

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7362976号
(P7362976)

(45)発行日 令和5年10月18日(2023.10.18)

(24)登録日 令和5年10月10日(2023.10.10)

(51)国際特許分類

F I

G 1 0 L 13/06 (2013.01)
G 1 0 L 13/10 (2013.01)
G 1 0 L 25/30 (2013.01)

G 1 0 L 13/06 1 2 0 Z
G 1 0 L 13/10 1 1 1 F
G 1 0 L 13/10 1 1 1 E
G 1 0 L 13/10 1 1 1 A
G 1 0 L 13/10 1 1 1 C

請求項の数 6 (全27頁) 最終頁に続く

(21)出願番号 特願2022-49374(P2022-49374)
(22)出願日 令和4年3月25日(2022.3.25)
(62)分割の表示 特願2018-113433(P2018-113433)
の分割
原出願日 平成30年6月14日(2018.6.14)
(65)公開番号 特開2022-81691(P2022-81691A)
(43)公開日 令和4年5月31日(2022.5.31)
審査請求日 令和4年3月25日(2022.3.25)

(73)特許権者 000004352
日本放送協会
東京都渋谷区神南2丁目2番1号
(73)特許権者 399060908
一般財団法人NHK財団
東京都世田谷区用賀四丁目10番1号
(74)代理人 100121119
弁理士 花村 泰伸
(72)発明者 清山 信正
東京都世田谷区砧一丁目10番11号
日本放送協会放送技術研究所内
(72)発明者 栗原 清
東京都世田谷区砧一丁目10番11号
日本放送協会放送技術研究所内
(72)発明者 熊野 正

最終頁に続く

(54)【発明の名称】 音声合成装置及びプログラム

(57)【特許請求の範囲】

【請求項1】

音声合成対象のテキストを言語分析し、言語特徴量を求める言語分析部と、
前記言語分析部により求めた前記言語特徴量に、音響の特徴を調整するための調整パラ
メータの調整量情報を追加する調整量追加部と、
前記調整量追加部により前記調整量情報が追加された前記言語特徴量に基づき、予め学
習された統計モデルを用いて、音響特徴量を推定する音響特徴量推定部と、
前記音響特徴量推定部により推定された前記音響特徴量に基づいて、音声信号を合成し
、前記テキストに対して前記調整パラメータによる調整が加えられた音声信号を出力する
音声生成部と、を備えた音声合成装置であって、
前記音響特徴量推定部が用いる統計モデルは、
予め設定されたテキストを言語分析し、学習言語特徴量を求める学習言語分析部と、
前記テキストに対応する音声信号を音響分析し、学習音響特徴量を求める音声分析部と、
前記学習言語分析部により求めた前記学習言語特徴量及び前記音声分析部により求めた
前記学習音響特徴量を時間的に対応付ける対応付け部と、
前記対応付け部により対応付けられた前記学習言語特徴量に、音響の特徴を調整するた
めの調整パラメータの調整量情報を追加する学習調整量追加部と、
前記対応付け部により対応付けられた前記学習音響特徴量を、前記調整パラメータの前
記調整量情報に従って調整する学習音響特徴量調整部と、
前記学習調整量追加部により前記調整量情報が追加された前記学習言語特徴量及び前記

学習音響特徴量調整部により調整された前記学習音響特徴量を用いて、統計モデルを学習する学習部と、

を備えた学習装置によって、予め学習された統計モデルであることを特徴とする音声合成装置。

【請求項 2】

請求項 1 に記載の音声合成装置において、
前記統計モデルは、ニューラルネットワークで構成された時間長モデル及び音響モデルからなり、

前記音響特徴量推定部は、
前記時間長モデルを用いて、音素毎の前記言語特徴量を前記時間長モデルの入力データとして、前記時間長モデルの出力データである音素毎の時間長を推定し、

音素毎の前記時間長からフレーム毎の時間長を生成し、
前記音響モデルを用いて、フレーム毎の前記言語特徴量及びフレーム毎の前記時間長を入力データとし、前記音響モデルの出力データであるフレーム毎の前記音響特徴量を推定する、ことを特徴とする音声合成装置。

10

【請求項 3】

請求項 1 または 2 に記載の音声合成装置において、
前記調整パラメータを、話速または時間長、パワー、ピッチ、及び抑揚の 4 つのパラメータのうちいずれか 1 つまたは 2 つ以上の組み合わせとする、ことを特徴とする音声合成装置。

20

【請求項 4】

請求項 1 または 2 に記載の音声合成装置において、
前記調整パラメータを、話速または時間長、パワー、ピッチ、及び抑揚の 4 つのパラメータとし、
当該 4 つのパラメータのうちいずれか 1 つのパラメータの調整量は、所定範囲内の任意の値が指定され、他の 3 つのパラメータの調整量は、固定値が用いられる、ことを特徴とする音声合成装置。

【請求項 5】

請求項 1 または 2 に記載の音声合成装置において、
前記調整パラメータを、話速または時間長、パワー、ピッチ、及び抑揚の 4 つのパラメータとし、
当該 4 つのパラメータにおけるそれぞれの調整量は、それぞれの所定範囲内の任意の値が指定される、ことを特徴とする音声合成装置。

30

【請求項 6】

コンピュータを、請求項 1 から 5 までのいずれか一項に記載の音声合成装置として機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、テキストから音声信号を合成するための統計モデルを用いて音声信号を合成する音声合成装置及びプログラムに関する。

40

【背景技術】

【0002】

従来、テキストとこれに対応する音声信号を用いて統計モデルを学習し、任意のテキストに対する合成音声を得る方法として、ディープニューラルネットワーク (DNN: Deep Neural Network) を用いた深層学習 (DL: Deep Learning) に基づく技術が知られている (例えば、非特許文献 1 を参照)。

【0003】

一方、音声信号の読み上げ方を調整する方法として、音声分析生成処理に基づく技術が知られている (例えば、非特許文献 2 を参照)。

50

【 0 0 0 4 】

図 1 5 は、非特許文献 1 に記載された従来の学習方法及び合成方法を示す説明図である。この学習方法を実現する学習装置は、事前に用意された音声コーパスのテキストとこれに対応する音声信号を用いて、テキストについては言語分析処理により言語特徴量を抽出する（ステップ S 1 5 0 1）。また、学習装置は、音声信号について音声分析処理により音響特徴量を抽出する（ステップ S 1 5 0 2）。

【 0 0 0 5 】

学習装置は、言語特徴量と音響特徴量の時間対応付けを行い（ステップ S 1 5 0 3）、言語特徴量と音響特徴量を用いて統計モデルを学習する（ステップ S 1 5 0 4）。

【 0 0 0 6 】

また、この合成方法を実現する音声合成装置は、任意のテキストを入力し、テキストの言語分析処理により言語特徴量を抽出する（ステップ S 1 5 0 5）。そして、音声合成装置は、学習装置により学習された統計モデルを用いて、言語特徴量から音響特徴量を推定し（ステップ S 1 5 0 6）、音声生成処理により、音響特徴量から音声信号波形を求める（ステップ S 1 5 0 7）。これにより、任意のテキストに対応する合成音声信号を得ることができる。

【 0 0 0 7 】

図 1 6 は、非特許文献 2 に記載された従来の音声信号調整方法を示す説明図である。この音声信号調整方法を実現する音声調整装置は、音声分析処理により、音声信号からフレーム毎の音響特徴量を抽出し（ステップ S 1 6 0 1）、調整パラメータに基づいて、音響特徴量の所望の部分に所望の調整を加える（ステップ S 1 6 0 2）。

【 0 0 0 8 】

音声調整装置は、音声生成処理により、調整が加えられたフレーム毎の音響特徴量から音声信号を生成する（ステップ S 1 6 0 3）。これにより、調整を加えた音声信号を得ることができる。

【 先行技術文献 】

【 非特許文献 】

【 0 0 0 9 】

【 文献 】 Zhizheng Wu, Oliver Watts, Simon King, “ Merlin : An Open Source Neural Network Speech Synthesis System ”, in Proc. 9th ISCA Speech Synthesis Workshop (SSW9), September 2016, Sunnyvale, CA, USA.

M. Morise, F. Yokomori, and K. Ozawa, “ WORLD : a vocoder-based high-quality speech synthesis system for real-time applications ”, IEICE transactions on information and systems, vol. E99-D, no, 7, pp. 1877-1884, 2016

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 1 0 】

例えば、放送番組等のコンテンツ制作に合成音声信号を利用する際に、演出効果として、テキストの特定部分の読み上げ方を調整した合成音声信号が求められることがある。

【 0 0 1 1 】

前述の非特許文献 1 の方法は、任意のテキストに対して合成音声信号を得るものであり、同一のテキストに対して常に同一の合成音声信号が得られる。また、前述の非特許文献 2 の方法は、音声信号の読み上げ方を調整するものである。

【 0 0 1 2 】

そこで、テキストの特定部分の読み上げ方を調整した合成音声信号を求める方法として、前述の非特許文献 1 , 2 を組み合わせることが想定される。

【 0 0 1 3 】

図 1 7 は、非特許文献 1 , 2 の従来技術を組み合わせた想定例を示す説明図である。この想定例の学習方法は、図 1 5 に示したステップ S 1 5 0 1 ~ S 1 5 0 4 と同様である（ステップ S 1 7 0 1 ~ S 1 7 0 4）。

10

20

30

40

50

【 0 0 1 4 】

この想定例の合成方法は、図 1 5 に示したステップ S 1 5 0 5 ~ S 1 5 0 7 の処理に、図 1 6 に示したステップ S 1 6 0 2 の処理を挿入したものである。具体的には、音声合成装置は、任意のテキストから言語特徴量を抽出し（ステップ S 1 7 0 5 ）、統計モデルを用いて言語特徴量から音響特徴量を推定する（ステップ S 1 7 0 6 ）。

【 0 0 1 5 】

音声合成装置は、調整パラメータに基づいて、音響特徴量の所望の部分に所望の調整を加える（ステップ S 1 7 0 7 ）。音声合成装置は、音声生成処理により、調整が加えられたフレーム毎の音響特徴量から音声信号を生成する（ステップ S 1 7 0 8 ）。これにより、任意のテキストに対応する合成音声信号を得ることができる。

10

【 0 0 1 6 】

しかしながら、この想定例では、ステップ S 1 7 0 6 にて統計モデルを用いて言語特徴量から推定した音響特徴量は、実際の音声信号から音声分析処理により抽出した音響特徴量とは異なり、時間的に平滑化された特性を持っている。このため、ステップ S 1 7 0 7 にて統計モデルを用いて推定した音響特徴量に調整を加え、ステップ S 1 7 0 8 にて調整後のフレーム毎の音響特徴量から合成音声信号を得ると、合成音声信号に音質劣化を生じてしまう。

【 0 0 1 7 】

このように、図 1 7 に示した想定例では、高品質の合成音声信号を得ることができないという問題があった。このため、テキストの特定部分の読み上げ方を調整した、高品質の合成音声信号を得るために、新たな手法が所望されていた。

20

【 0 0 1 8 】

そこで、本発明は前記課題を解決するためになされたものであり、その目的は、テキストの特定部分の読み上げ方を調整した合成音声信号を生成する際に、高品質の合成音声信号を得ることが可能な音声合成装置及びプログラムを提供することにある。

【課題を解決するための手段】

【 0 0 1 9 】

前記課題を解決するために、請求項 1 の音声合成装置は、音声合成対象のテキストを言語分析し、言語特徴量を求める言語分析部と、前記言語分析部により求めた前記言語特徴量に、音響の特徴を調整するための調整パラメータの調整量情報を追加する調整量追加部と、前記調整量追加部により前記調整量情報が追加された前記言語特徴量に基づき、予め学習された統計モデルを用いて、音響特徴量を推定する音響特徴量推定部と、前記音響特徴量推定部により推定された前記音響特徴量に基づいて、音声信号を合成し、前記テキストに対して前記調整パラメータによる調整が加えられた音声信号を出力する音声生成部と、を備えた音声合成装置であって、前記音響特徴量推定部が用いる統計モデルは、予め設定されたテキストを言語分析し、学習言語特徴量を求める学習言語分析部と、前記テキストに対応する音声信号を音響分析し、学習音響特徴量を求める音声分析部と、前記学習言語分析部により求めた前記学習言語特徴量及び前記音声分析部により求めた前記学習音響特徴量を時間的に対応付ける対応付け部と、前記対応付け部により対応付けられた前記学習言語特徴量に、音響の特徴を調整するための調整パラメータの調整量情報を追加する学習調整量追加部と、前記対応付け部により対応付けられた前記学習音響特徴量を、前記調整パラメータの前記調整量情報に従って調整する学習音響特徴量調整部と、前記学習調整量追加部により前記調整量情報が追加された前記学習言語特徴量及び前記学習音響特徴量調整部により調整された前記学習音響特徴量を用いて、統計モデルを学習する学習部と、を備えた学習装置によって、予め学習された統計モデルであることを特徴とする。

30

40

【 0 0 2 0 】

また、請求項 2 の音声合成装置は、請求項 1 に記載の音声合成装置において、前記統計モデルが、ニューラルネットワークで構成された時間長モデル及び音響モデルからなり、前記音響特徴量推定部が、前記時間長モデルを用いて、音素毎の前記言語特徴量を前記時間長モデルの入力データとして、前記時間長モデルの出力データである音素毎の時間長を

50

推定し、音素毎の前記時間長からフレーム毎の時間長を生成し、前記音響モデルを用いて、フレーム毎の前記言語特徴量及びフレーム毎の前記時間長を入力データとし、前記音響モデルの出力データであるフレーム毎の前記音響特徴量を推定する、ことを特徴とする。

【0021】

また、請求項3の音声合成装置は、請求項1または2に記載の音声合成装置において、前記調整パラメータを、話速または時間長、パワー、ピッチ、及び抑揚の4つのパラメータのうちのいずれか1つまたは2つ以上の組み合わせとする、ことを特徴とする。

【0022】

また、請求項4の音声合成装置は、請求項1または2に記載の音声合成装置において、前記調整パラメータを、話速または時間長、パワー、ピッチ、及び抑揚の4つのパラメータとし、当該4つのパラメータのうちのいずれか1つのパラメータの調整量は、所定範囲内の任意の値が指定され、他の3つのパラメータの調整量は、固定値が用いられる、ことを特徴とする。

10

【0023】

また、請求項5の音声合成装置は、請求項1または2に記載の音声合成装置において、前記調整パラメータを、話速または時間長、パワー、ピッチ、及び抑揚の4つのパラメータとし、当該4つのパラメータにおけるそれぞれの調整量は、それぞれの所定範囲内の任意の値が指定される、ことを特徴とする。

【0024】

また、請求項6のプログラムは、コンピュータを、請求項1から5までのいずれか一項に記載の音声合成装置として機能させることを特徴とする。

20

【発明の効果】

【0025】

以上のように、本発明によれば、テキストの特定部分の読み上げ方を調整した合成音声信号を生成する際に、高品質の合成音声信号を得ることが可能となる。

【図面の簡単な説明】

【0026】

【図1】本発明の実施形態による学習装置の構成例を示すブロック図である。

【図2】学習装置による事前学習処理例を示すフローチャートである。

【図3】言語特徴量のデータ構成例を説明する図である。

30

【図4】音声分析部による音声分析処理例を示すフローチャートである。

【図5】音響特徴量のデータ構成例を説明する図である。

【図6】時間情報が追加された言語特徴量のデータ構成例を説明する図である。

【図7】調整量情報が追加された言語特徴量のデータ構成例を説明する図である。

【図8】時間長モデルの学習処理例を説明する図である。

【図9】音響モデルの学習処理例を説明する図である。

【図10】本発明の実施形態による音声合成装置の構成例を示すブロック図である。

【図11】音声合成装置による音声合成処理例を示すフローチャートである。

【図12】時間長モデルを用いた時間長推定処理例を説明する図である。

【図13】音響モデルを用いた音響特徴量推定処理例を説明する図である。

40

【図14】音声生成部による音声合成処理例を説明する図である。

【図15】非特許文献1に記載された従来の学習方法及び合成方法を示す説明図である。

【図16】非特許文献2に記載された従来の音声信号調整方法を示す説明図である。

【図17】非特許文献1, 2の従来技術を組み合わせた想定例を示す説明図である。

【発明を実施するための形態】

【0027】

以下、本発明を実施するための形態について図面を用いて詳細に説明する。

〔学習装置〕

まず、本発明の実施形態による学習装置について説明する。図1は、学習装置の構成例を示すブロック図であり、図2は、学習装置による事前学習処理例を示すフローチャート

50

である。

【 0 0 2 8 】

この学習装置 1 は、記憶部 1 0 , 1 7、言語分析部 1 1、音声分析部 1 2、対応付け部 1 3、調整量追加部 1 4、音響特徴量調整部 1 5 及び学習部 1 6 を備えている。音声信号はモノラルであり、標本化周波数 4 8 k H z 及びビット数 1 6 で標本化されているものとする。

【 0 0 2 9 】

記憶部 1 0 には、予め設定された音声コーパスが格納されている。音声コーパスは、テキストと、これに対応する音声信号から構成される。例えば、A T R (株式会社国際電気通信基礎技術研究所) により作成された音素バランス 5 0 3 文を利用する場合、テキストと、これを読み上げた音声信号は、5 0 3 対からなる。音声コーパスについては、以下の文献を参照されたい。

磯健一、渡辺隆夫、桑原尚夫、「音声データベース用文セットの設計」、音講論(春)、pp.89-90(1988.3)

【 0 0 3 0 】

言語分析部 1 1 は、記憶部 1 0 から音声コーパスの各テキストを読み出し、テキストについて既知の言語分析処理を行い、音素毎の所定情報からなる言語特徴量を求める(ステップ S 2 0 1)。そして、言語分析部 1 1 は、音素毎の言語特徴量を対応付け部 1 3 に出力する。

【 0 0 3 1 】

具体的には、言語分析部 1 1 は、言語分析処理により、文を構成する音素毎に、音素情報、アクセント情報、品詞情報、アクセント句情報、呼気段落情報及び総数情報を求め、これらの情報からなる言語特徴量を求める。

【 0 0 3 2 】

言語分析処理としては、例えば以下に記載された形態素解析処理が用いられる。

“MeCab: Yet Another Part-of-Speech and Morphological Analyzer”, インターネット < U R L : <http://taku910.github.io/mecab/> >

また、言語分析処理としては、例えば以下に記載された係り受け解析処理が用いられる。

“CaboCha/南瓜: Yet Another Japanese Dependency Structure Analyzer”, インターネット < U R L : <https://taku910.github.io/cabocha/> >

【 0 0 3 3 】

図 3 は、言語特徴量のデータ構成例を説明する図である。図 3 に示すように、言語特徴量は、音素毎に、音素情報、アクセント情報、品詞情報、アクセント句情報、呼気段落情報及び総数情報から構成される。

【 0 0 3 4 】

図 1 及び図 2 に戻って、音声分析部 1 2 は、記憶部 1 0 から音声コーパスの各テキストに対応する各音声信号を読み出し、フレーム毎に音声信号を切り出し、フレーム毎の音声信号について既知の音響分析処理を行う。そして、音声分析部 1 2 は、フレーム毎の所定情報からなる音響特徴量を求め(ステップ S 2 0 2)、フレーム毎の音響特徴量を対応付け部 1 3 に出力する。音響特徴量は、後述するように、1 9 9 次元のデータから構成される。

【 0 0 3 5 】

音響分析処理としては、例えば以下に記載された音響分析処理が用いられる。

“A high-quality speech analysis, manipulation and synthesis system”, インターネット < U R L : <https://github.com/mmorise/World> >

また、音響分析処理としては、例えば以下に記載された音声信号処理が用いられる。

“Speech Signal Processing Toolkit(SPTK) Version 3.11 December 25, 2017”, インターネット < U R L : <http://sp-tk.sourceforge.net/> >

“REFERENCE MANUAL for Speech Signal Processing Toolkit Ver. 3.9”

【 0 0 3 6 】

10

20

30

40

50

図4は、音声分析部12による音声分析処理例を示すフローチャートである。音声分析部12は、記憶部10から音声コーパスの各音声信号を読み出し、フレーム長25msの音声信号をフレームシフト5ms毎に切り出す(ステップS401)。そして、音声分析部12は、フレーム毎の音声信号について音響分析処理を行い、スペクトル、ピッチ周波数及び非周期成分を求める(ステップS402)。

【0037】

音声分析部12は、スペクトルをメルケプストラム分析してメルケプストラム係数MGCを求める(ステップS403)。また、音声分析部12は、ピッチ周波数から有声/無声判定情報VUVを求め、ピッチ周波数の有声区間を対数化し、無声及び無音区間については前後の有声区間の情報を用いて補間することにより、対数ピッチ周波数LFOを求める(ステップS404)。また、音声分析部12は、非周期成分をメルケプストラム分析して帯域非周期成分BAPを求める(ステップS405)。

10

【0038】

これにより、静特性の音響特徴量として、フレーム毎に、メルケプストラム係数MGC、有声/無声判定情報VUV、対数ピッチ周波数LFO及び帯域非周期成分BAPが得られる。

【0039】

音声分析部12は、メルケプストラム係数MGCの1次差分 ΔMGC を算出して1次差分メルケプストラム係数 ΔMGC を求め(ステップS406)、2次差分 $\Delta^2 MGC$ を算出して2次差分メルケプストラム係数 $\Delta^2 MGC$ を求める(ステップS407)。

20

【0040】

音声分析部12は、対数ピッチ周波数LFOの1次差分 ΔLFO を算出して1次差分対数ピッチ周波数 ΔLFO を求め(ステップS408)、2次差分 $\Delta^2 LFO$ を算出して2次差分対数ピッチ周波数 $\Delta^2 LFO$ を求める(ステップS409)。

【0041】

音声分析部12は、帯域非周期成分BAPの1次差分 ΔBAP を算出して1次差分帯域非周期成分 ΔBAP を求め(ステップS410)、2次差分 $\Delta^2 BAP$ を算出して2次差分帯域非周期成分 $\Delta^2 BAP$ を求める(ステップS411)。

【0042】

これにより、動特性の音響特徴量として、フレーム毎に、1次差分メルケプストラム係数 ΔMGC 、2次差分メルケプストラム係数 $\Delta^2 MGC$ 、1次差分対数ピッチ周波数 ΔLFO 、2次差分対数ピッチ周波数 $\Delta^2 LFO$ 、1次差分帯域非周期成分 ΔBAP 及び2次差分帯域非周期成分 $\Delta^2 BAP$ が得られる。

30

【0043】

音声分析部12は、フレーム毎の静特性及び動特性の所定情報からなる音響特徴量を対応付け部13に出力する。

【0044】

図5は、音響特徴量のデータ構成例を説明する図である。図5に示すように、音響特徴量は、フレーム毎に、静特性のメルケプストラム係数MGC、対数ピッチ周波数LFO及び帯域非周期成分BAP、動特性の1次差分メルケプストラム係数 ΔMGC 、1次差分対数ピッチ周波数 ΔLFO 、1次差分帯域非周期成分 ΔBAP 、2次差分メルケプストラム係数 $\Delta^2 MGC$ 、2次差分対数ピッチ周波数 $\Delta^2 LFO$ 及び2次差分帯域非周期成分 $\Delta^2 BAP$ 、並びに静特性の有声/無声判定情報VUVから構成される。この音響特徴量は、後述するように、199次元のデータから構成される。

40

【0045】

図1及び図2に戻って、対応付け部13は、言語分析部11から音素毎の言語特徴量を入力すると共に、音声分析部12からフレーム毎の音響特徴量を入力する。そして、対応付け部13は、既知の音素アラインメントの技術を用いて、音素毎の言語特徴量とフレーム毎の音響特徴量とを時間的に対応付けることで、テキストの文を構成する各音素が音声信号のどの時刻に位置(対応)するのかを算出する(ステップS203)。

50

【 0 0 4 6 】

対応付け部 1 3 は、音素毎に、対応する開始フレームの番号及び終了フレームの番号からなる時間情報を生成し、言語特徴量を構成する音素毎の所定情報に時間情報を追加すると共に、音素の時間長（フレーム数）を求める。そして、対応付け部 1 3 は、対応付けた音素毎の時間情報を追加した言語特徴量を調整量追加部 1 4 に出力する。また、対応付け部 1 3 は、音素毎の時間長を音響特徴量に含め、対応付けたフレーム毎の音響特徴量（時間長については音素毎のデータ）を音響特徴量調整部 1 5 に出力する。

【 0 0 4 7 】

ここで、言語特徴量に追加される時間情報は、ミリ秒単位の情報である。また、音素毎の時間長は、後述する統計モデルにおける時間長モデルの出力データに用いられ、音素におけるミリ秒単位の時間の長さをフレームシフト 5 m s で除算した 5 m s フレーム単位の数値、すなわち音素のフレーム数が用いられる。

10

【 0 0 4 8 】

音素アラインメントの技術としては、例えば以下に記載された音声認識処理が用いられる。

“ The Hidden Markov Model Toolkit (HTK) ” , インターネット < URL : <http://htk.eng.cam.ac.uk> >

“ Speech Signal Processing Toolkit(SPTK) Version 3.11 December 25, 2017 ”

【 0 0 4 9 】

尚、対応付け部 1 3 は、言語特徴量及び音響特徴量の時間的な対応付け処理の後に、各文の文頭及び文末の無音区間を削除する。

20

【 0 0 5 0 】

図 6 は、時間情報が追加された言語特徴量のデータ構成例を説明する図である。図 6 に示すように、時間情報が追加された言語特徴量は、図 3 に示した言語特徴量に時間情報を追加して構成される。具体的には、この言語特徴量は、音素毎に、時間情報、音素情報、アクセント情報、品詞情報、アクセント句情報、呼気段落情報及び総数情報から構成される。

【 0 0 5 1 】

図 1 及び図 2 に戻って、調整量追加部 1 4 は、対応付け部 1 3 から音素毎の言語特徴量を入力すると共に、所定の調整パラメータを入力する。そして、調整量追加部 1 4 は、言語特徴量を構成する音素毎の所定情報に、調整パラメータの調整量情報を追加する（ステップ S 2 0 4）。調整量追加部 1 4 は、音素毎の調整量情報を追加した言語特徴量を学習部 1 6 に出力する。

30

【 0 0 5 2 】

所定の調整パラメータは、音声信号を調整する（音響の特徴を調整する）ためのパラメータであり、話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} のうちのいずれか、またはこれらの組み合わせとし、ユーザにより選択されるものとする。また、調整パラメータは、学習部 1 6 において学習データの一部として用いられる。

【 0 0 5 3 】

話速 R_{ST} は話速の調整量を示し、パワー R_{PW} はパワー（声の大きさ）の調整量を示し、 R_{PT} はピッチ（声の高さ）の調整量を示し、抑揚 R_{PD} は抑揚（声の高さの変化幅）の調整量を示す。尚、話速の代わりに、時間長を用いるようにしてもよい。

40

【 0 0 5 4 】

話速 R_{ST} の範囲（話速の調整量範囲）は、例えば以下のとおりとする。

（遅い） 0.5 = R_{ST} = 4.0（速い）

これは、話速 R_{ST} は 0.5 から 4.0 までの範囲において、0.5 に近いほど遅く、4.0 に近いほど速いことを意味する。

【 0 0 5 5 】

パワー R_{PW} の範囲（パワーの調整量範囲）は、例えば以下のとおりとする。

（小さい） 1.0E-5 = R_{PW} = 2.0（大きい）

50

これは、パワー R_{PW} は $1.0E-5$ から 2.0 までの範囲において、 $1.0E-5$ に近いほど小さく、 2.0 に近いほど大きいことを意味する。

【 0 0 5 6 】

ピッチ R_{PT} の範囲（ピッチの調整量範囲）は、例えば以下のとおりとする。

（低い） $0.5 = R_{PT} = 2.0$ （高い）

これは、ピッチ R_{PT} は 0.5 から 2.0 までの範囲において、 0.5 に近いほど低く、 2.0 に近いほど高いことを意味する。

【 0 0 5 7 】

抑揚 R_{PD} の範囲（抑揚の調整量範囲）は、例えば以下のとおりとする。

（小さい） $1.0E-5 = R_{PD} = 2.0$ （大きい）

これは、抑揚 R_{PD} は $1.0E-5$ から 2.0 までの範囲において、 $1.0E-5$ に近いほど小さく、 2.0 に近いほど大きいことを意味する。話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} の標準値は、いずれも 1.0 とする。

【 0 0 5 8 】

また、これらの調整パラメータのそれぞれは、例えば以下に示す 11 個のデータから選択されるものとする。すなわち、学習装置 1 における話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} の調整パラメータは、それぞれ 11 個のデータのいずれかが使用される。

[数 1]

$R_{ST} \in \{0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1.6, 2.2, 2.8, 3.4, 4.0\}$

$R_{PW} \in \{0.00001, 0.2, 0.4, 0.6, 0.8, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0\}$

$R_{PT} \in \{0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0\}$. . . (1)

$R_{PD} \in \{0.00001, 0.2, 0.4, 0.6, 0.8, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0\}$

【 0 0 5 9 】

ここで、4つの調整パラメータを以下の調整ベクトルで表現する。

$$R = (R_{ST}, R_{PW}, R_{PT}, R_{PD})$$

話速、パワー等の調整量を変化させないで元の話速、パワー等を維持する場合、調整ベクトルは以下のとおりである。

$$R = (1.0, 1.0, 1.0, 1.0)$$

【 0 0 6 0 】

4つの調整パラメータにおいて、それぞれ 11 個のデータから 1 個のデータが選択されるものとする、全ての組み合わせ数は、 $11^4 = 14,641$ となる。このため、統計モデルを学習するためには、膨大なデータ量が必要となることから、学習の負荷が高くなり、時間もかかってしまう。

【 0 0 6 1 】

そこで、本発明の実施形態では、ユーザは、4つの調整パラメータのうちの1つの調整パラメータについて、所定範囲の 11 個のデータから 1 個のデータを選択し、他の3つの調整パラメータについては、標準値 1.0 を固定値として用いるようにしてもよい。音響特徴量調整部 15 、及び後述する図 10 の音声合成装置 2 についても同様である。

【 0 0 6 2 】

例えば、ユーザは、話速 R_{ST} について 11 個のデータから 1 個のデータを選択し、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} について標準値 1.0 を固定値として用いるものとする、調整ベクトルは以下のとおりである。

$$R = (R_{ST}, 1.0, 1.0, 1.0)$$

10

20

30

40

50

この場合、調整量追加部 14 は、調整パラメータとして、ユーザにより 11 個のデータのうち 1 個のデータが選択された話速 R_{ST} 、並びに、標準値 1.0 を固定値としたパワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} を入力する。

【0063】

このように、4 つの調整パラメータのうちの 1 つの調整パラメータについては 11 個のデータから 1 個のデータが選択され、他の 3 つの調整パラメータについては標準値である 1.0 を固定値として用いることは、調整ベクトル R のいずれか 1 つの要素の軸方向のみに調整量をプロットしたと等価である。この場合の組み合わせ数は、 $10 \times 4 + 1 = 41$ となる。これにより、統計モデルを学習する際に、学習データの数を減らすことができるから、学習処理の負荷を低減し、学習処理の時間を短縮することができる。

10

【0064】

また、本発明の実施形態における他の例として、ユーザは、4 つの調整パラメータを 11 段階で連動させて選択するようにしてもよい。音響特徴量調整部 15、及び後述する図 10 の音声合成装置 2 についても同様である。

【0065】

この場合、調整量追加部 14 は、調整パラメータとして、予め設定された 11 種類のパターンのうち、ユーザにより選択されたいずれかのパターンの話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} を入力する。11 種類のパターンの調整ベクトルは以下のとおりである。

$$R = (a_1, b_1, c_1, d_1), (a_2, b_2, c_2, d_2), \dots, (a_{11}, b_{11}, c_{11}, d_{11})$$

20

$a_1, b_1, \dots, c_{11}, d_{11}$ は、対応する調整パラメータの調整量範囲に含まれる値とする。

【0066】

この場合の組み合わせ数は、11 となる。これにより、統計モデルを学習する際に、学習データの数を一層減らすことができるから、その負荷を一層低減し、その時間を一層短縮することができる。

【0067】

尚、調整量追加部 14 は、文章単位、呼気段落単位またはアクセント句単位で、異なる調整パラメータを入力するようにしてもよい。音響特徴量調整部 15、及び後述する音声合成装置 2 についても同様である。

30

【0068】

図 7 は、調整量情報が追加された言語特徴量のデータ構成例を説明する図である。図 7 に示すように、調整量情報が追加された言語特徴量は、図 6 に示した言語特徴量に、調整パラメータの調整量情報を追加して構成される。具体的には、この言語特徴量は、音素毎に、時間情報、音素情報、アクセント情報、品詞情報、アクセント句情報、呼気段落情報、総数情報及び調整量情報から構成される。

【0069】

調整量情報は、話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} の調整パラメータにおける調整量が反映された情報である。

40

【0070】

前述のとおり、調整量追加部 14 は、話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} のうちのどれか、またはこれらの組み合わせの調整パラメータを入力する。調整量追加部 14 は、例えば話速 R_{ST} のみの調整パラメータを入力した場合、言語特徴量に、入力した話速 R_{ST} 、並びに固定値である標準値 1.0 のパワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} の調整量情報を追加する。また、調整量追加部 14 は、例えば話速 R_{ST} 及びパワー R_{PW} の調整パラメータを入力した場合、言語特徴量に、入力した話速 R_{ST} 及びパワー R_{PW} 、並びに固定値である標準値 1.0 のピッチ R_{PT} 及び抑揚 R_{PD} の調整量情報を追加する。

【0071】

50

図 1 及び図 2 に戻って、音響特徴量調整部 1 5 は、対応付け部 1 3 から、調整量追加部 1 4 が入力する音素毎の言語特徴量に対応するフレーム毎の音響特徴量（時間長については音素毎のデータ）を入力する。また、音響特徴量調整部 1 5 は、調整量追加部 1 4 と同様の所定の調整パラメータを入力する。

【 0 0 7 2 】

音響特徴量調整部 1 5 は、調整パラメータに従ってフレーム毎の音響特徴量を調整し、調整後のフレーム毎の音響特徴量（時間長については音素毎のデータ）を学習部 1 6 に出力する。

【 0 0 7 3 】

話速 R_{ST} の調整パラメータに従い話速が調整される場合、音響特徴量調整部 1 5 は、以下の式のとおり、対応付け部 1 3 から入力した時間長 DUR に話速 R_{ST} の逆数を乗算し、乗算結果を整数化し、新たな時間長 DUR' を求めることで、時間長を調整する。

[数 2]

$$DUR' = \text{int}(DUR \times 1 / R_{ST}) \quad \dots (2)$$

対応付け部 1 3 から入力した時間長を DUR 、調整後の時間長を DUR' とする。

【 0 0 7 4 】

尚、話速 R_{ST} の代わりに時間長の調整パラメータ $R_{DR} (= 1 / R_{ST})$ に従い時間長が調整される場合、音響特徴量調整部 1 5 は、対応付け部 1 3 から入力した時間長 DUR に対し、話速 R_{ST} の逆数の代わりに、時間長の調整パラメータ R_{DR} を乗算し、乗算結果を整数化し、新たな時間長 DUR' を求めることで、時間長を調整する。

【 0 0 7 5 】

音響特徴量調整部 1 5 は、調整後の時間長に応じて、対応付け部 1 3 から入力したフレームの音響特徴量を繰り返しまたは間引きして、音響特徴量のフレーム数を揃えることで、音響特徴量を調整する。このように、音素毎の時間長の調整に応じて、音響特徴量のフレーム数が揃えられる。

【 0 0 7 6 】

尚、音響特徴量調整部 1 5 は、調整後の時間長に応じて、対応するフレームの音響特徴量を繰り返しまたは間引くことで音響特徴量を調整する際に、前後のフレームの音響特徴量を用いて補間を行うようにしてもよい。これにより、高品質の音響特徴量を得ることができる。また、話速 R_{ST} の調整パラメータ及び他の調整パラメータに従い話速等が調整される場合、音響特徴量調整部 1 5 は、話速を調整する前に、他の調整パラメータによる調整を行う。

【 0 0 7 7 】

また、パワー R_{PW} の調整パラメータに従い音声のパワーが調整される場合、音響特徴量調整部 1 5 は、対応付け部 1 3 から入力した音響特徴量に含まれる静特性のメルケプストラム係数 MGC における 0 次元目の値 $MGC[0]$ に、パワー R_{PW} を対数化した値を加算する。

【 0 0 7 8 】

音響特徴量調整部 1 5 は、以下の式のとおり、加算した値と 0 とを比較して大きい方を、新たな静特性のメルケプストラム係数 MGC における 0 次元目の値 $MGC[0]'$ として求めることで、音響特徴量を調整する。

[数 3]

$$MGC[0]' = \max(0, MGC[0] + \log R_{PW}) \quad \dots (3)$$

対応付け部 1 3 から入力した音響特徴量に含まれる静特性のメルケプストラム係数 MGC における 0 次元目の値を $MGC[0]$ 、調整後の値を $MGC[0]'$ とする。

【 0 0 7 9 】

また、ピッチ R_{PT} の調整パラメータに従い音声のピッチ周波数が調整される場合、音響特徴量調整部 1 5 は、対応付け部 1 3 から入力した音響特徴量に含まれる静特性の対数ピッチ周波数 $LF0$ における 0 次元目の値 $LF0[0]$ に、ピッチ R_{PT} を対数化した値を加算する。

10

20

30

40

50

【 0 0 8 0 】

音響特徴量調整部 1 5 は、以下の式のとおり、加算した値と 0 とを比較して大きい方を、新たな静特性の対数ピッチ周波数 L F 0 における 0 次元目の値 L F 0 [0] ' として求めることで、音響特徴量を調整する。

[数 4]

$$L F 0 [0] ' = \max (0 , L F 0 [0] + \log R_T) \quad \dots (4)$$

対応付け部 1 3 から入力した音響特徴量に含まれる静特性の対数ピッチ周波数 L F 0 における 0 次元目の値を L F 0 [0]、調整後の値を L F 0 [0] ' とする。

【 0 0 8 1 】

また、抑揚 R_{PD} の調整パラメータに従い音声の抑揚が調整される場合、音響特徴量調整部 1 5 は、対応付け部 1 3 から入力した音響特徴量に含まれる静特性の対数ピッチ周波数 L F 0 から、予め算出しておいた平均値 $\mu_{L F 0}$ を減算する。そして、音響特徴量調整部 1 5 は、減算結果を、予め算出しておいた標準偏差 $\sigma_{L F 0}$ で除算し、除算結果を求める。平均値 $\mu_{L F 0}$ は、対応付け部 1 3 から入力した音響特徴量に含まれる静特性の対数ピッチ周波数 L F 0 の平均値であり、標準偏差 $\sigma_{L F 0}$ はその標準偏差である。

10

【 0 0 8 2 】

音響特徴量調整部 1 5 は、以下の式のとおり、対応付け部 1 3 から入力した音響特徴量に含まれる静特性の対数ピッチ周波数 L F 0 について、その平均値 $\mu_{L F 0}$ 及び標準偏差 $\sigma_{L F 0}$ を文毎に算出しておくものとする。N は、文に対応するフレーム数である。

[数 5]

20

$$\mu_{L F 0} = \frac{1}{N} \sum_{i=1}^N L F 0_i \quad \dots (5)$$

[数 6]

$$\sigma_{L F 0} = \left(\frac{1}{N} \sum_{i=1}^N (L F 0_i - \mu_{L F 0})^2 \right)^{\frac{1}{2}} \quad \dots (6)$$

30

【 0 0 8 3 】

音響特徴量調整部 1 5 は、標準偏差 $\sigma_{L F 0}$ に、抑揚 R_{PD} を対数化した値を加算し、加算結果と 0 とを比較して大きい方を求める。そして、音響特徴量調整部 1 5 は、前記除算結果に、大きい方の値を乗算し、乗算結果に平均値 $\mu_{L F 0}$ を加算する。

【 0 0 8 4 】

音響特徴量調整部 1 5 は、加算した値と 0 とを比較して大きい方を、新たな静特性の対数ピッチ周波数 L F 0 ' として求める。音響特徴量調整部 1 5 による演算処理の式は以下のとおりである。

[数 7]

$$L F 0 ' = \max (0 , ((L F 0 - \mu_{L F 0}) / \sigma_{L F 0}) \times \max (0 , L F 0 + \log R_{PD}) + \mu_{L F 0}) \quad \dots (7)$$

40

対応付け部 1 3 から入力した音響特徴量に含まれる静特性の対数ピッチ周波数を L F 0、その平均値を $\mu_{L F 0}$ 、その標準偏差を $\sigma_{L F 0}$ 、調整後の静特性の対数ピッチ周波数を L F 0 ' とする。

【 0 0 8 5 】

音響特徴量調整部 1 5 は、前記のように各調整パラメータに従い算出された新たな静特性の 1 次差分 $\Delta L F 0$ を算出して新たな動特性の 1 次差分を求める。また、音響特徴量調整部 1 5 は、2 次差分 $\Delta^2 L F 0$ を算出して新たな動特性の 2 次差分を求める。このようにして、音響特徴量調整部 1 5 は、音響特徴量を調整する。

【 0 0 8 6 】

50

尚、音響特徴量調整部 15 による音響特徴量の調整処理は、調整量追加部 14 による調整量情報の言語特徴量への追加処理と連動するものとする。

【0087】

学習部 16 は、調整量追加部 14 から音素毎の言語特徴量を入力すると共に、音響特徴量調整部 15 からフレーム毎の音響特徴量（時間長については音素毎のデータ）を入力する。そして、学習部 16 は、これらのデータを標準化し、統計モデルである時間長モデル及び音響モデルを学習する。

【0088】

（時間長モデルの学習）

次に、学習部 16 による時間長モデルの学習処理について説明する。図 8 は、時間長モデルの学習処理例を説明する図である。学習部 16 は、調整量追加部 14 から入力した音素毎の言語特徴量に基づいて、言語特徴を表す 3 1 2 次元のバイナリ値及び 1 3 次元の数値データ、並びに 1 次元の調整データを生成する。1 次元の調整データは話速データであり、言語特徴量の次元数は 3 2 6 である。

10

【0089】

ここで、言語特徴量における 3 1 2 次元のバイナリ値及び 1 3 次元の数値データは、言語特徴量に含まれる音素情報、アクセント情報、品詞情報、アクセント句情報、呼気段落情報及び総数情報に基づいて生成される。言語特徴量における 1 次元の調整データは、言語特徴量に含まれる調整量情報（話速の調整量、パワーの調整量、ピッチの調整量及び抑揚の調整量）のうち、話速の調整量に基づいて生成される。

20

【0090】

学習部 16 は、言語特徴量の 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 1 次元の調整データ（話速データ）からなる 3 2 6 次元のデータを、時間長モデルの入力データとして扱う（ステップ S 8 0 1）。

【0091】

学習部 16 は、言語特徴量の 3 2 6 次元の全てのデータを用いて、次元毎に、最大値及び最小値を求めて記憶部 17 に格納すると共に、全てのデータのそれぞれについて、次元毎の最大値及び最小値を用いて標準化する（ステップ S 8 0 2）。

【0092】

また、学習部 16 は、音響特徴量調整部 15 から入力したフレーム毎の音響特徴量（時間長については音素毎のデータ）のうちの音素毎の時間長について、当該時間長の 1 次元のデータを、時間モデルの出力データとして扱う（ステップ S 8 0 3）。この時間長は、5 m s 単位のフレーム数であり、テキストを表現する音素毎に 1 次元の整数値からなる。

30

【0093】

学習部 16 は、時間長の 1 次元の全てのデータを用いて、平均値及び標準偏差を求めて記憶部 17 に格納すると共に、全てのデータのそれぞれについて、平均値及び標準偏差を用いて標準化する（ステップ S 8 0 4）。

【0094】

学習部 16 は、ステップ S 8 0 2、S 8 0 4 から移行して、音素毎に、言語特徴量の 3 2 6 次元の標準化されたデータを入力データとし、時間長の 1 次元の標準化されたデータを出力データとして時間長モデルを学習する（ステップ S 8 0 5）。そして、学習部 16 は、学習済みの時間長モデルを記憶部 17 に格納する。

40

【0095】

ステップ S 8 0 5 における時間長モデルの学習の際には、以下のサイトに記載された技術が用いられる。

“CSTR-Edinburgh/merlin”，インターネット<URL : <https://github.com/CSTR-Edinburgh/merlin>>

後述する図 9 のステップ S 9 0 5 における音響モデルの学習の場合も同様である。

【0096】

時間長モデルは、例えば入力層を 3 2 6 次元、隠れ層を 1 0 2 4 次元の 6 層、出力層を

50

1次元とした順伝播型のニューラルネットワークで構成される。隠れ層における活性化関数は双曲線正接関数が用いられ、損失誤差関数は平均二乗誤差関数が用いられる。また、ミニバッチ数を64、エポック数を100、dropout(ドロップアウト)率を0.5、学習係数の最適化方法として確率的勾配降下法、開始学習率を0.01、10エポックを過ぎてからエポック毎に学習率を指数減衰させ、誤差逆伝播法にて学習するものとする。尚、15エポックを過ぎてから、5エポック連続して評価誤差が減少しない場合は学習を早期終了するものとする。

【0097】

これにより、記憶部17には、統計モデルとして時間長モデルが格納される。また、記憶部17には、統計モデルとして、時間長モデルの入力データである言語特徴量の312次元のバイナリ値、13次元の数値データ及び1次元の調整データ(話速データ)からなる326次元のデータに関する次元毎の最大値及び最小値が格納される。また、記憶部17には、統計モデルとして、時間長モデルの出力データである時間長の1次元のデータに関する平均値及び標準偏差が格納される。

10

【0098】

(音響モデルの学習)

次に、学習部16による音響モデルの学習処理について説明する。図9は、音響モデルの学習処理例を説明する図である。学習部16は、調整量追加部14から入力した音素毎の言語特徴量に基づいて、言語特徴を表す312次元のバイナリ値、13次元の数値データ、4次元の時間データ及び3次元の調整データを生成する。

20

【0099】

4次元の時間データは、当該フレームに対応する音素のフレーム数(1次元のデータ)、及び当該フレームの音素内における位置(3次元のデータ)からなる。3次元の調整データは、パワーデータ、ピッチデータ及び抑揚データである。これらの調整データは、言語特徴量に含まれる調整量情報(話速の調整量、パワーの調整量、ピッチの調整量及び抑揚の調整量)のうち、パワーの調整量、ピッチの調整量及び抑揚の調整量に基づいて生成される。また、言語特徴量の次元数は332である。

【0100】

学習部16は、音素毎の言語特徴量における312次元のバイナリ値、13次元の数値データ、4次元の時間データ及び3次元の調整データ(パワーデータ、ピッチデータ及び抑揚データ)からなる332次元のデータから、フレーム毎の言語特徴量における332次元のデータを生成する。

30

【0101】

学習部16は、フレーム毎の言語特徴量について、言語特徴量の312次元のバイナリ値、13次元の数値データ、4次元の時間データ及び3次元の調整データ(パワーデータ、ピッチデータ及び抑揚データ)からなる332次元のデータを、音響モデルの入力データとして扱う(ステップS901)。

【0102】

学習部16は、言語特徴量の332次元の全てのデータを用いて、次元毎に、最大値及び最小値を求めて記憶部17に格納すると共に、全てのデータのそれぞれについて、次元毎の最大値及び最小値を用いて標準化する(ステップS902)。

40

【0103】

また、学習部16は、音響特徴量調整部15から入力したフレーム毎の音響特徴量(時間長については音素毎のデータ)のうちの時間長を除く音響特徴量について、199次元のデータを、音響モデルの出力データとして扱う(ステップS903)。

【0104】

ここで、前述のとおり、時間長を除く音響特徴量は、静特性のメルケプストラム係数MGC、対数ピッチ周波数LF0及び帯域非周期成分BAP、動特性の1次差分メルケプストラム係数MGC、1次差分対数ピッチ周波数LF0、1次差分帯域非周期成分BAP、2次差分メルケプストラム係数²MGC、2次差分対数ピッチ周波数²LF0及

50

び2次差分帯域非周期成分²BAP、並びに静特性の有声/無声判定情報VUVからなる。

【0105】

具体的には、時間長を除く音響特徴量は、静特性の60次元のメルケプストラム係数、1次元の対数ピッチ周波数及び5次元の帯域非周期成分を併せた静特性の66次元のデータと、これらの静特性のデータを1次差分及び2次差分して得られた動特性の132次元のデータと、1次元の有声/無声判定データとからなる。つまり、時間長を除く音響特徴量の次元数は199である。

【0106】

学習部16は、音響特徴量の199次元の全てのデータを用いて、次元毎に、平均値及び標準偏差を求めて記憶部17に格納すると共に、全てのデータのそれぞれについて、次元毎の平均値及び標準偏差を用いて標準化する(ステップS904)。

10

【0107】

学習部16は、ステップS902、S904から移行して、フレーム毎に、言語特徴量の332次元の標準化されたデータを入力データとし、音響特徴量の199次元の標準化されたデータを出力データとして音響モデルを学習する(ステップS905)。そして、学習部16は、学習済みの音響モデルを記憶部17に格納する。

【0108】

音響モデルは、例えば入力層を332次元、隠れ層を1024次元の6層、出力層を199次元とした順伝播型のニューラルネットワークで構成される。隠れ層における活性化関数は双曲線正接関数が用いられ、損失誤差関数は平均二乗誤差関数が用いられる。また、ミニバッチ数を256、エポック数を100、dropout(ドロップアウト)率を0.5学習係数の最適化方法として確率的勾配降下法、開始学習率を0.001、10エポックを過ぎてからエポック毎に学習率を指数減衰させ、誤差逆伝播法にて学習するものとする。尚、15エポックを過ぎてから、5エポック連続して評価誤差が減少しない場合は学習を早期終了するものとする。

20

【0109】

これにより、記憶部17には、統計モデルとして音響モデルが格納される。また、記憶部17には、統計モデルとして、音響モデルの入力データである言語特徴量の312次元のバイナリ値、13次元の数値データ、4次元の時間データ及び3次元の調整データ(パワーデータ、ピッチデータ及び抑揚データ)からなる332次元のデータに関する次元毎の最大値及び最小値が格納される。また、記憶部17には、統計モデルとして、音響モデルの出力データである音響特徴量の199次元のデータに関する次元毎の平均値及び標準偏差が格納される。

30

【0110】

以上のように、本発明の実施形態の学習装置1によれば、言語分析部11は、音声コーパスのテキストについて既知の言語分析処理を行い、音素毎の言語特徴量を求める。音声分析部12は、音声コーパスのテキストに対応する音声信号をフレーム毎に切り出し、フレーム毎の音声信号について既知の音響分析処理を行い、フレーム毎の音響特徴量を求める。

40

【0111】

対応付け部13は、既知の音素アラインメントの技術を用いて、音素毎の言語特徴量とフレーム毎の音響特徴量とを時間的に対応付け、音素毎の時間長を求める。そして、対応付け部13は、時間情報を追加した音素毎の言語特徴量を生成し、対応付けたフレーム毎の音響特徴量(時間長については音素毎のデータ)を生成する。

【0112】

調整量追加部14は、時間情報を追加した音素毎の言語特徴量に、調整パラメータの調整量情報を追加する。音響特徴量調整部15は、調整パラメータに従って、フレーム毎の音響特徴量(時間長については音素毎のデータ)を調整する。

【0113】

50

学習部 16 は、言語特徴量の 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 1 次元の調整データ（話速データ）からなる 3 2 6 次元のデータに基づいて、次元毎に、最大値及び最小値を求め、全てのデータのそれぞれを標準化する。また、学習部 16 は、時間長の 1 次元のデータに基づいて平均値及び標準偏差を求め、時間長の 1 次元のデータを標準化する。

【 0 1 1 4 】

学習部 16 は、音素毎に、言語特徴量の 3 2 6 次元の標準化されたデータを入力データとし、時間長の 1 次元の標準化されたデータを出力データとして時間長モデルを学習する。

【 0 1 1 5 】

学習部 16 は、言語特徴量の 3 1 2 次元のバイナリ値、1 3 次元の数値データ、4 次元の時間データ及び 3 次元の調整データ（パワーデータ、ピッチデータ及び抑揚データ）からなる 3 3 2 次元のデータに基づいて、次元毎に、最大値及び最小値を求め、全てのデータのそれぞれを標準化する。また、学習部 16 は、音響特徴量の 1 9 9 次元のデータに基づいて、次元毎に、平均値及び標準偏差を求め、全てのデータのそれぞれを標準化する。

10

【 0 1 1 6 】

学習部 16 は、フレーム毎に、言語特徴量の 3 3 2 次元の標準化されたデータを入力データとし、音響特徴量の 1 9 9 次元の標準化されたデータを出力データとして音響モデルを学習する。

【 0 1 1 7 】

これにより、記憶部 17 には、学習済みの統計モデルとして、調整パラメータの調整量情報が反映された時間長モデル、音響モデル及び最大値等が格納される。

20

【 0 1 1 8 】

そして、後述の音声合成装置 2 により、調整パラメータの調整量情報が反映された学習モデルを用いて、調整パラメータの調整量情報が追加された言語特徴量に基づき音響特徴量が推定され、フレーム毎の音響特徴量から合成音声信号が生成される。

【 0 1 1 9 】

図 17 に示した非特許文献 1, 2 の従来技術を組み合わせた想定例では、学習モデルを用いた推定により時間的に平滑化された特性を有する音響特徴量に調整を加え、調整後のフレーム毎の音響特徴量から合成音声信号を生成することから、合成音声信号に音質劣化が生じてしまう。さらに、入力文章の特定部分に対応する音響特徴量に調整を加え、調整後のフレーム毎の音響特徴量から合成音声信号を生成することから、調整を加えた部分と、これに隣接する調整を加えていない部分との間の接続部分において、合成音声信号に不連続を生じてしまう。

30

【 0 1 2 0 】

これに対し、本発明の実施形態による音声合成装置 2 は、調整パラメータの調整量情報が反映された学習モデルを用いて音響特徴量を推定し、合成音声信号を生成することから、学習モデルを用いた推定により時間的に平滑化された特性を有する音響特徴量に調整を加える必要がない。また、入力文章の特定部分に対応する言語特徴量を調整したものを学習モデルに入力して音響特徴量を求め、合成音声信号を生成することから、調整を加えた部分と、これに隣接する調整を加えていない部分との間の接続部分において、合成音声信号に不連続を生じることがない

40

【 0 1 2 1 】

したがって、テキストの特定部分の読み上げ方を調整した合成音声信号を生成する際に、高品質の合成音声信号を得ることができる。

【 0 1 2 2 】

また、本発明の実施形態では、調整パラメータは、話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} のうちのどれか、またはこれらの組み合わせであり、ユーザにより選択される。この場合、ユーザは、例えば 4 つの調整パラメータのうちの 1 つの調整パラメータについて、1 1 個のデータから 1 個のデータを選択し、他の 3 つの調整パラメータについては、標準値 1.0 を固定値として用いる。または、ユーザは、例えば 4 つの調整パラメータ

50

を 1 1 段階で連動させて選択する。

【 0 1 2 3 】

このように、調整パラメータの選択範囲を限定することにより、統計モデルを学習する際の学習データを少なくすることができ、低負荷かつ短時間で、統計モデルを学習することができる。

【 0 1 2 4 】

〔音声合成装置〕

次に、本発明の実施形態による音声合成装置について説明する。図 1 0 は、音声合成装置の構成例を示すブロック図であり、図 1 1 は、音声合成装置による音声合成処理例を示すフローチャートである。

【 0 1 2 5 】

この音声合成装置 2 は、言語分析部 2 0、調整量追加部 2 1、音響特徴量推定部 2 2、記憶部 1 7 及び音声生成部 2 3 を備えている。記憶部 1 7 は、図 1 に示した記憶部 1 7 に相当し、学習装置 1 により学習された統計モデルとして、時間長モデル、音響モデル及び最大値等が格納されている。

【 0 1 2 6 】

尚、学習装置 1 により学習された統計モデルは、学習装置 1 に備えた記憶部 1 7 から読み出され、音声合成装置 2 に備えた記憶部 1 7 に格納されるようにしてもよい。また、音声合成装置 2 は、インターネットを介して、学習装置 1 に備えた記憶部 1 7 へ直接アクセスするようにしてもよい。

【 0 1 2 7 】

言語分析部 2 0 は、音声合成対象のテキストを入力し、図 1 に示した言語分析部 1 1 と同様に、テキストについて既知の言語分析処理を行い、音素毎の所定情報からなる言語特徴量を求める（ステップ S 1 1 0 1）。そして、言語分析部 2 0 は、音素毎の言語特徴量を調整量追加部 2 1 に出力する。

【 0 1 2 8 】

調整量追加部 2 1 は、言語分析部 2 0 から音素毎の言語特徴量を入力すると共に、所定の調整パラメータを入力する。そして、調整量追加部 2 1 は、図 1 に示した調整量追加部 1 4 と同様に、言語特徴量を構成する音素毎の所定情報に、調整パラメータの調整量情報を追加する（ステップ S 1 1 0 2）。調整量追加部 2 1 は、音素毎の調整量情報を追加した言語特徴量を音響特徴量推定部 2 2 に出力する。

【 0 1 2 9 】

所定の調整パラメータは、前述と同様に、話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} のうちのどれか、またはこれらの組み合わせとし、ユーザにより指定されるものとする。調整パラメータの値は、前述した調整の範囲において任意の実数とする。つまり、所定の調整パラメータは、話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} のうちのいずれか 1 つまたは 2 つ以上の組み合わせとする。

【 0 1 3 0 】

尚、所定の調整パラメータは、話速 R_{ST} 、パワー R_{PW} 、ピッチ R_{PT} 及び抑揚 R_{PD} とし、これらの 4 つのパラメータのうちのいずれか 1 つのパラメータの調整量は、所定範囲内の任意の値が指定され、他の 3 つのパラメータの調整量は、固定値が用いられるようにしてもよい。また、所定の調整パラメータは、前述の 4 つのパラメータとし、それぞれの調整量は、それぞれの所定範囲内の任意の値が指定されるようにしてもよい。

【 0 1 3 1 】

尚、調整量追加部 2 1 は、図 1 に示した調整量追加部 1 4 と同様に、文章単位、呼気段落単位またはアクセント句単位で、異なる調整パラメータを入力するようにしてもよい。

【 0 1 3 2 】

音響特徴量推定部 2 2 は、調整量追加部 2 1 から音素毎の言語特徴量を入力し、記憶部 1 7 に格納された最大値等を用いて標準化及び逆標準化の処理を行い、時間長モデルを用いて音素毎の時間長を推定する。

10

20

30

40

50

【 0 1 3 3 】

音響特徴量推定部 2 2 は、記憶部 1 7 に格納された最大値等を用いて標準化及び逆標準化の処理を行い、音響モデルを用いてフレーム毎の音響特徴量を推定する（ステップ S 1 1 0 3）。音響特徴量推定部 2 2 は、フレーム毎の音響特徴量を音声生成部 2 3 に出力する。

【 0 1 3 4 】

（時間長モデルを用いた時間長の推定）

次に、音響特徴量推定部 2 2 による時間長モデルを用いた時間長の推定処理について説明する。図 1 2 は、時間長モデルを用いた時間長推定処理例を説明する図である。音響特徴量推定部 2 2 は、調整量追加部 2 1 から入力した音素毎の言語特徴量に基づいて、言語特徴を表す 3 1 2 次元のバイナリ値及び 1 3 次元の数値データ、並びに 1 次元の調整データ（話速データ）を生成する。言語特徴量の次元数は 3 2 6 である。

10

【 0 1 3 5 】

音響特徴量推定部 2 2 は、言語特徴量の 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 1 次元の調整データ（話速データ）からなる 3 2 6 次元のデータを、時間長モデルの入力データとして扱う（ステップ S 1 2 0 1）。

【 0 1 3 6 】

音響特徴量推定部 2 2 は、記憶部 1 7 から、時間長モデルの入力データである言語特徴量の 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 1 次元の調整データ（話速データ）からなる 3 2 6 次元のデータに関する次元毎の最大値及び最小値を読み出す。そして、音響特徴量推定部 2 2 は、言語特徴量の 3 2 6 次元のデータのそれぞれについて、次元毎に、最大値及び最小値を用いて標準化を行う（ステップ S 1 2 0 2）。

20

【 0 1 3 7 】

音響特徴量推定部 2 2 は、記憶部 1 7 に格納された時間長モデルを用いて、言語特徴量の 3 2 6 次元の標準化されたデータを時間長モデルの入力データとして、時間長モデルの出力データである時間長の 1 次元の標準化されたデータを推定する（ステップ S 1 2 0 3）。

【 0 1 3 8 】

音響特徴量推定部 2 2 は、記憶部 1 7 から、時間長モデルの出力データである時間長の 1 次元のデータに関する平均値及び標準偏差を読み出す。そして、音響特徴量推定部 2 2 は、ステップ S 1 2 0 3 にて推定した時間長の 1 次元の標準化されたデータについて、平均値及び標準偏差を用いて逆標準化を行い（ステップ S 1 2 0 4）、時間長の 1 次元のデータを求める（ステップ S 1 2 0 5）。

30

【 0 1 3 9 】

これにより、記憶部 1 7 に格納された時間長モデル、時間長モデルの入力データである言語特徴量の 3 2 6 次元のデータに関する次元毎の最大値及び最小値、並びに、時間長モデルの出力データである時間長の 1 次元のデータに関する平均値及び標準偏差を用いて、音素毎の言語特徴量の 3 2 6 次元のデータから、音素毎の時間長の 1 次元のデータを求めることができる。

【 0 1 4 0 】

（音響モデルを用いた音響特徴量の推定）

次に、音響特徴量推定部 2 2 による音響モデルを用いた音響特徴量の推定処理について説明する。図 1 3 は、音響モデルを用いた音響特徴量推定処理例を説明する図である。音響特徴量推定部 2 2 は、ステップ S 1 2 0 5 にて求めた音素毎の時間長の 1 次元のデータに基づいて、図 9 のステップ S 9 0 1 と同様に、音素に対応する複数フレームのそれぞれについて、時間データの 4 次元のデータを生成する（ステップ S 1 3 0 1）。

40

【 0 1 4 1 】

音響特徴量推定部 2 2 は、調整量追加部 2 1 から入力した音素毎の言語特徴量に基づいて、言語特徴を表す 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 3 次元の調整データ（パワーデータ、ピッチデータ及び抑揚データ）を生成する。そして、音響特徴量推

50

定部 2 2 は、音素毎の言語特徴量における 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 3 次元の調整データ（パワーデータ、ピッチデータ及び抑揚データ）からなる 3 2 8 次元のデータから、フレーム毎の言語特徴量における 3 2 8 次元のデータを生成する。

【 0 1 4 2 】

音響特徴量推定部 2 2 は、フレーム毎の言語特徴量の 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 3 次元の調整データ（パワーデータ、ピッチデータ及び抑揚データ）からなる 3 2 8 次元のデータ、並びにステップ S 1 3 0 1 にて生成した時間データの 4 次元のデータを、音響モデルの入力データとして扱う（ステップ S 1 3 0 2）。

【 0 1 4 3 】

音響特徴量推定部 2 2 は、記憶部 1 7 から、音響モデルの入力データである言語特徴量の 3 1 2 次元のバイナリ値、1 3 次元の数値データ、4 次元の時間データ及び 3 次元の調整データ（パワーデータ、ピッチデータ及び抑揚データ）からなる 3 3 2 次元のデータに関する次元毎の最大値及び最小値を読み出す。そして、音響特徴量推定部 2 2 は、言語特徴量の 3 2 8 次元のデータ及び時間データの 4 次元のデータからなる 3 3 2 次元のデータのそれぞれについて、次元毎に、最大値及び最小値を用いて標準化を行う（ステップ S 1 3 0 3）。

10

【 0 1 4 4 】

音響特徴量推定部 2 2 は、記憶部 1 7 に格納された音響モデルを用いて、言語特徴量の 3 2 8 次元の標準化されたデータ及び時間データの 4 次元の標準化されたデータからなる 3 3 2 次元の標準化されたデータを音響モデルの入力データとして、音響モデルの出力データである音響特徴量の 1 9 9 次元の標準化されたデータを推定する（ステップ S 1 3 0 4）。

20

【 0 1 4 5 】

音響特徴量推定部 2 2 は、記憶部 1 7 から、音響モデルの出力データである音響特徴量の 1 9 9 次元のデータに関する平均値及び標準偏差を読み出す。そして、音響特徴量推定部 2 2 は、ステップ S 1 3 0 4 にて推定した音響特徴量の 1 9 9 次元の標準化されたデータについて、次元毎に、平均値及び標準偏差を用いて逆標準化を行う（ステップ S 1 3 0 5）。音響特徴量推定部 2 2 は、フレーム毎の音響特徴量の 1 9 9 次元のデータを生成する（ステップ S 1 3 0 6）。

【 0 1 4 6 】

このようにして推定され逆標準化された音響特徴量は、フレーム毎に離散的な値をとる。そこで、音響特徴量推定部 2 2 は、連続するフレーム毎の音響特徴量の 1 9 9 次元のデータに対して、最尤推定または移動平均をとり、新たなフレーム毎の音響特徴量の 1 9 9 次元のデータを求める。これにより、フレーム毎の音響特徴量は滑らかな値となる。

30

【 0 1 4 7 】

これにより、記憶部 1 7 に格納された音響モデル、音響モデルの入力データである言語特徴量の 3 3 2 次元のデータに関する次元毎の最大値及び最小値、並びに、音響モデルの出力データである音響特徴量の 1 9 9 次元のデータに関する平均値及び標準偏差を用いて、フレーム毎の言語特徴量の 3 2 8 次元のデータ及び時間データの 4 次元のデータから、フレーム毎の音響特徴量の 1 9 9 次元のデータを得ることができる。

40

【 0 1 4 8 】

図 1 0 及び図 1 1 に戻って、音声生成部 2 3 は、音響特徴量推定部 2 2 からフレーム毎の音響特徴量を入力し、フレーム毎の音響特徴量に基づいて音声信号を合成する（ステップ S 1 1 0 4）。そして、音声生成部 2 3 は、音声合成対象のテキストに対して調整パラメータによる調整が加えられた音声信号を出力する。

【 0 1 4 9 】

図 1 4 は、音声生成部 2 3 による音声合成処理例を説明する図である。音声生成部 2 3 は、音響特徴量推定部 2 2 から入力したフレーム毎の音響特徴量のうち、フレーム毎のメルケプストラム係数 M G C、対数ピッチ周波数 L F 0 及び帯域非周期成分 B A P である静特性の音響特徴量を選択する（ステップ S 1 4 0 1）。

50

【 0 1 5 0 】

音声生成部 2 3 は、メルケプストラム係数 M G C をメルケプストラムスペクトル変換し、スペクトルを求める（ステップ S 1 4 0 2）。また、音声生成部 2 3 は、対数ピッチ周波数 L F 0 から有声 / 無声判定情報 V U V を求め、対数ピッチ周波数 L F 0 の有声区間を指数化し、無声及び無音区間についてはゼロとし、ピッチ周波数を求める（ステップ S 1 4 0 3）。また、音声生成部 2 3 は、帯域非周期成分 B A P をメルケプストラムスペクトル変換し、非周期成分を求める（ステップ S 1 4 0 4）。

【 0 1 5 1 】

音声生成部 2 3 は、ステップ S 1 4 0 2 にて求めたフレーム毎のスペクトル、ステップ S 1 4 0 3 にて求めたフレーム毎のピッチ周波数、及びステップ S 1 4 0 4 にて求めたフレーム毎の非周期成分を用いて連続的に音声波形を生成し（ステップ S 1 4 0 5）、音声信号を出力する（ステップ S 1 4 0 6）。

10

【 0 1 5 2 】

これにより、音声合成対象のテキストに対して所定の調整パラメータによる調整が加えられた音声信号を得ることができる。

【 0 1 5 3 】

以上のように、本発明の実施形態の音声合成装置 2 によれば、言語分析部 2 0 は、音声合成対象のテキストについて既知の言語分析処理を行い、音素毎の言語特徴量を求め、調整量追加部 2 1 は、音素毎の言語特徴量に、調整パラメータの調整量情報を追加する。

【 0 1 5 4 】

音響特徴量推定部 2 2 は、言語特徴量の 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 1 次元の調整データ（話速データ）からなる 3 2 6 次元のデータを、記憶部 1 7 に格納された最大値等を用いて標準化する。そして、音響特徴量推定部 2 2 は、記憶部 1 7 に格納された時間長モデルを用いて、これらの標準化されたデータを入力データとして、出力データである時間長の 1 次元の標準化されたデータを推定する。

20

【 0 1 5 5 】

音響特徴量推定部 2 2 は、時間長の 1 次元の標準化されたデータを、記憶部 1 7 に格納された平均値等を用いて逆標準化し、フレーム毎の時間データを求める。音響特徴量推定部 2 2 は、言語特徴量の 3 2 9 次元のデータのうち 3 1 2 次元のバイナリ値、1 3 次元の数値データ及び 3 次元の調整データ（パワーデータ、ピッチデータ及び抑揚データ）からなる 3 2 8 次元のデータ、並びに時間データの 4 次元のデータを、記憶部 1 7 に格納された最大値等を用いて標準化する。そして、音響特徴量推定部 2 2 は、記憶部 1 7 に格納された音響モデルを用いて、これらの標準化されたデータを入力データとして、出力データである音響特徴量の 1 9 9 次元の標準化されたデータを推定する。

30

【 0 1 5 6 】

音響特徴量推定部 2 2 は、音響特徴量の 1 9 9 次元の標準化されたデータを、記憶部 1 7 に格納された平均値等を用いて逆標準化し、フレーム毎の音響特徴量を求める。そして、音声生成部 2 3 は、フレーム毎の音響特徴量に基づいて音声信号を合成し、合成音声信号を生成する。

【 0 1 5 7 】

図 1 7 に示した非特許文献 1 , 2 の従来技術を組み合わせた想定例では、学習モデルを用いた推定により時間的に平滑化された特性を有する音響特徴量に調整を加え、調整後のフレーム毎の音響特徴量から合成音声信号を生成することから、合成音声信号に音質劣化を生じてしまう。さらに、入力文章の特定部分に対応する音響特徴量に調整を加え、調整後のフレーム毎の音響特徴量から合成音声信号を生成することから、調整を加えた部分と、これに隣接する調整を加えていない部分との間の接続部分において、合成音声信号に不連続を生じてしまう。

40

【 0 1 5 8 】

これに対し、本発明の実施形態による音声合成装置 2 は、調整パラメータの調整量情報が反映された学習モデルを用いて音響特徴量を推定し、合成音声信号を生成するから、学

50

習モデルを用いた推定により時間的に平滑化された特性を有する音響特徴量に調整を加える必要がない。また、入力文章の特定部分に対応する言語特徴量を調整したものを学習モデルに入力して音響特徴量を求め、合成音声信号を生成することから、調整を加えた部分と、これに隣接する調整を加えていない部分との間の接続部分において、合成音声信号に不連続を生じることがない。

【0159】

したがって、テキストの特定部分の読み上げ方を調整した合成音声信号を生成する際に、高品質の合成音声信号を得ることができる。

【0160】

以上、実施形態を挙げて本発明を説明したが、本発明は前記実施形態に限定されるものではなく、その技術思想を逸脱しない範囲で種々変形可能である。

10

【0161】

尚、本発明の実施形態による学習装置1及び音声合成装置2のハードウェア構成としては、通常のコンピュータを使用することができる。学習装置1及び音声合成装置2は、CPU、RAM等の揮発性の記憶媒体、ROM等の不揮発性の記憶媒体、及びインターフェース等を備えたコンピュータによって構成される。

【0162】

学習装置1に備えた記憶部10、17、言語分析部11、音声分析部12、対応付け部13、調整量追加部14、音響特徴量調整部15及び学習部16の各機能は、これらの機能を記述したプログラムをCPUに実行させることによりそれぞれ実現される。また、音声合成装置2に備えた言語分析部20、調整量追加部21、音響特徴量推定部22、記憶部17及び音声生成部23の各機能も、これらの機能を記述したプログラムをCPUに実行させることによりそれぞれ実現される。

20

【0163】

これらのプログラムは、前記記憶媒体に格納されており、CPUに読み出されて実行される。また、これらのプログラムは、磁気ディスク（フロッピー（登録商標）ディスク、ハードディスク等）、光ディスク（CD-ROM、DVD等）、半導体メモリ等の記憶媒体に格納して頒布することもでき、ネットワークを介して送受信することもできる。

【符号の説明】

【0164】

1 学習装置

2 音声合成装置

10, 17 記憶部

11, 20 言語分析部

12 音声分析部

13 対応付け部

14, 21 調整量追加部

15 音響特徴量調整部

16 学習部

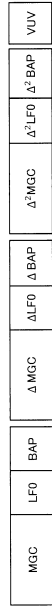
22 音響特徴量推定部

23 音声生成部

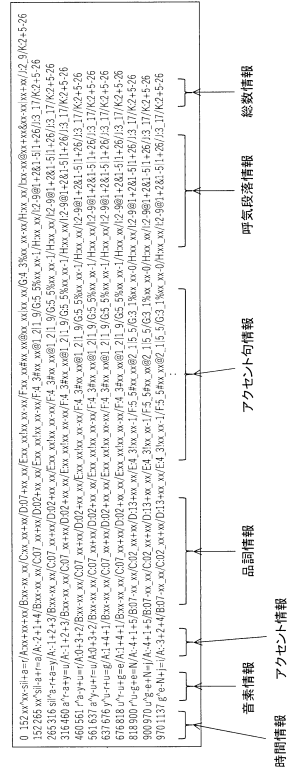
30

40

【図 5】



【図 6】

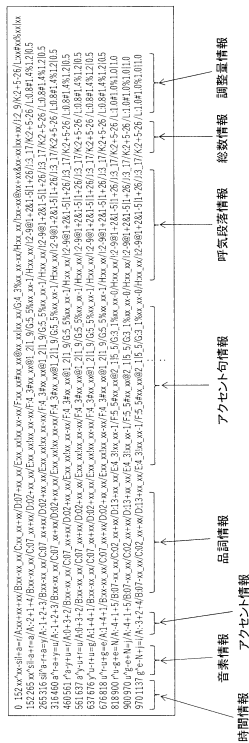


時間情報が追加された言語特徴量

10

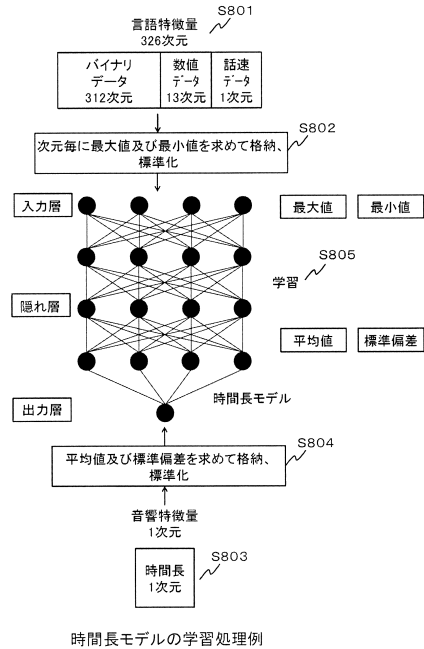
20

【図 7】



調整量情報が追加された言語特徴量

【図 8】



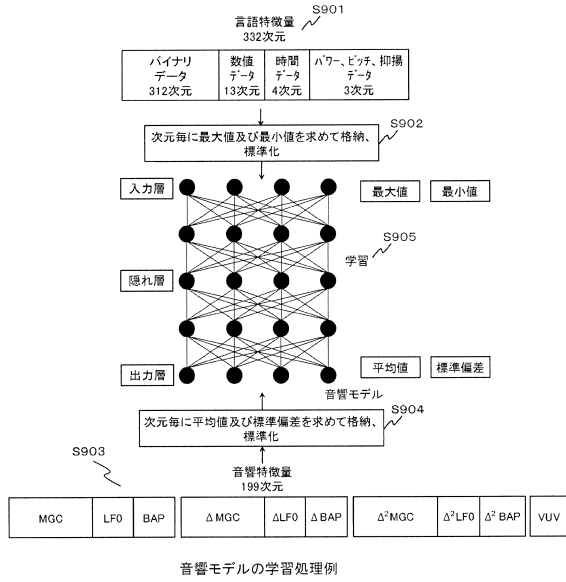
時間長モデルの学習処理例

30

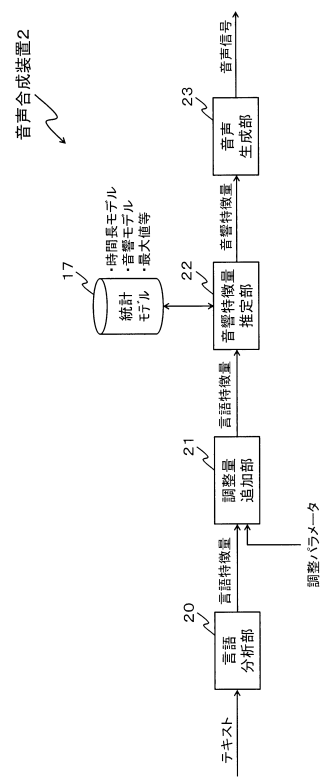
40

50

【図9】



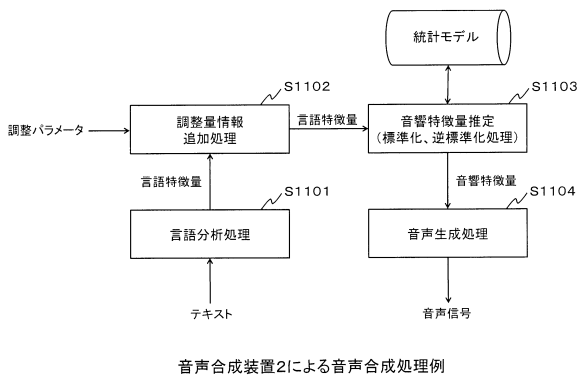
【図10】



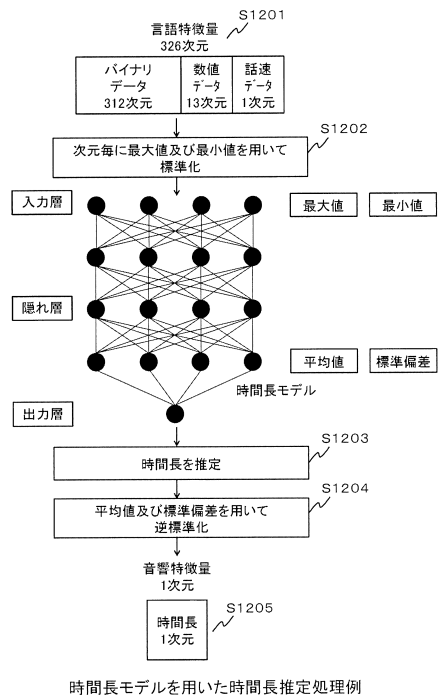
10

20

【図11】



【図12】

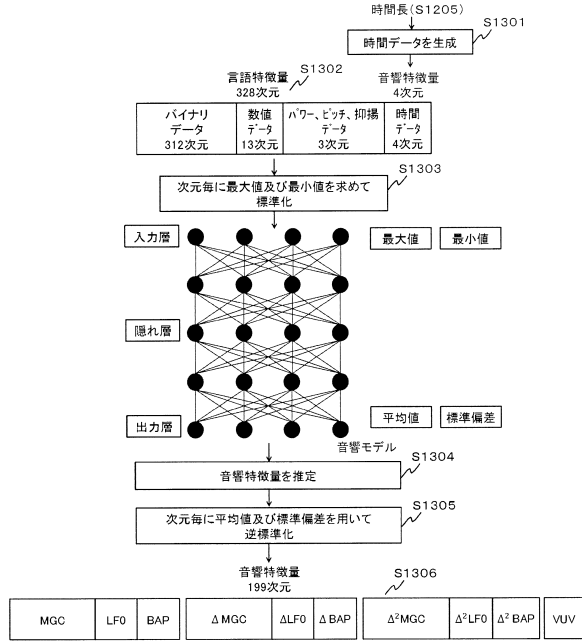


30

40

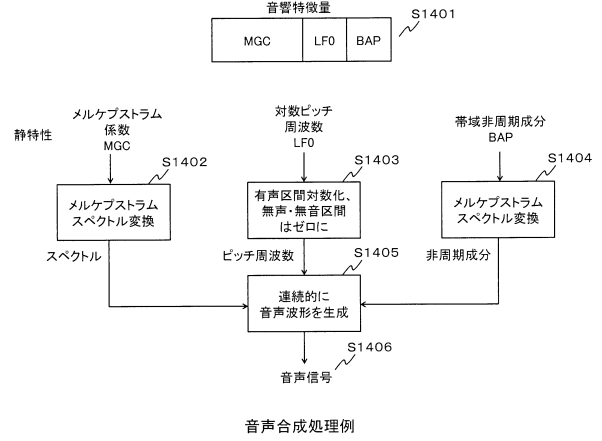
50

【図13】



音響モデルを用いた音響特徴量推定処理例

【図14】

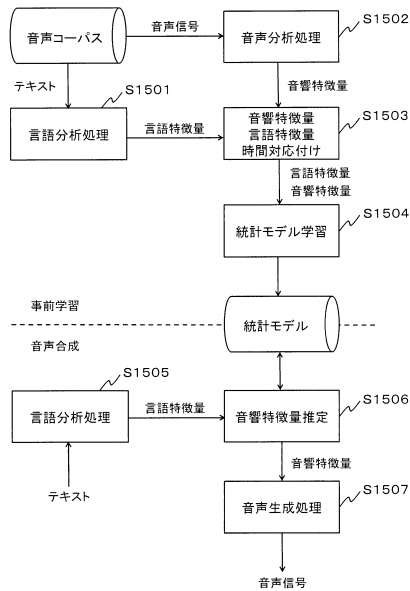


音声合成処理例

10

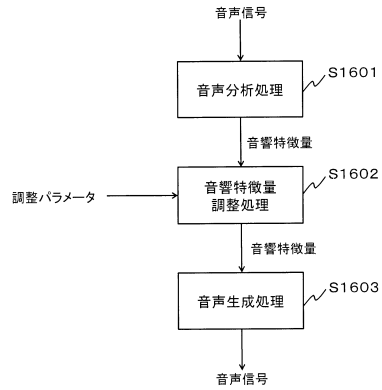
20

【図15】



従来技術(非特許文献1)

【図16】



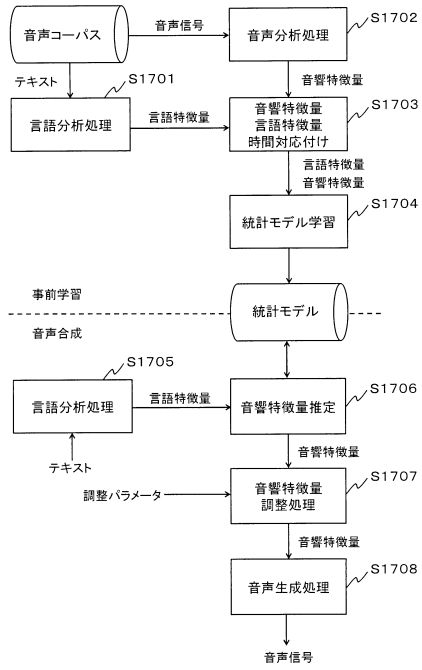
従来技術(非特許文献2)

30

40

50

【図 17】



従来技術(非特許文献1, 2)を組み合わせた想定例

10

20

30

40

50

フロントページの続き

(51)国際特許分類

F I
G 1 0 L 25/30

東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内

(72)発明者 今井 篤

東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内

(72)発明者 都木 徹

東京都世田谷区砧一丁目10番11号 一般財団法人NHKエンジニアリングシステム内

審査官 大野 弘

(56)参考文献

再公表特許第2009/107441(JP, A1)

特開2017-032839(JP, A)

山田 修平 Shuhei YAMADA, テーラーメイド音声合成のための差分特徴量を用いたDNN
に基づくF0制御, 日本音響学会 2017年 春季研究発表会講演論文集CD-ROM [CD-ROM], 日本, 2023年02月17日, PP271-274

(58)調査した分野 (Int.Cl., DB名)

G 1 0 L 13/06

G 1 0 L 13/10

G 1 0 L 25/30