

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6846237号
(P6846237)

(45) 発行日 令和3年3月24日(2021.3.24)

(24) 登録日 令和3年3月3日(2021.3.3)

(51) Int.Cl.	F I
G 1 O L 13/10 (2013.01)	G 1 O L 13/10 1 1 4
G 1 O L 13/033 (2013.01)	G 1 O L 13/10 1 1 3 B
	G 1 O L 13/10 1 1 2 C
	G 1 O L 13/10 1 1 1 F
	G 1 O L 13/10 1 1 3 C
請求項の数 4 (全 17 頁) 最終頁に続く	

(21) 出願番号	特願2017-42169 (P2017-42169)	(73) 特許権者	000004352
(22) 出願日	平成29年3月6日(2017.3.6)		日本放送協会
(65) 公開番号	特開2018-146803 (P2018-146803A)		東京都渋谷区神南2丁目2番1号
(43) 公開日	平成30年9月20日(2018.9.20)	(73) 特許権者	591053926
審査請求日	令和2年2月3日(2020.2.3)		一般財団法人NHKエンジニアリングシステム
			東京都世田谷区砧一丁目10番11号
		(74) 代理人	100121119
			弁理士 花村 泰伸
		(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	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
			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】

事前に学習されたDNN(ディープニューラルネットワーク)を用いて、音声波形を合成する音声合成装置において、

音素の言語特徴量、話者を識別するための話者ラベル、及び感情を識別するための感情ラベルが入力層に与えられ、音素の時間長が出力層に与えられることで学習された時間長DNNと、

音素フレームの言語特徴量、前記話者ラベル及び前記感情ラベルが入力層に与えられ、音素フレームの音響特徴量が出力層に与えられることで学習された音響特徴量DNNと、

テキスト、話者情報及び感情情報を入力し、前記時間長DNN及び前記音響特徴量DNNを用いて、前記テキスト、前記話者情報及び前記感情情報に対応する音声波形を合成する合成処理部と、を備え、

前記合成処理部は、

前記テキストをテキスト解析して音素の言語特徴量を生成し、

前記時間長DNNを用いて、前記音素の言語特徴量、前記話者情報に付与した話者ラベル、及び前記感情情報に付与した感情ラベルに基づいて、音素の時間長を生成し、

前記音素の言語特徴量及び前記音素の時間長に基づいて、音素フレームの言語特徴量を生成し、

前記音響特徴量DNNを用いて、前記音素フレームの言語特徴量、前記話者ラベル及び前記感情ラベルに基づいて、音素フレームの音響特徴量を生成し、

10

20

当該音素フレームの音響特徴量に基づいて、前記音声波形を合成する、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、

前記合成処理部は、

前記話者情報に前記話者ラベルを付与すると共に、前記感情情報に前記感情ラベルを付与する話者感情ラベル処理部と、

前記テキストをテキスト解析して前記音素の言語特徴量を生成し、前記音素の言語特徴量及び前記音素の時間長に基づいて、前記音素フレームの言語特徴量を生成するテキスト解析部と、

10

前記時間長 DNN を用いて、前記テキスト解析部により生成された前記音素の言語特徴量、前記話者感情ラベル処理部により付与された前記話者ラベル及び前記感情ラベルに基づいて、前記音素の時間長を生成し、

前記音響特徴量 DNN を用いて、前記テキスト解析部により生成された前記音素フレームの言語特徴量、前記話者ラベル及び前記感情ラベルに基づいて、前記音素フレームの音響特徴量を生成する時間長及び音響特徴量生成部と、

前記時間長及び音響特徴量生成部により生成された前記音素フレームの音響特徴量に基づいて、前記音声波形を合成する音声波形合成部と、
を備えたことを特徴とする音声合成装置。

【請求項 3】

20

請求項 1 または 2 に記載の音声合成装置において、

さらに、テキスト、話者情報、感情情報及び音声波形が格納された音声コーパスを用いて、前記時間長 DNN 及び前記音響特徴量 DNN を学習する学習部を備え、

前記学習部は、

前記音声コーパスから前記テキストを読み出し、当該テキストをテキスト解析して音素の言語特徴量を生成し、

前記音声コーパスから前記音声波形を読み出し、当該音声波形を音響分析して音素の区切り位置を求めると共に、音素の時間長を求め、

前記音声コーパスから前記話者情報及び前記感情情報を読み出し、話者ラベル及び感情ラベルをそれぞれ付与し、

30

前記音素の言語特徴量及び前記音素の時間長に基づいて、音素フレームの言語特徴量を生成し、

前記音素の言語特徴量、前記話者ラベル及び前記感情ラベル、並びに前記音素の時間長を用いて、前記時間長 DNN を学習し、

前記音素フレームの言語特徴量、前記話者ラベル及び前記感情ラベル、並びに前記音素フレームの音響特徴量を用いて、前記音響特徴量 DNN を学習する、ことを特徴とする音声合成装置。

【請求項 4】

コンピュータを、請求項 1 から 3 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。

40

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、DNN (Deep Neural Network: ディープニューラルネットワーク) を用いた音声合成装置及びプログラムに関する。

【背景技術】

【0002】

従来、統計モデルを用いた音声合成技術が進展し、スマートフォンまたはパソコンを通して、身近なサービスとして使用できるようになっている。放送分野では、音声合成を用いて解説放送を補完する音声ガイドシステムの研究が進められている(例えば、非特許文

50

献 1 を参照)。音声ガイドシステムでは、ガイド音声を表示する際に、演出効果として、多様な話者性及び感情表現が求められている。

【 0 0 0 3 】

一方、統計モデルを用いた音声合成技術の主流は、HMM (Hidden Markov Model : 隠れマルコフモデル) 方式とDNN方式とに分類される。HMM方式では、話者性及び感情表現を制御可能な音声合成を実現している (例えば、非特許文献 2 を参照)。

【 0 0 0 4 】

DNN方式は、一般にHMM方式よりも音質が良いと言われており、様々な手法で音声合成を実現しているが (例えば、非特許文献 3 ~ 5 を参照)、話者性及び感情表現の両方を制御可能な音声合成の実現に至っていない。

【先行技術文献】

【非特許文献】

【 0 0 0 5 】

【非特許文献 1】今井他, 電子情報通信学会総合大会講演論文集, H-4-11, Mar 2016

【非特許文献 2】J.Yamagishi et al, vol. E88-D, no. 3, pp. 503-509, Mar 2005

【非特許文献 3】Zhizheng Wu et al, ISCA SSW9, vol. PS2-13, pp. 218-223, Sep 2016

【非特許文献 4】H.Zen et al, IEICE Trans. Inf. & Syst., vol. E90-D, no. 5, pp. 825-834, May 2007

【非特許文献 5】北条他, 日本音響学会講演論文集, pp. 215-218, Sep 2015

【発明の概要】

【発明が解決しようとする課題】

【 0 0 0 6 】

DNN方式を用いた音声合成技術において、話者性及び感情表現を制御可能な音声合成を実現するためには、話者及び感情の組み合わせ毎に、DNNを用意する手法が想定される。例えば、怒りの感情を有する話者 a の DNN、喜びの感情を有する話者 a の DNN、・・・、怒りの感情を有する話者 b の DNN、喜びの感情の有する話者 b の DNN 等を用意する必要がある。

【 0 0 0 7 】

しかしながら、この手法では、話者と感情との組み合わせが膨大であり、用意すべき DNN の数が多くなり、実現が困難である。また、音声合成の際に、複数の異なる DNN を用いる場合には、連続的な自然な読み上げ音声を生成することが困難となる。

【 0 0 0 8 】

このように、DNN方式を用いた音声合成技術では、話者性及び感情表現を制御可能な音声合成を実現する際に、膨大な数の DNN を用意する必要のない新たな手法が所望されていた。

【 0 0 0 9 】

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、簡易な構成にて、話者性及び感情表現を同時に制御可能な音声合成を実現する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 0 】

前記課題を解決するために、請求項 1 の音声合成装置は、事前に学習された DNN (ディープニューラルネットワーク) を用いて、音声波形を合成する音声合成装置において、音素の言語特徴量、話者を識別するための話者ラベル、及び感情を識別するための感情ラベルが入力層に与えられ、音素の時間長が出力層に与えられることで学習された時間長 DNN と、音素フレームの言語特徴量、前記話者ラベル及び前記感情ラベルが入力層に与えられ、音素フレームの音響特徴量が出力層に与えられることで学習された音響特徴量 DNN と、テキスト、話者情報及び感情情報を入力し、前記時間長 DNN 及び前記音響特徴量 DNN を用いて、前記テキスト、前記話者情報及び前記感情情報に対応する音声波形を合成する合成処理部と、を備え、前記合成処理部が、前記テキストをテキスト解析して音素

10

20

30

40

50

の言語特徴量を生成し、前記時間長DNNを用いて、前記音素の言語特徴量、前記話者情報に付与した話者ラベル、及び前記感情情報に付与した感情ラベルに基づいて、音素の時間長を生成し、前記音素の言語特徴量及び前記音素の時間長に基づいて、音素フレームの言語特徴量を生成し、前記音響特徴量DNNを用いて、前記音素フレームの言語特徴量、前記話者ラベル及び前記感情ラベルに基づいて、音素フレームの音響特徴量を生成し、当該音素フレームの音響特徴量に基づいて、前記音声波形を合成する、ことを特徴とする。

【0011】

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、前記合成処理部が、前記話者情報に前記話者ラベルを付与すると共に、前記感情情報に前記感情ラベルを付与する話者感情ラベル処理部と、前記テキストをテキスト解析して前記音素の言語特徴量を生成し、前記音素の言語特徴量及び前記音素の時間長に基づいて、前記音素フレームの言語特徴量を生成するテキスト解析部と、前記時間長DNNを用いて、前記テキスト解析部により生成された前記音素の言語特徴量、前記話者感情ラベル処理部により付与された前記話者ラベル及び前記感情ラベルに基づいて、前記音素の時間長を生成し、前記音響特徴量DNNを用いて、前記テキスト解析部により生成された前記音素フレームの言語特徴量、前記話者ラベル及び前記感情ラベルに基づいて、前記音素フレームの音響特徴量を生成する時間長及び音響特徴量生成部と、前記時間長及び音響特徴量生成部により生成された前記音素フレームの音響特徴量に基づいて、前記音声波形を合成する音声波形合成部と、を備えたことを特徴とする。

【0012】

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、さらに、テキスト、話者情報、感情情報及び音声波形が格納された音声コーパスを用いて、前記時間長DNN及び前記音響特徴量DNNを学習する学習部を備え、前記学習部が、前記音声コーパスから前記テキストを読み出し、当該テキストをテキスト解析して音素の言語特徴量を生成し、前記音声コーパスから前記音声波形を読み出し、当該音声波形を音響分析して音素の区切り位置を求めると共に、音素の時間長を求め、前記音声コーパスから前記話者情報及び前記感情情報を読み出し、話者ラベル及び感情ラベルをそれぞれ付与し、前記音素の言語特徴量及び前記音素の時間長に基づいて、音素フレームの言語特徴量を生成し、前記音素の言語特徴量、前記話者ラベル及び前記感情ラベル、並びに前記音素の時間長を用いて、前記時間長DNNを学習し、前記音素フレームの言語特徴量、前記話者ラベル及び前記感情ラベル、並びに前記音素フレームの音響特徴量を用いて、前記音響特徴量DNNを学習する、ことを特徴とする。

【0013】

さらに、請求項4のプログラムは、コンピュータを、請求項1から3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。

【発明の効果】

【0014】

以上のように、本発明によれば、話者及び感情の組み合わせ毎の膨大な数のDNNを用意する必要がないから、簡易な構成にて、話者性及び感情表現を同時に制御可能な音声合成を実現することができる。

【図面の簡単な説明】

【0015】

【図1】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2】事前学習部の構成例を示すブロック図である。

【図3】テキスト解析部の構成例を示すブロック図である。

【図4】音響分析部の構成例を示すブロック図である。

【図5】合成処理部の構成例を示すブロック図である。

【図6】時間長DNNの構成の概要を説明する図である。

【図7】音響特徴量DNNの構成の概要を説明する図である。

【図8】言語特徴量及び音響特徴量の関係について説明する図である。

10

20

30

40

50

【図 9】音素の言語特徴量、話者ラベル及び感情ラベルの例を説明する図である。

【図 10】事前学習部の処理例を示すフローチャートである。

【図 11】合成処理部の処理例を示すフローチャートである。

【発明を実施するための形態】

【0016】

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。本発明は、話者性及び感情表現と音声波形とを関連付けて、時間長 DNN 及び音響特徴量 DNN を事前に学習する。また、本発明は、事前に学習した時間長 DNN 及び音響特徴量 DNN を用いて、話者性及び感情表現を反映した音声合成を実現する。

【0017】

時間長 DNN は、音素の言語特徴量、話者ラベル及び感情ラベルを入力層の各ユニットに与え、音素の時間長を出力層のユニットに与えることで、音素毎に学習されたモデルである。音響特徴量 DNN は、音素フレームの言語特徴量、話者ラベル及び感情ラベルを入力層の各ユニットに与え、音素フレームの音響特徴量を出力層の各ユニットに与えることで、音素フレーム毎に学習されたモデルである。

【0018】

これにより、時間長 DNN 及び音響特徴量 DNN の 2 つの DNN を用意すればよいから、話者及び感情の組み合わせ毎の膨大な数の DNN を用意する必要がなく、簡易な構成にて、話者性及び感情表現を同時に制御可能な音声合成を実現することができる。

【0019】

〔音声合成装置〕

まず、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 は、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。この音声合成装置 1 は、音声コーパスが格納された記憶部 2、事前学習部 3、時間長 DNN 及び音響特徴量 DNN が格納された記憶部 4、及び合成処理部 5 を備えている。

【0020】

記憶部 2 には、特定の文章が複数の話者と感情で読み上げられた音声に関する情報、すなわち、テキスト、話者情報、感情情報及び音声波形の各情報により構成された音声コーパスが格納されている。音声コーパスは、話者及び感情の組み合わせを単位としたデータベースである。

【0021】

話者情報は、個々の発話者を識別するための情報であり、感情情報は、例えば喜び、怒り、悲哀、平静等の発話表現を識別するための情報であり、音声波形は、テキストに対する音声波形情報である。

【0022】

事前学習部 3 は、記憶部 2 から、所定の音声コーパスのテキスト、話者情報、感情情報及び音声波形を読み出し、話者情報及び感情情報に話者ラベル及び感情ラベルをそれぞれ付与する。話者ラベルは、話者を識別するためのラベルであり、感情ラベルは、感情を識別するためのラベルである。

【0023】

事前学習部 3 は、テキストに対し、所定のテキスト解析を行うと共に、音声波形に対し、所定の音響解析を行うことで、時間長 DNN 及び音響特徴量 DNN を学習するための言語特徴量及び音響特徴量等の情報を生成する。事前学習部 3 は、言語特徴量及び音響特徴量等の情報、並びに話者ラベル及び感情ラベルを用いて、記憶部 4 に格納された時間長 DNN 及び音響特徴量 DNN を事前に学習する。

【0024】

テキスト解析の手法及び音響解析の手法は既知であるから、ここでは詳細な説明は省略する。時間長 DNN 及び音響特徴量 DNN の学習は、例えば LSTM (Long Short Term Memory: 長期短期記憶) 方式にて行われる。

【0025】

10

20

30

40

50

記憶部 4 には、事前学習部 3 により学習された時間長 DNN 及び音響特徴量 DNN が格納される。

【 0 0 2 6 】

図 6 は、時間長 DNN の構成の概要を説明する図である。時間長 DNN は、学習時に、音素の言語特徴量、話者（話者ラベル）及び感情（感情ラベル）が入力層の各ユニットに与えられ、音素の時間長が出力層のユニットに与えられることで、入力層、隠れ層及び出力層の各ユニットの重み等が計算され、音素単位の学習が行われる。

【 0 0 2 7 】

学習のための音素の言語特徴量は、例えば、音素ラベル、アクセントの位置、品詞情報、アクセント句の情報、呼気段落の情報等からなる。音素の時間長は、例えば音素を構成する音素フレームの数で表される。

10

【 0 0 2 8 】

また、後述する音声合成時には、時間長 DNN の入力層の各ユニットに、音素の言語特徴量、話者ラベル及び感情ラベルが与えられることで、出力層のユニットから、当該音素の言語特徴量、話者ラベル及び感情ラベルに対応する音素の時間長が出力される。

【 0 0 2 9 】

図 7 は、音響特徴量 DNN の構成の概要を説明する図である。音響特徴量 DNN は、学習時に、音素フレームの言語特徴量、話者（話者ラベル）及び感情（感情ラベル）が入力層の各ユニットに与えられ、音素フレームの音響特徴量が出力層の各ユニットに与えられることで、入力層、隠れ層及び出力層の各ユニットの重み等が計算され、音素フレーム単位の学習が行われる。

20

【 0 0 3 0 】

学習のための音素フレームの言語特徴量は、例えば、音素の言語特徴量と同様の音素ラベル、アクセントの位置、品詞情報、アクセント句の情報、呼気段落の情報等に加え、音素を構成するフレームの番号（フレームの位置）、音素を構成するフレームの総数等の情報からなる。音素フレームの音響特徴量は、例えば、スペクトル係数、雑音性係数、ピッチ、有声／無声判定等の情報からなる。

【 0 0 3 1 】

また、後述する音声合成時には、音響特徴量 DNN の入力層の各ユニットに、音素フレームの言語特徴量、話者ラベル及び感情ラベルが与えられることで、出力層の各ユニットから、当該音素フレームの言語特徴量、話者ラベル及び感情ラベルに対応する音素フレームの音響特徴量が出力される。

30

【 0 0 3 2 】

図 1 に戻って、合成処理部 5 は、合成対象の音声波形に対応するテキスト、話者情報及び感情情報を入力し、話者情報及び感情情報に基づいて、話者ラベル及び感情ラベルをそれぞれ付与する。そして、合成処理部 5 は、テキストに対し、所定のテキスト解析を行い、言語特徴量等の情報を生成し、記憶部 4 に格納された時間長 DNN 及び音響特徴量 DNN を用いて、時間長及び音響特徴量を生成し、音声波形を合成して出力する。

【 0 0 3 3 】

例えば、話者 A による喜びの感情を表現した音声波形を合成する場合、合成処理部 5 は、所定のテキスト、話者 A を示す話者情報、及び喜びの感情を示す感情情報を入力し、時間長 DNN 及び音響特徴量 DNN を用いて、所定のテキストに対応する音声波形を合成する。

40

【 0 0 3 4 】

〔事前学習部 3 / 構成〕

次に、図 1 に示した事前学習部 3 の構成について詳細に説明する。図 2 は、事前学習部 3 の構成例を示すブロック図である。この事前学習部 3 は、テキスト解析部 1 1、話者感情ラベル処理部 1 2 及び音響分析部 1 3 を備えている。

【 0 0 3 5 】

事前学習部 3 は、記憶部 2 から、事前学習対象の音声コーパスのテキスト、話者情報、

50

感情情報及び音声波形を読み出す。テキスト解析部 11 は、記憶部 2 の音声コーパスから読み出されたテキストに対し、テキスト解析を行い、音素毎に音素の言語特徴量を生成し、音素の言語特徴量に含まれる音素ラベルを音響分析部 13 に出力する。

【0036】

図 9 は、音素の言語特徴量、話者ラベル及び感情ラベルの例を説明する図である。この音素の言語特徴量は、テキスト解析により生成された情報である。図 9 に示すように、テキスト解析により生成された音素の言語特徴量は、音素毎に、「音素ラベル」「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」の各種情報からなる。「音素ラベル」は、テキストを構成する音素を特定するための情報（音素情報）であり、当該音素に加え、前後の音素も含まれる。「話者ラベル」は、話者情報に付与された情報であり、「感情ラベル」は、感情情報に付与された情報である。

10

【0037】

図 2 に戻って、テキスト解析部 11 は、テキスト解析にて生成した音素の言語特徴量に基づいて、事前学習のための音素の言語特徴量を生成し、話者感情ラベル処理部 12 から話者ラベル及び感情ラベルを入力する。そして、テキスト解析部 11 は、事前学習のための音素の言語特徴量、話者ラベル及び感情ラベルを、記憶部 4 の時間長 DNN における入力層の各ユニットに出力する。

【0038】

事前学習のための音素の言語特徴量は、テキスト解析により生成した音素の言語特徴量の一部の情報、及び、テキスト解析により生成した音素の言語特徴量を加工した情報からなる。事前学習のための音素の言語特徴量は、例えば「音素ラベル」「音素情報」「有声音の有無」「アクセントの位置」等の各種情報からなる。

20

【0039】

テキスト解析部 11 は、音響分析部 13 から音素の時間長を入力し、事前学習のための音素の言語特徴量及び音素の時間長に基づいて、音素の時間長が示す音素フレーム数分の音素フレームの言語特徴量を生成する。そして、テキスト解析部 11 は、音素フレームの言語特徴量、話者ラベル及び感情ラベルを、記憶部 4 の音響特徴量 DNN における入力層の各ユニットに出力する。

【0040】

音素フレームの言語特徴量は、事前学習のための音素の言語特徴量の各種情報に加え、音素フレームを特定するための情報からなる。音素フレームの言語特徴量は、例えば「音素ラベル」「有声音の有無」「アクセントの位置」「フレームの番号」「フレームの総数」等の情報からなる。テキスト解析部 11 の詳細については後述する。

30

【0041】

話者感情ラベル処理部 12 は、記憶部 2 の音声コーパスから読み出された話者情報及び感情情報に対し、話者ラベル及び感情ラベルをそれぞれ付与し、話者ラベル及び感情ラベルをテキスト解析部 11 に出力する。

【0042】

音響分析部 13 は、テキスト解析部 11 から音素ラベルを入力し、記憶部 2 の音声コーパスから読み出された音声波形に対し、所定の学習データを用いて音響分析を行うと共に、音素の区切り位置を求める。そして、音響分析部 13 は、音素の区切り位置から音素の時間長を求めると共に、音素フレームの音響特徴量を生成する。音響分析部 13 は、音素の時間長をテキスト解析部 11 に出力すると共に、記憶部 4 の時間長 DNN における出力層のユニットに出力する。音響分析部 13 は、音素フレームの音響特徴量を、記憶部 4 の音響特徴量 DNN における出力層の各ユニットに出力する。

40

【0043】

音響分析により音素の区切り位置及び音素の時間長を求め、音素フレームの音響特徴量を生成する手法は既知であるから、ここでは詳細な説明は省略する。

【0044】

音素フレームの音響特徴量は、例えば、「スペクトル係数」「雑音性係数」「ピッチ」

50

「有声／無声判定」等の情報からなる。音響分析部 1 3 の詳細については後述する。

【 0 0 4 5 】

このような事前学習により、時間長 D N N 及び音響特徴量 D N N は、話者性及び感情表現と音声波形とを関連付けたモデルとなる。

【 0 0 4 6 】

〔事前学習部 3 / 処理〕

次に、図 2 に示した事前学習部 3 の処理について説明する。図 1 0 は、事前学習部 3 の処理例を示すフローチャートである。事前学習部 3 のテキスト解析部 1 1 は、記憶部 2 の音声コーパスから読み出されたテキストに対し、テキスト解析を行い（ステップ S 1 0 0 1 ）、音素の言語特徴量を生成する。そして、テキスト解析部 1 1 は、テキスト解析にて生成した音素の言語特徴量に基づいて、事前学習のための音素の言語特徴量を生成する（ステップ S 1 0 0 2 ）。

10

【 0 0 4 7 】

話者感情ラベル処理部 1 2 は、記憶部 2 の音声コーパスから読み出された話者情報及び感情情報に対し、話者ラベル及び感情ラベルをそれぞれ付与する（ステップ S 1 0 0 3 ）。

【 0 0 4 8 】

音響分析部 1 3 は、記憶部 2 の音声コーパスから読み出された音声波形に対し、音響解析を行い（ステップ S 1 0 0 4 ）、音素の区切り位置を求め、音素の時間長を求める（ステップ S 1 0 0 5 ）。

20

【 0 0 4 9 】

テキスト解析部 1 1 は、事前学習のための音素の言語特徴量、話者ラベル及び感情ラベルを時間長 D N N に出力すると共に、音響分析部 1 3 は、音素の時間長を時間長 D N N に出力する。これにより、時間長 D N N の事前学習が行われる（ステップ S 1 0 0 6 ）。

【 0 0 5 0 】

音響分析部 1 3 は、音声波形を音響分析することで、ステップ S 1 0 0 5 にて求めた音素の区切り位置に基づいて、音素フレームの音響特徴量を生成する（ステップ S 1 0 0 7 ）。

【 0 0 5 1 】

テキスト解析部 1 1 は、事前学習のための音素の言語特徴量、及び音響分析部 1 3 により求めた音素の時間長に基づいて、音素フレームの言語特徴量を生成する（ステップ S 1 0 0 8 ）。

30

【 0 0 5 2 】

テキスト解析部 1 1 は、音素フレームの言語特徴量、話者ラベル及び感情ラベルを音響特徴量 D N N に出力すると共に、音響分析部 1 3 は、音素フレームの音響特徴量を音響特徴量 D N N に出力する。これにより、音響特徴量 D N N の事前学習が行われる（ステップ S 1 0 0 9 ）。

【 0 0 5 3 】

〔言語特徴量と音響特徴量との間の関係〕

次に、時間長 D N N 及び音響特徴量 D N N の入出力データである音素の言語特徴量、音素フレームの言語特徴量、音素の時間長及び音素フレームの音響特徴量の関係について説明する。図 8 は、言語特徴量及び音響特徴量の関係について説明する図である。

40

【 0 0 5 4 】

テキストを「い」「ま」とし、「い」の音素ラベルを「i 」、 「ま」の音素ラベルを「m 」「a ）」とする。また、音素ラベル「i 」「m 」「a ）」における音素の時間長をそれぞれ「1 2 」「8 」「1 5 ）」とする。音素の時間長は、1 音素あたりの音素フレームの数を示す。音素フレームの時間長は、例えば 5 msec である。

【 0 0 5 5 】

図 8 に示すように、音素ラベル「i ）」の時間区間において、この 1 音素に対応して、1 組の音素の言語特徴量（の各情報）が生成され、1 2 組の音素フレームの言語特徴量（の

50

各情報)が生成され、12組の音素フレームの音響特徴量(の各情報)が生成される。

【0056】

また、音素ラベル「m」の時間区間において、この1音素に対応して、1組の音素の言語特徴量が生成され、8組の音素フレームの言語特徴量が生成され、8組の音素フレームの音響特徴量が生成される。

【0057】

また、音素ラベル「a」の時間区間において、この1音素に対応して、1組の音素の言語特徴量が生成され、15組の音素フレームの言語特徴量が生成され、15組の音素フレームの音響特徴量が生成される。

【0058】

このように、事前学習において、時間長DNNの入力層の各ユニットには、音素の言語特徴量、話者ラベル及び感情ラベルが与えられ、出力層のユニットには、音素の時間長が与えられ、この事前学習は音素を単位として行われる。つまり、時間長DNNには、音素毎に、音素の言語特徴量、話者ラベル、感情ラベル及び音素の時間長が与えられ、事前学習が行われる。音声合成においては、音素毎に、時間長DNNを用いて、音素の言語特徴量、話者ラベル及び感情ラベルに基づいて、音素の時間長が生成され出力される。

【0059】

また、事前学習において、音響特徴量DNNの入力層の各ユニットには、音素フレームの言語特徴量、話者ラベル及び感情ラベルが与えられ、出力層の各ユニットには、音素フレームの音響特徴量が与えられ、この事前学習は音素フレームを単位として行われる。つまり、音響特徴量DNNには、音素フレーム毎に、音素フレームの言語特徴量、話者ラベル、感情ラベル及び音素フレームの音響特徴量が与えられ、事前学習が行われる。音声合成においては、音素フレーム毎に、音響特徴量DNNを用いて、音素フレームの言語特徴量、話者ラベル及び感情ラベルに基づいて、音素フレームの音響特徴量が生成され出力される。

【0060】

(テキスト解析部11)

次に、図2に示したテキスト解析部11について詳細に説明する。図3は、テキスト解析部11の構成例を示すブロック図である。このテキスト解析部11は、テキスト解析手段31、前処理手段32及びフレーム処理手段33を備えている。

【0061】

テキスト解析手段31は、記憶部2の音声コーパスから読み出されたテキストに対し、形態素解析等のテキスト解析を行い、音素毎に音素の言語特徴量を生成する。そして、テキスト解析手段31は、音素の言語特徴量を前処理手段32に出力する。

【0062】

前処理手段32は、テキスト解析手段31から、テキスト解析により生成された音素の言語特徴量を入力すると共に、話者感情ラベル処理部12から話者ラベル及び感情ラベルを入力する。そして、前処理手段32は、図9に示したように、テキスト解析により生成された音素の言語特徴量、話者ラベル及び感情ラベルからなる情報群を生成する。

【0063】

前処理手段32は、テキスト解析により生成された音素の言語特徴量(図9を参照)に基づいて、事前学習のための音素の言語特徴量を生成する。そして、前処理手段32は、事前学習のための音素の言語特徴量、話者ラベル及び感情ラベルをフレーム処理手段33に出力すると共に、記憶部4の時間長DNNにおける入力層の各ユニットに出力する。

【0064】

フレーム処理手段33は、前処理手段32から、事前学習のための音素の言語特徴量、話者ラベル及び感情ラベルを入力すると共に、音響分析部13から音素の時間長を入力する。そして、フレーム処理手段33は、事前学習のための音素の言語特徴量及び音素の時間長に基づいて、音素の時間長が示す音素フレーム数分の音素フレームの言語特徴量を生成する。

10

20

30

40

50

【 0 0 6 5 】

フレーム処理手段 3 3 は、音素フレームの言語特徴量、話者ラベル及び感情ラベルを、記憶部 4 の音響特徴量 D N N における入力層の各ユニットに出力する。

【 0 0 6 6 】

(音響分析部 1 3)

次に、図 2 に示した音響分析部 1 3 について詳細に説明する。図 4 は、音響分析部 1 3 の構成例を示すブロック図である。この音響分析部 1 3 は、音素区切り処理手段 3 4 及び音響分析手段 3 5 を備えている。

【 0 0 6 7 】

音素区切り処理手段 3 4 は、テキスト解析部 1 1 から音素ラベルを入力し、記憶部 2 の音声コーパスから読み出された音声波形に対し、所定の学習データを用いて音響分析を行う。そして、音素区切り処理手段 3 4 は、音素ラベルの示す音素が音声波形内でどの位置にあるかを特定し、音素の区切り位置を求める。また、音素区切り処理手段 3 4 は、音素の区切り位置に基づいて、音素ラベルの示す音素の時間長を求める。前述のとおり、音素の時間長は、音素を構成する音素フレームの数で表される。

【 0 0 6 8 】

音素区切り処理手段 3 4 は、音素の区切り位置を音響分析手段 3 5 に出力し、音素の時間長をテキスト解析部 1 1 に出力すると共に、記憶部 4 の時間長 D N N における出力層のユニットに出力する。

【 0 0 6 9 】

音響分析手段 3 5 は、音素区切り処理手段 3 4 から音素の区切り位置を入力し、記憶部 2 の音声コーパスから読み出された音声波形に対し、音響分析を行い、音素を構成する複数の音素フレームのそれぞれについて、音素フレームの音響特徴量を生成する。

【 0 0 7 0 】

音響分析手段 3 5 は、音素フレームの音響特徴量を、記憶部 4 の音響特徴量 D N N における出力層の各ユニットに出力する。

【 0 0 7 1 】

(合成処理部 5)

次に、図 1 に示した合成処理部 5 の構成について詳細に説明する。図 5 は、合成処理部 5 の構成例を示すブロック図である。この合成処理部 5 は、テキスト解析部 2 1、話者感情ラベル処理部 2 2、時間長及び音響特徴量生成部 2 3 及び音声波形合成部 2 4 を備えている。

【 0 0 7 2 】

テキスト解析部 2 1 は、図 2 に示したテキスト解析部 1 1 と同様の処理を行う。具体的には、テキスト解析部 2 1 は、合成対象の音声波形に対応するテキストを入力し、テキストに対してテキスト解析を行い、音素毎に音素の言語特徴量を生成する。

【 0 0 7 3 】

テキスト解析部 2 1 は、テキスト解析にて生成した音素の言語特徴量に基づいて、図 2 に示したテキスト解析部 1 1 により生成された事前学習のための音素の言語特徴量と同様の音素の言語特徴量を生成し、話者感情ラベル処理部 2 2 から話者ラベル及び感情ラベルを入力する。そして、テキスト解析部 2 1 は、音素の言語特徴量、話者ラベル及び感情ラベルを時間長及び音響特徴量生成部 2 3 に出力する。

【 0 0 7 4 】

テキスト解析部 2 1 は、時間長及び音響特徴量生成部 2 3 から、当該時間長及び音響特徴量生成部 2 3 に出力した音素の言語特徴量、話者ラベル及び感情ラベルに対応する音素の時間長を入力し、音素の言語特徴量及び音素の時間長に基づいて、音素の時間長が示す音素フレーム数分の音素フレームの言語特徴量を生成する。そして、テキスト解析部 2 1 は、音素フレームの言語特徴量、話者ラベル及び感情ラベルを、時間長及び音響特徴量生成部 2 3 に出力する。

【 0 0 7 5 】

10

20

30

40

50

話者感情ラベル処理部 2 2 は、図 2 に示した話者感情ラベル処理部 1 2 と同様の処理を行う。具体的には、話者感情ラベル処理部 2 2 は、話者情報及び感情情報を入力し、話者情報及び感情情報に対し、話者ラベル及び感情ラベルをそれぞれ付与し、話者ラベル及び感情ラベルをテキスト解析部 2 1 に出力する。

【 0 0 7 6 】

時間長及び音響特徴量生成部 2 3 は、テキスト解析部 2 1 から音素の言語特徴量、話者ラベル及び感情ラベルを入力し、記憶部 4 の時間長 DNN を用いて、音素の言語特徴量、話者ラベル及び感情ラベルに基づいて、音素の時間長を生成する。そして、時間長及び音響特徴量生成部 2 3 は、音素の時間長をテキスト解析部 2 1 に出力する。

【 0 0 7 7 】

時間長及び音響特徴量生成部 2 3 は、テキスト解析部 2 1 から音素フレームの言語特徴量、話者ラベル及び感情ラベルを入力し、記憶部 4 の音響特徴量 DNN を用いて、音素フレームの言語特徴量、話者ラベル及び感情ラベルに基づいて、音素フレームの音響特徴量を生成する。

【 0 0 7 8 】

時間長及び音響特徴量生成部 2 3 は、音素フレームの音響特徴量を音声波形合成部 2 4 に出力する。

【 0 0 7 9 】

音声波形合成部 2 4 は、時間長及び音響特徴量生成部 2 3 から音素フレームの音響特徴量を入力し、音素フレームの音響特徴量に基づいて、音声波形を合成し、合成した音声波形を出力する。

【 0 0 8 0 】

具体的には、音声波形合成部 2 4 は、音素フレームの音響特徴量に含まれるピッチ、雑音特性等の情報に基づいて、声帯音源波形を生成する。そして、音声波形合成部 2 4 は、声帯音源波形に対し、音素フレームの音響特徴量に含まれるスペクトル係数等の情報に基づいて声道フィルタ処理を施し、音声波形を合成する。

【 0 0 8 1 】

音素フレームの音響特徴量に基づいて音声波形を合成する手法は既知であるから、ここでは詳細な説明を省略する。

【 0 0 8 2 】

このような音声合成により、話者情報及び感情情報と音声波形とを関連付けた時間長 DNN 及び音響特徴量 DNN を用いることで、テキスト、話者情報及び感情情報に対応する音声波形が合成される。

【 0 0 8 3 】

〔合成処理部 5 / 処理〕

次に、図 5 に示した合成処理部 5 の処理について説明する。図 1 1 は、合成処理部 5 の処理例を示すフローチャートである。合成処理部 5 のテキスト解析部 2 1 は、合成対象の音声波形に対応するテキストに対し、テキスト解析を行い(ステップ S 1 1 0 1)、音素の言語特徴量を生成する(ステップ S 1 1 0 2)。

【 0 0 8 4 】

話者感情ラベル処理部 2 2 は、話者情報及び感情情報に対し、話者ラベル及び感情ラベルをそれぞれ付与する(ステップ S 1 1 0 3)。

【 0 0 8 5 】

時間長及び音響特徴量生成部 2 3 は、時間長 DNN を用いて、音素の言語特徴量、話者ラベル及び感情ラベルに基づき音素の時間長を生成する(ステップ S 1 1 0 4)。そして、テキスト解析部 2 1 は、音素の言語特徴量及び音素の時間長に基づいて、音素フレームの言語特徴量を生成する(ステップ S 1 1 0 5)。

【 0 0 8 6 】

時間長及び音響特徴量生成部 2 3 は、音響特徴量 DNN を用いて、音素フレームの言語特徴量、話者ラベル及び感情ラベルに基づき音素フレームの音響特徴量を生成する(ステ

10

20

30

40

50

ップS1106)。

【0087】

音声波形合成部24は、音素フレームの音響特徴量に基づいて、音声波形を合成し、合成した音声波形を出力する(ステップS1107)。

【0088】

以上のように、本発明の実施形態の音声合成装置1によれば、学習時に、事前学習部3のテキスト解析部11は、音声コーパスから読み出されたテキストに対しテキスト解析を行い、音素の言語特徴量を生成する。音響分析部13は、音声コーパスから読み出された音声波形に対して音響分析を行い、音素の区切り位置を求め、音素の時間長を求める。テキスト解析部11は、音素の言語特徴量、話者ラベル及び感情ラベルを、時間長DNNにおける入力層の各ユニットに出力すると共に、音響分析部13は、音素の時間長を、時間長DNNにおける出力層のユニットに出力する。これにより、時間長DNNの事前学習が行われる。

10

【0089】

また、音響分析部13は、音声波形を音響分析することで、音素の区切り位置に基づいて、音素フレームの音響特徴量を生成する。テキスト解析部11は、音素の言語特徴量及び音素の時間長に基づいて、音素フレームの言語特徴量を生成し、音素フレームの言語特徴量、話者ラベル及び感情ラベルを、音響特徴量DNNにおける入力層の各ユニットに出力すると共に、音響分析部13は、音素フレームの音響特徴量を、音響特徴量DNNにおける出力層の各ユニットに出力する。これにより、音響特徴量DNNの事前学習が行われる。

20

【0090】

さらに、本発明の実施形態の音声合成装置1によれば、音声合成時に、合成処理部5のテキスト解析部21は、対象のテキストに対しテキスト解析を行い、音素の言語特徴量を生成する。時間長及び音響特徴量生成部23は、時間長DNNを用いて、音素の言語特徴量、話者ラベル及び感情ラベルに基づき音素の時間長を生成する。

【0091】

テキスト解析部21は、音素の言語特徴量及び音素の時間長に基づいて、音素フレームの言語特徴量を生成する。時間長及び音響特徴量生成部23は、音響特徴量DNNを用いて、音素フレームの言語特徴量、話者ラベル及び感情ラベルに基づき音素フレームの音響特徴量を生成する。そして、音声波形合成部24は、音素フレームの音響特徴量に基づいて、音声波形を合成する。

30

【0092】

これにより、話者性及び感情表現と音声波形とを関連付けた時間長DNN及び音響特徴量DNNを用いるようにしたから、HMM方式よりも音質の良い音声合成を実現することができると共に、話者性及び感情表現を同時に制御することができる。このため、学習時には、話者性及び感情表現に対応した複雑な特徴抽出作業を行う必要がない。

【0093】

また、時間長DNN及び音響特徴量DNNからなる2つのDNNを用いて音声合成を行うようにしたから、話者及び感情の組み合わせ毎の膨大なDNNを用いる必要がなく、途切れることのない連続的なかつ自然な読み上げ音声を生成することができる。

40

【0094】

したがって、簡易な構成にて、話者性及び感情表現を同時に制御可能な音声合成を実現することができる。

【0095】

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。前記実施形態では、図1に示したように、音声合成装置1は、事前学習を行う事前学習部3と、音声合成を行う合成処理部5とを備えるようにした。これに対し、事前学習部3と合成処理部5とを、それぞれ異なる装置に備えるようにしてもよい。

50

【 0 0 9 6 】

具体的には、記憶部 2 を備えた記憶装置、事前学習部 3 を備えた学習装置、記憶部 4 を備えた記憶装置、及び合成処理部 5 を備えた合成装置により音声合成システムが構成される。この場合、学習装置と、記憶部 2 を備えた記憶装置及び記憶部 4 を備えた記憶装置とは、インターネットを介して接続されるようにしてもよい。また、合成装置と、記憶部 4 を備えた記憶装置とは、同様にインターネットを介して接続されるようにしてもよい。さらに、学習装置は、記憶部 2、事前学習部 3 及び記憶部 4 を備え、合成装置は、記憶媒体を介して可搬された記憶部 4、及び合成処理部 5 を備えるようにしてもよい。

【 0 0 9 7 】

また、前記実施形態では、時間長 DNN 及び音響特徴量 DNN におけるそれぞれの入力層のユニットに、話者情報及び感情情報を与えるようにした。これに対し、これらの入力層のユニットに、複数の話者情報及び複数の感情情報を与えるようにしてもよい。例えば、話者が複数の観点から分類され、話者に対して複数の話者情報が紐付けられ、同様に、感情が複数の観点から分類され、感情に対して複数の感情情報が紐付けられ、これらを入力層のユニットに与えるようにしてもよい。

10

【 0 0 9 8 】

尚、本発明の実施形態による音声合成装置 1 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。音声合成装置 1 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。音声合成装置 1 に備えた事前学習部 3 及び合成処理部 5 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、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。

20

【 符号の説明 】

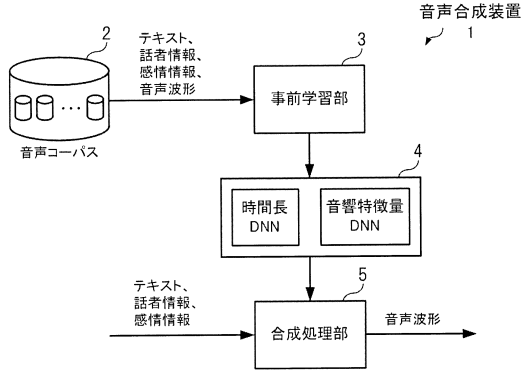
【 0 0 9 9 】

- 1 音声合成装置
- 2, 4 記憶部
- 3 事前学習部
- 5 合成処理部
- 11, 21 テキスト解析部
- 12, 22 話者感情ラベル処理部
- 13 音響分析部
- 23 時間長及び音響特徴量生成部
- 24 音声波形合成部
- 31 テキスト解析手段
- 32 前処理手段
- 33 フレーム処理手段
- 34 音素区切り処理手段
- 35 音響分析手段

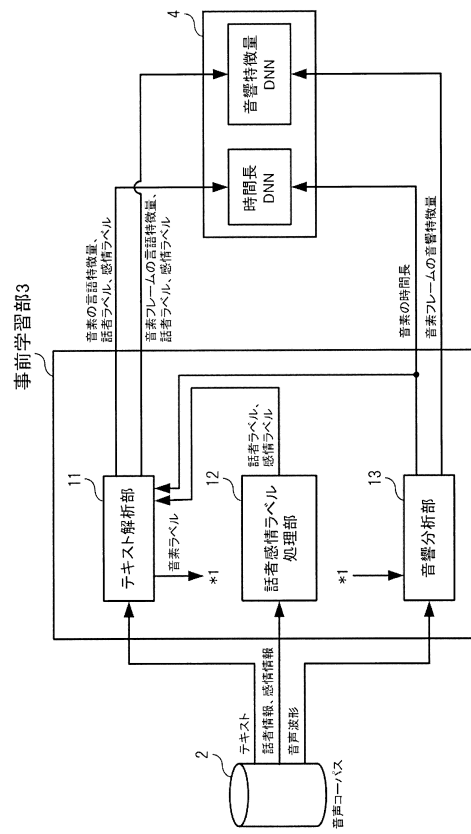
30

40

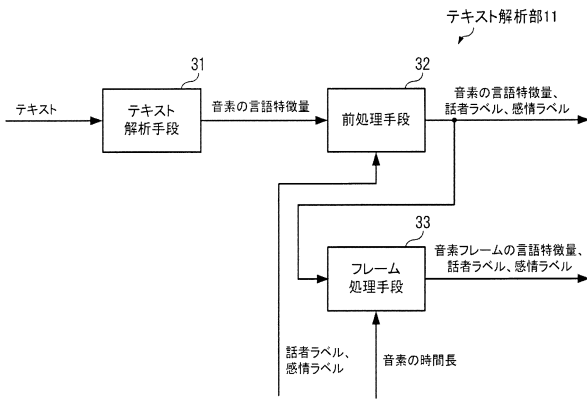
【図1】



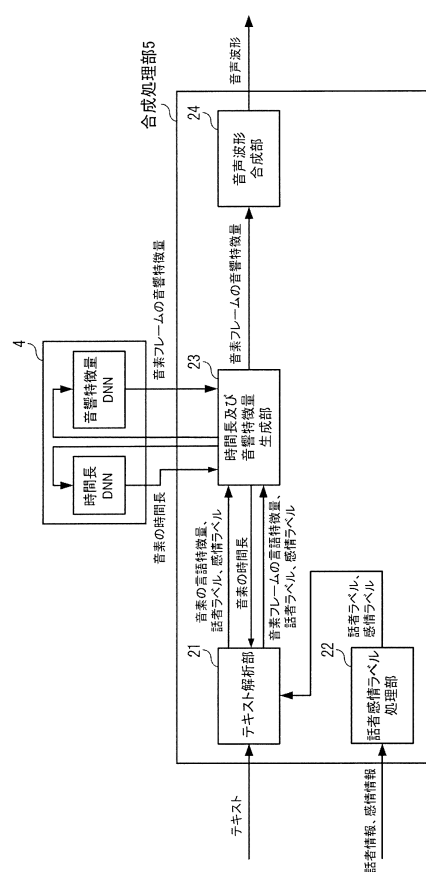
【図2】



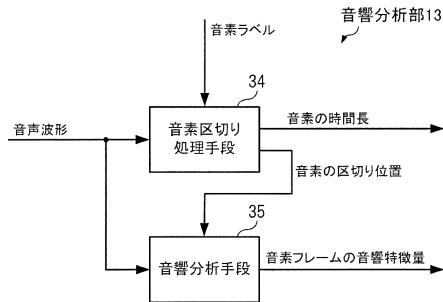
【図3】



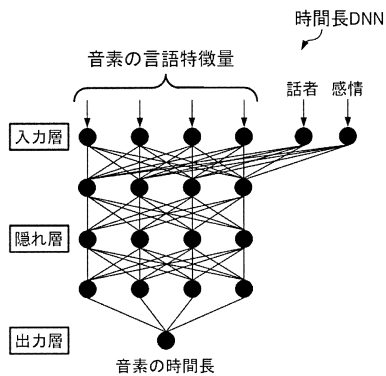
【図5】



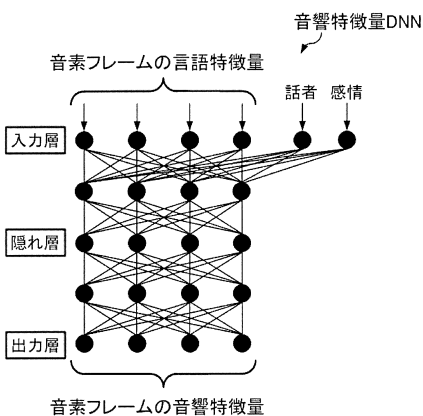
【図4】



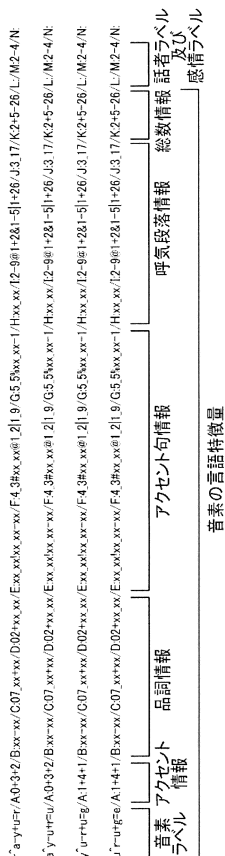
【図6】



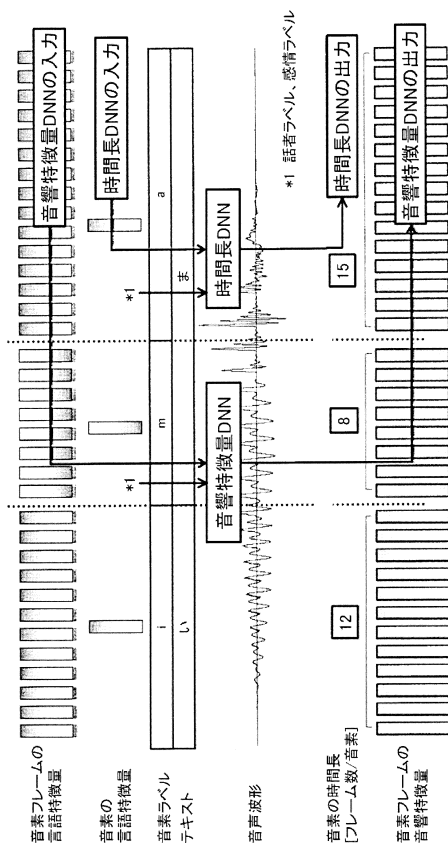
【図7】



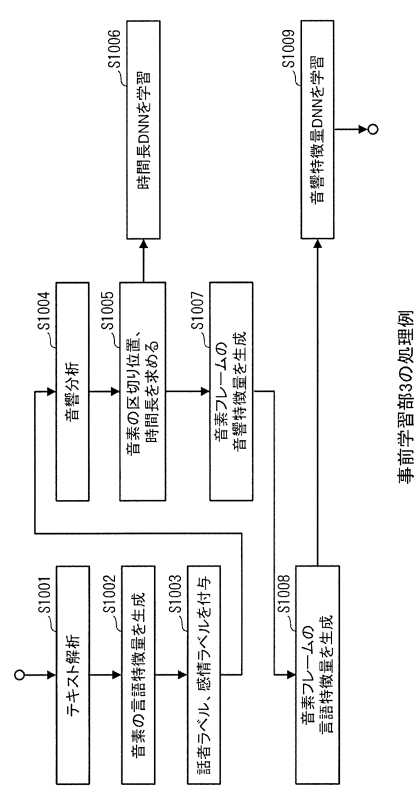
【図9】



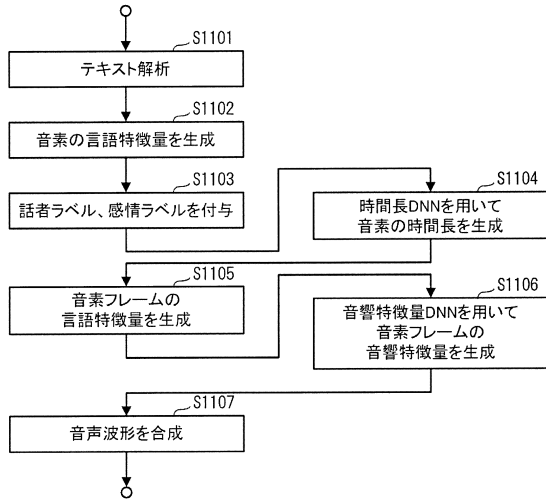
【図8】



【図10】



【図11】



合成処理部5の処理例

フロントページの続き

(51)Int.Cl. F I
G 1 0 L 13/033 1 0 2 B

(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 上田 雄

(56)参考文献 米国特許第08527276(US, B1)
特開平02-072399(JP, A)
中国特許出願公開第104538024(CN, A)
清山 信正, "解説02 音声合成技術の動向と放送・通信分野における応用展開", NHK技研
R&D, 2017年 1月15日, No. 161
高木 信二, "とてもDeepなテキスト音声合成", 電子情報通信学会技術研究報告, 2017
年 1月14日, Vol. 116, No. 414, pp. 41 - 46
LUONG, Hieu Thi、外3名, "DNNに基づくテキスト音声合成における話者・ジェンダー・年齢
コード利用の検討", 電子情報通信学会技術研究報告, 2016年10月20日, Vol. 11
6, No. 279

(58)調査した分野(Int.Cl., DB名)
G 1 0 L 13 / 0 0 - 1 3 / 1 0