



**【特許請求の範囲】****【請求項 1】**

目標音声に対応する音韻系列を合成単位で区切って得られる複数のセグメントを取得する第 1 の取得部と、

前記目標音声に対応する各々の前記セグメントの韻律情報を取得する第 2 の取得部と、

各々の前記セグメントごとに、当該セグメントに対し、当該セグメントの前記韻律情報に基づいて、予め用意された複数の音声素片のうちから、複数の音声素片を選択する選択部と、

各々の前記セグメントごとに、当該セグメントに対して選択された複数の前記音声素片を融合することによって、融合素片を生成する融合部と、

各々の前記セグメントごとに、前記選択部により選択された複数の前記音声素片に関する特徴量と、前記融合部により生成された前記融合素片に関する特徴量との少なくとも一方を用いて、当該セグメントに係る前記融合素片に対して行うべきフォルマント強調における強調度合いを推定する推定部と、

各々の前記セグメントごとに、当該セグメントに係る前記融合素片に対して、前記推定部が推定した前記強調度合いに基づくフォルマント強調を行うフォルマント強調フィルタ部とを備えたことを特徴とする音声処理装置。

**【請求項 2】**

各々の前記セグメントについて前記フォルマント強調フィルタ部によりそれぞれ得られたフォルマント強調された前記融合素片に係る音声波形をもとにして、合成音声を生成する生成部を更に備えたことを特徴とする請求項 1 に記載の音声処理装置。

**【請求項 3】**

各々の前記セグメントについて前記フォルマント強調フィルタ部によりそれぞれ得られたフォルマント強調された前記融合素片をそのまま出力する出力部を更に備えたことを特徴とする請求項 1 に記載の音声処理装置。

**【請求項 4】**

前記出力部は、前記融合素片を、テキスト音声合成に供するための音声素片を記憶する記憶部に出力することを特徴とする請求項 3 に記載の音声処理装置。

**【請求項 5】**

前記予め用意された複数の音声素片を記憶する音声素片記憶部を更に備えたことを特徴とする請求項 1 ないし 4 のいずれか 1 項に記載の音声処理装置。

**【請求項 6】**

前記推定部は、各々の前記セグメントごとに、前記融合部により生成された前記融合素片のスペクトル包絡が、前記選択部により選択された前記音声素片のスペクトル包絡から、どの程度鈍ったかを推定し、推定されたスペクトル包絡の鈍り具合が大きいセグメントほど、強めのフォルマント強調度合いを推定することを特徴とする請求項 1 ないし 5 のいずれか 1 項に記載の音声処理装置。

**【請求項 7】**

前記推定部は、各々の前記セグメントごとに、前記融合部により生成された前記融合素片のスペクトル包絡と、前記選択部により選択された前記音声素片のスペクトル包絡の形状との差を推定し、推定されたスペクトル包絡の形状の差が大きいセグメントほど、強めのフォルマント強調度合いを推定することを特徴とする請求項 1 ないし 5 のいずれか 1 項に記載の音声処理装置。

**【請求項 8】**

前記推定部は、各々の前記セグメントごとに、目標音声と前記融合部により生成された前記融合素片による音声との差を、当該セグメントの目標音声に対応する韻律情報と前記選択部により選択された前記音声素片の韻律情報とから推定し、推定された当該目標音声と融合素片による音声との差が大きいセグメントほど、強めのフォルマント強調度合いを推定することを特徴とする請求項 1 ないし 5 のいずれか 1 項に記載の音声処理装置。

**【請求項 9】**

前記推定部は、前記複数のセグメントのそれぞれに対して、フォルマントごと又は複数に分割した周波数帯域ごとにフォルマント強調度合いを推定し、

前記フォルマント強調フィルタ部は、それぞれのフォルマント又は周波数帯域に対して推定されたフォルマント強調度合いに従って、フォルマント又は周波数帯域間で異なる強さのフォルマント強調を行うことを特徴とする請求項 1 ないし 5 のいずれか 1 項に記載の音声処理装置。

【請求項 10】

第 1 の取得部、第 2 の取得部、選択部、融合部、推定部及びフォルマント強調フィルタ部を備えた音声処理装置の音声処理方法であって、

前記第 1 の取得部が、目標音声に対応する音韻系列を合成単位で区切って得られる複数のセグメントを取得するステップと、

前記第 2 の取得部が、前記目標音声に対応する各々の前記セグメントの韻律情報を取得するステップと、

前記選択部が、各々の前記セグメントごとに、当該セグメントに対し、当該セグメントの前記韻律情報に基づいて、予め用意された複数の音声素片のうちから、複数の音声素片を選択するステップと、

前記融合部が、各々の前記セグメントごとに、当該セグメントに対して選択された複数の前記音声素片を融合することによって、融合素片を生成するステップと、

前記推定部が、各々の前記セグメントごとに、前記選択部により選択された複数の前記音声素片に関する特徴量と、前記融合部により生成された前記融合素片に関する特徴量との少なくとも一方を用いて、当該セグメントに係る前記融合素片に対して行うべきフォルマント強調における強調度合いを推定するステップと、

前記フォルマント強調フィルタ部が、各々の前記セグメントごとに、当該セグメントに係る前記融合素片に対して、前記推定部が推定した前記強調度合いに基づくフォルマント強調を行うステップとを有することを特徴とする音声処理方法。

【請求項 11】

前記音声処理装置は、生成部を更に備えるものであり、

前記音声処理方法は、前記生成部が、各々の前記セグメントについて前記フォルマント強調フィルタ部によりそれぞれ得られたフォルマント強調された前記融合素片に係る音声波形をもとにして、合成音声を生成するステップを更に含むことを特徴とする請求項 10 に記載の音声処理方法。

【請求項 12】

前記音声処理装置は、出力部を更に備えるものであり、

前記音声処理方法は、前記出力部が、各々の前記セグメントについて前記フォルマント強調フィルタ部によりそれぞれ得られたフォルマント強調された前記融合素片をそのまま出力するステップを更に含むことを特徴とする請求項 10 に記載の音声処理方法。

【請求項 13】

第 1 の取得部、第 2 の取得部、選択部、融合部、推定部及びフォルマント強調フィルタ部を備えた音声処理装置としてコンピュータを機能させるためのプログラムであって、

前記第 1 の取得部が、目標音声に対応する音韻系列を合成単位で区切って得られる複数のセグメントを取得するステップと、

前記第 2 の取得部が、前記目標音声に対応する各々の前記セグメントの韻律情報を取得するステップと、

前記選択部が、各々の前記セグメントごとに、当該セグメントに対し、当該セグメントの前記韻律情報に基づいて、予め用意された複数の音声素片のうちから、複数の音声素片を選択するステップと、

前記融合部が、各々の前記セグメントごとに、当該セグメントに対して選択された複数の前記音声素片を融合することによって、融合素片を生成するステップと、

前記推定部が、各々の前記セグメントごとに、前記選択部により選択された複数の前記音声素片に関する特徴量と、前記融合部により生成された前記融合素片に関する特徴量

10

20

30

40

50

との少なくとも一方を用いて、当該セグメントに係る前記融合素片に対して行うべきフォルマント強調における強調度合いを推定するステップと、

前記フォルマント強調フィルタ部が、各々の前記セグメントごとに、当該セグメントに係る前記融合素片に対して、前記推定部が推定した前記強調度合いに基づくフォルマント強調を行うステップとをコンピュータに実行させるためのプログラム。

【請求項 14】

前記音声処理装置は、生成部を更に備えるものであり、

前記プログラムは、前記生成部が、各々の前記セグメントについて前記フォルマント強調フィルタ部によりそれぞれ得られたフォルマント強調された前記融合素片に係る音声波形をもとにして、合成音声を作成するステップを更にコンピュータに実行させることを特徴とする請求項 13 に記載の音声処理装置。

10

【請求項 15】

前記音声処理装置は、出力部を更に備えるものであり、

前記プログラムは、前記出力部が、各々の前記セグメントについて前記フォルマント強調フィルタ部によりそれぞれ得られたフォルマント強調された前記融合素片をそのまま出力するステップを更にコンピュータに実行させることを特徴とする請求項 13 に記載の音声処理装置。

【発明の詳細な説明】

【技術分野】

【0001】

20

本発明は、音声処理装置、音声処理方法及びプログラムに関する。

【背景技術】

【0002】

任意の文章から人工的に音声信号を作り出すことを、テキスト音声合成という。テキスト音声合成は、一般的に、言語処理部、韻律処理部及び音声合成部の 3 つ段階によって行われる。

【0003】

入力されたテキストは、まず言語処理部において、形態素解析や構文解析が行われ、次に韻律処理部において、アクセントやイントネーションの処理が行われて、音韻系列・韻律情報（基本周波数、音韻継続時間長、パワーなど）が出力される。最後に、音声合成部において、音韻系列・韻律情報から音声信号を合成する。そこで、音声合成部に用いる音声合成方法は、韻律処理部で生成される任意の音韻系列を、任意の韻律で音声合成することが可能な方法でなければならない。

30

【0004】

従来、このような音声合成方法として、入力の音韻系列を分割して得られる複数の合成単位（合成単位列）のそれぞれに対して、入力された音韻系列・韻律情報を目標にして、予め記憶された大量の音声素片の中から音声素片を選択し、選択した音声素片を合成単位間で接続することによって、音声を合成する、音声合成方法（素片選択型の音声合成方法）が知られている。例えば、特許文献 1 に開示された素片選択型の音声合成方法では、音声を合成することで生じる音声合成の劣化の度合いを、コストで表すこととし、予め定義されたコスト関数を用いて計算されるコストが小さくなるように、音声素片を選択する。例えば、音声素片を編集・接続することで生じる変形歪み及び接続歪みを、コストを用いて数値化し、このコストに基づいて、音声合成に使用する音声素片系列を選択し、選択した音声素片系列に基づいて、合成音声を生成する。

40

【0005】

特許文献 1 に開示された音声合成方法のように、音声を合成することで生じる音声合成の劣化の度合いを考慮して、大量の音声素片の中から適切な音声素片系列を選択することによって、音声素片の編集及び接続による音質の劣化を抑えた合成音声を生成することができる。

【0006】

50

しかしながら、特許文献 1 に開示された素片選択型の音声合成方法には、部分的に合成音の音質が劣化する問題点がある。この理由は次のようなものである。

【 0 0 0 7 】

第 1 の理由は、予め記憶された音声素片が非常に多い場合であっても、様々な音韻・韻律環境に対して適切な音声素片が存在するとは限らないことである。

【 0 0 0 8 】

第 2 の理由は、人が実際に感じる合成音声の劣化の度合いをコスト関数が完全に表現できないため、必ずしも最適な素片系列が選ばれない場合があるからである。

【 0 0 0 9 】

第 3 の理由は、音声素片が非常に多いために予め不良な音声素片を排除しておくことが困難であり、また不良な音声素片を取り除くためのコスト関数の設計も難しいため、選択された音声素片系列中に、突発的に不良な音声素片が混入する場合があるからである。

【 0 0 1 0 】

そこで、合成単位当たり 1 つずつの音声素片を選ぶのではなく、合成単位当たり複数個の音声素片を選択し、これを融合することによって新たな音声素片を生成し、こうして生成された音声素片を使って音声を合成する方法が開示されている（特許文献 2 参照）。以下、この方法を「複数素片選択融合型の音声合成方法」と呼ぶ。

【 0 0 1 1 】

特許文献 2 に開示された複数素片選択融合型の音声合成方法では、合成単位毎に複数の音声素片を融合することによって、目標とする音韻・韻律環境に合う適切な音声素片が存在しない場合や、最適な音声素片が選択されない場合、不良素片が選択されてしまった場合でも、高品質な音声素片を新たに生成することができ、さらに、この新たに生成した音声素片を使用して音声合成を行うことで、前述した素片選択型の音声合成方法の問題点を改善することができ、より安定性を増した高音質の音声合成を実現することができる。

【 0 0 1 2 】

この複数素片選択融合型の音声合成方法においては、音声素片の融合による平均化の副作用によってスペクトル包絡が原音に比べて若干鈍る傾向があり、その結果、こもり感やブザー感が生じる場合がある。こうしたこもり感やブザー感の主観的な改善には、音声符号化や音声合成でよく用いられるようなフォルマント強調フィルタを、融合された素片に対して適用することが効果的である。

【 0 0 1 3 】

フォルマント強調フィルタは、入力音声波形のスペクトル包絡のフォルマントによる山・谷を強調したような音声波形を出力するフィルタで、適度な度合いでフォルマントを強調できれば、スペクトル包絡が鈍ったことによって生じるこもり感やブザー感を改善できる。一般的に、フォルマント強調フィルタは入力波形のスペクトル特性に応じてフィルタ特性を変える点では適応的だが、どの程度フォルマントを強調するかについては、適切な強調度合いを決めるための客観尺度が存在しないため、主観評価などによって実験的に決めるしかなく、ハイパーパラメータなどの値を外部から指定することによって制御することが多い。

【 0 0 1 4 】

そのため、複数素片選択型の音声合成方法で用いる場合には、フォルマントの強調度合いは、合成音声の主観的な音質が総合的に良くなるように、主観評価などによって実験的に決める。すなわち、フォルマントの強調度合いは、融合されたあらゆる素片に対して共通のものが適用される。

【特許文献 1】特開 2 0 0 1 - 2 8 2 2 7 8 公報

【特許文献 2】特開 2 0 0 5 - 1 6 4 7 4 9 公報

【発明の開示】

【発明が解決しようとする課題】

【 0 0 1 5 】

しかしながら、音声素片の融合によるスペクトル包絡の鈍り具合は、通常、合成単位に

10

20

30

40

50

よって異なり、一様ではない。例えば、合成単位に対して選ばれた複数の素片が類似のスペクトル包絡を持つ場合は、融合してもさほどスペクトル包絡は鈍らないと考えられるが、フォルマントの位置が素片間で大きく異なるなど、選ばれた音声素片のスペクトル包絡がそれぞれ異なる特徴を持つ場合には、融合するとスペクトル包絡が鈍ってしまう可能性が高い。

#### 【 0 0 1 6 】

このような状況において、全音声素片に対して同じ強調度合いのフォルマント強調フィルタを一様に適用すると、融合によってスペクトル包絡が大きく鈍った箇所にはフォルマント強調の程度が不十分であるのに対し、逆に融合によるスペクトル包絡の鈍りが小さい箇所はフォルマントが強調されすぎて人工的な音になる問題がある。

10

#### 【 0 0 1 7 】

本発明は、上記事情を考慮してなされたもので、こもり感やブザー感が少なく、かつ人工的でない高音質な合成音声を生成できる音声処理装置、音声処理方法及びプログラムを提供することを目的とする。

#### 【課題を解決するための手段】

#### 【 0 0 1 8 】

本発明に係る音声処理装置は、目標音声に対応する音韻系列を合成単位で区切って得られる複数のセグメントを取得する第1の取得部と、前記目標音声に対応する各々の前記セグメントの韻律情報を取得する第2の取得部と、各々の前記セグメントごとに、当該セグメントに対し、当該セグメントの前記韻律情報に基づいて、予め用意された複数の音声素片のうちから、複数の音声素片を選択する選択部と、各々の前記セグメントごとに、当該セグメントに対して選択された複数の前記音声素片を融合することによって、融合素片を生成する融合部と、各々の前記セグメントごとに、前記選択部により選択された複数の前記音声素片に関する特徴量と、前記融合部により生成された前記融合素片に関する特徴量との少なくとも一方を用いて、当該セグメントに係る前記融合素片に対して行うべきフォルマント強調における強調度合いを推定する推定部と、各々の前記セグメントごとに、当該セグメントに係る前記融合素片に対して、前記推定部が推定した前記強調度合いに基づくフォルマント強調を行うフォルマント強調フィルタ部とを備えたことを特徴とする。

20

#### 【発明の効果】

30

#### 【 0 0 1 9 】

本発明によれば、こもり感やブザー感が少なく、かつ人工的でない高音質な合成音声を生成できる。

#### 【発明を実施するための最良の形態】

#### 【 0 0 2 0 】

以下、図面を参照しながら本発明の実施形態について説明する。

#### 【 0 0 2 1 】

(第1の実施形態)

本発明の第1の実施形態に係るテキスト音声合成装置(音声処理装置)について説明する。

40

#### 【 0 0 2 2 】

図1に、本実施形態に係るテキスト音声合成を行うテキスト音声合成装置(音声処理装置)の全体構成例を示す。

#### 【 0 0 2 3 】

図1に示されるように、本実施形態のテキスト音声合成装置は、テキスト入力部1、言語処理部2、韻律処理部3、音声合成部4を備えている。

#### 【 0 0 2 4 】

テキスト入力部1は、テキストを入力する。

#### 【 0 0 2 5 】

言語処理部2は、テキスト入力部1から入力されるテキストの形態素解析・構文解析を

50

行い、これら言語解析により得られた言語解析結果を韻律処理部 3 へ出力する。

【 0 0 2 6 】

韻律制御部 3 は、該言語解析結果を入力し、該言語解析結果からアクセントやイントネーションの処理を行って、音韻系列及び韻律情報を生成し、生成した音韻系列及び韻律情報を音声合成部へ出力する。

【 0 0 2 7 】

音声合成部 4 は、該音韻系列及び韻律情報を入力し、該音韻系列及び韻律情報から音声波形を生成して出力する。

【 0 0 2 8 】

以下、音声合成部 4 を中心に、その構成及び動作について詳細に説明する。

10

【 0 0 2 9 】

図 2 に、本実施形態の音声合成部 4 の構成例を示す。

【 0 0 3 0 】

図 2 に示されるように、音声合成部 4 は、音韻系列・韻律情報入力部 4 1、音声素片記憶部 4 2、素片選択部 4 3、素片融合部 4 4、フォルマント強調フィルタ部 4 5、フォルマント強調度合い推定部 4 6、素片編集・接続部 4 7、音声波形出力部 4 8 を備えている。

【 0 0 3 1 】

音韻系列・韻律情報入力部（以下、情報入力部と略記する。）4 1 は、音声合成部 4 への入力として、韻律制御部 3 から音韻系列・韻律情報を受理する。

20

【 0 0 3 2 】

音声素片記憶部（以下、素片記憶部と略記する。）4 2 は、大量の音声素片を蓄積している。また、素片記憶部 4 2 は、それら蓄積されている音声素片の全てについて、それぞれ、当該音声素片に対する音韻・韻律環境を併せて蓄積している。

【 0 0 3 3 】

素片選択部 4 3 は、素片記憶部 4 2 に蓄積された音声素片の中から、複数の音声素片を選択する。

【 0 0 3 4 】

素片融合部 4 4 は、素片選択部 4 3 により選択された複数の音声素片を融合して、新たな音声素片（以下、「融合素片」とも呼ぶ。）を生成する。

30

【 0 0 3 5 】

フォルマント強調フィルタ部 4 5 は、（次のフォルマント強調度合い推定部 4 6 により推定された、強調の程度に応じて）素片融合部 4 4 により生成された音声素片に対して、フォルマント強調を行う（すなわち、フォルマント強調された融合素片を生成する）。

【 0 0 3 6 】

フォルマント強調度合い推定部 4 6 は、フォルマント強調フィルタ部 4 5 においてフォルマントを強調する程度を推定する。

【 0 0 3 7 】

素片編集・接続部 4 7 は、フォルマント強調フィルタ部 4 5 から得られた音声素片を韻律変形及び接続して、合成音声の音声波形を生成する。

40

【 0 0 3 8 】

音声波形出力部 4 8 は、素片編集・接続部 4 7 で生成した音声波形を出力する。

【 0 0 3 9 】

なお、情報入力部 4 1 ~ 音声波形出力部 4 8 の各部の機能は、コンピュータに格納されたプログラムに実現できる。

【 0 0 4 0 】

次に、図 2 の音声合成部 4 の各ブロックについて詳しく説明する。

【 0 0 4 1 】

< 情報入力部 >

まず、情報入力部 4 1 は、韻律制御部 3 から入力された音韻系列・韻律情報を、素片選

50

択部 4 4 へ出力する。音韻系列は、例えば、音韻記号の系列である。また、韻律情報は、例えば、基本周波数、音韻継続時間長、パワーなどである。

【 0 0 4 2 】

以下、情報入力部 4 1 に入力される音韻系列、韻律情報を、それぞれ、入力音韻系列、入力韻律情報と呼ぶ。

【 0 0 4 3 】

< 素片記憶部 >

次に、素片記憶部 4 2 には、合成音声を生成するときに用いられる音声の単位（以下、「合成単位」と称する。）で、音声素片が大量に蓄積されている。

【 0 0 4 4 】

ここで、「合成単位」とは、音素あるいは音素を分割したもの（例えば、半音素など）の組み合わせ、例えば、半音素、音素（C、V）、ダイフオン（CV、VC、VV）、トライフオン（CVC、VCV）、音節（CV、V）、などであり（ここで、V は母音、C は子音を表す。）、また、これらが混在しているなど可変長であってもよい。

【 0 0 4 5 】

また、「音声素片」は、合成単位に対応する音声信号の波形もしくはその特徴を表すパラメータ系列などを表すものとする。

【 0 0 4 6 】

図 3 に、素片記憶部 4 2 に蓄積される音声素片の例を示す。図 3 に示すように、素片記憶部 4 2 には、各音素の音声信号の波形である音声素片が、当該音声素片を識別するための素片番号とともに記憶されている。これらの音声素片は、別途収録された多数の音声データに対して音素毎にラベル付けし、ラベルにしたがって音素毎に音声波形を切り出したものである。

【 0 0 4 7 】

また、素片記憶部 4 2 には、大量の音声素片とともに、各音声素片に対応した音韻・韻律環境が蓄積されている。

【 0 0 4 8 】

ここで、「音韻・韻律環境」とは、対応する音声素片にとって環境となる要因の組み合わせである。要因としては、例えば、当該音声素片の音素名、先行音素、後続音素、後々続音素、基本周波数、音韻継続時間長、パワー、ストレスの有無、アクセント核からの位置、息継ぎからの時間、発声速度、感情などがある。

【 0 0 4 9 】

また、素片記憶部 4 2 には、上記の他、音声素片の始端・終端でのケプストラム係数など、音声素片の音響特徴のうち音声素片の選択に用いる情報も蓄積されている。

【 0 0 5 0 】

以下では、素片記憶部 4 2 に蓄積される音声素片の音韻・韻律環境と音響特徴量とを総称して、「素片環境」と呼ぶ。

【 0 0 5 1 】

図 4 に、素片記憶部 4 2 に蓄積される素片環境の例を示す。図 4 に示す環境記憶部 4 3 には、素片記憶部 4 2 に蓄積される各音声素片の素片番号に対応して素片環境が記憶されている。ここでは、音韻・韻律環境として、音声素片に対応した音韻（音素名）、隣接音韻（この例では、当該音韻の前後それぞれ 2 音素ずつ）、基本周波数、音韻継続時間長が記憶され、音響特徴量として、音声素片始終端のケプストラム係数が記憶されている。

【 0 0 5 2 】

なお、これらの素片環境は、音声素片を切り出す元になった音声データを分析して抽出することによって得られる。また、図 4 では、音声素片の合成単位が音素である場合を示しているが、半音素、ダイフオン、トライフオン、音節、あるいはこれらの組み合わせや可変長であってもよい。

【 0 0 5 3 】

< 素片選択部 >

10

20

30

40

50



次に、図 2 の音声合成部 4 の動作を詳しく説明する。

【0054】

図 2 において、情報入力部 41 を介して素片選択部 43 に入力された音韻系列は、素片選択部 47 において、合成単位毎に区切られる。以下、この区切られた合成単位を、「セグメント」と呼ぶ。

【0055】

素片選択部 43 は、入力された入力音韻系列と入力韻律情報を基に、素片記憶部 42 を参照し、各セグメントに対して、それぞれ、融合する複数個の音声素片の組み合わせを選択する。

【0056】

このとき素片選択部 43 は、各音声素片候補を用いて音声を作成した場合の合成音声と目標音声との歪みができるだけ小さくなるように、融合する音声素片の組み合わせを選択する。ここでは、素片選択部 43 は、一般の素片選択型音声合成方法や従来の複数素片選択融合型音声合成方法と同様に、音声素片の選択の尺度として、各音声素片候補を用いて音声を作成した場合の合成音声と目標音声との歪みの大きさを間接的に表すコストを用い、このコストができるだけ小さくなるように、融合する音声素片の組み合わせを選択する。

【0057】

ここで、「目標音声」とは、音声を作成する際の目標となる（仮想的な）音声、すなわち、入力された音韻の並びと韻律を実現し、かつ、理想的に自然な音声をいう。

【0058】

最初に、コストについて説明する。

【0059】

合成音声の目標音声に対する歪みの度合いを表すコストには、大きく分けて、目標コストと接続コストの 2 種類のコストがある。

【0060】

目標コストは、コストの算出対象である音声素片（対象素片）を目標の音韻・韻律環境で使用するによって生じるコストである。

【0061】

接続コストは、対象素片を隣接する音声素片と接続したときに生じるコストである。

【0062】

具体的には、次の通りである。

【0063】

目標コストとしては、音声素片が持つ基本周波数と目標の基本周波数の違い（差）によって生じる歪み（基本周波数コスト）、音声素片の音韻継続時間長と目標の音韻継続時間長の違い（差）によって生じる歪み（継続時間長コスト）、音声素片が属していた音韻環境と目標の音韻環境の違いによって生じる歪み（音韻環境コスト）などがある。接続コストとしては、音声素片境界でのスペクトルの違い（差）によって生じる歪み（スペクトル接続コスト）や、音声素片境界での基本周波数の違い（差）によって生じる歪み（基本周波数接続コスト）などがある。

【0064】

コストを用いて、一セグメント当たり複数個の音声素片を選択する方法については、どのような方法を用いても構わない。

【0065】

例えば、特許文献 2 に開示された方法を用いても良い。ここでは、この選択方法の概要について、図 5 の処理手順例を参照しながら、一セグメント当たり M 個の音声素片を選ぶ場合について説明する。

【0066】

まず、ステップ S101 において、素片選択部 43 は、入力された音韻系列を、合成単位毎のセグメントに分割する。ここで、分割されたセグメントの数を N とする。

10

20

30

40

50

## 【 0 0 6 7 】

次に、ステップ S 1 0 2 において、素片記憶部 4 2 に記憶されている音声素片群の中から、各セグメントにつき 1 つずつの音声素片の系列を選択する。このときの選択においては、入力された目標の音韻系列・韻律情報と、素片記憶部 4 2 の音声素片環境の情報を基に、系列としてのコストの総和（トータルコスト）が最小となるような音声素片の系列（最適素片系列）を求める。この最適素片系列の探索は、動的計画法（DP(dynamic programming)）を用いることで、効率的に行うことができる。

## 【 0 0 6 8 】

次に、ステップ S 1 0 3 において、セグメント番号を表すカウンター i に、初期値「1」をセットする。

10

## 【 0 0 6 9 】

次に、ステップ S 1 0 4 において、セグメント i に対する複数の音声素片候補の各々に対してコストを算出する。このときに用いるコストには、当該音声素片候補での目標コストと、当該音声素片候補の前後のセグメントの最適音声素片（最適素片系列に含まれる音声素片）と当該音声素片候補との接続コストとの和を用いる。

## 【 0 0 7 0 】

次に、ステップ S 1 0 5 において、ステップ S 1 0 4 で算出したコストを用いて、セグメント i について、コストの小さい上位 M 個の音声素片を選択する。

## 【 0 0 7 1 】

次に、ステップ S 1 0 6 において、カウンター i が N 以下かどうかを判定する。

20

## 【 0 0 7 2 】

カウンター i が N 以下である場合（ステップ S 1 0 6 の YES）には、ステップ S 1 0 7 に進んで、カウンター i の値を 1 つ増やした後に、ステップ S 1 0 4 に進んで、次のセグメントに係る処理を行う。

## 【 0 0 7 3 】

カウンター i が N に達した場合（ステップ S 1 0 6 の NO）には、この素片選択の処理を終了する。

## 【 0 0 7 4 】

このように、素片選択部 4 4 は、各セグメントに対して M 個ずつの音声素片を選択し、選択した音声素片を分離部 4 5 に出力する。

30

## 【 0 0 7 5 】

素片選択部 4 4 においてセグメント当たり複数個の音声素片を選択する方法は、上記した方法に限定する必要はなく、コストであっても、コスト以外であっても、何らかの評価尺度の下で、適切な音声素片の組を選べる方法であれば、いかなる方法を用いても良い。

## 【 0 0 7 6 】

## &lt; 素片融合部 &gt;

素片融合部 4 4 は、それぞれのセグメント毎に、素片選択部 4 3 から入力された複数個の音声素片を融合して、新たな音声素片を生成する。

## 【 0 0 7 7 】

音声素片を融合する方法については、どのような方法を用いても構わない。

40

## 【 0 0 7 8 】

例えば、特許文献 2 に開示された方法を用いても良い。ここでは、この方法について図 6 及び図 7 を参照しながら説明する。

## 【 0 0 7 9 】

図 6 は、一つのセグメントに対する複数個の音声素片の波形を融合して、新たな音声波形を生成する手順を示すフローチャートである。図 7 は、あるセグメントに対して選択された 3 つの音声素片からなる素片組み合わせ候補（図中、6 0）を融合して、新たな音声素片（図中、6 3）を生成する例を示す図である。

## 【 0 0 8 0 】

まず、ステップ S 2 0 1 において、（ある一つのセグメントについて）選択されたそれ

50

それぞれの音声素片からピッチ波形を切り出す。

【 0 0 8 1 】

ここで、「ピッチ波形」とは、その長さが音声の基本周期の数倍程度で、それ自身は基本周期を持たない比較的短い波形であって、そのスペクトルが音声信号のスペクトル包絡を表すものである。

【 0 0 8 2 】

このようなピッチ波形を抽出する方法には、どのような方法が用いられても良いが、その一つの方法として、基本周期同期窓を用いる方法があり、ここでは、この方法が用いられる場合を例にとって説明する。

【 0 0 8 3 】

具体的には、それぞれの音声素片の音声波形に対して基本周期間隔毎にマーク（ピッチマーク）を付し、このピッチマークを中心にして、窓長が基本周期の2倍のハニング窓で窓掛けすることによって、ピッチ波形を切り出す。図7のピッチ波形系列61は、素片組み合わせ候補60の各音声素片から切り出して得られたピッチ波形の系列の例を示している。

【 0 0 8 4 】

次に、ステップS202において、それぞれの音声素片に対するピッチ波形の個数が、音声素片間で同一になるように、ピッチ波形の数を揃える。

【 0 0 8 5 】

このときに、揃える対象となるピッチ波形の数は、目標の音韻継続時間長の合成音声を生成するために必要なピッチ波形数とするが、例えば、最もピッチ波形数の多いものに揃えても良い。

【 0 0 8 6 】

ピッチ波形の少ない系列は、系列に含まれるいくつかのピッチ波形を複製することによってピッチ波形数を増やし、ピッチ波形の多い系列は、系列中のいくつかのピッチ波形を間引くことによってピッチ波形数を減らす。図7のピッチ波形系列62は、ピッチ波形の数を6つに揃えた例を示している。

【 0 0 8 7 】

次に、ステップS203において、ピッチ波形数を揃えた後、それぞれの音声素片に対応するピッチ波形系列中のピッチ波形を、その位置毎に融合することによって、新たなピッチ波形系列を生成する。

【 0 0 8 8 】

例えば、図7で生成された新たなピッチ波形63に含まれるピッチ波形63aは、ピッチ波形系列62のうち、6番目のピッチ波形62a, 62b, 62cを融合することによって得られる。このようにして生成された新たなピッチ波形系列63を、融合された音声素片とする。

【 0 0 8 9 】

ここで、ピッチ波形を融合する方法としては、例えば、次のような方法がある。

【 0 0 9 0 】

第1の方法は、単純にピッチ波形の平均を計算する方法である。

【 0 0 9 1 】

第2の方法は、ピッチ波形間の相関が最大になるよう時間方向に各ピッチ波形の位置を補正してから平均化する方法である。

【 0 0 9 2 】

第3の方法は、ピッチ波形を帯域分割して、帯域毎にピッチ波形間の相関が最大になるようピッチ波形の位置を補正して平均化した結果を、帯域間で足し合わせる方法である。

【 0 0 9 3 】

いずれの方法を用いても良いが、本実施形態では、最後に説明した第3の方法を用いる場合を例にとって説明する。

【 0 0 9 4 】

10

20

30

40

50

素片融合部 4 4 は、上記した方法を用いて、各セグメントについて、複数の音声素片を融合して新たな音声素片を生成し、フォルマント強調フィルタ部 4 5 に出力する。

【 0 0 9 5 】

< フォルマント強調フィルタ部 >

さて、上記のように融合によって生成された音声素片の音声波形は、融合の影響によって、融合元の音声素片の波形よりもスペクトル包絡がなまってしまい、いくつかのフォルマントが弱められてしまった結果、明瞭感が下がってしまうことが多い。そこで、フォルマント強調フィルタ部 4 5 は、素片融合部 4 4 から入力された融合素片に対して、フォルマントを強調するためのフィルタリングを行い、素片編集・接続部 4 7 に出力する。

【 0 0 9 6 】

ここで用いるフォルマント強調フィルタとしては、例えば、J. Chenらの文献(J. Chen, etc., 「Adaptive Postfiltering for Quality Enhancement of Coded Speech」, IEEE Trans. Speech and Audio Processing, vol. 3, Jan 1995) (以下、文献 3 と呼ぶ。) によって開示されているものを、用いることができる。

【 0 0 9 7 】

こうしたフォルマント強調フィルタを、融合素片の音声波形に対して適用することによって、スペクトル包絡中のフォルマントを強調し、融合による明瞭性の低下を補償することが可能である。

【 0 0 9 8 】

フォルマント強調フィルタの概要を、文献 3 で開示されているフォルマント強調フィルタを例に用いて説明する。文献 3 で開示されているフォルマント強調フィルタは、数式 ( 1 ) のような伝達関数を持つフィルタである。

【 数 1 】

$$H(z) = G \frac{1 - P(z/\beta)}{1 - P(z/\alpha)} (1 - \mu z^{-1}), \quad 0 < \beta < \alpha < 1 \quad \dots (1)$$

【 0 0 9 9 】

ただし、 $P(z)$  は、数式 ( 2 ) で表される。ここで、 $a_i$  は入力波形を線形予測分析したときの  $i$  番目の線形予測係数 ( L P C ) を表し、 $M$  は線形予測次数である。

【 数 2 】

$$P(z) = \sum_{i=1}^M a_i z^{-i} \quad \dots (2)$$

【 0 1 0 0 】

数式 ( 1 ) における  $1 / [ 1 - P(z/\alpha) ]$  は、 $\alpha = 1$  の場合は、線形予測フィルタを表し、入力波形の L P C スペクトルと同じ周波数応答を持つ。 $\alpha$  を小さくすると、L P C スペクトルを鈍らせたような周波数応答になり、0 に近づくにつれ、フラットな周波数応答になる。よって、入力波形のスペクトル中のパワーの大きい周波数成分は、より大きくなり、パワーの小さい周波数成分は、より小さくなるため、スペクトル中の山・谷を強調する効果を持つ。また、一般的な音声のスペクトル包絡には、低域から高域に向かって負の傾斜が見られるため、 $1 / [ 1 - P(z/\alpha) ]$  の周波数応答は、全体的に、同様の負の傾斜を持つ。すなわち、スペクトルの山・谷を強調する効果に加え、副作用としてローパス特性を持っている。そこで、 $[ 1 - P(z/\alpha) ]$  および  $[ 1 - \mu z^{-1} ]$  の項によって、このローパス特性を補正する。 $[ 1 - P(z/\alpha) ]$  は、L P C スペクトルの極と同じ周波数に零点を持つフィルタであり、 $1 / [ 1 - P(z/\alpha) ]$  でのスペクトルの傾斜を補償する効果を持つ。一方、 $[ 1 - \mu z^{-1} ]$  は、単純なハイパスフィルタで、残っているスペクトルの傾きを無くすよう調整するための項である。なお、 $G$  は、フィルタリング前後でパワーが変化するのを防ぐためのパワー調整用のゲインであり、文献 3 で開示されている方法により、入力波形に応じて自動で決めることができる。

10

20

30

40

50

## 【 0 1 0 1 】

このフォルマント強調フィルタでは、パラメータ  $\mu$  を変えることによって、フォルマント強調の度合いを変えることができる（ただし、 $\mu$  の値に応じて、ローパス特性を補償するような適切な  $\mu$  も決める必要がある）。 $\mu$  が 1 に近いほど強調の度合いが強く、 $\mu$  が小さくなるにつれ強調の度合いが弱まり、 $\mu$  が 0.5 以下になるとほとんど強調されない。どの程度フォルマントを強調すべきかは音声波形の特徴によって異なるが、これを決めるための客観尺度が存在しないため、通常、音声符号化や音声合成においてフォルマント強調フィルタを用いる場合には、フォルマント強調の度合いは主観評価などによって実験的に求める。

## 【 0 1 0 2 】

しかしながら、複数素片選択融合型の合成方法においては、融合によるスペクトル包絡の鈍り具合がセグメントごとに大きく変わり得るため、1 文など全体に対して同じパラメータを適用すると、融合によってスペクトル包絡が大きく鈍った箇所にはフォルマント強調の程度が不十分であるのに対し、逆に融合によるスペクトル包絡の鈍りが小さい箇所はフォルマントが強調されすぎて人工的な音になるという問題がある。

## 【 0 1 0 3 】

そこで、本実施形態では、融合されてできたそれぞれの音声素片に対し（あるいは、それぞれの融合素片の各ピッチ波形に対し）、適切なフォルマント強調の度合いをフォルマント強調度合い推定部 46 で推定し、フォルマント強調フィルタ部 45 は、推定されたフォルマント強調度合いに応じてフォルマント強調フィルタの係数を変える。すなわち、融合素片ごとに（あるいは、融合素片のピッチ波形ごとに）、フォルマント強調度合いを適応的に制御する。ここで、フォルマント強調度合い推定部 46 から与えられるフォルマント強調度合いは、例えば、0（強調無し）から 1.00（フォルマント強調フィルタの制御可能な範囲で最も強い強調）まで連続的に変化するようなものでもよいし、また、例えば、0（強調無し）から 4（非常に強く強調）までの 5 段階で指定できるような離散的なものであってもよい。上述の文献 3 で開示されているフォルマント強調フィルタを用いる場合は、フォルマント強調度合い推定部 46 で推定されたフォルマント強調度合いが大きい場合は  $\mu$  の値を 1 に近づけ、逆にフォルマント強調度合いが小さい場合は  $\mu$  を 0.5 に近づける。おおよび  $\mu$  の値も  $\mu$  の値に応じて変えるが、各  $\mu$  の値に対して適切な  $\mu$  と  $\mu$  の値は、実験的に求めることが可能である。

## 【 0 1 0 4 】

また、フォルマント強調度合い推定部 46 で推定されたフォルマント強調度合いを、フィルタ係数に具体的に反映するためのマッピングは、主観評価などによって実験的に得ることができる。

## 【 0 1 0 5 】

本実施形態においては、文献 3 で開示されているフォルマント強調フィルタを用いる場合について説明したが、フォルマントが強調でき、フォルマントの強調度合いがパラメータなどで制御できるフォルマント強調フィルタであれば、いかなるものでも用いることができる。

## 【 0 1 0 6 】

< フォルマント強調度合い推定部 >

フォルマント強調度合い推定部 46 は、素片選択部 43 や素片融合部 44 から与えられた融合素片や融合元の複数の音声素片の情報を元に、融合素片に対して適切なフォルマント強調度合いを推定し、推定したフォルマント強調度合いをフォルマント強調フィルタ部 45 に出力する。

## 【 0 1 0 7 】

前述のように、ある波形に対して適切なフォルマント強調度合いを決めるような客観尺度は存在しないが、融合素片と融合元の複数の音声素片の間でスペクトル包絡に関する特徴量を比較することによって、音声素片の融合によってどの程度スペクトル包絡が鈍ったかをある程度見積もることは可能である。そこで、フォルマント強調度合い推定部 46 で

10

20

30

40

50

は、融合によるスペクトル包絡の鈍り具合を以下のような方法で推定し、これに基づいてフォルマントの強調度合いを決める。

【 0 1 0 8 】

融合によるスペクトル包絡の鈍りが大きいほど、融合元の各音声素片と融合素片との間でスペクトル包絡の形状の差が大きくなると考えられる。そこで、融合元の各音声素片と融合素片との間でのスペクトル包絡の形状の差を見積もることができれば、音声素片の融合によるスペクトル包絡の鈍り具合を推定できると考えられる。

【 0 1 0 9 】

スペクトル包絡の特徴を表すパラメータとしては、ケプストラムや L S P (線スペクトル対) などがある。以下では、ケプストラムの一つであるメル周波数ケプストラム係数 (M F C C) を用いて、融合元の各音声素片と融合素片の間でのスペクトル包絡の形状の差を間接的に見積もる場合を例によって説明する。

【 0 1 1 0 】

M F C C は、音声認識の分野で広く用いられている音響特徴量で、音声合成においても上述の「スペクトル接続コスト」の評価尺度としてよく用いられる。M F C C は、人間の聴覚特性を考慮した特徴量で、低い次元でもスペクトル包絡の特徴を良く表せる利点も持つ。M F C C の低次の係数はスペクトル包絡の概形を、高次の係数はスペクトル包絡の細部を表現する。素片 1 と素片 2 の  $i$  次の M F C C をそれぞれ  $c_{1i}$ 、 $c_{2i}$  とすると、数式 (3) により、素片 1 と素片 2 との間の M F C C 距離が算出できる。

【 数 3 】

$$D_{MFCC} = \sum_{i=1}^p (c_{1i} - c_{2i})^2 \quad \dots (3)$$

【 0 1 1 1 】

ただし、 $p$  は M F C C の次元を表す。

【 0 1 1 2 】

なお、本例においては、M F C C の次元は 20 次程度とする。

【 0 1 1 3 】

次に、この M F C C 距離を使って、音声素片の融合によるスペクトル包絡の鈍り具合を推定する方法について説明する。

【 0 1 1 4 】

図 8 に、この場合の処理手順の一例を示す。

【 0 1 1 5 】

ここで、融合素片の元になった融合元の素片数は  $N$  とする。

【 0 1 1 6 】

まず、融合素片の M F C C ( $c_0$ ) を算出する (ステップ S 3 0 1)。

【 0 1 1 7 】

次に、カウンター  $i$  を 1 に、 $D_{sum}$  を 0 に初期化して (ステップ S 3 0 2、ステップ S 3 0 3)、ステップ S 3 0 4 に進む。

【 0 1 1 8 】

ステップ S 3 0 4 では、融合元の  $N$  個の音声素片のうち、 $i$  番目の音声素片の M F C C ( $c_i$ ) を算出する。

【 0 1 1 9 】

次に、 $c_0$  と  $c_i$  との間の M F C C 距離 ( $D_i$ ) を、数式 (3) を用いて算出する (ステップ S 3 0 4)。

【 0 1 2 0 】

次のステップ S 3 0 5 では、算出された  $D_i$  を  $D_{sum}$  に加算して、ステップ S 3 0 7 に進む。

【 0 1 2 1 】

ステップ S 3 0 7 では、カウンタ  $i$  が  $N$  以下であるかを判定する。

【 0 1 2 2 】

カウンタ  $i$  が  $N$  以下である場合（ステップ S 3 0 7 の Y E S）には、ステップ S 3 0 8 に進んで、カウンタ  $i$  の値を 1 つ増やした後に、ステップ S 3 0 4 に進んで、次の音声素片に係る処理を行う。

【 0 1 2 3 】

カウンタ  $i$  が  $N$  に達した場合（ステップ S 3 0 7 の N O）には、ループ処理を終了し、ステップ S 3 0 9 に進む。

【 0 1 2 4 】

ステップ S 3 0 9 では、 $D_{sum}$  を  $N$  で割ることによって、平均 M F C C 距離（ $D_{mean}$ ）を求め、全ての処理を終了する。

【 0 1 2 5 】

本実施形態では、このようにして求めた平均 M F C C 距離を、融合によるスペクトル包絡の鈍り具合を反映する評価尺度として用いる。すなわち、平均 M F C C 距離が小さいほどスペクトル包絡の鈍り具合が小さく、平均 M F C C 距離が大きいほどスペクトル包絡の鈍り具合が大きいとして、平均 M F C C 距離をそのままスペクトル包絡の鈍り具合とするか、平均 M F C C 距離の分布などに基づいて何らかの変換を行って得た値をスペクトル包絡の鈍り具合とする。

【 0 1 2 6 】

次に、このようにして得たスペクトル包絡の鈍り具合に基づいて、フォルマント強調度合いを求める必要があるが、スペクトル包絡の鈍り具合が大きいほど強いフォルマント強調を施すべきと考えられるため、ここでは、スペクトル包絡の鈍り具合が増すとともに単調増加するような関数（ただし、フォルマント強調度合いが離散値の場合は、階段状に変化）を用いてフォルマント強調度合いに変換する。関数の形状については、例えば、スペクトル包絡の鈍り具合に対して途中まで線形に増加し、ある閾値を超えるとフォルマント強調度合いの上限値をとるようなものであっても良いし、シグモイド関数のように増加率がスペクトル包絡の鈍り具合に応じて変化するような形状のものであっても良く、それらの関数のパラメータ（傾き、など）は実験的に適切なものを得れば良い。

【 0 1 2 7 】

なお、本実施形態においては、融合によるスペクトル包絡の鈍り具合を推定する方法の一例として、上記の M F C C を用いる方法を例にとりて説明したが、スペクトル包絡の形状の差を適切に見積もれる音響パラメータであれば、どのようなものを用いてもよい。例えば、L S P 係数の二乗誤差を用いても良いし、F F T（高速フーリエ変換）によって得られた F F T スペクトルを確率分布のように見なすことによって、確率分布の差を計算するのによく用いられる K L 距離（Kullback-Leibler 距離）を算出して、これを用いても良い。

【 0 1 2 8 】

また、融合によるスペクトル包絡の鈍り具合を推定する方法として、素片選択部 4 3 で算出された目標コストを用いる方法も考えられる。融合元の複数の音声素片がいずれも適切な音韻・韻律環境から選ばれた場合、目標コストは小さくなり、かつ、融合によるスペクトル包絡の鈍り具合も小さくなると考えられる。逆に、目標の音韻・韻律環境と異なる音声素片ばかりが選ばれた場合、目標コストは大きくなり、融合によるスペクトル包絡の鈍り具合も大きくなると考えられる。そこで、融合によるスペクトル包絡の鈍り具合を表す一つの指標として、融合元の音声素片が選ばれた際の目標コストを用いてもよいと考えられる。この方法は、前述の音響パラメータを用いる方法よりは間接的だが、非常に単純である。

【 0 1 2 9 】

フォルマント強調度合い推定部 4 6 は、上述のようにして推定した、融合によるスペクトル包絡の鈍り具合を、フォルマント強調フィルタ部 4 5 に出力する。

【 0 1 3 0 】

10

20

30

40

50

### < 素片編集・接続部 >

素片編集・接続部 47 は、フォルマント強調部 45 から渡されたセグメント毎の音声素片を、入力韻律情報に従って変形して接続することによって、合成音声の音声波形を生成する。

#### 【0131】

図 9 は、素片編集・接続部 47 での処理を説明するための図である。図 9 には、フォルマント強調部 45 から入力された、音素「a」「N」「s」「a」「a」の各合成単位に対する音声素片を、変形・接続して、「aNsaa」という音声波形を生成する場合を示している。

#### 【0132】

この例では、有声音の音声素片はピッチ波形の系列で表現されている。一方、無声音の音声素片は、フレーム毎の波形として表現されている。

#### 【0133】

図 9 の点線は、目標の音韻継続時間長に従って分割した音素毎のセグメントの境界を表し、白い三角は、目標の基本周波数に従って配置した各ピッチ波形を重畳する位置（ピッチマーク）を示している。

#### 【0134】

図 9 のように、有声音については音声素片のそれぞれのピッチ波形を対応するピッチマーク上の重畳し、無声音については各フレームの波形をセグメント中の各フレームに対応する部分に貼り付けることによって、所望の韻律（ここでは、基本周波数、音韻継続時間長）を持った音声波形を生成する。

#### 【0135】

以上のように本実施形態によれば、素片融合によるフォルマントの鈍り具合に応じて、セグメントごとに適切な強さのフォルマント強調を行うので、こもり感やブザー感が少なく、かつ人工的でない高音質な合成音声を生成できる。

#### 【0136】

##### （第 2 の実施形態）

本発明の第 2 の実施形態に係るテキスト音声合成を行うテキスト音声合成装置（音声処理装置）について説明する。

#### 【0137】

第 1 の実施形態では、音声素片の融合処理およびフォルマント強調の処理に大きな計算量を要するため、CPU スペックが比較的低いミドルウェア向けの応用などには適用が向かないこともあり得る。

#### 【0138】

そこで、本実施形態では、音声素片の融合およびフォルマント強調の処理を予め行った音声素片をオフラインで作成しておき、実際の動作時には、こうして作成された音声素片から適切な音声素片を選択して接続するだけの処理で合成波形を生成する。

#### 【0139】

本実施形態に係るテキスト音声合成装置の全体構成例は、図 1 と同様であり、テキスト入力部 1、言語処理部 2、韻律処理部 3、音声合成部 4 を備えている。

#### 【0140】

図 10 に、本実施形態の音声合成部 4 の構成例を示す。

#### 【0141】

以下、図 10 を参照しながら、本実施形態について、第 1 の実施形態と相違する点を中心に説明する。

#### 【0142】

図 10 に示されるように、本実施形態の音声合成部 4 は、情報入力部 41、素片記憶部 42、素片選択部 43、素片編集・接続部 47、音声波形出力部 48 を備えている。

#### 【0143】

第 1 の実施形態（図 2）と比較すると、本実施形態の音声合成部 4 は、図 2 の素片融合

10

20

30

40

50



部 4 4、フォルマント強調フィルタ部 4 5、フォルマント強調度合い推定部 4 6 が省かれている。

【 0 1 4 4 】

また、本実施形態の素片記憶部 4 2 には、後述の方法によって生成された融合済みの音声素片が格納されている。

【 0 1 4 5 】

第 1 の実施形態の素片選択部 4 4 が各セグメントに対して複数個ずつの音声素片を選択するのに対し、本実施形態の素片選択部 4 4 は、各セグメントに対して 1 つずつの融合済み音声素片の最適系列を選択する。

【 0 1 4 6 】

素片選択部 4 4 の動作としては、例えば第 1 の実施形態で図 5 のフローチャートを用いる場合と比較すると、本実施形態では、図 5 のフローチャートのうち、ステップ S 1 0 1 とステップ S 1 0 2 だけを実行すればよい。もちろん、各セグメントに対して 1 つずつの融合済み音声素片の最適系列を選択する方法は、これに限られるものではなく、種々の方法が可能である。

【 0 1 4 7 】

なお、素片編集・接続部 4 7 および音声波形出力部 4 8 の動作は、第 1 の実施形態のものと同様である。

【 0 1 4 8 】

次に、音声素片記憶部 4 2 に格納する融合済みの音声素片を学習する方法について、図 1 1 及び図 1 2 を参照しながら説明する。

【 0 1 4 9 】

本実施形態では、融合済みの音声素片を作成する融合済み音声素片作成部 5 を用いる。融合済み音声素片作成部 5 は、図 1 0 のテキスト音声合成装置に含まれても良い。この場合、テキスト音声合成に供するための「フォルマント強調された融合素片」の生成時には、図 1 の音声合成部 4 を融合済み音声素片作成部 5 に置き換えた構成で用いれば良い。

【 0 1 5 0 】

また、融合済み音声素片作成部 5 は、テキスト音声合成装置に含まれなくても良い。この場合、例えば、融合済み音声素片作成部 5 を、独立した音声処理装置（テキスト音声合成に供するための「フォルマント強調された融合素片」を生成する音声処理装置）として構成しても良い。この場合、独立した音声処理装置は、図 1 の音声合成部 4 を融合済み音声素片作成部 5 に置き換えた構成にすれば良い。

【 0 1 5 1 】

図 1 1 に、融合済み音声素片作成部 5 の構成例を示す。

【 0 1 5 2 】

融合済み音声素片作成部 5 の構成は、第 1 の実施形態の音声合成部 4 の構成とほとんど同じであるため、ここでは相違する点について説明する。

【 0 1 5 3 】

融合済み音声素片作成部 5 は、第 1 の実施形態の音声合成部 4 の素片編集・接続部 4 7 および音声波形出力部 4 8 の代わりに、音声素片出力部 4 9 を持つ。第 1 の実施形態の音声合成部 4 の素片編集・接続部 4 7 および音声波形出力部 4 8 は、フォルマント強調部 4 5 から入力された各セグメントに対する音声素片を接続して、入力テキストに対する合成波形を生成するのに対し、音声素片出力部 4 9 は、フォルマント強調部 4 5 から入力された音声素片をそのまま出力する。

【 0 1 5 4 】

すなわち、融合済み音声素片作成部 5 は、音声素片（フォルマント強調された融合素片）を、図 1 0 のテキスト音声合成装置の音声素片記憶部 4 2 へ出力し、音声素片（フォルマント強調された融合素片）は、音声素片記憶部 4 2 に記憶される。

【 0 1 5 5 】

次に、音声素片記憶部 4 2 に格納する融合済みの音声素片を学習する手順について説明

10

20

30

40

50

する。

【0156】

図12に、この場合の処理手順の一例を示す。

【0157】

まず、ステップS501において、融合済み音声素片作成部5を備えたテキスト音声合成装置又は独立した音声処理装置に対して、大量の文を入力する。

【0158】

次に、ステップS502において、入力された各文の各セグメントに対して生成された融合済み音声素片が、融合済み音声素片生成部5から出力される。

【0159】

次に、ステップS503において、外部から指定された音声素片記憶部42に格納する音声素片の総数のうち、それぞれの素片種別に対して幾つずつ配分するかを決める。

【0160】

ここで、素片種別とは、音声素片の音韻環境などで分類された種別を指す。例えば、素片種別/a/は、音素/a/に対応する音声素片のこととする。

【0161】

各素片種別に何個ずつ素片を配分するかは、各素片種別の音声素片の出現頻度などに応じて決める。例えば、素片種別/a/の素片が素片種別/u/の素片よりも出現頻度が高い場合は、素片種別/a/に多めの素片を配分することとする。

【0162】

素片種別*i*に配分する音声素片の個数を $N_i$ とする。

【0163】

次に、ステップS504において、素片種別番号*i*に初期値1をセットする。

【0164】

次に、ステップS505において、素片種別*i*の融合済み音声素片を、ステップS502で出力された素片種別*i*の音声素片の中から、出現頻度が上位のものを $N_i$ ずつ抽出する。

【0165】

次に、ステップS506において、*i*と素片種別数を比較する。

【0166】

*i*が素片種別数以下である場合(ステップS506のYES)には、ステップS507に進んで、*i*の値を1つ増やし、そして、ステップS505～ステップS506を繰り返す。

【0167】

*i*が素片種別数を超過している場合(すなわち、全ての素片種別に対する処理が完了している場合)(ステップS506のNO)には、全ての処理を終了する。

【0168】

上記のようにして抽出した融合済み音声素片を、音声素片記憶部42に格納する。

【0169】

ここで、音声素片記憶部42に格納するために選択する音声素片の個数は、トータルでの音声素片サイズと合成音声の音質とのトレードオフで、任意に決めることができる。より多くの音声素片を選択して格納すれば、サイズは大きくなるが、合成音声の音質を高くすることができ、音声素片の数を減らせば、合成音声の音質は犠牲になるが、サイズを小さくすることができる。

【0170】

なお、上記では出現頻度の高い素片を抽出する方法を説明したが、音声素片の両端で算出したメルケプストラムなどの音声素片の特徴量を用いて抽出しても良い。

【0171】

この場合、各素片種別に対して出力された融合済み音声素片をそれぞれ、音声素片の特徴量を用いてクラスタリングし、分割された各クラスタの中心(セントロイド)に最も近

10

20

30

40

50

い素片を抽出する。クラスタリングにおけるクラスタ数は、各素片種別に配分する素片数に応じて決める。

【0172】

出現頻度に基づいて素片を抽出する場合は、出現頻度が低いコンテキストに対して適切な素片が抽出されない可能性があり、入力テキストによっては音質が大きく劣化してしまう可能性があるが、本方法によって素片を抽出した場合、特徴量空間をできるだけ広く覆うような音声素片のセットが抽出できるため、出現頻度に基づいて抽出した場合より安定した合成音が生成できる。

【0173】

以上のように本実施形態によれば、複数の音声素片を融合する処理とフォルマント強調の処理を予めオフラインで行うので、第1の実施形態よりも少ない計算量で実現でき、CPUスペックが比較的低いミドルウェア向けなどの応用にも適用可能である。

【0174】

また、合成音声の音質とのトレードオフで、格納する音声素片のトータルのサイズもスケラブルに決めることができる。

【0175】

(第3の実施形態)

本発明の第3の実施形態に係るテキスト音声合成装置について説明する。

【0176】

本実施形態は、フォルマント強調度合い推定部46の推定方法が、第1の実施形態で説明した例とは相違するものであり、以下、この相違点を中心に説明する。

【0177】

第1の実施形態では、フォルマント強調度合い推定部46でフォルマント強調度合いを推定する方法として、融合元の各音声素片と融合素片の間でのスペクトル包絡の差を算出することによって推定する方法を説明したが、融合元の各音声素片と融合素片との間でのスペクトル包絡の差と、融合によるスペクトル包絡の鈍り具合の間には、高い相関はあると考えられるものの、直接的な関係があるわけではない。そこで、スペクトル包絡の鈍り具合を、より直接的に求められる方法があれば、より確度の高い推定を行うことが可能と考えられる。

【0178】

その一つの方法として、線形予測極(LP極)を用いる方法が考えられる。LP極は、数式(2)の $P(z)$ について $(1 - P(z)) = 0$ とおいたときに得られる解(複素数)のことで、この解の $z$ 平面上での位置と単位円との関係から、各フォルマントの周波数とバンド幅を推定することができる。それぞれの極が各フォルマントに対応すると考えられ、 $i$ 番目の極に関して、極と原点を結ぶ線の角度を $\theta_i$ 、極と原点の距離を $r_i$ とした場合、 $i$ 番目のフォルマントの周波数 $F_i$ およびバンド幅 $BW_i$ は、数式(4)のように推定できる。

【数4】

$$F_i = \frac{\theta_i}{2\pi T_s}, BW_i = \frac{-\ln(r_i)}{\pi T_s} \quad \dots (4)$$

【0179】

このようにして推定した各フォルマントの周波数とバンド幅を用いれば、スペクトル包絡のうち、特にフォルマントに関する鈍り具合がより正確に推定できると考えられる。

【0180】

以下、LP極から推定される各フォルマントのバンド幅を用いて、スペクトル包絡の鈍り具合を推定する方法の一例を、図13を参照しながら説明する。

【0181】

図 13 に、LP 極から推定される各フォルマントのバンド幅を用いてスペクトル包絡の鈍り具合を推定する手順の一例を示す。

【0182】

まず、ステップ S601 において、融合素片の LP 極を算出する。具体的には、融合素片の音声波形に対して LPC 分析を行い、得られた線形予測係数を係数に持つ数式 (2) の  $P(z)$  について、 $(1 - P(z)) = 0$  とおいたときの解を得る。

【0183】

次のステップ S602 では、融合元の音声素片それぞれに対する LP 極を、ステップ S601 と同様の方法で算出する。

【0184】

次に、ステップ S603 では、フォルマントバンド幅比率の和  $R_{sum}$  を 0 に、ステップ S604 では、用いた LP 極の個数  $N_{LP}$  を 0 に、ステップ S605 では、カウンター  $i$  を 1 に、それぞれ初期化して、ステップ S606 に進む。

【0185】

ステップ S606 では、融合素片の  $i$  番目の LP 極が実軸上 (すなわち虚数項が 0) かどうかを判定し、実軸上である場合 (ステップ S506 の YES) には、ステップ S620 に進んで、カウンター  $i$  の値を 1 つ増やした後に、再び S606 に進む。

【0186】

これは、実軸上の LP 極がフォルマントには対応しない (スペクトル包絡全体の形状に寄与) ため、実軸上である場合については、ステップ S607 以降の処理をスキップし、フォルマントに対応した LP 極のみを考慮するためのものである。

【0187】

LP 極が実軸上でない場合 (ステップ S606 の NO) には、ステップ S607 に進む。

【0188】

ステップ S607 では、 $N_{LP}$  の値を 1 つ増やした後に、ステップ S608 に進む。

【0189】

ステップ S608 では、融合素片の  $i$  番目の LP 極に対するフォルマントのバンド幅  $B_{Wi}$  を、数式 (4) を用いて算出する。

【0190】

次のステップ S609 では、融合元の音声素片のフォルマントに関するバンド幅の和  $B_{W_{iorg\_sum}}$  を 0 に初期化し、ステップ S610 に進む。

【0191】

ステップ S610 では、カウンター  $k$  を 1 に初期化して、ステップ S611 に進む。

【0192】

ステップ S611 では、融合元の音声素片 (計  $N_{fused}$  個) のうち  $k$  番目の音声素片 (「音声素片  $k$ 」と呼ぶ。) について、この音声素片の LP 極の中で、融合素片の  $i$  番目の LP 極が表すフォルマントに対応するような LP 極を選択する。具体的には、音声素片  $k$  の LP 極の中で、融合素片の  $i$  番目の LP 極に最も近いものを選択する。LP 極の間の距離については、例えば数式 (5) (文献 “Goncharoff, etc., 「Interplation of LP C spectra via pole shifting.」, IEEE ICASSP, Detroit, MI, Vol.1, pp.780-783, 1995” 参照) を用いて算出できる。ただし、 $p_i$  は LP 極の複素数表現、 $r_i$  は LP 極と原点の距離を表し、 $D(p_0, p_1)$  が LP 極  $p_0$  と  $p_1$  の距離を表す。

10

20

30

40

【数 5】

$$D(p_0, p_1) = \begin{cases} \left| \ln\left(\frac{p_1}{p_0}\right) \right| \left\{ \frac{\ln((1-r_0^2)/(1-r_1^2))}{\ln(r_1/r_0)} \right\}, & r_0 \neq r_1 \\ \left| \ln\left(\frac{p_1}{p_0}\right) \right| \{2r^2/(1-r^2)\}, & r = r_0 = r_1 \end{cases} \quad \dots (5)$$

10

【0193】

この数式(5)によって、融合素片の*i*番目のLP極との距離を、融合元の音声素片のLP極のそれぞれについて算出し、最も距離が小さいLP極を選択する。

【0194】

次のステップS612では、ステップS611で選択されたLP極に対するバンド幅  $BW_{i\_org\_k}$  を、数式(4)を用いて算出する。

【0195】

次に、ステップS613において、ステップS612で算出した  $BW_{i\_org\_k}$  を  $BW_{i\_org\_sum}$  に加算する。

【0196】

続いて、ステップS613において、カウンター*k*が融合元の音声素片数  $N_{fuse\_d}$  以下かどうかを判定する。

20

【0197】

カウンター*k*が  $N_{fuse\_d}$  以下である場合(ステップS613のYES)には、ステップS619に進んで、カウンター*k*の値を1つ増やした後に、ステップS611からのステップを繰り返す。一方、カウンター*k*が  $N_{fuse\_d}$  を超える場合(ステップS613のNO)には、ステップS615に進む。

【0198】

ステップS615では、 $BW_{i\_org\_sum}$  を  $N_{fuse\_d}$  で割ることによって、融合素片の*i*番目のLP極に対応するような、融合元の各音声素片のLP極についての、フォルマントのバンド幅の平均値  $BW_{i\_org\_mean}$  を算出する。

30

【0199】

次のステップS616では、ステップS615で算出した  $BW_{i\_org\_mean}$  に対する、融合素片の*i*番目のLP極のバンド幅  $BW_i$  の比率を、フォルマントバンド幅比率の和  $R_{sum}$  に加算する。

【0200】

続いて、ステップS617では、カウンター*i*が、 $N_{max\_LP}$  という設定値以下かどうかを判定する。

【0201】

ここで、 $N_{max\_LP}$  は、フォルマントの鈍り具合を推定するのに用いるLP極の個数の最大値を表す。

40

【0202】

この値は、例えば、LPC分析での分析次数の1/2などに設定する。

【0203】

カウンター*i*が  $N_{max\_LP}$  以下である場合(ステップS617のYES)には、ステップS620に進んで、カウンター*i*の値を1つ増やした後に、ステップS606からの処理を繰り返す。一方、カウンター*i*が  $N_{max\_LP}$  を越える場合(ステップS617のNO)には、ステップS618に進む。

【0204】

ステップS618では、フォルマントバンド幅比率の和  $R_{sum}$  を、用いたNP極の個

50

数  $N_{LP}$  で割ることによって、フォルマントバンド幅比率の平均値  $R_{mean}$  を算出し、全ての処理を終了する。

【0205】

本実施形態では、上記のような方法で算出したフォルマントバンド幅比率の平均値  $R_{mean}$  を、音声素片の融合によるスペクトル包絡の鈍り具合を表す尺度として用いる。この値は、フォルマントのバンド幅がほぼ変わらずスペクトル包絡がほとんど鈍らなかった場合には 1 に近い値、フォルマントのバンド幅が融合元の音声素片より広がってスペクトル包絡が鈍った場合には 1 より大きい値となり、スペクトル包絡の鈍り具合が大きければ大きいほど大きな値になると考えられる。そこで、本実施形態においては、 $R_{mean}$  が 1 以下の場合には強調無しで、1 より大きい場合は強調度合いが単調増加するような何らかの関数を用いることによって、この  $R_{mean}$  からフォルマント強調度合いを算出することとする。

10

【0206】

このように、融合素片と融合元の各音声素片に対して推定されたフォルマントのバンド幅を用いることによって、スペクトル包絡形状の差を用いる場合よりも、融合によるスペクトル包絡の鈍り具合をより直接的に求められるので、フォルマント強調度合いをより高い確度で推定することが可能である。

【0207】

(第4の実施形態)

本発明の第4の実施形態に係るテキスト音声合成装置について説明する。

20

【0208】

本実施形態は、フォルマント強調度合い推定部46の推定方法が、第1、第3の実施形態で説明した例とは相違するものであり、以下、この相違点を中心に説明する。

【0209】

第3の実施形態においては、融合素片の各LP極に対して求めたフォルマントバンド幅比率を平均化することによって、スペクトル包絡全体での鈍り具合を推定しているが、実際にはスペクトルの鈍り具合がフォルマント毎で異なる場合も考えられる。そこで、各LP極に対して求めたフォルマントバンド幅比率(以下、 $R_i$ とする。)をそのまま用いることによって、フォルマントごとに強調度合いが異なるようなフォルマント強調を行うことも可能である。

30

【0210】

ここで、融合素片の*i*番目のLP極を  $p_i$  とすると、数式(2)の  $P(z)$  に関して、数式(6)のように表せる。

【数6】

$$1 - P(z) = \prod_{i=1}^M (1 - p_i z^{-i}) \quad \dots (6)$$

【0211】

$H(z) = 1 / (1 - P(z))$  という伝達関数を持つフィルタ(線形予測フィルタ)に、LPC分析したときの予測残差を入力すると完全に元の波形が再現できるが、上記の  $p_i$  を  $z$  平面上の単位円に近づくように変更したフィルタに予測残差を入力すると、*i* 番目のLP極に対応するフォルマントのバンド幅が狭まり、結果的に、このフォルマントを強調することができる。すなわち、 $R_i$  に応じて適切に  $p_i$  を変更したフィルタをフォルマント強調フィルタとして用いれば、フォルマントごとに適切なフォルマント強調を行うことができる。

40

【0212】

図14に、本実施形態のフォルマント強調フィルタ部45の構成例を示す。

【0213】

50

L P C 分析部 4 5 1 は、入力された波形に対して L P C 分析を行い、算出された L P C を L P C 変形部 4 5 2 に、予測残差を線形予測フィルタ部 4 5 3 に出力する。

【 0 2 1 4 】

L P C 変形部 4 5 2 では、フォルマント強調度合い推定部 4 6 から入力された各 L P 極に対するフォルマントバンド幅比率  $R_i$  に応じて L P C 係数を変形し、この変形された L P C 係数を線形予測フィルタ部 4 5 3 に与える。

【 0 2 1 5 】

線形予測フィルタ部 4 5 3 では、L P C 変形部 4 5 2 から与えられた L P C 係数をフィルタ係数に用いて、L P C 分析部 4 5 1 から入力された予測残差をフィルタリングすることによって、フォルマント強調された波形を出力する。

10

【 0 2 1 6 】

なお、L P C 変形部 4 5 2 においては、まず、入力された L P C 係数から数式  $P(z)$  を得た後、 $(1 - P(z))$  を数式 (6) のように因数分解することによって L P 極  $p_i$  を得る。

【 0 2 1 7 】

次に、L P 極  $p_i$  を  $R_i$  に応じて変更する。

【 0 2 1 8 】

例えば、数式 (7) のように変更すれば、フォルマントのバンド幅は  $1 / R_i$  倍となり、融合元の音声素片での平均的なフォルマントのバンド幅に近づくようバンド幅を狭めることが可能である。

20

【 数 7 】

$$p_i = |p_i|^{\frac{1}{R_i}-1} p_i \quad \dots (7)$$

【 0 2 1 9 】

このような方法で  $R_i$  に応じて変更した L P 極  $p_i$  を数式 (6) に代入して、この数式を展開することによって、変形された L P C 係数を得ることができる。

【 0 2 2 0 】

本実施形態においては、融合素片および融合元の各音声素片に対して求めた L P 極を用いてフォルマントごとに強調度合いを変える方法を説明したが、この方法以外にも、フォルマントごとあるいは周波数帯域によって強調度合いが変わるようなフォルマント強調を行うことも可能である。

30

【 0 2 2 1 】

例えば、フォルマント強調度合い推定部 4 6 において、融合素片および融合元の各音声素片の波形を複数の周波数帯域に分割し、それぞれの帯域においてスペクトル包絡の鈍り具合を推定することによって、それぞれの帯域でのフォルマント強調度合いを推定する。そして、フォルマント強調フィルタ部 4 5 において、融合素片の波形を帯域分割して得た各周波数帯域の波形に対し、フォルマント強調度合い推定部 4 6 から入力された各帯域の強調度合いに従ってフォルマント強調した後、周波数帯域間で波形を足し合わせれば、各周波数帯域でのスペクトル包絡の鈍り具合に応じたスペクトル強調を行うことが可能である。

40

【 0 2 2 2 】

なお、以上の各機能は、ソフトウェアとして記述し適当な機構をもったコンピュータに処理させても実現可能である。

また、本実施形態は、コンピュータに所定の手順を実行させるための、あるいはコンピュータを所定の手段として機能させるための、あるいはコンピュータに所定の機能を実現させるためのプログラムとして実施することもできる。加えて該プログラムを記録したコンピュータ読取り可能な記録媒体として実施することもできる。

【 0 2 2 3 】

50

なお、本発明は上記実施形態そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化できる。また、上記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。

【図面の簡単な説明】

【0224】

【図1】本発明の一実施形態に係るテキスト音声合成装置の構成例を示すブロック図

【図2】同実施形態に係る音声合成部の構成例を示すブロック図

【図3】同実施形態に係る音声素片記憶部に蓄積される音声素片の例を示す図

10

【図4】同実施形態に係る音声素片記憶部に蓄積される素片属性情報の例を示す図

【図5】音声素片の選択手順の一例を示すフローチャート

【図6】音声波形を融合して新たな音声波形を生成する手順の一例を示すフローチャート

【図7】選択された3つの音声素片からなる素片組み合わせ候補を融合して新たな音声素片を生成する例について説明するための図

【図8】音声素片の融合によるスペクトル包絡の鈍り具合を推定する手順の一例を示すフローチャート

【図9】同実施形態に係る素片編集・接続部での処理を説明するための図

【図10】同実施形態に係る音声合成部の他の構成例を示すブロック図

【図11】同実施形態に係る融合済み音声素片作成部の構成例を示すブロック図

20

【図12】融合済みの音声素片を学習する手順の一例を示すフローチャート

【図13】フォルマント強調度合いを推定する手順の一例を示すフローチャート

【図14】同実施形態に係るフォルマント強調フィルタ部の構成例を示すブロック図

【符号の説明】

【0225】

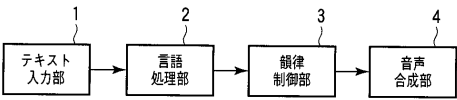
1...テキスト入力部、2...言語処理部、3...韻律処理部、4...音声合成部、41...音韻系列・韻律情報入力部、42...音声素片記憶部、43...素片選択部、44...素片融合部、45...フォルマント強調フィルタ部、46...フォルマント強調度合い推定部、47...素片編集・接続部、48...音声波形出力部、49...音声素片出力部、5...融合済み音声素片作成部、451...LPC分析部、452...LPC変形部、453...線形予測フィルタ部

30



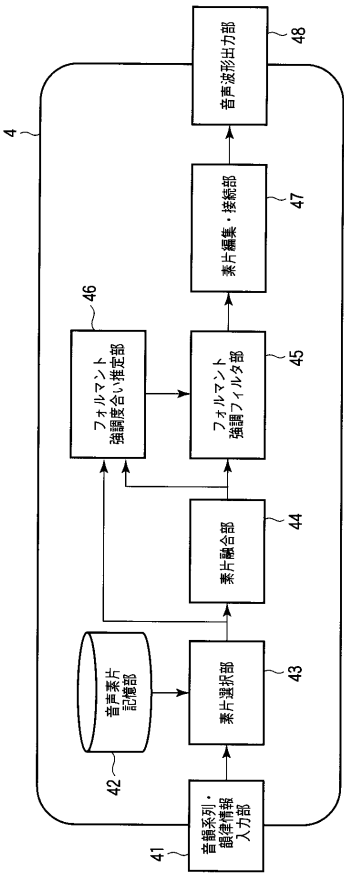
【図 1】

図 1



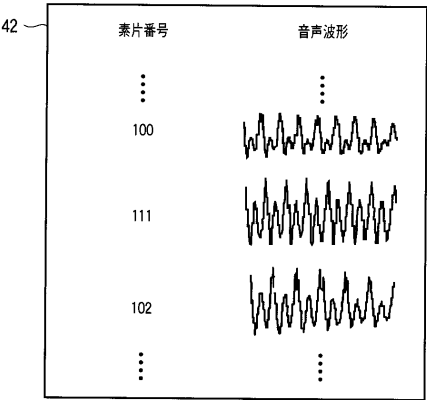
【図 2】

図 2



【図 3】

図 3



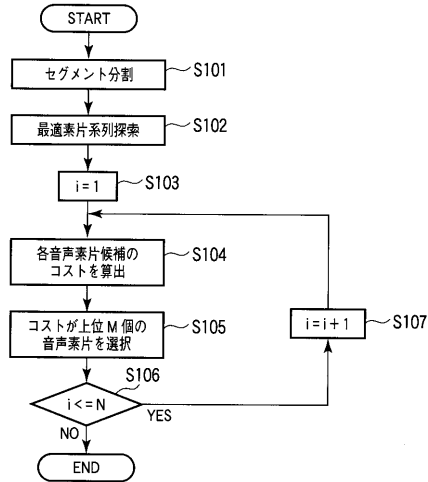
【図 4】

図 4

素片番号	当該音韻	隣接音韻 [2音韻ずつ]	基本 周波数	音韻 継続 時間長	ケプストラム係数	
					始端	終端
...	...	...	...	...	...	...
100	/a/	/-//k/, /i//m/	221Hz	83msec	254.024, ...	249.018, ...
101	/a/	//a//m/, /k//e/	296Hz	125 msec	233.028, ...	2.55.022, ...
102	/i/	/o//k/, /r//u/	240Hz	61msec	254.-0.35, ...	223.0.09, ...
...	...	...	...	...	...	...

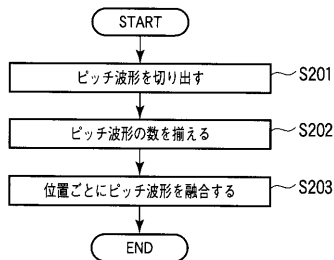
【図 5】

図 5



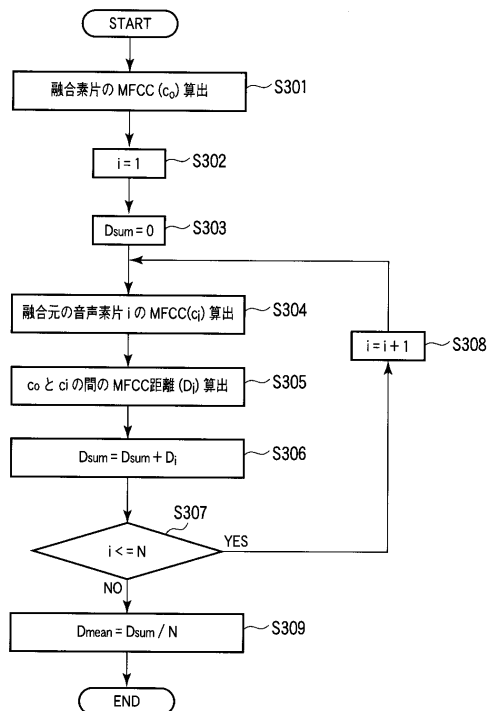
【図 6】

図 6



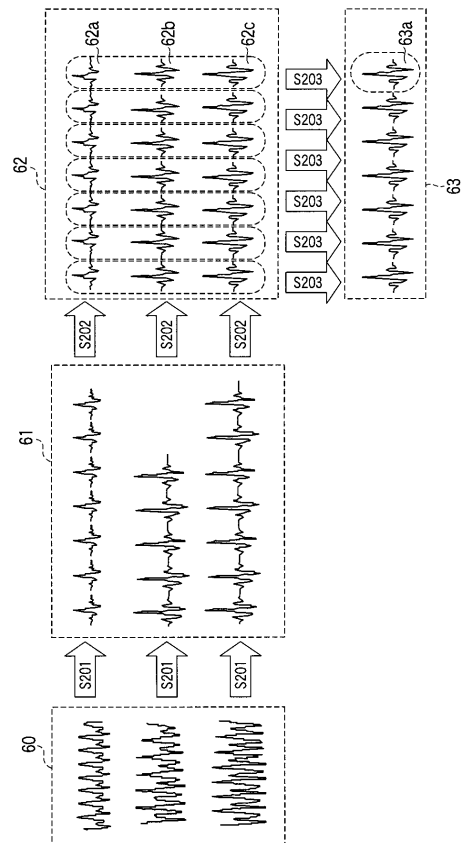
【図 8】

図 8



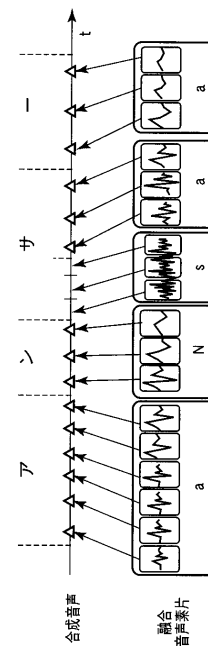
【図 7】

図 7



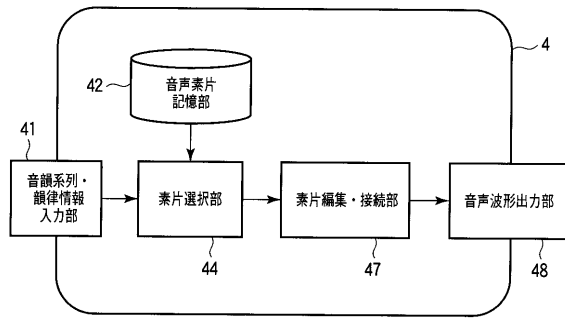
【図 9】

図 9



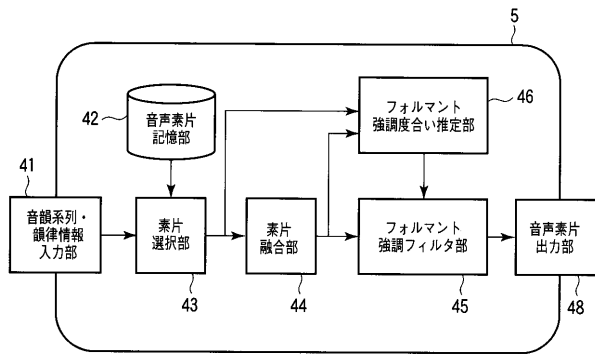
【図 10】

図 10



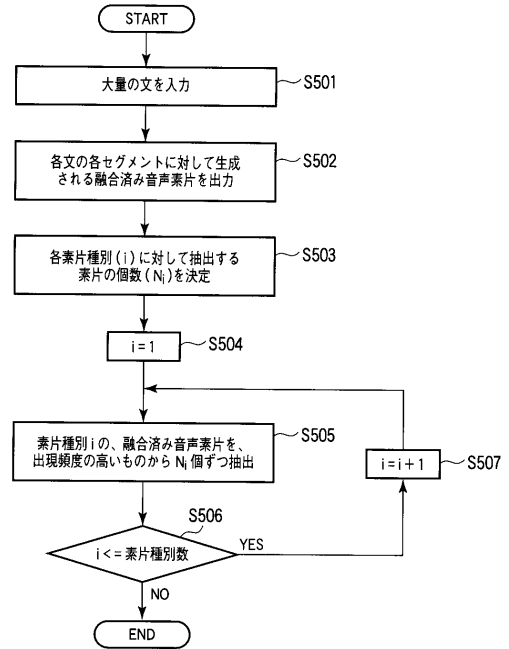
【図 11】

図 11



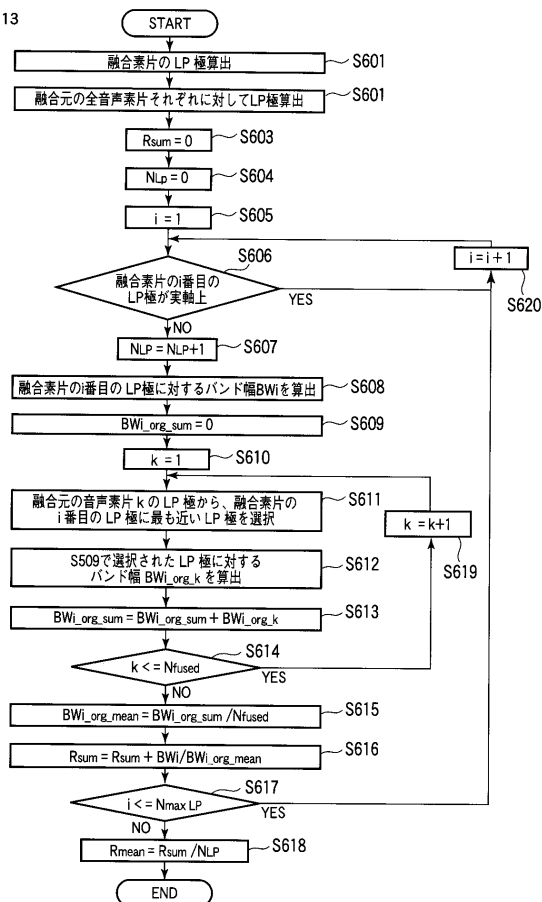
【図 12】

図 12



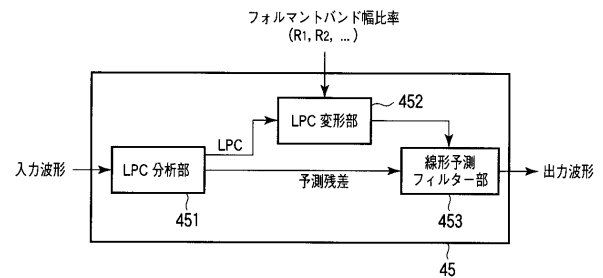
【図 13】

図 13



【図 14】

図 14



---

フロントページの続き

- (74)代理人 100095441  
弁理士 白根 俊郎
- (74)代理人 100084618  
弁理士 村松 貞男
- (74)代理人 100103034  
弁理士 野河 信久
- (74)代理人 100119976  
弁理士 幸長 保次郎
- (74)代理人 100153051  
弁理士 河野 直樹
- (74)代理人 100140176  
弁理士 砂川 克
- (74)代理人 100100952  
弁理士 風間 鉄也
- (74)代理人 100101812  
弁理士 勝村 紘
- (74)代理人 100070437  
弁理士 河井 将次
- (74)代理人 100124394  
弁理士 佐藤 立志
- (74)代理人 100112807  
弁理士 岡田 貴志
- (74)代理人 100111073  
弁理士 堀内 美保子
- (74)代理人 100134290  
弁理士 竹内 将訓
- (74)代理人 100127144  
弁理士 市原 卓三
- (74)代理人 100141933  
弁理士 山下 元
- (72)発明者 森田 眞弘  
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 籠嶋 岳彦  
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 平林 剛  
東京都港区芝浦一丁目1番1号 株式会社東芝内