

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2006-251375
(P2006-251375A)

(43) 公開日 平成18年9月21日(2006.9.21)

(51) Int. Cl. F I テーマコード (参考)
G 1 O L 21/04 (2006.01) G 1 O L 21/04 1 2 O C
G 1 O L 21/02 (2006.01) G 1 O L 21/02 4 0 0

審査請求 未請求 請求項の数 8 O L (全 15 頁)

(21) 出願番号	特願2005-67907 (P2005-67907)	(71) 出願人	000004075 ヤマハ株式会社 静岡県浜松市中沢町10番1号
(22) 出願日	平成17年3月10日 (2005.3.10)	(74) 代理人	100098084 弁理士 川▲崎▼ 研二
		(72) 発明者	劔持 秀紀 静岡県浜松市中沢町10番1号 ヤマハ株式会社内
		(72) 発明者	吉岡 靖雄 静岡県浜松市中沢町10番1号 ヤマハ株式会社内
		(72) 発明者	ジョルディ ボナダ スペイン バルセロナ パセイ デ シル コンバルーラシオ 8

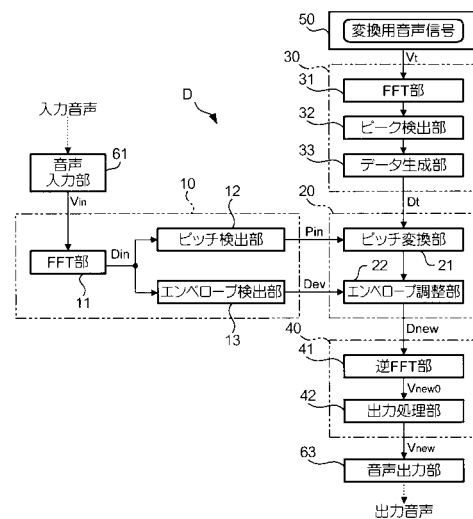
(54) 【発明の名称】 音声処理装置およびプログラム

(57) 【要約】

【課題】 入力音声を簡易な構成によって多人数での合唱音や合奏音に変換する。

【解決手段】 ピッチ検出部12は、音声入力部61から供給される入力音声信号V_{in}のピッチP_{in}を検出する。エンベロープ検出部13は、入力音声信号V_{in}のスペクトルエンベロープを検出する。スペクトル取得手段30は、並列に発生された複数の音声を含む変換用音声の周波数スペクトルを取得する。ピッチ変換部21は、スペクトル取得手段30が取得した周波数スペクトルの各ピークの周波数をピッチP_{in}に応じて変化させる。エンベロープ調整部22は、ピッチ変換部21による処理後の周波数スペクトルのスペクトルエンベロープをエンベロープ検出部13が検出したスペクトルエンベロープと略一致するように調整する。音声生成手段40は、エンベロープ調整部22による調整後の周波数スペクトルから出力音声信号V_{new}を生成する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出手段と、
並列に発生した複数の音声を含む変換用音声の周波数スペクトルである変換用スペクトルを取得するスペクトル取得手段と、

前記スペクトル取得手段が取得した変換用スペクトルのスペクトルエンベロープを前記エンベロープ検出手段が検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整手段と、

前記エンベロープ調整手段による調整後の変換用スペクトルから出力音声信号を生成する音声生成手段と

を具備する音声処理装置。

10

【請求項 2】

前記入力音声のピッチを検出するピッチ検出手段と、

前記スペクトル取得手段が取得した変換用スペクトルに含まれる各ピークの周波数を前記ピッチ検出手段が検出したピッチに応じて変化させるピッチ変換手段と

を具備し、

前記エンベロープ調整手段は、前記ピッチ変換手段による処理後の変換用スペクトルのスペクトルエンベロープを調整する

請求項 1 に記載の音声処理装置。

【請求項 3】

前記ピッチ変換手段は、前記ピッチ検出手段が検出したピッチに応じて変換用スペクトルを周波数軸の方向に伸長または縮小する

請求項 2 に記載の音声処理装置。

20

【請求項 4】

前記ピッチ変換手段は、前記変換用スペクトルにおける各ピークの周波数を含むスペクトル分布領域の各々を前記ピッチ検出手段が検出したピッチに応じて周波数軸の方向に移動させる

請求項 2 に記載の音声処理装置。

【請求項 5】

前記入力音声のピッチを検出するピッチ検出手段を具備し、

前記スペクトル取得手段は、各々のピッチが相違する複数の変換用音声のうち前記ピッチ検出手段が検出したピッチに近似するピッチの変換用音声の変換用スペクトルを取得する

請求項 1 に記載の音声処理装置。

30

【請求項 6】

入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出手段と、

変換用音声の周波数スペクトルである第 1 変換用スペクトル、および、前記第 1 変換用スペクトルが示す変換用音声と略同一ピッチの音声の周波数スペクトルであり前記第 1 変換用スペクトルよりも各ピークの帯域幅が広い第 2 変換用スペクトルの何れかを取得するスペクトル取得手段と、

40

前記スペクトル取得手段が取得した変換用スペクトルのスペクトルエンベロープを前記エンベロープ検出手段が検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整手段と、

前記エンベロープ調整手段による調整後の変換用スペクトルから出力音声信号を生成する音声生成手段と

を具備する音声処理装置。

【請求項 7】

コンピュータに、

入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出処理と、

並列に発生した複数の音声を含む変換用音声の周波数スペクトルである変換用スペクトル

50

ルを取得するスペクトル取得処理と、

前記スペクトル取得処理にて取得した変換用スペクトルのスペクトルエンベロープを前記エンベロープ検出処理にて検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整処理と、

前記エンベロープ調整処理後の変換用スペクトルから出力音声信号を生成する音声生成処理と

を実行させるプログラム。

【請求項 8】

コンピュータに、

入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出処理と、

変換用音声の周波数スペクトルである第 1 変換用スペクトル、および、前記第 1 変換用スペクトルが示す変換用音声と略同一ピッチの音声の周波数スペクトルであり前記第 1 変換用スペクトルよりも各ピークの帯域幅が広い第 2 変換用スペクトルの何れかを取得するスペクトル取得処理と、

前記スペクトル取得処理にて取得した変換用スペクトルのスペクトルエンベロープを前記エンベロープ検出処理にて検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整処理と、

前記エンベロープ調整処理後の変換用スペクトルから出力音声信号を生成する音声生成処理と

を実行させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声の特性を変化させる技術に関する。

【背景技術】

【0002】

利用者が発声した音声（以下「入力音声」という）に音楽的な効果を付与するための種々の技術が従来から提案されている。例えば特許文献 1 には、入力音声のピッチを変換することによって生成された協和音（入力音声と和音を構成する音声）を入力音声と加算して出力する技術が開示されている。この構成によれば、実際の発声者がひとりであっても、恰も複数人にて別個の旋律を合唱しているかのような音声を出力することができる。また、例えば入力音声を楽器の演奏音とすれば、複数の楽器によって別個の旋律を合奏しているかのような音声生成される。

【特許文献 1】特開平 10 - 78776 号公報（段落 0013 および図 1 参照）

【発明の開示】

【発明が解決しようとする課題】

【0003】

ところで、合唱や合奏の形態としては、各歌唱者や演奏者が別個の旋律を歌唱または演奏する形態（いわゆるコーラス）のほか、複数の歌唱者や演奏者が同一の旋律を歌唱または演奏するユニゾンと呼ばれる形態がある。特許文献 1 に記載された構成においては、入力音声のピッチを変換することによって協和音が生じられるため、複数人が別個の旋律を歌唱ないし演奏したときの音声を生成することはできるものの、複数人が共通の旋律を歌唱または演奏するユニゾンの効果を入力音声に付与することはできない。なお、特許文献 1 に記載された構成においても、例えば入力音声のピッチを変更せずに音響的な特性（声質）のみを変換した音声を入力音声とともに出力すれば、複数人が共通の旋律を歌唱または演奏しているかのような効果を付与することも一応は可能である。しかしながら、この場合には、ユニゾンを構成する音声ごとに入力音声の特性を変換するための仕組みを用意することが不可欠となる。したがって、多人数によるユニゾンを実現しようとするれば、DSP（Digital Signal Processor）などのハードウェアによって入力音声の特性が変換される構成においてはその回路規模が肥大化し、この変換がソフトウェアによって実現され

10

20

30

40

50

る構成においては演算装置の処理負荷が過大になるといった問題がある。本発明は、このような事情に鑑みてなされたものであり、入力音声を簡易な構成によって多人数での合唱音や合奏音に変換することを目的としている。

【課題を解決するための手段】

【0004】

この課題を解決するために、本発明に係る音声処理装置は、入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出手段と、並列に発生した複数の音声を含む変換用音声の周波数スペクトルである変換用スペクトルを取得するスペクトル取得手段と、スペクトル取得手段が取得した変換用スペクトルのスペクトルエンベロープをエンベロープ検出手段が検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整手段と、エンベロープ調整手段による調整後の変換用スペクトルから出力音声信号を生成する音声生成手段とを具備することを特徴としている。なお、本発明にいう「音声」には、人間が発声した音声や楽器の演奏音といった種々の音響が含まれる。

10

この構成によれば、並列に発生した複数の音声を含む変換用音声の変換用スペクトルのエンベロープが入力音声のスペクトルエンベロープと略一致するように調整されるから、入力音声と同様の音韻をもった複数の音声（すなわち合奏音や合唱音）を示す出力音声信号を生成することができる。しかも、複数の音声の各々について入力音声の特性を変換するための仕組みは原理的に不要であるから、音声処理装置の構成は特許文献1の構成と比較して大幅に簡素化される。

なお、エンベロープ検出手段が検出したスペクトルエンベロープと変換用スペクトルのスペクトルエンベロープとが「略一致する」とは、エンベロープ調整手段による調整後の周波数スペクトルから生成された出力音声信号に基づいて実際に音声が発音されたときに、その音声の音韻が聴感上において入力音声の音韻と同一であると知覚される程度に近似（理想的には一致）していることを意味する。したがって、入力音声のスペクトルエンベロープとエンベロープ調整手段による調整後のスペクトルエンベロープとは厳密な意味で完全に一致している必要は必ずしもない。

20

【0005】

本発明に係る音声処理装置において、音声生成手段が生成した出力音声信号は、例えばスピーカやイヤホンなどの放音機器に供給されて音声（以下「出力音声」という）として出力される。ただし、この出力音声信号が利用の態様は任意である。例えば、出力音声信号が記録媒体に記憶されたうえで、当該記憶手段を再生する他の装置にて出力音声が発音される態様としてもよいし、出力音声信号が通信回線を介して他の装置に送信されて当該装置にて音声として再生される態様としてもよい。

30

【0006】

音声生成手段が生成する出力音声信号のピッチ（換言すると出力音声のピッチ）は、入力音声のピッチとは無関係なピッチであってもよいが、より好適には入力音声に応じたピッチ（例えば入力音声と略一致するピッチや入力音声と協和音を構成するピッチ）とされる。この好適な態様においては、例えば、入力音声のピッチを検出するピッチ検出手段と、スペクトル取得手段が取得した変換用スペクトルに含まれる各ピークの周波数をピッチ検出手段が検出したピッチに応じて変化させるピッチ変換手段とが更に設けられ、エンベロープ調整手段は、ピッチ変換手段による処理後の変換用スペクトルのスペクトルエンベロープを調整する。この態様によれば、出力音声信号が入力音声に対応したピッチとされるから、この出力音声信号に基づいて放音される音声を聴感上において心地よい音声とすることができる。

40

より具体的な態様において、ピッチ変換手段は、ピッチ検出手段が検出したピッチに応じて変換用スペクトルを周波数軸の方向に伸長または縮小する。この態様によれば、変換用スペクトルの各周波数に対して入力音声のピッチに応じた数値を乗算するという簡易な処理によって変換用スペクトルのピッチを調整することができる。また、他の態様において、ピッチ変換手段は、変換用スペクトルにおける各ピークの周波数を含むスペクトル分布領域（例えばピークの周波数を中心とする所定幅の周波数帯域）の各々をピッチ検出手

50

段が検出したピッチに応じた周波数軸の方向に移動させる（図8参照）。この態様によれば、変換用スペクトルの各ピークの周波数を所期の周波数に合致させることができるから、変換用スペクトルのピッチを所望のピッチに精度よく調整することができる。

【0007】

もっとも、出力音声を入力音声に応じたピッチとするための構成は任意である。例えば、入力音声のピッチを検出するピッチ検出手段を設けたうえで、スペクトル取得手段が、各々のピッチが相違する複数の変換用音声のうちピッチ検出手段が検出したピッチに近似（理想的には一致）するピッチの変換用音声の変換用スペクトルを取得する態様としてもよい（図6参照）。この態様によれば、変換用音声のピッチを変換するための仕組みを不要とすることができる。ただし、変換用スペクトルのピッチを変換する構成と、各々のピッチが異なる複数の変換用音声の何れかを選択する構成とを組み合わせてもよい。例えば、各々が異なるピッチに対応する複数の変換用スペクトルのうち入力音声のピッチに近似するピッチに対応した変換用スペクトルをスペクトル取得手段が取得し、この選択した変換用スペクトルのピッチをピッチ変換手段がピッチデータに応じて変換する構成も採用される。

10

【0008】

ところで、複数の歌唱者や演奏者から略同一のピッチにて同時に（並列に）発せられた音声の周波数スペクトルは、その各ピークの帯域幅（例えば図3に示される帯域幅 W_2 ）が、単一の歌唱者や演奏者から発せられた音声の周波数スペクトルにおける各ピークの帯域幅（例えば図2に示される帯域幅 W_1 ）よりも広い場合が多い。いわゆるユニゾンにおいては、各歌唱者や各演奏者の音声のピッチが厳密には一致していないからである。このような観点から、本発明に係る音声処理装置は、入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出手段と、変換用音声の周波数スペクトルである第1変換用スペクトル、および、第1変換用スペクトルが示す変換用音声と略同一ピッチの音声の周波数スペクトルであり第1変換用スペクトルよりも各ピークの帯域幅が広い第2変換用スペクトルの何れかを取得するスペクトル取得手段と、スペクトル取得手段が取得した変換用スペクトルのスペクトルエンベロープをエンベロープ検出手段が検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整手段と、エンベロープ調整手段による調整後の変換用スペクトルから出力音声信号を生成する音声生成手段とを具備する構成としても特定される。なお、この構成の具体例は第2実施形態（図5）として後述される。

20

30

この構成によれば、出力音声信号を生成するための周波数スペクトルとして第1変換用スペクトルおよび第2変換用スペクトルの何れかが選択されるから、第1変換用スペクトルに応じた特性の出力音声信号と第2変換用スペクトルに応じた特性の出力音声信号とを選択的に生成することができる。例えば、第1変換用スペクトルが選択された場合には単一の歌唱者または演奏者から発せられた出力音声信号を生成することができ、第2変換用スペクトルが選択された場合には複数の歌唱者や演奏者から発せられた出力音声信号を生成することができる。なお、ここでは第1変換用スペクトルと第2変換用スペクトルとが特定されているが、更に他の変換用スペクトルが選択手段による選択の対象とされる構成としてもよい。例えば、それぞれ各ピークの帯域幅が相違する複数の変換用スペクトルを記憶手段に記憶させておき、このうちの何れかを選択手段が選択して出力音声信号の生成に利用するといった構成も採用される。

40

【0009】

本発明に係る音声処理装置は、音声処理に専用されるDSPなどのハードウェアによって実現されるほか、パーソナルコンピュータなどのコンピュータとプログラムとの協働によっても実現される。このプログラムは、入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出処理と、並列に発生した複数の音声を含む変換用音声の周波数スペクトルである変換用スペクトルを取得するスペクトル取得処理と、スペクトル取得処理にて取得した変換用スペクトルのスペクトルエンベロープをエンベロープ検出処理にて検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整処理と

50

、エンベロープ調整処理後の変換用スペクトルから出力音声信号を生成する音声生成処理とをコンピュータに実行させる内容となる。また、他の態様に係るプログラムは、入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出処理と、変換用音声の周波数スペクトルである第1変換用スペクトル、および、第1変換用スペクトルが示す変換用音声と略同一ピッチの音声の周波数スペクトルであり第1変換用スペクトルよりも各ピークの帯域幅が広い第2変換用スペクトルの何れかを取得するスペクトル取得処理と、スペクトル取得処理にて取得した変換用スペクトルのスペクトルエンベロープをエンベロープ検出処理にて検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整処理と、エンベロープ調整処理後の変換用スペクトルから出力音声信号を生成する音声生成処理とをコンピュータに実行させる内容となる。これらのプログラムは、コンピュータ読取り可能な記録媒体（例えばCD-ROM）に格納された態様にて利用者に提供されてコンピュータにインストールされるほか、ネットワークを介した配信の形態にてサーバ装置から提供されてコンピュータにインストールされる。

10

20

30

40

50

【0010】

また、本発明は、入力音声进行处理するための方法としても特定される。この方法は、入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出過程と、並列に発生した複数の音声を含む変換用音声の周波数スペクトルである変換用スペクトルを取得するスペクトル取得過程と、スペクトル取得過程にて取得した変換用スペクトルのスペクトルエンベロープをエンベロープ検出過程にて検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整過程と、エンベロープ調整過程における調整後の変換用スペクトルから出力音声信号を生成する音声生成過程とを有する。他の観点に基づく音声処理方法は、入力された入力音声のスペクトルエンベロープを検出するエンベロープ検出過程と、変換用音声の周波数スペクトルである第1変換用スペクトル、および、第1変換用スペクトルが示す変換用音声と略同一ピッチの音声の周波数スペクトルであり第1変換用スペクトルよりも各ピークの帯域幅が広い第2変換用スペクトルの何れかを取得するスペクトル取得過程と、スペクトル取得過程にて取得した変換用スペクトルのスペクトルエンベロープをエンベロープ検出過程にて検出したスペクトルエンベロープと略一致するように調整するエンベロープ調整過程と、エンベロープ調整過程における調整後の変換用スペクトルから出力音声信号を生成する音声生成過程とを有する。

【発明の効果】

【0011】

以上のように、本発明によれば、簡易な構成によって多人数での合唱や合奏を実現することができる。

【発明を実施するための最良の形態】

【0012】

< A : 第1実施形態 >

まず、図1を参照して、本発明の第1実施形態に係る音声処理装置の構成および動作を説明する。同図に示される音声処理装置の各部は、例えばCPU (Central Processing Unit) などの演算回路がプログラムを実行することによって実現されてもよいし、DSP など音声処理に専用されるハードウェアによって実現されてもよい。後述する各実施形態においても同様である。

【0013】

図1に示されるように、音声処理装置Dは、周波数分析手段10と、スペクトル変換手段20と、スペクトル取得手段30と、音声生成手段40と、記憶手段50とを有する。このうち周波数分析手段10には音声入力部61が接続される。この音声入力部61は、利用者が発する入力音声に応じた信号（以下「入力音声信号」という）Vinを出力する手段であり、例えば、入力音声の時間軸上における波形を表わすアナログの電気信号を出力する収音機器（マイクロホン）と、この電気信号をデジタルの入力音声信号Vinに変換するA/D変換器とを有する。

【0014】

周波数分析手段 10 は、音声入力部 61 から供給される入力音声信号 V_{in} のピッチ P_{in} およびスペクトルエンベロープ $E_{V_{in}}$ を特定する手段であり、FFT (Fast Fourier Transform) 部 11 とピッチ検出部 12 とエンベロープ検出部 13 とを有する。このうち FFT 部 11 は、音声入力部 61 から供給される入力音声信号 V_{in} を所定の時間長 (例えば 5ms ないし 10ms) のフレームに切り出し、各フレームの入力音声信号 V_{in} に対して FFT 処理を含む周波数分析を実行して周波数スペクトル (以下「入力スペクトル」という) $S_{P_{in}}$ を検出する。入力音声信号 V_{in} の各フレームは時間軸上において相互に重なり合うように選定される。これらのフレームは簡易的には同一の時間長とされるが、入力音声信号 V_{in} のピッチ P_{in} (後述するようにピッチ検出部 12 によって検出される) に応じて時間長が変化する構成としてもよい。図 2 には、ひとりの利用者が発声した入力音声のうちひとつのフレームについて特定された入力スペクトル $S_{P_{in}}$ が例示されている。この場合の入力スペクトル $S_{P_{in}}$ は、基音および倍音に相当する各周波数においてスペクトル強度 M の局所的なピーク p が極めて狭い帯域幅 W_1 にて現れる。FFT 部 11 は、入力音声信号 V_{in} の入力スペクトル $S_{P_{in}}$ を表わすデータ (以下「入力スペクトルデータ」という) D_{in} をフレームごとにピッチ検出部 12 とエンベロープ検出部 13 とに出力する。入力スペクトルデータ D_{in} は複数の単位データを含む。各単位データは、周波数軸上に所定の間隔ごとに選定された複数の周波数 F_{in} の各々と当該周波数における入力スペクトル $S_{P_{in}}$ のスペクトル強度 M_{in} とが組み合わされたデータである。

10

【0015】

図 1 に示されるピッチ検出部 12 は、FFT 部 11 から供給される入力スペクトルデータ D_{in} に基づいて入力音声のピッチ P_{in} を検出する手段である。更に詳述すると、ピッチ検出部 12 は、図 2 に示されるように、入力スペクトルデータ D_{in} が示す入力スペクトル $S_{P_{in}}$ のうち基音に相当するピーク p (すなわち周波数が最小であるピーク p) の周波数をピッチ P_{in} として検出する。一方、エンベロープ検出部 13 は、入力音声のスペクトルエンベロープ (スペクトル包絡) $E_{V_{in}}$ を検出する手段である。スペクトルエンベロープ $E_{V_{in}}$ は、図 2 に示されるように、入力スペクトル $S_{P_{in}}$ のピーク p を連結した包絡線である。このスペクトルエンベロープ $E_{V_{in}}$ を検出する方法としては、例えば、入力スペクトル $S_{P_{in}}$ のうち周波数軸上において相互に隣接するピーク p の間隙を直線的に補間することによってスペクトルエンベロープ $E_{V_{in}}$ を折線として検出する方法や、各ピーク p を通過する曲線を 3 次のスプライン補間など各種の補間処理によって算定してスペクトルエンベロープ $E_{V_{in}}$ を検出する方法などが採用される。エンベロープ検出部 13 は、図 2 に示されるように、こうして検出したスペクトルエンベロープ $E_{V_{in}}$ を示すデータ (以下「エンベロープデータ」という) D_{ev} を出力する。エンベロープデータ D_{ev} は、入力スペクトルデータ D_{in} と同様に複数の単位データ U_{ev} を含む。各単位データ U_{ev} は、周波数軸上に所定の間隔ごとに選定された複数の周波数 F_{in} (F_{in1}, F_{in2}, \dots) の各々と当該周波数 F_{in} におけるスペクトルエンベロープ $E_{V_{in}}$ のスペクトル強度 M_{ev} (M_{ev1}, M_{ev2}, \dots) とが組み合わされたデータである。

20

30

【0016】

次に、図 1 に示されるスペクトル変換手段 20 は、入力音声の特性を変化させた出力音声の周波数スペクトル (以下「出力スペクトル」という) $S_{P_{new}}$ を示すデータ (以下「新規スペクトルデータ」という) D_{new} を生成する手段である。本実施形態におけるスペクトル変換手段 20 は、予め用意された特定の音声 (以下「変換用音声」という) の周波数スペクトル (以下「変換用スペクトル」という) S_{P_t} と入力音声のスペクトルエンベロープ $E_{V_{in}}$ とに基づいて出力音声の周波数スペクトル $S_{P_{new}}$ を特定する。なお、周波数スペクトル $S_{P_{new}}$ を生成する手順については後述する。

40

【0017】

一方、スペクトル取得手段 30 は、変換用スペクトル S_{P_t} を取得するための手段であり、FFT 部 31 とピーク検出部 32 とデータ生成部 33 とを有する。このうち FFT 部 31 には、記憶手段 50 (例えばハードディスク装置) から読み出された変換用音声信号 V_t が供給される。この変換用音声信号 V_t は、変換用音声の波形を特定の区間にわたって

50

表わす時間領域の信号であり、予め記憶手段50に格納されている。FFT部31は、入力音声に係る手順と同様に、記憶手段50から順次に供給される変換用音声信号Vtを所定の時間長のフレームに切り出し、各フレームの変換用音声信号Vtに対してFFT処理を含む周波数分析を実行することによって変換用スペクトルSPtを検出する。一方、ピーク検出部32は、FFT部31によって特定された変換用スペクトルSPtのピークptを検出してその周波数を特定する。ピークptを検出する方法としては、例えば、周波数軸上において近接する所定数のピークのうちスペクトル強度が最大となるものをピークptとして検出する方法が採用される。

【0018】

本実施形態においては、多数の発声者が略同一のピッチPtにて発声した音声（すなわち合唱や合奏といったユニゾンの音声）をマイクロホンなどの収音機器によって収音した信号が変換用音声信号Vinとして記憶手段50に記憶されている場合を想定する。このような変換用音声信号VtにFFT処理を施して得られる変換用スペクトルSPtは、図3に示されるように、変換用音声のピッチPtに応じた基音および倍音に相当する各周波数においてスペクトル強度Mの局所的なピークptが現れる点で図1の入力スペクトルSPinと共通するが、各ピークptの帯域幅W2が入力スペクトルSPinの各ピークpの帯域幅W1よりも広いという特性を有する。このようにピークptの帯域幅W2が広いのは、多数の発声者によって発声された各音声のピッチが完全には一致しないからである。

10

【0019】

図1に示されるデータ生成部33は、変換用スペクトルSPtを示すデータ（以下「変換用スペクトルデータ」という）Dtを生成するための手段である。変換用スペクトルデータDtは、図3に示されるように、複数の単位データUtと指示子Aとを含む。各単位データUtは、エンベロープデータDevと同様に、周波数軸上に所定の間隔ごとに選定された複数の周波数Ft（Ft1, Ft2, ...）の各々と当該周波数Ftにおける変換用スペクトルSPtのスペクトル強度Mt（Mt1, Mt2, ...）とが組み合わせられたデータ構造となっている。一方、指示子Aは、変換用スペクトルSPtのピークptを指示するためのデータ（例えばフラグ）であり、変換用スペクトルデータDtに含まれる総ての単位データUtのうちピーク検出部32によって検出されたピークptに対応する単位データUtに対して選択的に付加される。例えば、ピーク検出部32が周波数Ft3にピークptを検出した場合、図3に示されるように、周波数Ft3を含む単位データUtに指示子Aが付加され、これ以外の単位データUt（つまりピークpt以外の周波数に対応する単位データUt）に指示子Aは付加されない。

20

30

【0020】

図1に示されるように、スペクトル変換手段20は、ピッチ変換部21とエンベロープ調整部22とを有する。スペクトル取得手段30から出力された変換用スペクトルデータDtはピッチ変換部21に入力される。このピッチ変換部21は、変換用スペクトルデータDtが示す変換用スペクトルSPtの各ピークptの周波数を、ピッチ検出部12が検出したピッチPinに応じて変化させる手段である。本実施形態におけるピッチ変換部21は、変換用スペクトルデータDtが示す変換用音声のピッチPtがピッチ検出部12によって検出されたピッチPinと略一致するように変換用スペクトルSPtを变形する。この変換の具体的な手順について図4を参照して説明する。

40

【0021】

図4の部分(b)には、図3に示した変換用スペクトルSPtが図示されている。また、図4の部分(a)には、入力スペクトルSPin（図2に示したもの）が変換用スペクトルSPtとの対比のために併記されている。入力音声のピッチPinは利用者の発声に応じて変動するから、図4の部分(a)および部分(b)に示されるように、入力スペクトルSPinの各ピークpの周波数と変換用スペクトルSPtの各ピークptの周波数とは必ずしも一致しない。そこで、ピッチ変換部21は、変換用スペクトルSPtを周波数軸の方向に伸長または縮小することによって当該変換用スペクトルSPtの各ピークptの周波数を入力スペクトルSPinの各ピークpの周波数に合致させる。更に詳述すると、ピッチ変換

50

部 2 1 は、ピッチ検出部 1 2 が検出したピッチ P_{in} と変換用音声のピッチ P_t との比「 P_{in} / P_t 」を算定し、変換用スペクトルデータ D_t を構成する各単位データ U_t の周波数 F_t に対して当該比を乗算する。なお、変換用音声のピッチ P_t は、例えば、変換用スペクトル S_{P_t} の多数のピーク p_t のうち基音に相当するピーク p_t (すなわち周波数が最小であるピーク p_t) の周波数として特定される。この処理により、図 4 の部分 (c) に示されるように、変換用スペクトル S_{P_t} の各ピーク p_t は入力スペクトル $S_{P_{in}}$ の各ピーク p の周波数まで移動し、この結果として変換用音声のピッチ P_t は入力音声のピッチ P_{in} に略一致することになる。ピッチ変換部 2 1 は、こうしてピッチを変換した変換用スペクトル S_{P_t} を示す変換用スペクトルデータ D_t をエンベロープ調整部 2 2 に出力する。

【 0 0 2 2 】

エンベロープ調整部 2 2 は、この変換用スペクトルデータ D_t が示す変換用スペクトル S_{P_t} のスペクトル強度 M (換言すればスペクトルエンベロープ E_{V_t}) を調整することによって新規スペクトル $S_{P_{new}}$ を生成する手段である。更に詳述すると、エンベロープ調整部 2 2 は、図 4 の部分 (d) に示されるように、新規スペクトル $S_{P_{new}}$ のスペクトルエンベロープが、エンベロープ検出部 1 3 によって検出されたスペクトルエンベロープ $E_{V_{in}}$ と略一致するように、変換用スペクトル S_{P_t} のスペクトル強度 M を調整する。スペクトル強度 M を調整する方法の具体例は以下の通りである。

【 0 0 2 3 】

エンベロープ調整部 2 2 は、まず、変換用スペクトルデータ D_t のうち指示子 A が付加されたひとつの単位データ U_t を選定する。この単位データ U_t は、変換用スペクトル S_{P_t} の何れかのピーク p_t (以下では特に「注目ピーク p_t 」という) の周波数 F_t およびスペクトル強度 M_t を含む (図 3 参照)。次いで、エンベロープ調整部 2 2 は、エンベロープ検出部 1 3 から供給されるエンベロープデータ D_{ev} のうち注目ピーク p_t の周波数 F_t に近似または一致する周波数 F_{in} を含む単位データ U_{ev} を選定する。そして、エンベロープ調整部 2 2 は、この選定した単位データ U_{ev} に含まれるスペクトル強度 M_{ev} と注目ピーク p_t のスペクトル強度 M_t との比「 M_{ev} / M_t 」を算定し、注目ピーク p_t を中心とした所定の帯域に属する変換用スペクトル S_{P_t} の各単位データ U_t のスペクトル強度 M_t に対して当該比を乗算する。この一連の処理を変換用スペクトル S_{P_t} の総てのピーク p_t について繰り返すことにより、新規スペクトル $S_{P_{new}}$ は、図 4 の部分 (d) に示されるように、各ピークの頂点がスペクトルエンベロープ $E_{V_{in}}$ 上に位置する形状となる。エンベロープ調整部 2 2 は、この新規スペクトル $S_{P_{new}}$ を示す新規スペクトルデータ D_{new} を出力する。

【 0 0 2 4 】

ピッチ変換部 2 1 やエンベロープ調整部 2 2 による処理は入力音声信号 V_{in} を区分したフレームごとに行われる。ところで、変換用音声のフレーム数は記憶手段 5 0 に記憶された変換用音声信号 V_t の時間長に応じて制約されるのに対して入力音声のフレーム数は利用者による発声の期間に応じて変化するため、入力音声のフレーム数と変換用音声のフレーム数とは一致しない場合が多い。変換用音声のフレーム数が入力音声のフレーム数よりも多い場合には、余ったフレームに対応する変換用スペクトルデータ D_t を破棄すれば足りる。一方、変換用音声のフレーム数が入力音声のフレーム数よりも少ない場合には、総てのフレームに対応する変換用スペクトルデータ D_t の使用後に最初のフレームの変換用スペクトルデータ D_t を使用するという具合に、変換用スペクトルデータ D_t をループさせて使用すればよい。

【 0 0 2 5 】

次に、図 1 に示される音声生成手段 4 0 は、新規スペクトル $S_{P_{new}}$ に基づいて時間領域の出力音声信号 V_{new} を生成する手段であり、逆 FFT 部 4 1 と出力処理部 4 2 とを有する。このうち逆 FFT 部 4 1 は、エンベロープ調整部 2 2 からフレームごとに出力される新規スペクトルデータ D_{new} に対して逆 FFT 処理を施して時間領域の出力音声信号 V_{new0} を生成する。出力処理部 4 2 は、こうして生成されたフレームごとの出力音声信号 V_{new0} に時間窓関数を乗算し、これらを時間軸上において相互に重なり合うように連結して出力音声信号 V_{new} を生成する。この出力音声信号 V_{new} は音声出力部 6 3 に供給される。

10

20

30

40

50

音声出力部 63 は、出力音声信号 V_{new} をアナログの電気信号に変換する D/A 変換器と、この D/A 変換器からの出力信号に基づいて放音する放音機器（例えばスピーカやヘッドフォン）とを有する。

【0026】

以上に説明したように、本実施形態においては、多数の発声者によって並列に発せられた複数の音声を含む変換用音声のスペクトルエンベロープ E_{Vt} が入力音声のスペクトルエンベロープ E_{Vin} と略一致するように調整されるから、入力音声と同様の音韻をもった複数の音声（すなわち合唱音や合奏音）を示す出力音声信号 V_{new} を生成することができる。したがって、ひとりの利用者による音声や演奏音が入力音声とされた場合であっても、恰も多数の発声者や演奏者によって合唱や合奏が行なわれているかのような出力音声を音声出力部 63 から出力することができる。しかも、複数の音声の各々について入力音声の特性を変化させるための仕組みは原理的に不要である。したがって、音声処理装置 D の構成は特許文献 1 の構成と比較して大幅に簡素化される。さらに、本実施形態においては、入力音声のピッチ P_{in} に応じて変換用音声のピッチ P_t が変換されるから、任意のピッチの合唱音や合奏音を生成することができる。また、このピッチの変換が、変換用スペクトル S_{Pt} を周波数軸の方向に伸長するという簡素な処理（乗算処理）によって実現されるという利点もある。

10

【0027】

< B : 第 2 実施形態 >

次に、本発明の第 2 実施形態に係る音声処理装置について説明する。なお、本実施形態のうち第 1 実施形態と同様の要素については共通の符号を付してその説明を適宜に省略する。

20

【0028】

図 5 は、本実施形態に係る音声処理装置 D の構成を示すブロック図である。同図に示されるように、この音声処理装置 D は、記憶手段 50 の記憶内容およびスペクトル取得手段 30 の構成が第 1 実施形態の音声処理装置 D とは相違するが、他の要素は同様の構成である。本実施形態においては、第 1 変換用音声信号 V_{t1} と第 2 変換用音声信号 V_{t2} とが記憶手段 50 に記憶される。第 1 変換用音声信号 V_{t1} と第 2 変換用音声信号 V_{t2} とは、互いに略同一のピッチ P_t にて発せられた変換用音声を収録した信号である。ただし、第 1 変換用音声信号 V_{t1} は、図 2 に示した入力音声信号 V_{in} と同様に、単一の音声（ひとりの発声者からの音声やひとつの楽器からの演奏音）の波形を示す信号であるのに対し、第 2 変換用音声信号 V_{t2} は、第 1 実施形態の変換用音声信号 V_t と同様に、各々が並列に発せられた複数の音声（多数の発声者からの音声や多数の楽器からの演奏音）からなる変換用音声を収録した信号である。したがって、第 2 変換用音声信号 V_{t2} から特定される変換用スペクトル S_{Pt} の各ピークの帯域幅（図 3 に示す帯域幅 W_2 ）は、第 1 変換用音声信号 V_{t1} から特定される変換用スペクトル S_{Pt} の各ピークの帯域幅（図 1 に示す帯域幅 W_1 ）よりも広い。

30

【0029】

また、本実施形態におけるスペクトル取得手段 30 は FFT 部 31 の前段に選択部 34 を有する。この選択部 34 は、外部から供給される選択信号に基づいて、第 1 変換用音声信号 V_{t1} および第 2 変換用音声信号 V_{t2} の何れかを選択して記憶手段 50 から読み出す手段である。選択信号は、例えば、入力機器 67 に対する操作に応じて供給される。この選択部 34 によって読み出された変換用音声信号 V_t が FFT 部 31 に供給される。これ以後の構成および動作は第 1 実施形態と同様である。

40

【0030】

このように、本実施形態においては、第 1 変換用音声信号 V_{t1} および第 2 変換用音声信号 V_{t2} の何れかが選択的に新規スペクトル $S_{P_{new}}$ の生成に利用される。そして、第 1 変換用音声信号 V_{t1} が選択された場合には、入力音声の音韻と変換用音声の周波数特性とを兼ね備えた単一の出力音声が出力される一方、第 2 変換用音声信号 V_{t2} が選択された場合には、第 1 実施形態と同様に、入力音声の音韻を維持した多数の音声からなる出力音声が

50

出力される。すなわち、本実施形態においては、出力音声を単一の音声とするか複数の音声とするかを利用者が任意に選択することができる。

【0031】

なお、本実施形態においては入力機器67への操作に応じて変換用音声信号Vtが選択される構成を例示したが、この選択の基準となる要素は任意に変更される。例えば、所定の時間間隔にて発生するタイマ割込を契機として第1変換用音声信号Vt1および第2変換用音声信号Vt2の一方から他方に切り替える構成としてもよい。さらに、本実施形態に係る音声処理装置Dをカラオケ装置に適用した場合には、カラオケ演奏される楽曲の進行に同期して第1変換用音声信号Vt1および第2変換用音声信号Vt2の一方から他方に切り替える構成も採用される。また、本実施形態においては、単一の音声を示す第1変換用音声信号Vt1と複数の音声を示す第2変換用音声信号Vt2とが記憶手段50に記憶された構成を例示したが、各変換用音声信号Vtが示す音声数はこれに限られない。例えば、各々が並列に発生された所定数の音声からなる変換用音声を示す第1変換用音声信号Vt1と、これよりも多数の音声からなる変換用音声を示す第2変換用音声信号Vt2とを利用してよい。

10

【0032】

< C : 変形例 >

各実施形態に対しては種々の変形が加えられる。具体的な変形の態様は以下の通りである。なお、以下の各態様を適宜に組み合わせてもよい。

【0033】

(1) 各実施形態においてはひとつのピッチPtの変換用音声信号Vt(またはVt1, Vt2)が記憶手段50に記憶された構成を例示したが、図6に示されるように、各々のピッチPt(Pt1, Pt2, ...)が相違する複数の変換用音声信号Vtを記憶手段50に記憶させた構成も採用される。各変換用音声信号Vtは、並列に発生した多数の音声を含む変換用音声を収録したものである。図6の構成においては、ピッチ検出部12によって検出されたピッチPinがスペクトル取得手段30の選択部34にも供給されるようになっている。この選択部34は、入力音声のピッチPinに近似または一致するピッチPtの変換用音声信号Vtを選択的に記憶手段50から読み出してFFT部31に出力する手段である。この構成によれば、新規スペクトルSPnewの生成に利用される変換用音声信号VtのピッチPtを入力音声信号VinのピッチPinに近づけることができるから、ピッチ変換部21による処理にて変換用スペクトルSPtの各ピークptの周波数を変化させる量が低減される。したがって、自然な形状の新規スペクトルSPnewを生成することができるという利点がある。なお、ここでは変換用音声信号Vtの選択に加えてピッチ変換部21による処理も実行する構成としたが、多数のピッチPtの変換用音声信号Vtが記憶されていれば変換用音声信号Vtの選択のみによって所望のピッチの出力音声を生成することができるから、ピッチ変換部21は必ずしも必要ではない。

20

30

【0034】

(2) 各実施形態においては、変換用スペクトルデータDtの各単位データUtに含まれる周波数Ftに特定の数値(Pin/Pt)を乗算することによって変換用スペクトルSPtを周波数軸の方向に伸長または縮小する構成を例示したが、変換用スペクトルSPtのピッチPtを変換する方法は任意に変更される。例えば、各実施形態に示した方法においては、変換用スペクトルSPtが全帯域にわたって同率に伸長または縮小されるため、各ピークptの帯域幅が元のピークptの帯域幅よりも著しく広がってしまう場合が生じ得る。例えば、図7の部分(a)に示される変換用スペクトルSPtのピッチPtを第1実施形態の方法によって2倍のピッチに変換した場合、図7の部分(b)に示されるように各ピークptの帯域幅は2倍となる。このように各ピークptのスペクトル形状が大幅に変化すると変換用音声の特性とは著しく相違する出力音声が生成されることになる。このような問題を解消するために、ピッチ変換部21が、特定の数値(Pin/Pt)を乗算して得られた変換用スペクトルSPt(図7の部分(b)に示される周波数スペクトル)の各ピークptについて、図7の部分(c)に矢印Bにて示されるように、当該ピークptの帯域幅をピ

40

50

ッチ変換前のピーク p_t の帯域幅まで狭めるための演算処理を各单位データ U_t の周波数 F_t に施してもよい。この構成によれば、変換用音声の特性を忠実に再現した出力音声を生成することができる。

【0035】

また、ここでは各单位データ U_t の周波数 F_t に対する乗算処理によってピッチ P_t を変換する場合を例示したが、図8の部分(a)に示されるように、変換用スペクトル S_{Pt} を周波数軸上にて複数の帯域(以下「スペクトル分布領域」という) R に区分し、各スペクトル分布領域 R を周波数軸の方向に移動させることによってピッチ P_t を変化させてもよい。各スペクトル分布領域 R は、ひとつのピーク p_t とその前後の帯域とを含むように選定される。ピッチ変換部21は、図8の部分(b)に示されるように、各スペクトル分布領域 R に属するピーク p_t の周波数が、入力スペクトル S_{Pin} (図8の部分(c))に現れる各ピーク p の周波数と略一致するように、各スペクトル分布領域 R を周波数軸の方向に移動させる。なお、図8の部分(b)に示されるように、相互に隣接するスペクトル分布領域 R の間隙には周波数スペクトルが存在しない帯域が生じ得るが、この帯域についてはスペクトル強度 M を所定値(例えばゼロ)に選定すればよい。この処理によれば、変換用スペクトル S_{Pt} の各ピーク p_t の周波数を確実に入力音声のピーク p_t の周波数に一致させることができるから、所望のピッチの出力音声を精度よく生成することができるという利点がある。

10

【0036】

(3) 各実施形態においては、記憶手段50に記憶された変換用音声信号 V_t から変換用スペクトル S_{Pt} が特定される構成を例示したが、変換用スペクトル S_{Pt} を示す変換用スペクトルデータ D_t が予めフレームごとに記憶手段50に記憶された構成も採用される。この構成におけるスペクトル取得手段30は、記憶手段50から変換用スペクトルデータ D_t を読み出してスペクトル変換手段20に出力する構成であれば足り、FFT部31やピーク検出部32やデータ生成部33を備えている必要はない。また、ここでは記憶手段50に変換用スペクトルデータ D_t が記憶された構成を例示したが、スペクトル取得手段30は、例えば通信回線を介して接続された通信装置から変換用スペクトルデータ D_t を取得する手段であってもよい。このように、本発明におけるスペクトル取得手段30は、変換用スペクトル S_{Pt} を取得する手段であれば足り、その取得の方法や取得先の如何は不問である。

20

30

【0037】

(4) 各実施形態においては入力音声の周波数スペクトル S_{Pin} からピッチ P_{in} を検出する構成を例示したが、このピッチ P_{in} を検出する方法は任意に変更される。例えば、音声入力部61から入力された時間領域の入力音声信号 V_{in} からピッチ P_{in} を検出する構成としてもよい。ピッチ P_{in} を検出する方法としては、公知である各種の方法が採用される。

【0038】

(5) 各実施形態においては変換用音声のピッチ P_t を入力音声のピッチ P_{in} に一致させる構成を例示したが、変換用音声のピッチ P_t をこれ以外のピッチに変換してもよい。例えば、ピッチ変換部21が、入力音声のピッチ P_{in} と協和音を構成するピッチとなるように変換用音声のピッチ P_t を変換する構成も採用される。この構成に加え、出力処理部42から出力された出力音声信号 V_{new} と音声入力部61から入力された入力音声信号 V_{in} とを加算したうえで音声出力部63から放音する構成を採用すれば、利用者が発声した入力音声とともにコーラス音を出力することができる。このように、本発明のうちピッチ変換部21を備えた態様においては、このピッチ変換部21が変換用音声のピッチ P_t を入力音声のピッチ P_{in} に応じて(すなわちピッチ P_{in} の変化に伴って変換用音声のピッチ P_t が変化するように)変化させる構成であれば足りる。

40

【図面の簡単な説明】

【0039】

【図1】第1実施形態に係る音声処理装置の構成を示すブロック図である。

【図2】入力音声に関する処理を説明するための図である。

50

- 【図3】変換用音声信号に関する処理を説明するための図である。
- 【図4】スペクトル変換手段による処理の内容を説明するための図である。
- 【図5】第2実施形態に係る音声処理装置の構成を示すブロック図である。
- 【図6】変形例に係る音声処理装置の構成を示すブロック図である。
- 【図7】変形例に係る音声処理装置におけるピッチ変換について説明するための図である。
- 【図8】変形例に係る音声処理装置におけるピッチ変換について説明するための図である。

【符号の説明】

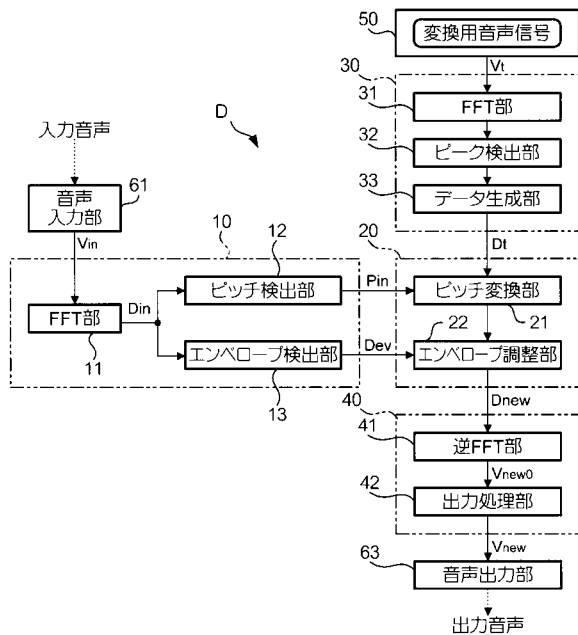
【0040】

D ... 音声処理装置、10 ... 周波数分析手段、11 ... FFT部、12 ... ピッチ検出部、13 ... エンベロープ検出部、20 ... スペクトル変換手段、21 ... ピッチ変換部、22 ... エンベロープ調整部、30 ... スペクトル取得手段、31 ... FFT部、32 ... ピーク検出部、33 ... データ生成部、34 ... 選択部、40 ... 音声生成手段、41 ... 逆FFT部、42 ... 出力処理部、50 ... 記憶手段、61 ... 音声入力部、63 ... 音声出力部、67 ... 入力機器、Vin ... 入力音声信号、Vt, Vt1, Vt2 ... 変換用音声信号、Vnew ... 出力音声信号、SPin ... 入力スペクトル、SPt ... 変換用スペクトル、SPnew ... 新規スペクトル、EVin ... スペクトルエンベロープ、Din ... 入力スペクトルデータ、Dt ... 変換用スペクトルデータ、Dnew ... 新規スペクトルデータ、Dev ... エンベロープデータ、R ... スペクトル分布領域。

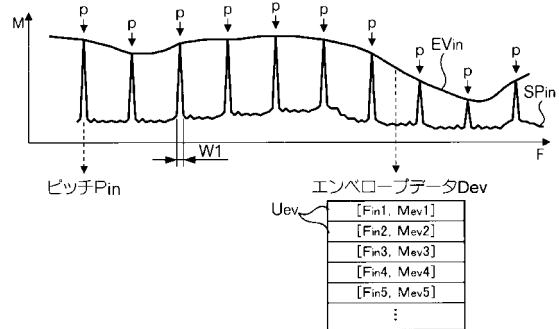
10

20

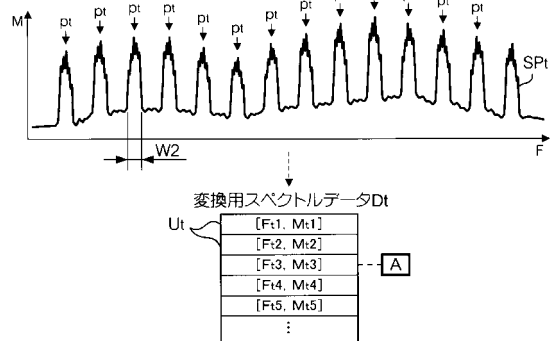
【図1】



【図2】



【図3】



【 図 8 】

