

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5360489号
(P5360489)

(45) 発行日 平成25年12月4日(2013.12.4)

(24) 登録日 平成25年9月13日(2013.9.13)

(51) Int. Cl. F I
G 1 O L 13/06 (2013.01) G 1 O L 13/06 1 4 O
G 1 O L 13/033 (2013.01) G 1 O L 13/02 1 2 2 B
G 1 O L 19/00 (2013.01) G 1 O L 19/00 2 2 O Z

請求項の数 15 (全 22 頁)

(21) 出願番号	特願2009-244698 (P2009-244698)	(73) 特許権者	000002897
(22) 出願日	平成21年10月23日(2009.10.23)		大日本印刷株式会社
(65) 公開番号	特開2011-90218 (P2011-90218A)		東京都新宿区市谷加賀町一丁目1番1号
(43) 公開日	平成23年5月6日(2011.5.6)	(74) 代理人	100111659
審査請求日	平成24年8月8日(2012.8.8)		弁理士 金山 聡
		(74) 代理人	100135954
			弁理士 深町 圭子
		(74) 代理人	100119057
			弁理士 伊藤 英生
		(74) 代理人	100122529
			弁理士 藤枿 裕実
		(74) 代理人	100131369
			弁理士 後藤 直樹

最終頁に続く

(54) 【発明の名称】 音素符号変換装置および音声合成装置

(57) 【特許請求の範囲】

【請求項1】

発声された1つの音節を所定のサンプリング周波数でサンプル数Jの音節波形信号として取得する音節波形取得手段と、

前記音節波形信号を当該サンプリング周波数を維持しながら時間軸上でQ(Qは2以上の整数)倍拡大し、J×Qサンプルの拡大波形信号に変換する音節波形拡大手段と、

前記拡大波形信号に対して所定の周波数解析を行い、発音開始時刻、発音終了時刻、音高、強さのデータを備える複数の符号コードで表現した符号コード群を作成する符号コード群作成手段と、

前記符号コード群を構成する全ての符号コードの音高をQ倍に補正し、全ての符号コードの発音開始時刻と発音終了時刻を1/Q倍に補正し、補正された符号コード群で構成される音節符号を作成する音節符号作成手段と、

前記音節符号を構成する各符号コードについて、各音高ごとに発音開始時刻と発音終了時刻との時間差と符号コードの強さとの積で与えられるエネルギー値の総和であるエネルギー総和値を算出する音高別エネルギー算出手段と、

複数の音節符号間で、各音高ごとに全てのエネルギー総和値を乗算して統合エネルギー値を算出する音高別エネルギー統合手段と、

前記統合エネルギー値が高い上位の音高を所定の個数だけ抽出し、抽出された各音高に対応する符号コードに、所定の強さ、所定の発音開始時刻、所定の発音終了時刻のパラメータを設定し、複数の符号コードで構成される音素符号に変換する符号変換手段と、

10

20

を有することを特徴とする音素符号変換装置。

【請求項 2】

請求項 1 において、

前記音高別エネルギー統合手段は、日本語カナ文字の「ア」に対応する複数個の「カ」「サ」「タ」「ナ」「ハ」「マ」などの複数の子音音節に対応する音節符号間に対応するエネルギー総和値を音高別に乗算して統合エネルギー値を算出し、

前記符号変換手段は、「A」などの共通する母音音素の音素符号に変換することを特徴とする音素符号変換装置。

【請求項 3】

請求項 1 において、

前記音高別エネルギー統合手段は、日本語カナ文字の「ア」「イ」「ウ」「エ」「オ」からなる 5 種の母音音節に対応する「カ」「キ」「ク」「ケ」「コ」などの複数の子音音節に対応する音節符号間に対応するエネルギー総和値を音高別に乗算して統合エネルギー値を算出し、

前記符号変換手段は、「K」などの共通する子音音素の音素符号に変換することを特徴とする音素符号変換装置。

【請求項 4】

請求項 3 において、

前記音高別エネルギー統合手段は、日本語カナ文字の「ア」「イ」「ウ」「エ」「オ」からなる 5 種の母音音節に対応する「カ」「キ」「ク」「ケ」「コ」などの複数の子音音節に対応する音節符号間に対応するエネルギー総和値を音高別に乗算する際、

あらかじめ、「ア」「イ」「ウ」「エ」「オ」からなる 5 種の母音音節に対応する音素符号を決定しておき、

各子音音節のエネルギー総和値の中で、各々対応する母音の前記決定された音素符号を構成する全ての音高に対応するエネルギー総和値に所定の 1 未満の実数値を乗算することにより縮小するようにしていることを特徴とする音素符号変換装置。

【請求項 5】

請求項 1 から請求項 4 のいずれかに記載の音素符号変換装置により作成され、日本語カナ文字の各音節を構成する音素に対応して、所定の種類以下の音高を同時にもち、音の強さおよび音の長さが均一の複数の符号コードで成される音素符号を、音素符号を識別する音素符号識別情報と対応付けて記録した音素符号データベースと、

与えられた合成指示データに記載されている音節識別情報をもとに音素符号識別情報に変換し、対応する音素符号を前記音素符号データベースから抽出し、前記音節識別情報に従って、発音の開始および終了を特定する時刻を設定し、母音音素に対応する音素の発音の終了を特定する時刻より所定の無音区間を加えた時刻を後続する音節の発音の開始を特定する時刻として設定することにより合成音声データを生成する音素編集処理手段と、

を有することを特徴とする音声合成装置。

【請求項 6】

請求項 5 において、

前記音素編集処理手段により生成された合成音声データを音声として出力する音声出力手段をさらに有することを特徴とする音声合成装置。

【請求項 7】

請求項 5 または請求項 6 において、

前記音素編集処理手段により生成された合成音声データを五線譜に変換し、印刷する印刷手段をさらに有することを特徴とする音声合成装置。

【請求項 8】

請求項 5 から請求項 7 のいずれかにおいて、

前記音素編集処理手段は、前記合成指示データに記載されている音節識別情報が母音音節で、日本語カナ文字の長音であるとき、その音節全体の発音時間を、所定の値だけ増加させることを特徴とする音声合成装置。

10

20

30

40

50

【請求項 9】

請求項 5 から請求項 7 のいずれかにおいて、

前記音素編集処理手段は、前記合成指示データに記載されている音節識別情報が、日本語カナ文字の促音であるとき、当該促音の直後に配置される音節に対応する第 1 の音素と同一の音素を、当該第 1 の音素の直前に配置して、各音素の発音の開始を特定する時刻、発音の終了を特定する時刻を設定することを特徴とする音声合成装置。

【請求項 10】

請求項 5 から請求項 7 のいずれかにおいて、

前記音素編集処理手段は、前記合成指示データに記載されている音節識別情報が、日本語カナ文字の「ヤ」「ユ」「ヨ」の拗音であるとき、前記拗音の直前の音節については、第 1 の音素のみを設定し、当該第 1 の音素の直後に、前記拗音に対応する「ヤ」「ユ」「ヨ」いずれかの音節に対応する 2 つの音素を配置して、各音素の発音の開始を特定する時刻、発音の終了を特定する時刻を設定することを特徴とする音声合成装置。

10

【請求項 11】

請求項 5 から請求項 10 のいずれかにおいて、

前記音素編集処理手段が、前記音節識別情報より変換された音素符号識別情報に対応する音素符号を前記音素符号データベースから抽出し、前記音節識別情報に従って、発音の開始および終了を特定する時刻を設定する際、前記無音区間に対して、設定された時間伸縮率を乗算し、前記発音の開始および終了を特定する時刻に対して所定の改変を施すようにしていることを特徴とする音声合成装置。

20

【請求項 12】

請求項 5 から請求項 11 のいずれかにおいて、

前記音素編集処理手段が、前記音節識別情報より変換された音素符号識別情報に対応する音素符号を前記音素符号データベースから抽出し、前記音節識別情報に従って、発音の開始および終了を特定する時刻を設定する際、設定された音高オフセットパラメータに基づいて、前記音素符号データベースに記録されている前記音素符号が母音の場合、当該音素符号を構成する各符号コードの音高に対して、前記音高オフセットパラメータを加算し、前記合成音声データを構成する全ての母音音素に対応する符号コードの音高に対して所定の改変を施すようにしていることを特徴とする音声合成装置。

30

【請求項 13】

請求項 12 において、

前記合成指示データには各音節ごとに音節識別情報とともに前記音高オフセットパラメータが定義されており、前記音素編集処理手段が、与えられた音節識別情報より変換された音素符号識別情報に対応する音素符号を前記音素符号データベースから抽出し、前記音節情報に従って、発音の開始および終了を特定する時刻を設定する際、前記各音節ごとに定義された音高オフセットパラメータに基づいて、前記音素符号データベースに記録されている前記音素符号が母音の場合、当該音素符号を構成する各符号コードの音高に対して、前記音高オフセットパラメータを加算し、前記合成音声データを構成する全ての母音音素に対応する符号コードの音高に対して、改変を施すようにしていることを特徴とする音声合成装置。

40

【請求項 14】

請求項 1 から請求項 4 のいずれかに記載の音素符号変換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。

【請求項 15】

請求項 5 から請求項 13 のいずれかに記載の音声合成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。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、本発明は人間の声を基にして、電子楽器、楽譜等で利用可能な符号データを

50

得るための技術に関する。

【背景技術】

【0002】

従来、人間の声を擬似的に合成する手法は、種々の分野で利用されている。出願人は、人間の声をPCMでデジタル化した後、フーリエ変換を用い、実効強度の大きい周波数に対応する符号コードを取得することにより音声合成を行う技術を提案している（特許文献1～5参照）。

【0003】

また、出願人は、玩具などに搭載されている性能の低いMIDI音源でも再生可能とし、既存の楽譜編集ツールに読み込ませて五線譜に変換すると、判読性のある譜面が得られるようにするために、各音節ごとの符号コード群を音素ごとの符号コード群に変換する技術を提案している（特許文献6参照）。

10

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特許第3795201号公報

【特許文献2】特許第4037542号公報

【特許文献3】特許第4156268号公報

【特許文献4】特許第4132362号公報

【特許文献5】特許第4061070号公報

20

【特許文献6】特願2009-143825号

【発明の概要】

【発明が解決しようとする課題】

【0005】

上記特許文献6に記載の発明では、MIDI等の符号コード群として構成される音節符号に対して、各音高ごとに発音開始時刻と発音終了時刻との時間差と符号コードの強さとの積で与えられるエネルギー値の総和であるエネルギー総和値を算出し、複数の音節符号間で、各音高ごとに全てのエネルギー総和値を乗算して統合エネルギー値を算出し、統合エネルギー値が高い上位の音高を抽出して、音素符号を得るようにしたので、実際に人間が発音することにより得られた音節符号を利用して、より少ない数で全ての音節を表現可能な音素符号を得ることが可能となった。

30

【0006】

しかしながら上記特許文献6に記載の発明では、同一話者が発した日本語音節71音を録音した波形音声データ一式を高精細なMIDIデータに変換し、変換された複数の音節MIDIデータどうしを掛け合わせることで、日本語音素20音に対応する音素MIDIデータを抽出してデータベース化する方法をとっていた。この場合、変換された高精細なMIDIデータの品質、特に音節における子音音素と母音音素の時間軸上における分離性能が生成される音素MIDIデータの品質を左右する。しかし、既提案のMIDI符号化方式（特許文献1～5）を今回対象とする短い音声の音節信号に適用すると、子音音素成分と母音音素成分が混在してしまい、先提案の音節MIDIデータどうしの掛け合わせでは、明瞭性の良い音素MIDIデータを生成することが難しいことが判明した。原因としては、周波数解析における時間分解能が不十分で音声の急峻な周波数変動に追従できず、音節成分より子音音素を分離・符号化することが実現できていなかった。時間分解能を向上させる方法として、周波数解析時のフレーム長を小さく設定する方法もあるが、返って周波数解析精度が低下し妥当な符号化データを得ることが難しかった。また、限られた和音の中に低音側の声帯基本振動音など子音音素と母音音素に共通する成分が重複して含まれ、明瞭性の低下につながっていた。

40

【0007】

そこで、本発明では、従来と同等な周波数解析精度を維持しながら、解析における時間分解能を向上させ、周波数変動を高精度に抽出した音節符号が得られるとともに、符号デ

50

ータを基本として音声合成機能を実現する場合において音声再生品質の明瞭性を改善することが可能な音素符号変換装置を提供することを課題とする。

【課題を解決するための手段】

【0008】

上記課題を解決するため、本発明では、発声された1つの音節を所定のサンプリング周波数でサンプル数 J の音節波形信号として取得する音節波形取得手段と、前記音節波形信号を当該サンプリング周波数を維持しながら時間軸上で Q (Q は2以上の整数)倍拡大し、 $J \times Q$ サンプルの拡大波形信号に変換する音節波形拡大手段と、前記拡大波形信号に対して所定の周波数解析を行い、発音開始時刻、発音終了時刻、音高、強さのデータを備える複数の符号コードで表現した符号コード群を作成する符号コード群作成手段と、前記符号コード群を構成する全ての符号コードの音高を Q 倍に補正し、全ての符号コードの発音開始時刻と発音終了時刻を $1/Q$ 倍に補正し、補正された符号コード群で構成される音節符号を作成する音節符号作成手段と、前記音節符号を構成する各符号コードについて、各音高ごとに発音開始時刻と発音終了時刻との時間差と符号コードの強さとの積で与えられるエネルギー値の総和であるエネルギー総和値を算出する音高別エネルギー算出手段と、複数の音節符号間で、各音高ごとに全てのエネルギー総和値を乗算して統合エネルギー値を算出する音高別エネルギー統合手段と、前記統合エネルギー値が高い上位の音高を所定の個数だけ抽出し、抽出された各音高に対応する符号コードに、所定の強さ、所定の発音開始時刻、所定の発音終了時刻のパラメータを設定し、複数の符号コードで構成される音素符号に変換する符号変換手段を有する音素符号変換装置を提供する。

10

20

【0009】

本発明によれば、デジタル化された音響信号の各強度配列を時間軸方向に所定の倍率だけ拡大した後、所定数 T 個の強度配列で構成される単位区間ごとに、所定数 P 種類の周波数に対応したスペクトル強度を算出し、周波数、時刻を含む P 個の符号コードを得て、 P 個の符号コードの周波数を Q 倍、時刻を $1/Q$ 倍に補正するようにしたので、従来と同様な周波数解析精度を維持しながら、解析における時間分解能を向上させ、主として音声信号における周波数変動を高精度に抽出した音節符号が得られ、さらに、得られた音節符号に対して、各音高ごとに発音開始時刻と発音終了時刻との時間差と符号コードの強さとの積で与えられるエネルギー値の総和であるエネルギー総和値を算出し、複数の音節符号間で、各音高ごとに全てのエネルギー総和値を乗算して統合エネルギー値を算出し、統合エネルギー値が高い上位の音高を抽出するようにしたので、実際に人間が発音することにより得られた音節符号を利用して、より少ない数で全ての音節を表現可能な音素符号を得ることが可能となる。

30

【0012】

また、本発明では、前記音素符号変換装置により作成され、日本語カナ文字の各音節を構成する音素に対応して、所定の種類以下の音高を同時にもち、音の強さおよび音の長さが均一の複数の符号コードで構成される音素符号を、音素符号を識別する音素符号識別情報と対応付けて記録した音素符号データベースと、与えられた合成指示データに記載されている音節情報をもとに音素符号識別情報に変換し、対応する音素符号を前記音素符号データベースから抽出し、前記音節情報に従って、発音の開始および終了を特定する時刻を設定し、母音音素に対応する音素の発音の終了を特定する時刻より所定の無音区間を加えた時刻を後続する音節の発音の開始を特定する時刻として設定することにより合成音声データを生成する音素編集処理手段を有する音声合成装置を提供する。

40

【0013】

本発明の音声合成装置によれば、日本語カナ文字の各音節を構成する音素を音素符号として記録した音素符号データベースを有し、与えられた合成指示データに記載されている音節情報をもとに、対応する音素符号を音素符号データベースから抽出し、音節情報に従って、発音の開始および終了を特定する時刻を設定し、母音音素に対応する音素の発音の

50

終了を特定する時刻より所定の無音区間を加えた時刻を後続する音節の発音の開始を特定する時刻として設定するようにしたので、音声合成を行うことにより作成される符号コード群は玩具などに搭載されている性能の低いMIDI音源で再生可能であるとともに、既存の楽譜編集ツールにより演奏者が楽器演奏で再生可能な判読性のある五線譜に変換することが可能となる。

【発明の効果】

【0014】

本発明によれば、従来と同等な周波数解析精度を維持しながら、解析における時間分解能を向上させ、周波数変動を高精度に抽出した音節符号が得られるとともに、符号データを基本として音声合成機能を実現する場合において音声再生品質の明瞭性を改善することが可能となるという効果を奏する。

10

【図面の簡単な説明】

【0015】

【図1】本発明における音節と音素の基本概念を示す図である。

【図2】本発明における音節と音素の基本概念を示す図である。

【図3】本発明に係る音素符号変換装置の一実施形態を示す構成図である。

【図4】音節録音信号から音節符号を作成する処理の概要を示すフローチャートである。

【図5】時間軸方向の拡大、周波数の増加・時間情報の縮小の概念を示す図である。

【図6】音節と音素の対応関係を記録した変換テーブルを示す図である。

【図7】音節符号から音素符号へ変換する際における符号コード群の変化の様子を示す図である。

20

【図8】音節符号と、音素符号の構成例を示す図である。

【図9】音素符号記憶部13に格納された男声の音素符号の例を示す図である。

【図10】音素符号記憶部13に格納された男声の音素符号の例を示す図である。

【図11】音素符号記憶部13に格納された女声の音素符号の例を示す図である。

【図12】音素符号記憶部13に格納された女声の音素符号の例を示す図である。

【図13】図9、図10に示した男声の音素符号を五線譜化した例を示す図である。

【図14】図11、図12に示した女声の音素符号を五線譜化した例を示す図である。

【図15】本発明に係る音声合成装置の一実施形態を示す構成図である。

【発明を実施するための形態】

30

【0016】

(1. 本発明における音節と音素の基本概念)

以下、本発明の好適な実施形態について図面を参照して詳細に説明する。最初に、本発明における音節と音素の基本概念について説明する。日本語の母音は、2つの特徴的な音声フォルマント成分を含む4つ以上の重音による和音で近似できることが知られている。子音は母音成分に加えて、摩擦音など雑音を表現する和音と母音への経過音を表現する和音の3種の連結された和音が理論上必要であるが、経過音は人間の聴覚上の補間現象(音脈分凝)に委ねれば、最初の雑音と母音の2つの和音に近似できる。そこで、出願人は、日本語の母音は単一の4和音、子音は2連の4和音を基本にして全音素を表現し、これらを時間軸上につなぎ合わせれば音声合成を実現できると考えた。

40

【0017】

しかしながら、様々な要因により、MIDI音源で種々の楽器音色を設定して再生しても音声の明瞭性に欠けるという問題があった。本発明では、母音、子音等の音節の構成を従来とは根本的に異なるものとした。具体的には、71音節それぞれ固有の音節符号とするのではなく、日本語5母音と15子音に対応する20の音素符号より1つまたは2つの音素符号を選択し組み合わせることにより、71音節を表現することとした。

【0018】

まず、母音音節と子音音節の違いについて説明する。本実施形態では、母音音節は第1音素だけの単独音素とする。子音音節は第1音素と第2音素の2音素構成で第1音素は第2音素に比べ短めにする。なお、本実施形態では、母音音節の第1音素、子音音節の第1

50

音素、第2音素は、いずれも8和音となっている。

【0019】

子音音節の第1音素は、「K、S、T、N、H、M、R、G、Z、D、B、P、Y、W」の14音素のいずれか1つが選択され、子音音節第2音素の、母音音節の第1音素は、「A、I、U、E、O、n」の6音素のいずれか1つが選択される。子音音節には、濁音（「バ」など）、半濁音（「パ」など）を含み、撥音（はつおん「ン」）は第1音素は音素「N」、第2音素は音素「n」とする。

【0020】

本発明では、上述のような構成の子音音節と母音音節を用いて、音声の特徴に応じてさらに多種の態様で合成することを特徴とする。音声の特徴に応じた合成態様の分類については、様々なものが考えられるが、本実施形態では、母音、子音（濁音、半濁音、撥音含む）、長音の母音、長音の子音、促音、拗音の6通りに分類している。

10

【0021】

図1、図2は、本実施形態における音節と音素の基本概念を示す図である。図1(a)~(d)、図2(a)(b)は、上記各分類に対応したものとなっている。図1、図2において、各矩形は、音節または音素を示しており、矩形と矩形の間の空白は無音部分を示している。また、左右方向は時間軸であり、矩形の横幅は、音節の発音時間に対応している。音節の発音時間は、基本的な母音、子音の発音時間を単位区間として設定し、他の分類の音節についても、この単位区間を基準にして定められている。単位区間の具体的な長さは適宜設定することが可能であるが、本実施形態では0.25秒としている。また、詳しくは後述するが、本実施形態では、音節は、2つのパターンで符号化され、1つの音素で構成される音節と、2つの連続する音素（以下、2連音素という）で構成される音節がある。

20

【0022】

図1(a)は、母音の場合の合成パターンを示す例である。母音の場合、音節は1つの音素（第1音素）で構成される。図1(a)に示すように、第1音素を単位区間分発音した後、単位区間分無音とし、その後、他の音節が続く。図1(b)は、子音の場合の合成パターンを示す例である。子音の場合、音節は第1音素と第2音素の2つの音素で構成される。図1(b)に示すように、単位区間の1/4だけ第1音素を発音し、続いて、単位区間の3/4だけ第2音素を発音する。すなわち、第1音素の発音終了と同時に第2音素の発音を開始し、第1音素の発音開始から第2音素の発音終了までがちょうど単位区間となる。その後、単位区間分無音とし、その後、他の音節が続く。

30

【0023】

図1(c)は、長音の母音の場合の合成パターンを示す例である。長音の母音の場合、1つの音素で構成されるが、図1(c)に示すように、通常の母音に比べ、単位区間分発音時間を延ばす。すなわち、第1音素を2単位区間分発音する。その後、単位区間分無音とし、その後、他の音節が続く。図1(d)は、長音の子音の場合の合成パターンを示す例である。長音の子音の場合、2つの音素で構成されるが、図1(d)に示すように、通常の子音に比べ、第2音素の発音時間を単位区間分だけ延ばす。第1音素の発音時間は、通常の子音の場合と同じく、単位区間の1/4である。すなわち、第1音素の発音終了と同時に第2音素の発音を開始し、第1音素の発音開始から第2音素の発音終了までが2単位区間分となる。その後、単位区間分無音とし、その後、他の音節が続く。

40

【0024】

図2(a)は、促音の場合の合成パターンを示す例である。促音の場合、図2(a)に示すように、後続音節である子音の第1音素を、後続音節の直前に発音する。例えば、“ッタ”における“ッ”を合成する場合、後続音節“タ”の第1音素を、“タ”の第1音素の直前に配置する。このとき、促音の発音開始は、先行音素の発音終了から単位区間の3/4だけ経過した時点とする。すなわち、先行音素の発音終了から単位区間の3/4だけ経過した時点から、後続音節の第1音素が単位区間の1/4だけ2回繰り返して発音され、その後、後続音節の第2音素が単位区間の3/4だけ発音されることになる。

50

【 0 0 2 5 】

図 2 (b) は、拗音の場合の合成パターンを示す例である。拗音の場合、図 2 (b) に示すように、直前の子音音節を第 1 音素のみとする。そして、直前の子音音節の第 1 音素の発音終了時刻を、拗音に対応する子音の発音開始時刻として、その拗音に対応する子音の第 1 音素と第 2 音素を連続して発音する。例えば、“キャ”における“ヤ”を合成する場合、先行音節“キ”の第 1 音素の直後に、拗音“ヤ”に対応する子音“ヤ”の第 1 音素を配置する。このとき、直前の子音音節の第 1 音素の発音時間は単位区間の $1/4$ とし、拗音に対応する子音の発音時間は、通常の子音の場合と同様、第 1 音素が単位区間の $1/4$ 、第 2 音素が単位区間の $3/4$ となる。

【 0 0 2 6 】

図 1 (a) ~ 図 1 (d)、図 2 (b) においては、先行音節を省略し、図 2 (a)、(b) においては、後続音節を省略したが、基本的に全ての音素の間には単位区間分の無音区間が設定される。この単位区間の具体的長さは前述の通り 0.25 秒を標準として設定しているが、適宜変更可能である。また、図 1、2 に示したものは、基本様式であるため、各音節における各音素の発音時間の比率、発音時間や無音時間の長さは適宜変更することが可能である。例えば、隣接音節との間隔を変動させれば話速を変更できる。

【 0 0 2 7 】

(2 . 前準備)

次に、従来技術を利用した前準備について説明する。前準備として、人間の声を音節単位でデジタル化する。これは、各音節を人間が実際に発声し、これを録音したものを P C M 等の手法により行う。このとき、話者はネイティブの日本人男性または女性とし、同一人が 7 1 の音節に対してできる限り、ピッチ (音高) と発声区間を揃えて、発声することが望まれる。一般人の話者ではこのように発声を揃えるのは困難であるため、専門のボイストレーニングを受けたアナウンサー・声楽家などに発声してもらうことが望ましい。更に、男性または女性を含む複数の話者により複数のセットの 7 1 音節の録音信号を収集することが望ましい。前準備により、日本語 7 1 音節の録音信号が得られる。この日本語 7 1 音節とは、いわゆる五十音に加え、撥音、濁音、半濁音を含む計 7 1 音である。

【 0 0 2 8 】

(3 . 音素符号への変換)

次に、音節ごとの録音信号から音素符号への変換について説明する。図 3 は、本発明に係る音素符号変換装置の一実施形態を示す構成図である。記憶手段 1 0 は、音節録音データ記憶部 1 1、音節符号記憶部 1 2、音素符号記憶部 1 3 を有しており、コンピュータに接続されたハードディスク等の外部記憶装置により実現される。音節符号記憶部 1 1 には、上述の前処理により作成された録音信号が音節を特定する音節識別情報と対応付けて記憶されている。音節録音データ記憶部 1 1 には、音節と音素の対応関係を示した変換テーブルが記憶されており、符号化された音節符号が音節識別情報と対応付けて記憶される。音節符号記憶部 1 2 には、音節と音素の対応関係を示した変換テーブルが記憶されており、符号化された音節符号が音節識別情報と対応付けて記憶される。音素符号記憶部 1 3 には、符号化された音素符号が音素符号識別情報と対応付けて記憶される。

【 0 0 2 9 】

処理制御手段 2 0 は、音素符号変換装置全体の処理を統括するものであり、音節波形拡大手段 2 1、符号コード群作成手段 2 2、音節符号作成手段 2 3、音高別エネルギー算出手段 2 4、音高別エネルギー統合手段 2 5、符号変換手段 2 6 を有している。処理制御手段 2 0 は、CPU、メモリを含むコンピュータ本体であり、音節波形拡大手段 2 1、符号コード群作成手段 2 2、音節符号作成手段 2 3、音高別エネルギー算出手段 2 4、音高別エネルギー統合手段 2 5、符号変換手段 2 6 は、専用のプログラムを CPU が実行することにより実現される。符号表示手段 3 0 は、処理制御手段 2 0 により処理される音節符号、音素符号を表示するものであり、液晶ディスプレイ等の表示装置により実現される。

【 0 0 3 0 】

続いて、本実施形態に係る音素符号変換装置の処理動作について説明する。図 4 は、本

10

20

30

40

50

実施形態に係る音素符号変換装置において、音節符号を生成するまでの処理概要を示すフローチャートである。

【0031】

まず、処理制御手段20は、処理対象である音節録音信号を、図示しないデータ入力機器から読み込む。音節録音信号は、アナログの録音信号を所定のサンプリング周波数、量子化ビット数でサンプリングしたものであり、本実施形態では、サンプリング周波数44.1kHz、量子化ビット数16ビットでサンプリングした場合を例にとって以下説明していく。サンプリング周波数44.1kHzでサンプリングした場合、音節録音信号は、1秒間に44100個のサンプルを有するサンプル列(サンプルの配列)として構成されることになる。

10

【0032】

音節録音信号を読み込んだら、音節波形拡大手段21が、音節録音信号を構成するサンプルを時間軸方向に所定の倍率Q(Qは整数)だけ拡大する(S1)。具体的には、音節録音信号を構成するサンプルの数をQ倍にする。そして、Q個ごとに、元のサンプルと同じ値のものを配置し、その間の(Q-1)個のサンプルの値としては、両側に位置する元のサンプルの値を用いて線形補間したものを与える。原音節録音信号の各サンプルj(j=0...J-1)についてのサンプル値をx(j)とすると、コンピュータは、以下の〔数式1〕に従った処理を実行することにより、拡大後の音節録音信号の各サンプルj・Q+k(0≤k<Q-1)についてのサンプル値x'(j・Q+k)を算出する。以下の〔数式1〕において、wはk/(Q-1)で与えられる0≤w<1の値をとる実数値とする。

20

【0033】

〔数式1〕

$$x'(j \cdot Q + k) = (1 - w) \cdot x(j) + w \cdot x(j + 1)$$

【0034】

S1における処理の結果、デジタル音響信号を構成するJ個のサンプルは、J×Q個に拡大される。図5(a)にS1における拡大処理による波形の変化を示す。図5(a)における波形は、サンプルの値をプロットしたものを線で結んだものであるが、サンプル数が多いため、曲線状に表現されるものである。上記〔数式1〕に従った処理を実行することにより、左側に示したような波形が右側に示したような波形に変化することになる。なお、図5の例では、説明の便宜上Q=2の場合を示している。

30

【0035】

次に、符号コード群作成手段22が、時間軸方向に拡大されたサンプル上に単位区間を設定する(S2)。単位区間の長さ(サンプル数T)は、サンプリング周波数との関係で設定されるが、サンプリング周波数が44.1kHzの場合、低域部まで忠実に解析するためには、4096サンプル以上必要である。そこで、本実施形態では、1単位区間のサンプル数T=4096として単位区間を設定している。

【0036】

単位区間の設定は、特許文献1~5に開示されているように、デジタル音響信号の先頭から順次サンプルを抽出することにより行われる。単位区間は、全てのサンプルを漏らさず設定し、好ましくは、連続する単位区間においてサンプルが重複するように設定する。この場合、各単位区間の先頭の間隔(シフト幅という)は、様々な規則で設定することができる。最も単純なのは、シフト幅を固定、すなわち重複させるサンプル数を一定として設定する手法である。例えば、T=4096の場合、先頭の単位区間をj=0~4095、2番目の単位区間をj=2048~6143、2番目の単位区間をj=4096~8191というように、2048(=T/2)個のサンプルを重複させながら、設定することになる。しかし、時間分解能を向上させるためには、シフト幅を小さくしたいという要望があり、一方、シフト幅を小さくするほど、計算時間が増大するという問題がある。また、シフト幅を必要以上に小さくすると、後述する図4・S4の単音成分の連結処理において連結条件が満足されなくなり連結処理が適切に機能しなくなる。そこで、音響信号の状

40

50

態に合わせて最適なシフト幅を設定するため、本実施形態では、特許文献5に開示したような、ゼロ交差点間隔の粗密または自己相関解析により周波数変化が顕著なゼロ交差点を選別し、このゼロ交差点に位置するサンプルを先頭とする。

【0037】

ゼロ交差点とは、正負両極性の音響信号と、信号の0レベルとの交差点のことであり、ここでは、音響信号の信号の強度値（振幅）が0となる時刻を示す。ただし、デジタル化した音響信号は、アナログ信号におけるゼロ交差点をサンプルするとは限らない。そのため、実際には、強度値が丁度0になる場合に加え、サンプリング点の強度値が正から負、または負から正に変化した場合に、その前後のサンプリング点のどちらかをゼロ交差点とみなす処理を行う。なお、ゼロ交差点検出のためには、解析対象となる音響信号が正負両極性となっている必要がある。そのため、直流成分を含む音響信号については、直流成分を除去しておく必要がある。直流成分の除去については、周知の種々の手法を適用することができるので、ここでは詳細な説明は省略する。基本的には、ゼロ交差点に位置するサンプルを先頭として単位区間を設定するが、連続する単位区間のシフト幅が一定の範囲に収まるように、ゼロ交差点以外の位置を先頭として単位区間を設定する場合もある。具体的には、最大シフト幅（例えば $T/2$ ）を上回る場合は、ゼロ交差点以外の位置でも最大シフト幅となる位置を先頭にして単位区間を設定する。逆に、最小シフト幅（例えば $T/8$ ）を下回る場合は、最小シフト幅を上回るように幾つかのゼロ交差点を飛ばした位置を先頭にして単位区間を設定し、最小シフト幅を上回りかつ最大シフト幅の範囲で該当するゼロ交差点が存在しない場合は、上記と同様に最大シフト幅となる位置を先頭にして単位区間を設定するような補正を行う。

【0038】

符号コード群作成手段22は、続いて、設定された各単位区間を対象として周波数解析を実行し、各単位区間のスペクトルを算出する（S3）。各単位区間のスペクトルの算出は、特許文献1～5に開示されているように、MIDIのノートナンバー n （ $0 \leq n < 127$ ）に対応する128種の解析周波数 $f(n) = 440 \cdot 2^{(n-69)/12}$ の要素信号（要素関数）を基本にした一般化調和解析により、128個の成分を抽出することにより行う。ノートナンバー n に対応して解析周波数を設定した場合、周波数が高くなるにつれ、ノートナンバー間の周波数間隔が広がるため、特に、 n が60を超えると解析精度が低下してしまう。そこで、本実施形態では、特許文献3に開示したように、ノートナンバー間を M 個の微分音に分割した128 M 個の要素信号 $f(n, m) = 440 \cdot 2^{(n-69+m/M)/12}$ を用いて解析を行い、128 M 個の成分を抽出する。後述する図4・S4においてピッチベンド符号の付加など特殊な符号化を行わない限り、各ノートナンバーにおける M 個の微分音の情報は不要であるため、 M 個の微分音の成分の合算値を当該ノートナンバーにおける成分として代表させ、結果的に128個の成分を抽出する。

【0039】

符号コード群作成手段22による具体的な処理手順としては、まず、ノートナンバー分の強度配列 $E(n)$ （ $0 \leq n < 127$ ）と副周波数配列 $S(n)$ を設定し、初期値を全て0とする。続いて、 $0 \leq n < 127$ および $0 \leq m < M-1$ に対して以下の〔数式2〕に従った処理を実行し、 $E(n, m)$ を最大にする（ n_{max}, m_{max} ）を求める。

【0040】

〔数式2〕

$$A(n, m) = (1/T(n)) \cdot \sum_{i=0, T(n)-1} x(i) \sin(2\pi f(n, m) i / f_s)$$

$$B(n, m) = (1/T(n)) \cdot \sum_{i=0, T(n)-1} x(i) \cos(2\pi f(n, m) i / f_s)$$

$$\{E(n, m)\}^2 = \{A(n, m)\}^2 + \{B(n, m)\}^2$$

【0041】

上記〔数式2〕において $T(n)$ は解析フレーム長であり、単位区間 T を超えない範囲で要素信号の周期の最大の整数倍になるように設定し、 k を適当な整数値として、 $T(n) = k / f(n, m)$ で与える。 $E(n, m)$ を最大にする（ n_{max}, m_{max} ）を用いた $f(n_{max}, m_{max})$ が調和信号として選出されることになる。（ n_{max}, m_{max} ）

\max) が求められたら、コンピュータは、 $A(n_{\max}, m_{\max})$ および $B(n_{\max}, m_{\max})$ を用いて、以下の〔数式 3〕に従った処理を実行し、サンプル配列 $x(i)$ の全ての要素 ($0 \leq i < T - 1$) を更新する。

【0042】

〔数式 3〕

$$x(i) = x(i) - A(n_{\max}, m_{\max}) \cdot \sin(2\pi f(n_{\max}, m_{\max})i / fs) - B(n_{\max}, m_{\max}) \cdot \cos(2\pi f(n_{\max}, m_{\max})i / fs)$$

【0043】

上記〔数式 3〕においては、 $x(i)$ から含有信号を減じる処理を行っている。さらに、以下の〔数式 4〕に従った処理を実行し、強度配列 $E(n)$ 、副周波数配列 $S(n)$ を更新する。

10

【0044】

〔数式 4〕

$$E(n_{\max}) = E(n_{\max}) + E(n_{\max}, m_{\max})$$

$$S(n_{\max}) = m_{\max}$$

【0045】

コンピュータは、上記〔数式 2〕～〔数式 4〕の処理を全ての n ($0 \leq n < 127$) に対して実行し、全ての $E(n)$ および $S(n)$ の値を決定する。

【0046】

本実施形態では、処理負荷を軽減するため、 M の値については、ノートナンバーに基づいて可変に設定し、例えば解析する周波数間隔が 100 Hz 程度になるようにしている。そして、ノートナンバー 60 以下は分割せず $M = 1$ にする。また、精度は若干落ちるが、初回の〔数式 2〕の処理で $S(n)$ を決定し、2 回目以降の〔数式 2〕の処理は、 $m = S(n)$ に固定して行い、微分音解析を省略するようにしても良い。また、〔数式 2〕の処理で、既に同一ノートナンバーに対して副周波数が異なる信号成分が複数回に渡って解析される可能性があるが、 $E(n)$ と $S(n)$ に既に値がセットされている場合は、 $E(n, m)$ の最大値の選定候補から除外するようにしても良い。

20

【0047】

各单位区間について、スペクトル (128 個の周波数成分) が算出されたら、音節符号作成手段 23 が、周波数情報と、各周波数に対応するスペクトル強度、および単位区間の開始と終了を特定可能な時間情報で構成される符号コードを作成する (S4)。符号コードの作成にあたり、まず、算出したスペクトルに、各ノートナンバー n の時刻、時間長の情報を追加し、[開始時刻, 時間長, 主周波数 n , 副周波数 $S(n)$, 強度 $E(n)$] で構成される単音成分を作成する。「開始時刻」としては単位区間の先頭の時刻を、デジタル音響信号全体において特定できる情報であれば良く、本実施形態では、単位区間の先頭サンプル ($i = 0$) に付されたデジタル音響信号全体におけるサンプル番号 (絶対サンプルアドレス: j に対応) を記録している。この絶対サンプルアドレスをサンプリング周波数 (44100) で除算することにより、音響信号先頭からの時刻が得られる。時間長は、本実施形態では単位区間ごとに可変で与えられることを特徴とし、直後に後続する単位区間の開始時刻までの差分 (後続する単位区間の開始時刻 - 当該単位区間の開始時刻) で与えられる。

30

40

【0048】

S2 で設定された単位区間ごとに、128 個の単音成分が作成されるが、さらに、S4 においては、連続する単位区間において単音成分を連結する処理を行う。具体的には、連続する単位区間における同一ノートナンバーの単音成分が、所定の連結条件を満たす場合、2 つの単音成分を連結する。連結条件としては、同一の音として連続性を有する状態を適宜設定することができるが、本実施形態では、副周波数を考慮した周波数 (主周波数 + 副周波数) の差が所定の閾値 N_{diff} 未満で、双方の強度が所定の閾値 L_{min} 以上で、かつ双方の強度の差が所定の閾値 L_{diff} 未満である場合に、連続性を有するとして、後続の単音成分を前方の単音成分に連結する。ただし、連結後の主周波数、副周波数、強度

50

は大きい方の単音成分の各値を採用し、時間長は双方の和で与える。連結条件としての具体的な閾値は、本実施形態では、 $N d i f = 8 / 2 5$ [単位：ノートナンバー換算]、 $L m i n = 1$ [単位：128段階ベロシティ換算]、 $L d i f = 1 0$ [単位：128段階ベロシティ換算]としている。連結処理は、符号コードへの変換前に行うものであるため、各閾値は、ノートナンバー、ベロシティに換算したものである。

【0049】

同一ノートナンバーの単音成分の連結は、連結条件を満たす限り、後続する単位区間の単音成分に対して繰り返し行い、最終的に得られた[開始時刻，時間長，主周波数 n ，副周波数 $S(n)$ ，強度 $E(n)$]の単音成分を、符号コードに変換する。符号コードの形式としては、周波数情報と、各周波数に対応するスペクトル強度、および単位区間の開始と終了を特定可能な時間情報を有するものであれば、どのような形式のものであっても良いが、本実施形態では、MIDI形式に変換する。MIDIでは、発音開始と、発音終了を別のイベントとして発生するため、したがって、本実施形態では、1つの単音成分を2つのMIDIノートイベントに変換する。具体的には、「開始時刻」で、ノートナンバー n のノートオンイベントを発行し、ベロシティ値は強度 $E(n)$ の最大値を E_{max} として、 $128 \cdot \{E(n) / E_{max}\}^{1/4}$ で与える。時刻については、Standard MIDI Fileでは、直前イベントとの相対時刻(デルタタイム)で与える必要があり、その時刻単位は任意の整数値で定義でき、例えば、 $1 / 1536$ [秒]の単位に変換して与える。そして、「開始時刻」+「時間長」で特定される終了時刻で、ノートナンバー n のノートオフイベントを発行する。この際、時間長には、0以上1以下の実数を乗じる。これは、使用するMIDI音源の音色にも依存するが、MIDI音源の余韻を考慮して早めにノートオフ指示をするためである。時間長をそのまま用いてもMIDI音源の処理上問題はないが、発音の際、後続音と部分的に重なる場合がある。

【0050】

MIDI符号に変換する際、MIDI音源で処理可能な同時発音数についても考慮するため、同時発音数の調整を行う必要がある。MIDI音源で処理可能な同時発音数が32である場合、時間軸方向に発音期間中(ノートオン状態)のノートイベントの個数を連続的にカウントし、同時に32個のノートイベントが存在する箇所が見つかった場合は、各々対になるノートオフイベントを近傍区間内で探索し、各ノートイベント対のベロシティ値とデュレーション値(ノートオフ時刻-ノートオン時刻)の積(エネルギー値)で優先度を評価し、指定和音数(この場合“32”)以下になるように優先度の低い(エネルギー値の小さい)ノートイベント対を局所的に削除する処理を行う。“局所的に”とは、32を超えるノートイベントが存在する部分に限りという意味である。この際、ベロシティ値またはデュレーション値のいずれかが所定の下限值より低い場合、優先度に関係なく削除する処理も行う。

【0051】

さらに、MIDI符号に変換する際、MIDI音源で処理可能なビットレートについても考慮するため、ビットレートの調整を行う必要がある。時間軸方向に、1秒間隔にノートオンまたはノートオフイベントの個数をカウントし、各々の符号長を平均5バイト(40ビット)とし、MIDI音源で処理可能な最大ビットレートを9000 [bps(ビット/秒)]とすると、1秒間あたりイベント数が $9000 / 40 = 225$ 個を超えている区間が見つかった場合は、その区間に存在するノートオンまたはノートオフイベントと各々対になるノートオフまたはノートオンイベントを近傍区間内で探索し、各ノートイベント対のベロシティ値とデュレーション値(ノートオフ時刻-ノートオン時刻)の積(エネルギー値)で優先度を評価し、指定イベント個数(この場合“225”)以下になるように優先度の低い(エネルギー値の小さい)ノートイベント対を局所的に削除する処理を行う。この際、ベロシティ値またはデュレーション値のいずれかが所定の下限值より低い場合、優先度に関係なく削除する処理も行う。

【0052】

符号コードの作成が行われたら、音節符号作成手段23は、時間軸方向に拡大して処理

10

20

30

40

50

されたことによる変動を是正するため、各符号コードを補正する処理を行う（S5）。具体的には、まず、全てのノートイベント（ノートオンイベントまたはノートオフイベント）のノートナンバー値に $12 \cdot \log_2 Q$ だけ加算する処理を行う。例えば、 $Q = 4$ の場合、24半音（2オクターブ）だけ全体的に音高を上げる。この処理は、S1においてサンプル数をQ倍したことにより周波数が $1/Q$ になっているため、周波数をQ倍にして元の状態に戻すために行う。この補正によりノートナンバーが規格値上限の127を超えるノートナンバーをもつ符号コードは削除する。具体的には補正前のノートナンバーが $128 - 12 \cdot \log_2 Q$ 以上の符号コードが削除される。

【0053】

続いて、全てのノートイベントの時刻（ノートオン時刻またはノートオフ時刻）に $1/Q$ を乗算する。これにより、MIDI符号全体の演奏時間、および各ノートイベントの発音時間が $1/Q$ に縮小される。この処理は、S1においてサンプル数をQ倍したことにより全体の演奏時間がQ倍になっているため、時刻を $1/Q$ にして元の状態に戻すために行う。この処理を行うと、時間あたりのノートイベント数がQ倍に増大するため、上記S4で実行したビットレートの調整を再度実行する。

【0054】

S5における処理の結果、周波数（音高）はQ倍になるとともに、時間情報は $1/Q$ になる。S5の補正処理によるMIDIイベント（MIDI符号のノートイベント）の変化の様子を図5（b）に示す。図5（b）においては、 $Q = 2$ の場合のMIDIイベントの変化を、音符により示している。S5の補正処理により左側の“ミ”の音符は、右側では1オクターブ高い（周波数が2倍）“ミ”の音符に変化している。一方、左側の四分音符が、右側では時間的に $1/2$ の八分音符に変化している。このようにして得られた音節符号は、音節識別情報と対応付けられて音節符号記憶部12に記憶される。

【0055】

また、上記の例では、ビットレートの調整をS4、S5の2回行っているが、これらについては、最終的に少なくとも1回ずつ行われていれば良い。また、上記実施形態では、S2～S4の処理について好ましい処理例について具体的に説明したが、これらの処理については、本発明の趣旨を逸脱しない範囲で、公知の特許文献1～5に開示した技術を用いることができる。

【0056】

音節符号が得られたら、71の音節単位で構成される音節符号を基に、20の音素符号に変換する。ここで、音節と音素の対応関係を示した変換テーブルを図6に示す。図6において、カタカナで示す「ア」「イ」・・・の71音は音節であり、アルファベットで示す「A」「I」は音素である。母音音素「A」・・・「O」は水平方向の13音節のAND演算で決定され、子音音素「K」・・・「P」は垂直方向の5音節のAND演算で決定される。図6に示した横長の枠は、母音音素「A」が13個の音節「ア」・・・「パ」で決定されることを示しており、縦長の枠は、子音音素「K」が5個の音節「カ」・・・「コ」で決定されることを示している。なお、子音音素「Y」は3個の音節「ヤ」「ユ」「ヨ」で決定され、子音音素「W」は2個の音節「ワ」「ヲ」で決定され、音素「n」は1個の音節「ン」で決定される。

【0057】

処理制御手段20は、図6に示した変換テーブルを参照し、関連する複数の音節符号を読み込む。例えば、音素符号「K」を得るために、音節符号記憶部11から「カ」「キ」「ク」「ケ」「コ」の5個の音節符号を読み込む。

【0058】

処理制御手段20が、5個の音節符号を読み込んだら、音高別エネルギー算出手段24は、各音節符号単位で、構成する全ての符号コードを対象にして、音高（MIDIの場合、ノートナンバー）別に、エネルギー総和値を算出する。エネルギー総和値は、各音高における音の強度（MIDIの場合、ベロシティ）×発音時間（MIDIの場合、デュレーション：ノートオフ時刻 - ノートオン時刻）により算出する。ここで、エネルギー総和値

10

20

30

40

50

の“総和”とは、1つの音節符号内において、同一音高で2つ以上の符号コードが存在する場合、各符号コードについて総和をとることを意味している。

【0059】

各符号単位で、全音高についてエネルギー総和値が算出されたら、次に、音高別エネルギー統合手段25が、5個の音節符号間で、同一の音高別にエネルギー総和値を乗算し、統合エネルギー値を得る。乗算であるため、5個の音節符号のうち、いずれか1つの音節符号に存在しない音高成分については、“0”となる。したがって、統合エネルギー値を求めることにより、5個の音節符号に共通の成分のみが残ることになる。例えば、音素符号「K」を構成する音高成分は、音節符号「カ」「キ」「ク」「ケ」「コ」に共通に含まれている音高成分でなければならない。

10

【0060】

音高別に統合エネルギー値が算出されたら、符号変換手段26が、統合エネルギー値が上位の音高を指定和音数(例：8個)を超えないように選出する。

【0061】

指定和音数は、事前に設定されるものであり、自由に設定することができるが、本実施形態では、上述のように“8”としている。したがって、本実施形態では、32個の符号コードから8個の符号コードが、符号変換手段23により抽出されることになる。さらに、符号変換手段26は、元の各音節符号を構成する複数の符号コードのうち、最も早い区間開始時刻、最も遅い区間終了時刻を、それぞれ選出された全ての音高の区間開始時刻、区間終了時刻として設定し、選出された音高のベロシティを規定値に設定する。ベロシティの規定値については、ベロシティが“0”～“127”の値を取り得るため、本実施形態では、その最大の“127”としている。

20

【0062】

音高別エネルギー算出手段24、音高別エネルギー統合手段25、符号変換手段26による処理前と処理後の音素符号の変化の様子を図7に示す。図7において、横軸は時間、縦軸は周波数(ノートナンバー)に対応している。グラフ内に配置された矩形は符号コードを示しており、横方向の長さは横軸に従って時間的長さを示しているが、縦方向の長さは縦軸とは異なり、周波数ではなく強度(ベロシティ)を示している。

【0063】

実際には、図6に示したように、1つの音素は、13個の音節、または5個の音節により得られる(例外として、音素Yは3個の音節、音素Wは2個の音節、音素nは1個の音節)が、図7の例では、説明の便宜上2つの音節符号、を用いて、音素を抽出する場合を示している。図7(a)、(b)は、それぞれ音高別エネルギー算出手段21による処理前の音節符号、を構成する符号コード群を示したものである。上述のように、本実施形態では、同一時刻において32個の符号コードで音節符号を構成し、指定和音数は8に設定するのが一般的であるが、図7(a)では、説明の都合上、同一時刻において符号コードは最大6個となっており、指定和音数は4に設定している場合を示している。また、各符号コードを示す矩形の横方向および縦方向の長さからわかるように、各符号コードの再生時間(終了時刻-開始時刻)および強度も異なっている。

30

【0064】

音高別エネルギー算出手段24、音高別エネルギー統合手段25による処理後は、全ての音高についての統合エネルギー値が得られる。統合エネルギー値は、図7(a)(b)に示されるような、音高別エネルギー算出手段24により音高ごとに算出される発音開始時刻と発音終了時刻との時間差と符号コードの強さ(ベロシティ)との積の同一音高における総和値であるエネルギー総和値に対して、音節符号ととの間で対応する音高におけるエネルギー総和値同士を乗算したものである。図で表現するのは難しいが、イメージ的には、図7(c)に示すように、音高に対応して、その統合エネルギー値(図7(c)では、図面上下方向の幅で表現)が得られる。そして、符号変換手段26により、図7(d)に示すように、統合エネルギー値が大きい音高が4つ選出される。さらに、この後、符号変換手段26により、各音高に対応する符号コードの強度値は、上述のような規

40

50

定値に設定される。「カ」「キ」「ク」「ケ」「コ」の5個の音節符号(最大32の音高)と、「A」「I」「U」「E」「O」「K」の6個の音素符号(最大8の音高)の構成例を図8に示す。

【0065】

図6に示したように、1つの母音の音素符号は、13個の音節符号の各音高におけるエネルギー総和値同士を乗算することにより決定され、1つの子音の音素符号は一部の例外を除き5個の各音高におけるエネルギー総和値同士を乗算することにより決定され、基本的にはこれら20種の音素符号を決定する順位は問わない。そうすると、母音の音素符号と子音の音素符号との間で、ある音高が重複して選出される場合が生じる。即ち、8音など限られた音高で構成される子音の音素と母音の音素が音声合成時に時間的に連続して同一音高が再生される場合が生じ、明瞭性の点で好ましくない。そこで、本願では、子音の音素と母音の音素の符号間で、できるだけ同一の音高が重複して含まれないようにする方法を提案する。そのためには、先に、13個の音節符号の各音高におけるエネルギー総和値同士を乗算することにより5種類の母音の音素符号を決定する。続いて、5個の母音「ア」「イ」「ウ」「エ」「オ」のいずれかを含む5個の子音の音節符号(例えば、「カ」「キ」「ク」「ケ」「コ」)の各音高におけるエネルギー総和値同士を乗算する前に、各音節符号に含まれる母音の既に決定された音素符号の全ての音高に対応するエネルギー総和値を一定の割合で縮小させる。例えば、エネルギー総和値に1/1000を乗算させて他の音高のエネルギー総和値に対して相対的に小さな値に改変する。乗算する値は、必ずしも1/1000とする必要はなく、1未満であれば良い。その後、5個の音節符号の各音高における改変されたエネルギー総和値同士を乗算することにより5種類の母音の音素符号を決定する。このような補正処理を施すことにより、子音の音素と母音の音素の符号間で、同一の音高が重複して選択される確率が低くなる。(ただし、エネルギー総和値の高い音高の種類が元来少ない場合、多少重複して選択されることもある。)

【0066】

音高別エネルギー算出手段24、音高別エネルギー統合手段25、符号変換手段26は、音節符号記憶部12に記憶されている各音節符号について処理を行い、得られた各音素符号を音素符号記憶部13に格納する。音素符号記憶部13に格納された音素符号の例を図9~図12に示す。このうち、図9、図10は男声を符号化したものであり、図11、図12は女声を符号化したものである。図9~図12中、“C、C#、D、D#、E、F、F#、G、G#、A、A#、B”は、“ド、ド#、レ、レ#、ミ、ファ、ファ#、ソ、ソ#、ラ、ラ#、シ”の音名の英語表記で、列記されている数字はオクターブ番号を示し、音名とオクターブ番号の対記号でMIDI規格のノートナンバーを特定でき、本願ではMIDI規格ノートナンバーの69をA3と表記する(国際的にはA4をMIDI規格ノートナンバーの69を示す表記も多数存在する)。音素符号を構成する符号コードが、MIDI規格で定義されている場合、市販の楽譜編集ツールにより五線譜に変換することができる。図9、図10の男声の音素符号を五線譜化した例を図13に、図11、図12の女声の音素符号を五線譜化した例を図14にそれぞれ示す。

【0067】

(4. 音声の合成)

次に、得られた音素符号を利用した音声の合成について説明する。図15は、本発明に係る音声合成装置の一実施形態を示す構成図である。図15において、音素符号データベース13aは、得られた音素符号を、音素符号識別情報と対応付けて記録したものである。音素符号データベース13aに格納されている音素符号は、上述の音素符号変換装置により変換され、音素符号記憶部13に格納されたものと同じである。したがって、上述の音素符号変換装置は、この音素符号データベース13aを作成するためのものであるとも言える。また、音素符号データベース13aには、各音素符号識別情報と、音節を特定する音節識別情報との対応関係を示した変換テーブルが記録されている。この変換テーブルは、図6に示したのと同じである。合成音声データ記憶手段14は、音素編集処理手段50により合成された合成音声データを記憶するものであり、ハードディスク等の記憶装

10

20

30

40

50

置により実現される。

【 0 0 6 8 】

音素編集処理手段 5 0 は、合成指示データの内容に従って、音素符号データベース 1 3 a から対応する音素符号を抽出し、所定の加工を施して合成音声データを生成し、所定の出力先に出力する処理を行う。生成された合成音声データは、設定に従って合成音声データ記憶手段 1 4、音声出力手段 6 0、印刷手段 7 0 のうち、1 つ以上に出力される。音声出力手段 6 0 は、音素編集処理手段 5 0 から受け取った合成音声データを実際の音声として発音するものであり、M I D I 音源を備えた M I D I 再生装置により実現される。印刷手段 7 0 は、音素編集処理手段 5 0 から受け取った合成音声データを五線譜に変換し、印刷するものであり、五線譜への変換は、公知の変換ソフトウェアを実行することにより実現され、印刷機能は、公知のプリンタ等により実現される。図 1 5 に示した音声合成装置は、現実には、入力機器、外部記憶装置を備え、M I D I 再生装置を接続したコンピュータに専用のプログラムを組み込むことにより実現される。

10

【 0 0 6 9 】

音声合成装置に入力される合成指示データは、音節識別情報を所定の順序で配置したものであり、この音節識別情報は、音節を識別することができるものであれば、どのような形式であっても良い。本実施形態では、音節識別情報として、音節に対応する文字コードを記録したテキストデータを用いている。この場合、音素符号データベース 1 3 a 内の変換テーブルには、音節識別情報に対応する文字コードと音素符号識別情報が対応付けて記録されている必要がある。

20

【 0 0 7 0 】

続いて、図 1 5 に示した音声合成装置の処理動作について説明する。まず、合成指示データを音声合成装置に入力する。音声合成装置は、合成指示データを読み込むと、音素編集処理手段 5 0 が合成指示データ内を先頭の音節識別情報から順に合成処理していく。具体的には、音素編集処理手段 5 0 は、合成指示データ内の音節識別情報で音素符号データベース 1 3 a 内の変換テーブルを参照して、音素符号識別情報を取得し、その音素符号識別情報に対応する音素符号を抽出する。

【 0 0 7 1 】

そして、抽出した音素符号が母音音素 1 つだけである場合は、母音音節であるので、先行する音節のノートオフ時刻の 0 . 2 5 秒後をノートオン時刻として設定し、その 0 . 2 5 秒後をノートオフ時刻とし、ノートナンバー、ペロシティは音素符号データベース 1 2 a に記録されていた値そのものとする M I D I イベントを作成する。ただし、ノートナンバーについてはオプション的に別途ユーザにより指示される音高オフセットパラメータに基づいて適宜上下され、ピッチ変換を行えるようにしてある。

30

【 0 0 7 2 】

抽出した音素符号が 2 つであり、それが子音音素と母音音素である場合は、子音音節であるので、先の子音音素符号について、前の音節のノートオフ時刻の 0 . 2 5 秒後をノートオン時刻として設定し、単位区間の 1 / 4、すなわち 0 . 0 6 2 5 秒後をノートオフ時刻とする。そして、後の母音音素符号について、先の子音音素符号のノートオフ時刻をノートオン時刻として設定し、単位区間の 3 / 4、すなわち 0 . 1 8 7 5 秒後をノートオフ時刻とする。子音音節の場合も、母音音節の場合と同様、ノートナンバー、ペロシティは音素符号データベース 1 2 a に記録されていた値そのものとするが、ノートナンバーについてはオプション的に別途ユーザにより指示される音高オフセットパラメータを加算することにより適宜上下され、ピッチ変換を実現することができる。また、上記の 0 . 2 5 秒、0 . 0 6 2 5 秒、0 . 1 8 7 5 秒という時間数値はあくまで基準値であり、別途ユーザにより指示される時間伸縮パラメータを乗算することにより適宜伸縮され、話速変換を実現することができる。

40

【 0 0 7 3 】

音節識別情報が長音を示すものであった場合（音節識別情報を文字コードで記録したときは、“ー”に対応する文字コードであった場合）、その直前の音節識別情報とともに 2

50

つの音節識別情報で1つの長音の音節を特定する。例えば、例えば、音節識別情報が“ア”と“ー”が連続した場合、2つの音節識別情報“アー”により、長音の母音音節であると判断する。音節識別情報が“カ”と“ー”が連続した場合、2つの音節識別情報“カー”により、長音の子音音節であると判断する。長音の場合、長音の母音音節と長音の子音音節で若干異なる。長音の母音音節の場合、ノートオン時刻からノートオフ時刻の間隔を0.5秒に増加して設定する。長音の子音音節の場合、先の子音音素符号については、長音でない通常の場合と同様、ノートオン時刻からノートオフ時刻の間隔を0.0625秒にして設定し、後の母音音素符号についてノートオン時刻からノートオフ時刻の間隔を0.4375秒にして設定する。したがって、長音の場合、音節全体の発音時間は、母音音節、子音音節ともに同じ0.5秒となる。子音音節については、第1音素の発音時間は、長音でない通常の場合と同じ0.0625秒であるが、第2音素の発音時間が、長音でない通常の場合と比べて長くなる。尚、上記の0.5秒、0.25秒、0.4375秒、0.5秒、0.0625という時間数値も同様にあくまで基準値であり、別途ユーザにより指示される時間伸縮パラメータを乗算することにより適宜伸縮され、話速変換を実現することができる。

10

【0074】

促音の場合、その直後の音節の第1音素と同じものを、直後の音節の第1音素の直前に加える。第1音素の発音時間は0.0625秒であるため、先行する音節の発音終了時刻から0.1875秒後に促音のノートオン時刻を設定することになり、促音のノートオフ時刻と、直後の音節の第1音素のノートオン時刻が同一となる。尚、上記の0.0625秒、0.1875秒という時間数値も同様にあくまで基準値であり、別途ユーザにより指示される時間伸縮パラメータを乗算することにより適宜伸縮され、話速変換を実現することができる。

20

【0075】

拗音の場合、直前の子音の第1音素の直後に加える。したがって、直前の子音の第1音素のノートオフ時刻と、拗音のノートオン時刻が同一となるように設定する。拗音の音節の構成自体は子音と同じであるので、拗音の第1音素のノートオフ時刻および第2音素のノートオン時刻は、第1音素のノートオン時刻の0.0625秒後であり、拗音の第2音素のノートオフ時刻は、そのノートオン時刻の0.1875秒後となる。尚、上記の0.0625秒、0.1875秒という時間数値も同様にあくまで基準値であり、別途ユーザにより指示される時間伸縮パラメータを乗算することにより適宜伸縮され、話速変換を実現することができる。

30

【0076】

音素編集処理手段50は、読み込んだ合成指示データ内の音節識別情報単位で音素の合成処理を行っていき、処理が終わった音節単位で順に、合成音声データ(MIDIデータ)を、音声出力手段60に渡していく。音声出力手段60は、音素編集処理手段50から受け取ったMIDIデータを順に再生していく。以上のようにして、音声合成装置は、読み込んだ合成指示データに従って音声の再生が可能となる。

【0077】

五線譜として出力する場合は、合成音声データを印刷手段70により五線譜データに変換した後、印刷出力する。また、上記の例のように、合成指示データに従って音声合成をリアルタイムで行い、音声再生したり、五線譜出力することも可能であるが、この音声合成装置では、音素編集処理手段50による処理結果であるMIDIデータを合成音声データ記憶手段13に蓄積し、別途このMIDIデータをMIDI再生装置により音声再生するようにしても良い。MIDIデータを記憶装置に蓄積する方法としては、SMF(Standard MIDI File)形式ファイルを用いると、市販の種々の音楽関係ソフトウェアに渡すことができ、作成されたMIDIデータからは、市販の楽譜作成ツールを用いて、楽譜を作成することができる。この場合、楽譜は、SMF形式に記録されていた音素符号を基にして作成される。そして、作成された楽譜を印刷装置から出力すれば、読みやすい楽譜として、楽器演奏の際に利用することができる。

40

50

【0078】

上述の通り、音素編集処理手段50は、合成指示データ内の音節識別情報で音素符号データベース13aから対応する音素符号を抽出し、MIDIイベントを作成する際、そのノートナンバーについては音素符号データベース13aに収録されている当該音素符号を構成する各音符のノートナンバーに対して、オプション的に別途ユーザにより指示される音高オフセットパラメータを加算し適宜上下させ、ピッチ変換を行えるようにしてある。この場合は、合成音声データ全体のピッチを上下させるのではなく、母音音素に限定して上下させるようにする。また、合成指示データ内の音節識別情報とともに音高オフセットパラメータを音節ごとに定義すれば、各音節ごとに構成される母音音素のピッチを個別に上下させることもできる。すなわち、あらかじめ作成した旋律の隣接音符間での音高変化（音程情報）を、合成指示データ内の音節識別情報とともに定義される音高オフセットパラメータとして与えれば、歌声合成を実現することができる。

10

【産業上の利用可能性】

【0079】

本発明は、イベントや余興目的に行われる人間の音声再生を模倣した音楽作品制作・作曲の支援産業に利用することができる。また、エンターテインメント分野において、電子楽器を主体とした玩具（ロボット、ぬいぐるみを含む）、玩具型のアコースティック楽器（室内装飾用のミニチュアピアノ）、オルゴール、携帯電話の着信メロディ等の音階再生媒体に対して音声合成機能を付加する産業に利用することができる。また、SMF（Standard MIDI File）等によるMIDI音楽コンテンツ配布時における著作権保護等の産業に利用することができる。

20

【符号の説明】

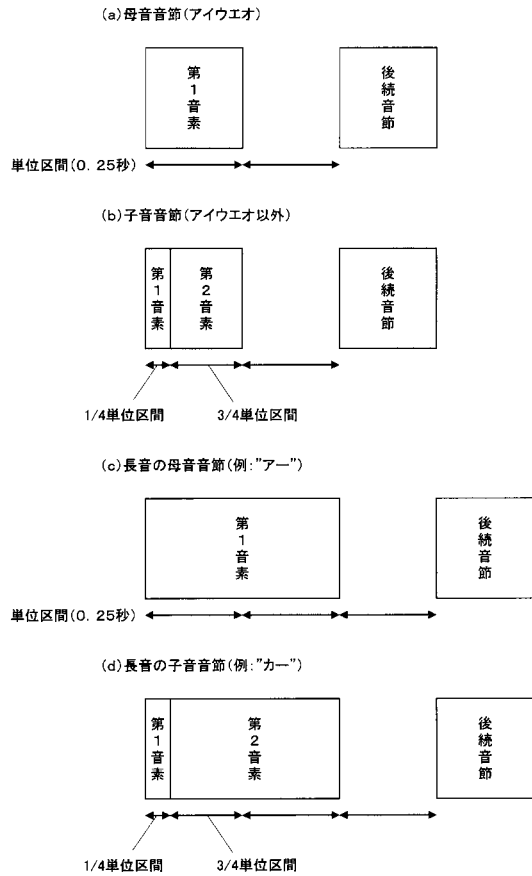
【0080】

- 10・・・記憶手段
- 11・・・音節録音データ記憶部
- 12・・・音節符号記憶部
- 13・・・音素符号記憶部
- 13a・・・音素符号データベース
- 14・・・合成音声データ記憶手段
- 20・・・処理制御手段
- 21・・・音節波形拡大手段
- 22・・・符号コード群作成手段
- 23・・・音節符号作成手段
- 24・・・音高別エネルギー算出手段
- 25・・・音高別エネルギー統合手段
- 26・・・符号変換手段
- 30・・・符号表示手段
- 50・・・音素編集処理手段
- 60・・・音声出力手段
- 70・・・印刷手段

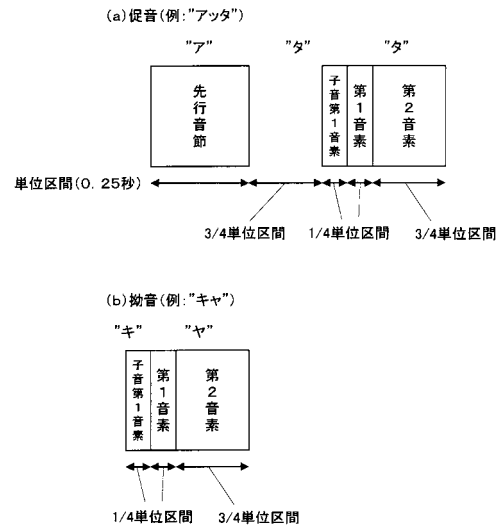
30

40

【図1】

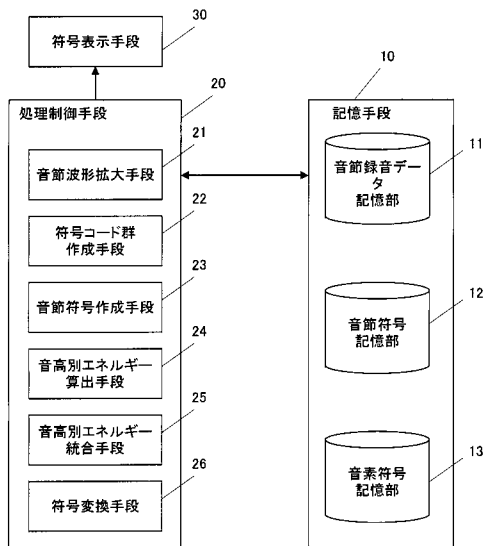


【図2】



【図3】

音素符号変換装置



【図4】

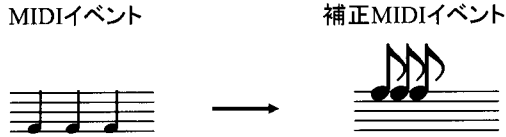


【 図 5 】

(a) 時間軸方向にQ倍に拡大



(b) 周波数をQ倍、時間軸を1/Q



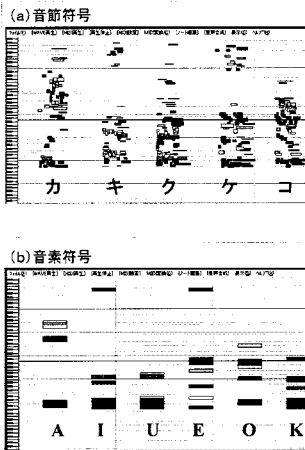
【 図 6 】

変換テーブル

	K	S	T	N	H	M	R	G	Z	D	B	P	Y	W
A	カ	サ	タ	ナ	ハ	マ	ラ	ガ	ザ	ダ	バ	パ	ヤ	ワ
I	キ	シ	チ	ニ	ヒ	ミ	リ	ギ	ジ	チ	ビ	ピ		
U	ク	ス	ツ	ヌ	フ	ム	ル	グ	ズ	ツ	ブ	ユ		
E	ケ	セ	テ	ネ	ヘ	メ	レ	ゲ	ゼ	デ	ベ			
O	コ	ソ	ト	ノ	ホ	モ	ロ	ゴ	ゾ	ド	ボ	ポ	ヨ	ヲ
n				ン										

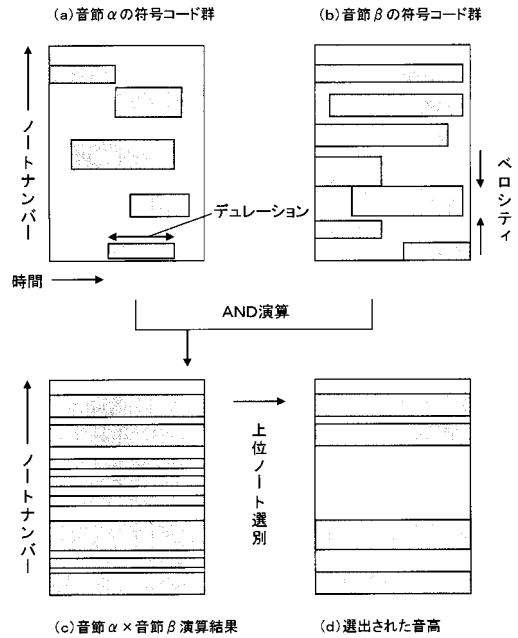
音節: カタカナ (71個)
音素: アルファベット (20個)

【 図 8 】



【 図 7 】

符号コード群の変化の様子



【 図 9 】

日本語音素MIDIコード・データベース例(男声)1

[A]	[I]	[U]	[E]	[O]	[K]	[S]	[T]	[N]	[H]
F5	A#6	D3	A#3	F#4	F4	B6	D#4	C4	D4
E5	C#3	C#3	A3	F4	A#3	A6	F3	D#3	A3
D#5	A#2	A#2	G#3	E4	C3	C4	D#3	G#2	F3
A4	D2	A2	E3	G#3	B2	B3	B2	G2	E3
G#4	C#2	D2	G#2	G3	A#2	F3	G#2	F#2	D#3
C2	C2	C#2	D2	F#3	G2	F2	F#2	F2	F#2
B1	B1	C2	C#2	C2	C#2	E2	E2	E2	E2
A#1	A#1	B1	A1	B1	A1	D#2	A1	A1	A1

【 図 10 】

日本語音素MIDIコード・データベース例(男声)2

[M]	[R]	[G]	[Z]	[D]	[B]	[P]	[Y]	[W]	[n]
B3	B3	C#4	C#4	D4	C4	C4	C#4	E4	D3
D#3	D#3	B3	C4	F3	B3	A#3	B3	D4	B2
C3	C3	F3	C3	C3	F3	A3	A3	C#4	A2
G2	B2	C3	B2	G#2	D#3	F3	F3	B3	G#2
F#2	G#2	B2	G2	G2	C3	G2	C3	G#3	D2
F2	G2	G2	F#2	F2	B2	E2	G#2	G3	C#2
E2	F2	F#2	E2	E2	G2	D#2	D#2	C2	C2
D#2	A1	E2	D#2	A1	E2	A1	A1	A#1	B1

【図 1 1】

日本語音素MIDIコード・データベース例(女声)1

[A]	[I]	[U]	[E]	[O]	[K]	[S]	[T]	[N]	[H]
B5	A6	E5	G6	B4	D#5	G#6	G#6	C3	F4
A#5	G#6	D#5	D#6	A#4	B4	E4	D6	B2	E4
F#5	G6	D#4	B4	A4	D#4	D#4	G#5	A#2	D#4
D#5	F#6	C#4	F4	G#4	C4	B3	F4	A2	B3
C#5	G3	D#3	E4	D4	A3	E3	E4	G#2	A#3
C5	F3	D3	D4	C#4	D#3	D3	D#4	F2	E3
B4	E3	C#3	C#4	A#3	D3	C#3	E3	C2	D#3
G#4	D#3	A#2	C4	A3	C#3	B2	G#2	A#1	B2

【図 1 2】

日本語音素MIDIコード・データベース例(女声)2

[M]	[R]	[G]	[Z]	[D]	[B]	[P]	[Y]	[W]	[n]
F#4	G#6	G6	A6	G#6	G4	E5	D#6	D#5	B5
C3	D#4	F#6	G#6	D#4	D#4	E4	F#4	D5	A#5
B2	C4	E4	G#5	C4	B3	D#4	E4	B4	C4
A#2	B3	D4	D3	B3	B2	D4	C4	A#4	A#3
G#2	C3	C4	C3	D3	A2	B3	B3	E4	D#3
F#2	B2	B3	A2	C3	G#2	A#3	E3	D#4	D3
F2	A2	D#3	G#2	G2	G2	E3	A2	D4	C3
C#2	G#2	B2	C2	C2	B1	D#3	G#2	C4	B2

【図 1 3】

日本語音素の五線譜化の例(男声)

Figure 13 shows two staves of musical notation for male voice. The first staff is for the syllables 'A I U E O K S T' and the second for 'N H M R G Z D B P Y W n'. The notation includes pitch contours and chord symbols above the notes. A small note '(作曲者)' is present in the top right corner.

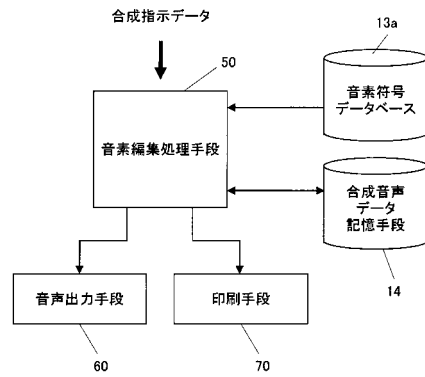
【図 1 4】

日本語音素の五線譜化の例(女声)

Figure 14 shows three staves of musical notation for female voice. The first staff is for 'A I U E O K S T', the second for 'N H M R G Z D B', and the third for 'P Y W n [Women]'. The notation includes pitch contours and chord symbols. A small note '(作曲者)' is present in the top right corner.

【図 1 5】

音声合成装置



フロントページの続き

(72)発明者 茂出木 敏雄
東京都新宿区市谷加賀町一丁目1番1号 大日本印刷株式会社内

審査官 山下 剛史

(56)参考文献 特開平4 - 349497 (JP, A)
特開2004 - 294816 (JP, A)
特開2006 - 139158 (JP, A)
特開平10 - 247099 (JP, A)
茂出木敏雄, "音声MIDIコードを用いたSMFファイルへの情報埋め込み手法", 電子情報通信学会2009年総合大会講演論文集, 2009年 3月, 情報・システム2, pp.S-35~S-36

(58)調査した分野(Int.Cl., DB名)

G10L 13/00 - 13/10, 19/00, 25/00 - 25/93
G10G 3/04
G10H 1/00