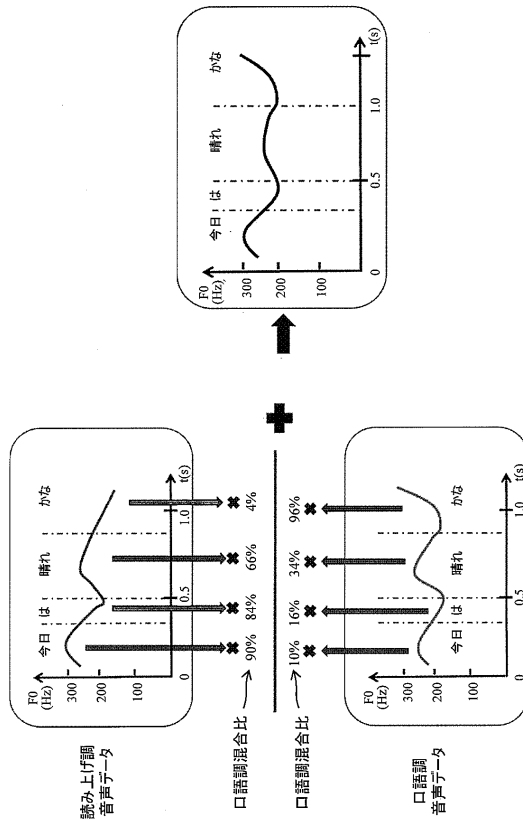
| ID | 表記  | 口語調度 | 口語調混合比 |
|----|-----|------|--------|
| 1  | 今日  | 5    | 10%    |
| 2  | は   | 8    | 16%    |
| 3  | 晴れ  | 17   | 34%    |
| 4  | かなあ | 48   | 96%    |

图 20

| ID | 音素 | 口語調混合比 | 予測された<br>口語調韻律<br>(F0) | 予測された<br>読み上げ調韻律<br>(F0) | 生成され<br>た韻律<br>(F0) |
|----|----|--------|------------------------|--------------------------|---------------------|
| 1  | KY | 10%    | 232                    | 232                      | 232                 |
| 2  | O  | 10%    | 248                    | 248                      | 248                 |
| 3  | O  | 10%    | 274                    | 274                      | 274                 |
| 4  | W  | 16%    | 262                    | 262                      | 262                 |
| 5  | A  | 16%    | 239                    | 239                      | 239                 |
| 6  | H  | 34%    | 235                    | 235                      | 235                 |
| 7  | A  | 34%    | 230                    | 230                      | 230                 |
| 8  | R  | 34%    | 239                    | 239                      | 239                 |
| 9  | E  | 34%    | 263                    | 263                      | 263                 |
| 10 | K  | 96%    | 252                    | 252                      | 252                 |
| 11 | A  | 96%    | 207                    | 237                      | 208                 |
| 12 | N  | 96%    | 215                    | 225                      | 215                 |
| 13 | A  | 96%    | 248                    | 200                      | 246                 |
| 14 | A  | 96%    | 259                    | 180                      | 256                 |

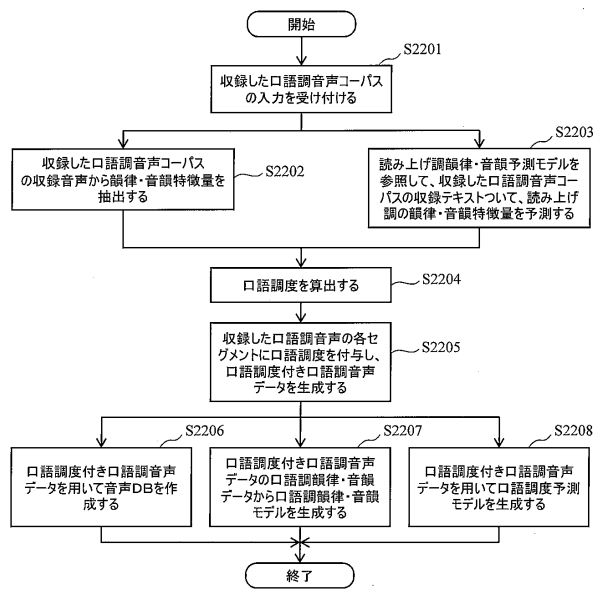
【図 21】

図 21



【図 22】

図 22



---

フロントページの続き

(56)参考文献 特開 2 0 1 3 - 2 4 2 5 1 5 ( J P , A )  
特開 2 0 1 3 - 0 1 5 6 9 3 ( J P , A )  
特開 2 0 0 3 - 3 3 7 5 9 2 ( J P , A )  
特開 2 0 1 4 - 0 6 2 9 7 0 ( J P , A )  
特開 2 0 1 2 - 1 9 8 2 7 7 ( J P , A )  
特開 2 0 0 3 - 3 0 2 9 9 2 ( J P , A )  
特開 2 0 0 4 - 2 2 6 5 0 5 ( J P , A )  
特開 2 0 1 4 - 1 6 3 9 7 8 ( J P , A )

(58)調査した分野(Int.Cl. , D B 名)

G 1 0 L 1 3 / 0 0 - 1 3 / 1 0