

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5582915号
(P5582915)

(45) 発行日 平成26年9月3日 (2014.9.3)

(24) 登録日 平成26年7月25日 (2014.7.25)

(51) Int.Cl.

F I

G 1 O G 3/04 (2006.01)

G 1 O G 3/04

G 1 O L 25/51 (2013.01)

G 1 O L 25/51 3 0 0

請求項の数 11 (全 30 頁)

(21) 出願番号	特願2010-177968 (P2010-177968)	(73) 特許権者	000005326
(22) 出願日	平成22年8月6日 (2010.8.6)		本田技研工業株式会社
(65) 公開番号	特開2011-39511 (P2011-39511A)		東京都港区南青山二丁目1番1号
(43) 公開日	平成23年2月24日 (2011.2.24)	(74) 代理人	100064908
審査請求日	平成24年11月27日 (2012.11.27)		弁理士 志賀 正武
(31) 優先権主張番号	61/234, 076	(74) 代理人	100108578
(32) 優先日	平成21年8月14日 (2009.8.14)		弁理士 高橋 詔男
(33) 優先権主張国	米国 (US)	(74) 代理人	100146835
			弁理士 佐伯 義文
		(74) 代理人	100094400
			弁理士 鈴木 三義
		(74) 代理人	100107836
			弁理士 西 和哉
		(74) 代理人	100108453
			弁理士 村山 靖彦

最終頁に続く

(54) 【発明の名称】 楽譜位置推定装置、楽譜位置推定方法および楽譜位置推定ロボット

(57) 【特許請求の範囲】

【請求項 1】

音響信号取得部と、

前記音響信号取得部が取得する音響信号に対応する楽譜情報を取得する楽譜情報取得部と、

前記音響信号に含まれる音階を構成する楽音のうちの1つの楽音に隣接する他の楽音のパワーを減じ、さらに前のフレーム時のパワーを減じて当該1つの楽音を強調し、前記強調された楽音を用いて当該音響信号の特徴量を抽出する音響信号の特徴量抽出部と、

前記楽譜情報の特徴量を抽出する楽譜情報の特徴量抽出部と、

前記音響信号のビート位置を推定するビート位置推定部と、

前記推定されたビート位置を用いて、前記音響信号の特徴量と前記楽譜情報の特徴量とのマッチングを行うことで、前記音響信号が対応する前記楽譜情報における位置を推定するマッチング部と、

を備えることを特徴とする楽譜位置推定装置。

【請求項 2】

前記音響信号の特徴量抽出部は、

フレーム時刻 t 毎にバンドパスフィルタによって前記1つの楽音 $c(i, t)$ (i は 1 ~ 12 の整数) を抽出し、前記抽出した前記1つの楽音 $c(i, t)$ に対して、次式

【数 1】

$$\begin{aligned}
c'(i,t) = & -c(i+1,t-1) - 2c(i+1,t) - c(i+1,t+1) \\
& -c(i,t-1) + 6c(i,t) + 3c(i,t+1) \\
& -c(i-1,t-1) - 2c(i-1,t) - c(i-1,t+1)
\end{aligned} \quad \dots(1)$$

の畳み込みを周期的に行い、前記畳み込みが行われた $c'(i, t)$ に基づいて音響クロマベクトルを算出する請求項 1 に記載の楽譜位置推定装置。

10

【請求項 3】

前記楽譜情報の特徴量抽出部は、前記楽譜情報から音符の出現頻度であるレアネスを算出し、

前記マッチング部は、前記レアネスを用いてマッチングを行う

ことを特徴とする請求項 1 または請求項 2 に記載の楽譜位置推定装置。

【請求項 4】

前記マッチング部は、前記算出されたレアネスと前記抽出された音響信号の特徴量と楽譜情報の特徴量との積に基づきマッチングを行う

ことを特徴とする請求項 3 に記載の楽譜位置推定装置。

20

【請求項 5】

前記レアネスは、フレーム内の所定の区間における前記音符の出現頻度である

ことを特徴とする請求項 3 または請求項 4 に記載の楽譜位置推定装置。

【請求項 6】

前記マッチング部は、

前記楽譜情報において進むべきフレーム数 F に対して次式のように重み付けを行い（ただし f_m は前記楽譜情報の m 番目のオンセットのフレーム、 f_{m+k} は前記楽譜情報の $m+k$ 番目のオンセット時刻、 k は進むべき前記楽譜情報のオンセット時刻、 σ は重み付けの分散値）、

【数 2】

$$W(k) = \exp\left(-\frac{(f_{m+k} - f_m - F)^2}{2\sigma^2}\right) \quad \dots(2)$$

30

前記楽譜情報における m 番目のオンセットのフレームと、前記音響信号における n 番目のオンセット時刻との類似性 $S(n, m)$ を算出し、

前記重み付けした値 $W(k)$ と前記算出された類似性 $S(n, m)$ を用いて、次式の範囲内で探索を行うことで前記進むべき楽譜情報のオンセット時刻 k を算出する

【数 3】

$$k = \arg\max_l W(l)S(n+1, m+l) \quad \dots(3)$$

40

ことを特徴とする請求項 1 から請求項 5 のいずれか 1 項に記載の楽譜位置推定装置。

【請求項 7】

前記音響信号の特徴量抽出部は、前記音響信号の特徴量を、クロマベクトルを用いて抽出し、

前記楽譜情報の特徴量抽出部は、前記楽譜情報の特徴量を、クロマベクトルを用いて抽出する

50

ことを特徴とする請求項 1 から請求項 6 のいずれか 1 項に記載の楽譜位置推定装置。

【請求項 8】

前記音響信号の特徴量抽出部は、抽出した音響信号の特徴量において高周波成分に重み付けを行い、重み付けした特徴量に基づき音符の出だしのタイミング時刻を算出し、

前記マッチング部は、算出された音符の出だしのタイミング時刻を用いて、マッチングを行う

ことを特徴とする請求項 1 から請求項 7 のいずれか 1 項に記載の楽譜位置推定装置。

【請求項 9】

前記ビート位置推定部は、異なる複数の観測誤差モデルを、スイッチング・カルマン・フィルタにより切り替えることでビート位置の推定を行う

10

ことを特徴とする請求項 1 から請求項 8 のいずれか 1 項に記載の楽譜位置推定装置。

【請求項 10】

楽譜位置推定装置の楽譜位置推定方法において、

音響信号取得部が、音響信号を取得する音響信号取得工程と、

楽譜情報取得部が、前記音響信号に対応する楽譜情報を取得する楽譜情報取得工程と、

音響信号の特徴量抽出部が、前記音響信号に含まれる音階を構成する楽音のうちの 1 つの楽音に隣接する他の楽音のパワーを減じ、さらに前のフレーム時のパワーを減じて当該 1 つの楽音を強調し、前記強調された楽音を用いて当該音響信号の特徴量を抽出する音響信号の特徴量抽出工程と、

楽譜情報の特徴量抽出部が、前記楽譜情報の特徴量を抽出する楽譜情報の特徴量抽出工程と、

20

ビート位置推定部が、前記音響信号のビート位置を推定するビート位置推定工程と、

マッチング部が、前記推定されたビート位置を用いて、前記音響信号の特徴量と前記楽譜情報の特徴量とのマッチングを行うことで、前記音響信号に対応する前記楽譜情報における位置を推定するマッチング工程と、

を含むことを特徴とする楽譜位置推定方法。

【請求項 11】

音響信号取得部と、

前記音響信号取得部が取得した音響信号に対して抑圧処理を行うことで、演奏に対応する音響信号を抽出する音響信号分離部と、

30

前記音響信号分離部が抽出した音響信号に対応する楽譜情報を取得する楽譜情報取得部と、

前記音響信号分離部が抽出した前記音響信号に含まれる音階を構成する楽音のうちの 1 つの楽音に隣接する他の楽音のパワーを減じ、さらに前のフレーム時のパワーを減じて当該 1 つの楽音を強調し、前記強調された楽音を用いて当該音響信号の特徴量を抽出する音響信号の特徴量抽出部と、

前記楽譜情報の特徴量を抽出する楽譜情報の特徴量抽出部と、

前記音響信号分離部が抽出した音響信号のビート位置を推定するビート位置推定部と、

前記推定されたビート位置を用いて、前記音響信号の特徴量と前記楽譜情報の特徴量とのマッチングを行うことで、前記音響信号に対応する前記楽譜情報における位置を推定するマッチング部と、

40

を備えることを特徴とする楽譜位置推定ロボット。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、楽譜位置推定装置、楽譜位置推定方法および楽譜位置推定ロボットに関する。

【背景技術】

【0002】

近年、ロボットの身体的な機能の顕著な向上により、人間社会に関わる家事や看護を行

50

う人の援助などが試みられている。このように、日々の場面で人間とロボットとが共存していく上で、ロボットには人間と自然な相互作用ができるようにしていくことが求められている。

【0003】

人間とロボットとの相互作用におけるコミュニケーションとして、音楽を通したコミュニケーションがある。音楽は、人間同士のコミュニケーションにおいても重要な役割を果たし、例えば言葉が通じ合わない人間同士でも音楽を通じて友好的で楽しい時間を共有することができる場合もある。このため、ロボットには、音楽を通して人間と相互作用を行うことができることが、人間と調和して共存していく上で重要になってくる。

【0004】

ロボットが音楽を通して人間とコミュニケーションを行う場面として、例えば、ロボットが、伴奏または歌声に合わせて歌ったり、音楽に合わせて自身の胴体を動かしたりすることが考えられる。

【0005】

このようなロボットにおいて、楽譜情報を解析し、解析結果に基づいて動作することが知られている。

楽譜に記載されている音符が、どの音符かを認識する技術として、楽譜の画像データを音符データに変換して、楽譜を自動認識する技術が提案されている（例えば、特許文献1参照）。また、楽譜データと予めグルーピングされた構造分析データにもとづき、楽曲データの拍節構造を分析し、演奏されている音響信号からテンポを推定する技術として、ビートトラッキング法が提案されている（例えば、特許文献2参照）。

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特許第3147846号公報

【特許文献2】特開2006-201278号公報

【発明の概要】

【発明が解決しようとする課題】

【0007】

特許文献2の拍節構造を分析する技術では、楽譜に基づく構造のみを分析していた。このため、ロボット自身が集音した音響信号に合わせてロボットに歌唱を行わせようとしても、音楽が途中から開始されると、楽譜のどの部分が不明であるので、演奏されている曲のビート時間やテンポの抽出に失敗する場合もあるという問題点があった。また、人間が行う演奏の場合、演奏のテンポが変動することもあり、この結果、ロボットが、演奏されている曲のビート時間やテンポの抽出に失敗する場合もあるという問題点があった。

以上のように、従来技術では、楽譜データに基づいて曲の拍節構造やビート時間やテンポを抽出していたので、実際の演奏が行われている場合、楽譜のどの位置が演奏されているかを精度良く検出することができなかった。

【0008】

本発明は、上記の問題点に鑑みてなされたものであって、演奏に対する楽譜位置の推定を行う楽譜位置推定装置、楽譜位置推定方法および楽譜位置推定ロボットを提供することを課題としている。

【課題を解決するための手段】

【0009】

上記目的を達成するため、本発明に係る楽譜位置推定装置は、音響信号取得部と、前記音響信号取得部が取得する音響信号に対応する楽譜情報を取得する楽譜情報取得部と、前記音響信号に含まれる音階を構成する楽音のうちの1つの楽音に隣接する他の楽音のパワーを減じ、さらに前のフレーム時のパワーを減じて当該1つの楽音を強調し、前記強調された楽音を用いて当該音響信号の特徴量を抽出する音響信号の特徴量抽出部と、前記楽譜情報の特徴量を抽出する楽譜情報の特徴量抽出部と、前記音響信号のビート位置を推定す

10

20

30

40

50

るビート位置推定部と、前記推定されたビート位置を用いて、前記音響信号の特徴量と前記楽譜情報の特徴量とのマッチングを行うことで、前記音響信号が対応する前記楽譜情報における位置を推定するマッチング部と、を備えることを特徴としている。

また、本発明に係る楽譜位置推定装置において、前記音響信号の特徴量抽出部は、フレーム時刻 t 毎にバンドパスフィルタによって前記 1 つの楽音 $c(i, t)$ (i は 1 ~ 12 の整数) を抽出し、前記抽出した前記 1 つの楽音 $c(i, t)$ に対して、次式

【数 1】

$$\begin{aligned} c'(i, t) = & -c(i+1, t-1) - 2c(i+1, t) - c(i+1, t+1) \\ & -c(i, t-1) + 6c(i, t) + 3c(i, t+1) \\ & -c(i-1, t-1) - 2c(i-1, t) - c(i-1, t+1) \end{aligned} \quad \dots(1)$$

10

の畳み込みを周期的に行い、前記畳み込みが行われた $c'(i, t)$ に基づいて音響クロマベクトルを算出するようにしてもよい。

【0010】

また、本発明に係る楽譜位置推定装置において、前記楽譜情報の特徴量抽出部は、前記楽譜情報から音符の出現頻度であるレアネスを算出し、前記マッチング部は、前記レアネスを用いてマッチングを行うようにしてもよい。

20

【0011】

また、本発明に係る楽譜位置推定装置において、前記マッチング部は、前記算出されたレアネスと前記抽出された音響信号の特徴量と楽譜情報の特徴量との積に基づきマッチングを行うようにしてもよい。

【0012】

また、本発明に係る楽譜位置推定装置において、前記レアネスは、フレーム内の所定の区間における前記音符の出現頻度であるようにしてもよい。

また、本発明に係る楽譜位置推定装置において、前記マッチング部は、前記楽譜情報において進むべきフレーム数 F に対して次式のように重み付けを行い（ただし f_m は前記楽譜情報の m 番目のオンセットのフレーム、 f_{m+k} は前記楽譜情報の $m+k$ 番目のオンセット時刻、 k は進むべき前記楽譜情報のオンセット時刻、 σ は重み付けの分散値）、

30

【数 2】

$$W(k) = \exp\left(-\frac{(f_{m+k} - f_m - F)^2}{2\sigma^2}\right) \quad \dots(2)$$

前記楽譜情報における m 番目のオンセットのフレームと、前記音響信号における n 番目のオンセット時刻との類似性 $S(n, m)$ を算出し、

前記重み付けした値 $W(k)$ と前記算出された類似性 $S(n, m)$ を用いて、次式の範囲内で探索を行うことで前記進むべき楽譜情報のオンセット時刻 k を算出するようにしてもよい

40

【数 3】

$$k = \underset{l}{\operatorname{argmax}} W(l) S(n+1, m+l) \quad \dots(3)$$

【0013】

また、本発明に係る楽譜位置推定装置において、前記音響信号の特徴量抽出部は、前記

50

音響信号の特徴量を、クロマベクトルを用いて抽出し、前記楽譜情報の特徴量抽出部は、前記楽譜情報の特徴量を、クロマベクトルを用いて抽出するようにしてもよい。

【0014】

また、本発明に係る楽譜位置推定装置において、前記音響信号の特徴量抽出部は、抽出した音響信号の特徴量において高周波成分に重み付けを行い、重み付けした特徴量に基づき音符の出だしのタイミング時刻を算出し、前記マッチング部は、算出された音符の出だしのタイミング時刻を用いて、マッチングを行うようにしてもよい。

【0015】

また、本発明に係る楽譜位置推定装置において、前記ビート位置推定部は、異なる複数の観測誤差モデルを、スイッチング・カルマン・フィルタにより切り替えることでビート位置の推定を行うようにしてもよい。

10

【0016】

上記目的を達成するため、本発明に係る楽譜位置推定装置における楽譜位置推定方法は、音響信号取得部が、音響信号を取得する音響信号取得工程と、楽譜情報取得部が、前記音響信号に対応する楽譜情報を取得する楽譜情報取得工程と、音響信号の特徴量抽出部が、前記音響信号に含まれる音階を構成する楽音のうちの1つの楽音に隣接する他の楽音のパワーを減じ、さらに前のフレーム時のパワーを減じて当該1つの楽音を強調し、前記強調された楽音を用いて当該音響信号の特徴量を抽出する音響信号の特徴量抽出工程と、楽譜情報の特徴量抽出部が、前記楽譜情報の特徴量を抽出する楽譜情報の特徴量抽出工程と、ビート位置推定部が、前記音響信号のビート位置を推定するビート位置推定工程と、マッチング部が、前記推定されたビート位置を用いて、前記音響信号の特徴量と前記楽譜情報の特徴量とのマッチングを行うことで、前記音響信号が対応する前記楽譜情報における位置を推定するマッチング工程と、を含むことを特徴としている。

20

【0017】

上記目的を達成するため、本発明に係る楽譜位置推定ロボットは、音響信号取得部と、前記音響信号取得部が取得した音響信号に対して抑圧処理を行うことで、演奏に対応する音響信号を抽出する音響信号分離部と、前記音響信号分離部が抽出した音響信号に対応する楽譜情報を取得する楽譜情報取得部と、前記音響信号分離部が抽出した前記音響信号に含まれる音階を構成する楽音のうちの1つの楽音に隣接する他の楽音のパワーを減じ、さらに前のフレーム時のパワーを減じて当該1つの楽音を強調し、前記強調された楽音を用いて当該音響信号の特徴量を抽出する音響信号の特徴量抽出部と、前記楽譜情報の特徴量を抽出する楽譜情報の特徴量抽出部と、前記音響信号分離部が抽出した音響信号のビート位置を推定するビート位置推定部と、前記推定されたビート位置を用いて、前記音響信号の特徴量と前記楽譜情報の特徴量とのマッチングを行うことで、前記音響信号が対応する前記楽譜情報における位置を推定するマッチング部と、を備えることを特徴としている。

30

【発明の効果】

【0018】

本発明によれば、取得した音響信号から特徴量とビート位置とを抽出し、取得した楽譜情報から特徴量を抽出する。そして、抽出したビート位置を用いて、音響信号の特徴量と楽譜情報の特徴量とをマッチングすることで、音響信号が対応する前記楽譜情報における位置を推定するようにした。この結果、音響信号に基づいて、楽譜位置を正確に推定することが可能になる。

40

本発明によれば、楽譜情報から音符の出現頻度であるレアネスを算出し、算出したレアネスを用いてマッチングを行うようにしたので、精度良く音響信号に基づいて、楽譜位置を正確に推定することが可能になる。

本発明によれば、レアネスと音響信号の特徴量と楽譜情報の特徴量との積に基づいてマッチングを行うようにしたので、精度良く音響信号に基づいて、楽譜位置を正確に推定することが可能になる。

本発明によれば、音符の出現頻度の低さをレアネスとして用いるようにしたので、精度良く音響信号に基づいて、楽譜位置を正確に推定することが可能になる。

50

本発明によれば、音響信号の特徴量と楽譜情報の特徴量とをクロマベクトルを用いて抽出するようにしたので、精度良く音響信号に基づいて、楽譜位置を正確に推定することが可能になる。

本発明によれば、音響信号の特徴量において高周波成分に重み付けを行い、重み付けした特徴量に基づき音符の出だしのタイミング時刻を用いて、マッチングを行うようにしたので、精度良く音響信号に基づいて、楽譜位置を正確に推定することが可能になる。

本発明によれば、異なる複数の観測誤差モデルを、スイッチング・カルマン・フィルタにより切り替えることでビート位置の推定を行うようにしたので、演奏が楽譜通りのテンポから外れた場合においても、精度良く音響信号に基づいて、楽譜位置を正確に推定することが可能になる。

10

【図面の簡単な説明】

【0019】

【図1】本実施形態に係る楽譜位置推定装置100を組み込んだロボット1の一例を説明する図である。

【図2】同実施形態に係る楽譜位置推定装置100のブロック図の一例を示す図である。

【図3】楽器演奏時の音響信号のスペクトラムの一例を示す図である。

【図4】楽器演奏時の音響信号の残響波形（パワーエンベロープ）の一例を示す図である。

。

【図5】実際の演奏に基づく音響信号と楽譜のクロマベクトルの一例を示す図である。

【図6】音楽演奏におけるスピード、又はテンポの変化を示したものである。

20

【図7】同実施形態に係る楽譜位置推定部120の構成を説明するブロック図である。

【図8】同実施形態に係る音響信号からの特徴量抽出部410がクロマベクトルとオンセット時刻と抽出する際に用いる式における記号を説明するリストである。

【図9】同実施形態に係る音響信号と楽譜からクロマベクトルを算出する過程を説明する図である。

【図10】同実施形態に係るオンセット時刻抽出手順の概略を説明する図である。

【図11】同実施形態に係るレアネスを説明する図である。

【図12】同実施形態に係るカルマン・フィルタを適用したビートトラッキングを説明する図である。

【図13】同実施形態に係る楽譜位置推定処理のフローチャートである。

30

【図14】楽譜位置推定装置を備えるロボット1と音源の設置関係を説明する図である。

【図15】2種類の音楽信号（（v）と（vi））と4つの手法（（i）～（iv））の結果を示している。

【図16】クリーン信号時の各手法の累積絶対値誤差平均値で分類された楽曲数を示している。

【図17】残響あり信号時の各手法の累積絶対値誤差平均値で分類された楽曲数を示している。

【発明を実施するための形態】

【0020】

以下、図面を用いて本発明の実施形態について詳細に説明する。なお、本発明は斯かる実施形態に限定されず、その技術思想の範囲内で種々の変更が可能である。

40

【0021】

図1は、本実施形態における楽譜位置推定装置100を組み込んだロボット1の一例を説明する図である。ロボット1は、図1に示すように、基体部11と、基体部11にそれぞれ可動連結される頭部12（可動部）と、脚部13（可動部）と、腕部14（可動部）とを備えている。また、ロボット1は、背負う格好で基体部11に収納部15を装着している。なお、基体部11には、スピーカ20が収納され、頭部12にはマイクロホン30が収納されている。なお、図1は、ロボット1を側面から見た図であり、マイクロホン30およびスピーカ20は、例えば正面から左右対称にそれぞれ複数収納されている。

【0022】

50

図 2 は、本実施形態における楽譜位置推定装置 100 のブロック図の一例を示す図である。図 2 のように、楽譜位置推定装置 100 にはマイクロホン 30、スピーカ 20 が接続されている。また、楽譜位置推定装置 100 は、音響信号分離部 110 と、楽譜位置推定部 120 と、歌声生成部 130 とを備えている。また、音響信号分離部 110 は、自己生成音抑制フィルタ部 111 を備え、楽譜位置推定部 120 は、楽譜データベース 121 と楽曲位置推定部 122 を備え、歌声生成部 130 は、歌詞とメロディーのデータベース 131 と音声生成部 132 を備えている。

【0023】

マイクロホン 30 は、演奏（伴奏）の音と、ロボット 1 自身のスピーカ 20 を介して出力した音声信号（歌声）とが混合された音を集音し、集音した音を音響信号に変換して音響信号分離部 110 に出力する。

10

音響信号分離部 110 には、マイクロホン 30 から集音された音響信号と、歌声生成部 130 から生成された音声信号とが入力される。音響信号分離部 110 の自己生成音抑制フィルタ部 111 は、入力された音響信号に対して、独立成分分析（ICA；Independent Component Analysis；）を行って、生成された音声信号と音響信号に含まれる残響音を抑圧する。これにより、音響信号分離部 110 は、演奏に関わる音響信号を分離して抽出する。音響信号分離部 110 は、抽出した音響信号を楽譜位置推定部 120 に出力する。

【0024】

楽譜位置推定部 120（楽譜情報取得部、音響信号の特徴量抽出部、楽譜情報の特徴量抽出部、ビート位置推定部、マッチング部）には、音響信号分離部 110 から分離された音響信号が入力される。楽譜位置推定部 120 の楽曲位置推定部 122 は、入力された音響信号から特徴量である音響クロマベクトルとオンセット時刻を算出する。楽曲位置推定部 122 は、楽譜データベース 121 から演奏されている曲の楽譜データを読み出し、楽譜データから特徴量である楽譜クロマベクトルと音符の出現頻度であるレアネスを算出する。楽曲位置推定部 122 は、入力された音響信号からビートトラッキングを行い、リズム間隔（テンポ）を検出する。楽曲位置推定部 122 は、抽出したリズム間隔（テンポ）に基づき、スイッチング・カルマン・フィルタ（SKF；Switching Kalman Filter）を用いて、テンポの外れ値やノイズ分を推定し、安定したリズム間隔（テンポ）を抽出する。楽曲位置推定部 122（音響信号の特徴量抽出部、楽譜情報の特徴量抽出部、ビート位置推定部、マッチング部）は、抽出したリズム間隔（テンポ）と、算出した音響クロマベクトル、オンセット時刻情報、楽譜クロマベクトルおよびレアネスとを用いて、演奏による音響信号と楽譜とのマッチングを行う。つまり、楽曲位置推定部 122 は、演奏されている曲が、楽譜のどの位置であるかを推定する。楽譜位置推定部 120 は、推定した楽譜位置を示す楽譜位置情報を歌声生成部 130 に出力する。

20

なお、楽譜データベース 121 に予め楽譜データが記憶されている例を説明したが、楽譜位置推定部 120 は、入力された楽譜データを楽譜データベース 121 に書き込んで記憶させるようにしてもよい。

30

【0025】

歌声生成部 130 には、推定された楽譜位置情報が入力される。歌声生成部 130 の音声生成部 132 は、入力された楽譜位置情報に基づき、歌詞とメロディーのデータベース 131 に記憶されている情報を用いて、公知の手法により演奏に合わせた歌声の音声信号を生成する。歌声生成部 130 は、生成した歌声の音声信号を、スピーカ 20 を介して出力する。

40

【0026】

次に、音響信号分離部 110 が、独立成分分析を用いて、生成された音声信号と音響信号に含まれる残響音を抑圧する概要について説明する。独立成分分析は、音源同士の独立性（確率密度）を仮定して分離を行う。ロボット 1 がマイクロホン 30 を介して取得した音響信号は、演奏されている音とロボット 1 がスピーカ 20 から出力した音声信号とが混合された信号である。この混合された信号のうち、ロボット 1 がスピーカ 20 から出力し

50

た音声信号は、音声生成部 1 3 2 で生成した信号であるため既知である。このため、音響信号分離部 1 1 0 は、周波数領域で独立成分分析を行い、混合された信号に含まれるロボット 1 の音声信号を抑圧することで、演奏されている音を分離する。

【 0 0 2 7 】

次に、本実施形態における、楽譜位置推定装置 1 0 0 に採用した技術の概略を説明する。演奏されている音楽（伴奏）からビートやテンポを抽出して、楽譜のどの位置が演奏されているのかを推定する場合、大きく分けると 3 つの技術がある。

【 0 0 2 8 】

第 1 の技術は、演奏されている音響信号に含まれている様々な楽器音の差異をどのようにに判別するかである。図 3 は、楽器演奏時の音響信号のスペクトラムの一例を示す図である。図 3 (a) は、ピアノで A 4 の音 (4 4 0 [H z]) を鳴らしたときの音響信号のスペクトラムを表し、図 3 (b) は、フルートで A 4 の音を鳴らしたときの音響信号のスペクトラムを表している。縦軸は信号の大きさを表し、横軸は周波数を表している。図 3 (a) と図 3 (b) のように、同じ周波数範囲で分析した各スペクトラムは、同じ基本周波数 4 4 0 [H z] の A 4 の音でも、楽器によってスペクトラムの形状や成分が異なっている。

【 0 0 2 9 】

図 4 は、楽器演奏時の音響信号の残響波形（パワーエンベロープ）の一例を示す図である。図 4 (a) は、ピアノにおける音響信号の残響波形を表し、図 4 (b) は、フルートにおける音響信号のスペクトラムを表す。縦軸は信号の大きさを表し、横軸は時間を表している。通常、楽器の残響波形は、アタック(出だし)部分 (2 0 1 , 2 1 1)、減衰部分 (2 0 2 , 2 1 2)、持続部分 (2 0 3 , 2 1 3)、リリース(消滅)部分 (2 0 4 , 2 1 4) から構成されている。図 4 (a) のように、ピアノやギターのような楽器の残響波形は、下降的な持続部分 2 0 3 を有し、図 4 (b) のように、フルートやヴァイオリン、サキソフォンなどの楽器の残響波形は、永続性の持続部分 2 1 3 を有している。

複合的な音符が様々な楽器から同時に演奏された場合、言い換えれば、和音音響を扱う場合、各音符の基本周波数を検出すること、又は持続音を認識することは、更に難しくなる。

【 0 0 3 0 】

このため、本実施形態では、演奏における波形の出だしであるオンセット時刻 (2 0 5 , 2 1 5) に着目する。

楽譜位置推定部 1 2 0 は、1 2 段階のクロマベクトル（音響特徴量）を用いて周波数領域の特徴量を抽出する。そして、楽譜位置推定部 1 2 0 は、抽出した周波数領域の特徴量に基づき、時間領域の特徴量であるオンセット時刻を算出する。クロマベクトルの利点としては、様々な楽器のスペクトル形状の変化への頑健性と和音音響信号に対する有効性が挙げられる。クロマベクトルは、基本周波数の代わりに各 1 2 音名、つまり、C、C #、...、B などのパワーを抽出する。本実施形態では、図 4 (a) の出だし部分 2 0 5、図 4 (b) の出だし部分 2 1 5 のように、パワーの急激な上昇周辺の頂点を「オンセット時刻」と定義する。オンセット時刻の抽出は、楽譜同期を行う上で、各音符のスタート時間を得るために必要である。さらに和音音響信号では、オンセット時刻の抽出は、時間領域におけるパワーの上昇部分として、持続部分やリリース部分より簡単に抽出できる。

【 0 0 3 1 】

次に、第 2 の技術は、演奏されている音響信号と楽譜の相違を推定するかである。図 5 は、実際の演奏に基づく音響信号と楽譜のクロマベクトルの一例を示す図である。図 5 (a) は、楽譜のクロマベクトルを表し、図 5 (b) は、実際の演奏に基づく音響信号のクロマベクトルを表している。図 5 (a) と図 5 (b) における縦軸は、1 2 段階の音の種類を表し、図 5 (a) における横軸は楽譜におけるビートを表し、図 5 (b) における横軸は時間を表している。また、図 5 (a) と図 5 (b) において、縦実線 3 1 1 は、各音（音符）のオンセット時刻を表している。楽譜中のオンセット時刻は、各音符フレームの開始部分として定義される。

図5(a)と図5(b)のように、実際の演奏による音響信号に基づくクロマベクトルと楽譜に基づくクロマベクトルには差異が見られる。実線で囲まれている符号301の領域では、図5(a)ではクロマベクトルが存在せず、図5(b)ではクロマベクトルが存在している。すなわち、楽譜中には音符が無い部分にもかかわらず、実際の演奏においては、前の音のパワーが持続している。点線で囲まれている符号302の領域では、逆に、図5(a)ではクロマベクトルが存在しているのに、図5(b)ではクロマベクトルがほとんど検出できない。

さらに、楽譜において、各音符の音量は明示されていない。

【0032】

以上により、本実施形態において、ほとんど使用されない音名の音符は、音響信号において、時として顕著に表されるという考えに基づき、音響信号と楽譜との相違を軽減する。まず、演奏される曲の楽譜を予め取得して、楽譜データベース121に登録しておく。そして、楽曲位置推定部122は、演奏される曲の楽譜を解析し、各音符の使用頻度を各々算出する。この楽譜中の各音名の出現頻度をレアネス(rareness)と定義する。レアネスの定義は情報エントロピーに類似している。図5(a)において、音名Bの数は他の音名の数より少ないため、音名Bのレアネスは高い。対照的に、音名Cや音名Eは楽譜中で頻繁に使用されているためレアネスは低い。

さらに、楽曲位置推定部122は、このように算出した各音名に対して、算出したレアネスに基づき重み付けを行う。

このように重み付けを行うことで、低頻出音符は高頻出音符に比べて和音音響信号からより簡単に抽出される可能性がある。

【0033】

次に、第3の技術は、演奏されている音響信号のテンポの変動を推定するかである。安定したテンポ推定は、ロボット1が正確に楽譜に同期して歌唱を実行するだけでなく、演奏されている曲に合わせてロボット1が滑らかで心地よい歌声を出力することにとっても不可欠である。人間が行う演奏においては、楽譜で指示されているテンポから外れる場合もある。さらに、公知のビートトラッキングを用いたテンポ推定時にも発生する。

図6は、音楽演奏におけるスピード、又はテンポの変化を示したものである。図6(a)は、人間の演奏に厳密に一致させたMIDI(登録商標(Musical Instrument Digital Interface; 電子楽器デジタルインタフェース))データから算出したビートの時間変動を示す図である。各テンポは楽譜中の音符の長さをその時間の長さで分割して得られる。図6(b)は、ビートトラッキングにおけるビートの時間変動を示す図である。テンポ列は相当数の外れ値を含む。外れ値は一般的にドラムのパターンの変化によって引き起こされる。図6において、縦軸は、時間あたりのビート数を表し、横軸は時間を表している。

このため、本実施形態では、楽曲位置推定部122は、テンポ推定にスイッチング・カルマン・フィルタ(SKF)を用いる。SKFは、誤りを含む一連のテンポから、次のテンポ推定を可能にする。

【0034】

次に、楽譜位置推定部120が行う処理について、図7~図12を用いて、詳細に説明する。図7は、楽譜位置推定部120の構成を説明するブロック図である。図7のように、楽譜位置推定部120は、楽譜データベース121と楽曲位置推定部122とを備えている。また、楽曲位置推定部122は、音響信号からの特徴量抽出部410(音響信号の特徴量抽出部)と、楽譜からの特徴量抽出部420(楽譜情報の特徴量抽出部)と、ビート間隔(テンポ)算出部430と、マッチング部440と、テンポ推定部450(ビート位置推定部)を備えている。また、マッチング部440は、類似度計算部441と重み付け計算部442を備えている。また、テンポ推定部450は、小さな観測誤差モデル451と外れ値となる大きな観測誤差モデル452を備えている。

【0035】

[音響信号からの特徴量抽出]

10

20

30

40

50

音響信号からの特徴量抽出部 410 には、音響信号分離部 110 により分離された音響信号が入力される。音響信号からの特徴量抽出部 410 は、入力された音響信号から、音響クロマベクトルとオンセット時刻とを抽出し、抽出したクロマベクトルとオンセット時刻情報をビート間隔(テンポ)算出 430 に出力する。

図 8 は、音響信号からの特徴量抽出部 410 がクロマベクトルとオンセット時刻情報と抽出する際に用いる式における記号を説明するリストである。図 8 において、 i は、西洋音階における 12 音 (C、C#、D、D#、E、F、F#、G、G#、A、A#、B) の名前のインデックスである。 t は、音響信号のフレーム時間である。 n は、音響信号におけるオンセット時刻のためのインデックスである。 t_n は、音響信号における n 番目のオンセット時刻である。 f は、楽譜のフレーム・インデックスである。 m は、楽譜におけるオンセット時刻のためのインデックスである。 f_m は、楽譜における m 番目のオンセット時刻である。

【0036】

音響信号からの特徴量抽出部 410 は、短時間フーリエ変換 (STFT; short-time Fourier transformation) を用いて、入力された音響信号からスペクトラムを算出する。短時間フーリエ変換は、入力された音響信号にハニング等の窓関数を音声信号に乗じて有限期間内で、解析位置をシフトしながらスペクトラムを算出する技術である。本実施形態では、ハニング窓が 4096 [ポイント]、変位間隔が 512 [ポイント]、サンプリングレートが 44.1 [kHz] の設定を用いた。ここでフレーム時間 t 、周波数 ω の時の $p(t, \omega)$ をパワーとする。

クロマベクトル $c(t) = [c(1, t), c(2, t), \dots, c(12, t)]^T$ (T はベクトルの転置を意味する) はフレーム時間 t 毎に生成される。図 9 のように、音響信号からの特徴量抽出部 410 が各音名のバンドパス・フィルタにより 12 音名のうちの 1 つに対応した各成分を抽出し、抽出した 12 音名のうちの 1 つに対応した各成分は数式 (1) のように表される。図 9 は、音響信号と楽譜からクロマベクトルを算出する過程を説明する図であり、図 9 (a) は、音響信号からクロマベクトルを算出する過程を説明する図である。

【0037】

【数 1】

$$c(i, t) = \sum_{h=Oct_L}^{Oct_H} \int_{-\infty}^{\infty} BPF_{i,h}(\omega) p(t, \omega) d\omega \quad \dots (1)$$

【0038】

式 (1) において、 $BPF_{i,h}$ は、 h 番目のオクターブにおける音名 i のバンドパス・フィルタである。また、 Oct_L と Oct_H は、それぞれ考慮される下限オクターブ及び上限オクターブである。周波帯のピークは、音の基本周波数である。周波帯の端は、隣接する音の周波数である。例えば、基本周波数 440 [Hz] である音 “A4” (第 4 オクターブの音 “A”) の BPF は、440 [Hz] にその周波帯のピークがある。その周波帯の一端は、“G#” (第 4 オクターブの音 “G#”) の 415 [Hz] であり、“A#” の 466 [Hz] である。本実施形態では、 $Oct_L = 3$ 及び $Oct_H = 7$ とする。言い換えれば、下限音は “C3” の 131 [Hz] とし、上限音は “B7”、3951 [Hz] とした。

次に、音響信号からの特徴量抽出部 410 は、音名を強調するために、式 (1) に対して、次式 (2) の畳み込みを行う。

【0039】

【数 2】

$$\begin{aligned}
 c'(i,t) = & -c(i+1,t-1) - 2c(i+1,t) - c(i+1,t+1) \\
 & -c(i,t-1) + 6c(i,t) + 3c(i,t+1) \\
 & -c(i-1,t-1) - 2c(i-1,t) - c(i-1,t+1) \quad \dots (2)
 \end{aligned}$$

【0040】

音響信号からの特徴量抽出部 410 は、式 (2) の畳み込みを、 i に対して周期的に処理を行う。例えば、 $i = 1$ (音名は“C”) のとき、 $c(i-1, t)$ は、 $c(12, t)$ (音名は“B”) に置き換えられる。

10

式 (2) の畳み込みにより、隣接する音名のパワーを減じ、他よりパワーを持つ成分が強調され、画像処理におけるエッジ抽出に類似する可能性がある。前のフレーム時間のパワーが減じられることで、パワーの増加量は強調される。

次に、音響信号からの特徴量抽出部 410 は、次式 (3) により、音響信号から音響クロマベクトル $c_{sig}(i, t)$ を算出することで特徴量を抽出する。

【0041】

【数 3】

20

$$c_{sig}(i,t) = \begin{cases} c'(i,t) & (c'(i,t) > 0) \\ 0 & otherwise \end{cases} \quad \dots (3)$$

【0042】

次に、音響信号からの特徴量抽出部 410 は、入力された音響信号から Rodet 他により提案されたオンセット抽出手法 (手法 1) を使用して、オンセット時刻を抽出する。

【0043】

30

文献 1 (手法 1) X. Rodet and F. Jaillet. Detection and modeling of fast attack transients. In International Computer Music Conference, pages 30-33, 2001.

【0044】

オンセット抽出において、特に高周波領域に位置するオンセット時刻のパワー増加量を利用する。音階のある楽器の音のオンセット時刻は、ドラムのような打楽器のオンセット時刻に比べ、より高い周波領域に重心がある。このように、この手法は、音階のある楽器のオンセット時刻検出に特に効果的である。

まず、音響信号からの特徴量抽出部 410 は、高周波成分と呼ばれるパワーを次式 (4) により算出する。

【0045】

40

【数 4】

$$h(t) = \sum_{\omega} \omega p(t, \omega) \quad \dots (4)$$

【0046】

高周波成分は重み付けされたパワーであり、その重みは周波数に対して直線的に増加する。音響信号からの特徴量抽出部 410 は、図 10 のように、オンセット時刻 t_n を、中央値フィルタを用いて $h(t)$ のピークを選択することにより判断する。図 10 は、オンセット時刻抽出手順の概略を説明する図である。図 10 のように、入力された音響信号の

50

スペクトラムを算出した後（図 10（a））、音響信号からの特徴量抽出部 410 は、高周波成分に重み付けしたパワーを算出する（図 10（b））。そして、音響信号からの特徴量抽出部 410 は、重み付けしたパワーに対して中央値フィルタを適用し、パワーのピーク部分の時刻をオンセット時刻として算出する（図 10（c））。

【0047】

音響信号からの特徴量抽出部 410 は、抽出した音響クロマベクトルとオンセット時刻情報とをマッチング部 440 に出力する。

【0048】

〔楽譜からの特徴量抽出〕

楽譜からの特徴量抽出 420 は、楽譜データベース 121 に記憶されている楽譜の中から必要な楽譜データを読み出す。なお、本実施形態では、予め演奏される曲名がロボット 1 に入力されているとし、楽譜からの特徴量抽出 420 は、指定されている曲の楽譜データを選択して読み出す。

次に、楽譜からの特徴量抽出 420 は、読み出した楽譜データを、図 9（b）のように、1 小節の 48 分の 1 と同等の長さのフレームに分割する。このフレーム解決法では、6 分音符や 3 連音符の処理が可能である。本実施形態では、楽譜のクロマベクトルを、次式（5）を用いて算出することで特徴量を抽出する。図 9（b）は、楽譜からクロマベクトルを算出する過程を説明する図である。

【0049】

【数 5】

$$c_{sco}(i, m) = \begin{cases} 1 & \text{pitch name } i \text{ starts at frame } f_m \\ 0 & \text{otherwise} \end{cases} \quad \dots (5)$$

【0050】

式（5）において、 f_m は、楽譜中の m 番目のオンセット時刻を表している。次に、楽譜からの特徴量抽出 420 は、抽出したクロマベクトルから、フレーム f_m における各音名（ i ）のレアネス $r(i, m)$ を、次式（7）により算出する。

【0051】

【数 6】

$$n(i, m) = \frac{\sum_{p \in M} c_{sco}(i, p)}{\sum_{i=1}^{12} \sum_{p \in M} c_{sco}(i, p)} \quad \dots (6)$$

【0052】

【数 7】

$$r(i, m) = \begin{cases} -\log_2 n(i, m) & (n(i, m) > 0) \\ \max_i (-\log_2 n(i, m)) & (n(i, m) = 0) \end{cases} \quad \dots (7)$$

【0053】

M はフレーム f_m に中心があり、長さが 2 小節分のフレーム範囲を意味している。したがって、 $n(i, m)$ は、フレーム f_m 周辺の各音名の分布を表している。

楽譜からの特徴量抽出 420 は、抽出した楽譜のクロマベクトルとレアネスとをマッチ

10

20

30

40

50

ング部 4 4 0 に出力する。

【 0 0 5 4 】

図 1 1 は、レアネスを説明する図である。図 1 1 (a) ~ 図 1 1 (c) において、縦軸は音名を表し、横軸は時間を表している。図 1 1 (a) は楽譜のクロマベクトルを表す図であり、図 1 1 (b) は演奏された音響信号のクロマベクトルを表す図である。図 1 1 (c) ~ 図 1 1 (e) は、レアネスの算出方法を説明する図である。

図 1 1 (c) のように、楽譜からの特徴量抽出 4 2 0 は、図 1 1 (a) の楽譜クロマベクトルについて、フレーム毎に前後 2 小節区間で、各音の出現頻度（使用頻度）を計算する。そして、図 1 1 (d) のように、楽譜からの特徴量抽出 4 2 0 は、前後 2 小節区間における各音名 i の使用頻度 p_i を算出する。次に、図 1 1 (e) のように、楽譜からの特徴量抽出 4 2 0 は、式 (7) を用いて算出した各音名 i の使用頻度 p_i の対数を取ってレアネス r_i を算出する。式 (7) および図 1 1 (e) のように、 $-\log p_i$ は、使用頻度が低い音名 i を抽出することを意味している。

【 0 0 5 5 】

楽譜からの特徴量抽出 4 2 0 は、抽出した楽譜クロマベクトルとレアネスとをマッチング部 4 4 0 に出力する。

【 0 0 5 6 】

[ビートトラッキング]

ビート間隔(テンポ)算出 4 3 0 は、村田らにより開発されたビートトラッキング手法（手法 2）を用いて、入力された音響信号からビート間隔(テンポ)を算出する。

【 0 0 5 7 】

文献 2（手法 2） K. Murata, K. Nakadai, K. Yoshii, R. Takeda, T. Torii, H. G. Okuno, Y. Hasegawa, and H. Tsujino. A robot uses its own microphone to synchronize its steps to musical beats while scatting and singing. In IROS, pages 2459-2464, 2008.

【 0 0 5 8 】

まず、ビート間隔(テンポ)算出 4 3 0 は、周波数が直線的音階にあるスペクトログラム $p(t, \varphi)$ は、6 4 段階のメル尺度に周波数がある $p_{mel}(t, \varphi)$ に次式 (9) を用いて変換する。ビート間隔(テンポ)算出 4 3 0 は、オンセットのベクトル $d(t, \varphi)$ を、次式 (8) を用いて算出する。なお、式 (8) で式 (9) 用いているファイ、すなわち、 $d(t, \text{ファイ})$ のファイと、 $p_{mel}(t, \varphi)$ および $d(t, \varphi)$ で用いている (ファイ) は同じである。

【 0 0 5 9 】

【数 8】

$$d(t, \varphi) = \begin{cases} p_{mel}^{sobel}(t, \varphi) & (p_{mel}^{sobel}(t, \varphi) > 0) \\ 0 & otherwise \end{cases} \quad \dots (8)$$

【 0 0 6 0 】

【数 9】

$$\begin{aligned} p_{mel}^{sobel}(t, \varphi) = & -p_{mel}(t-1, \varphi+1) + p_{mel}(t+1, \varphi+1) \\ & -2p_{mel}(t-1, \varphi) + 2p_{mel}(t+1, \varphi) \\ & -p_{mel}(t-1, \varphi-1) + p_{mel}(t+1, \varphi-1) \end{aligned} \quad \dots (9)$$

【 0 0 6 1 】

式(9)は、ゾーベル・フィルタによるオンセット強調を意味する。

次に、ビート間隔(テンポ)算出430は、ビート間隔(テンポ)推定を行う。ビート間隔(テンポ)算出430は、ビート間隔の信頼性 $R(t, k)$ を、正規化相互相関を用いて次式(10)により算出する。

【0062】

【数10】

$$R(t, k) = \frac{\sum_j \sum_{l=0}^{P_w-1} d(t-l, j) d(t-k-l, j)}{\sqrt{\sum_j \sum_{k=l}^{P_w-1} d(t-l, j)^2 \sum_j \sum_{l=0}^{P_w-1} d(t-k-l, j)^2}} \quad \dots (10)$$

【0063】

式(10)において、 P_w は、信頼性算出におけるウィンドウの長さであり、 k は時間シフトパラメータである。ビート間隔(テンポ)算出430は、ビート間隔 $I(t)$ を時間シフト値 k に基づいて判断する。また、ビート間隔の信頼性 $R(t, k)$ は、局所的なピークの値をとる。

【0064】

ビート間隔(テンポ)算出430は、このように算出したビート間隔(テンポ)情報をテンポ推定部450に出力する。

【0065】

[音響信号と楽譜のマッチング]

マッチング部440には、音響信号からの特徴量抽出部410が抽出した音響クロマベクトルとオンセット時刻情報と、楽譜からの特徴量抽出420が抽出した楽譜クロマベクトルとレアネスと、テンポ推定部450が推定した安定化したテンポ情報とが入力される。マッチング部440は、 (t_n, f_m) を最終マッチング対とする。 t_n は音響信号における時間、 f_m は楽譜のフレーム・インデックスである。 t_{n+1} で検出された音響信号の新しいオンセット時刻及びその時間のテンポを考える場合、楽譜中の進むべきフレームの数 F は、マッチング部440により次式(11)のように推定される。

【0066】

【数11】

$$F = A(t_{n+1} - t_n) \quad \dots (11)$$

【0067】

式(11)において、係数 A は、テンポに対応し、音楽が速く進むと、係数 A は大きくなる。また、楽譜フレーム f_{m+k} の重み付けを次式(12)のように定義する。

【0068】

10

20

30

40

【数 1 2】

$$W(k) = \exp\left(-\frac{(f_{m+k} - f_m - F)^2}{2\sigma^2}\right) \quad \dots (12)$$

【0 0 6 9】

式(12)、 k は進むべき楽譜中のオンセット時刻数であり、 σ は重み付けの分散値である。本実施形態では、 $\sigma=24$ として実行したが、これは音符の長さの半分に相当する。ここで k は、負数となる可能性もあることに留意する。負数 k の場合、 (t_{n+1}, f_{m-1}) のようなマッチングを考慮することになるが、それはマッチングが楽譜内で逆行することを意味する。

10

【0 0 7 0】

マッチング部440は、対 (t_n, f_m) の類似性を、次式(13)を用いて算出する。

【0 0 7 1】

【数 1 3】

$$s(n, m) = \sum_{i=1}^{12} \sum_{\tau=t_n}^{t_{n+1}} r(i, m) c_{sco}(i, m) c_{sig}(i, \tau) \quad \dots (13)$$

20

【0 0 7 2】

式(13)において、 i は音名、 $r(i, m)$ はレアネス、 c_{sco} 及び c_{sig} はそれぞれ楽譜及び音響信号から生成されたクロマベクトルである。すなわち、マッチング部440は、対 (t_n, f_m) の類似性を、レアネスと音響クロマベクトルと楽譜クロマベクトルの積に基づいて算出する。

最終マッチングが (t_n, f_m) の時、新しいマッチングは (t_{n+1}, f_{m+k}) となり、そのとき進むべき楽譜中のオンセット時刻数 k は、次式(14)である。

30

【0 0 7 3】

【数 1 4】

$$k = \underset{l}{\operatorname{argmax}} W(l) S(n+1, m+l) \quad \dots (14)$$

【0 0 7 4】

本実施形態では、マッチング部440が行う処理の実行中、各マッチングステップの進むべき楽譜中のオンセット時刻数 k の探索範囲は、計算コスト削減のため2小節内に制限する。

40

【0 0 7 5】

マッチング部440は、式(11)～(14)を用いて、最終マッチング対 (t_n, f_m) を算出し、算出した最終マッチング対 (t_n, f_m) を歌声生成部130に出力する。

【0 0 7 6】

[スイッチング・カルマン・フィルタを用いたテンポ推定]

マッチング結果とビートトラッキング手法によるテンポ推定における2種類の誤差に対処するため、テンポ推定部450は、スイッチング・カルマン・フィルタ(手法3)を使用してテンポ推定を行う。

50

【 0 0 7 7 】

文献 3 (手法 3) K. P. Murphy. Switching kalman filters. Technical report, 1998.

【 0 0 7 8 】

テンポ推定部 4 5 0 が対応すべき 2 つの誤差とは、「演奏スピードのわずかな変化による小さな誤差」と「ビートトラッキングによるテンポ推定の外れ値による誤差」である。テンポ推定部 4 5 0 は、スイッチング・カルマン・フィルタで構成され、小さな観測誤差モデル 4 5 1 と外れ値となる大きな観測誤差モデル 4 5 2 の 2 つのモデルを備える。スイッチング・カルマン・フィルタとは、カルマン・フィルタ (K F) を拡張したものである。カルマン・フィルタは、状態遷移モデルと観測モデルを有する直線的予測フィルタであり、状態が観測不能の時、離散時間系列内の誤差を含む観測値からその状態を推定する。スイッチング・カルマン・フィルタは、複合的な状態遷移モデル及び観測モデルを有する。スイッチング・カルマン・フィルタが観測値を得る毎に、モデルを各モデルの可能性に基づき自動的に切り替える。

10

本実施形態において、スイッチング・カルマン・フィルタが備える小さな観測誤差モデル 4 5 1 と外れ値となる大きな観測誤差モデル 4 5 2 の 2 つのモデルにおいて、状態遷移などの他のモデリング成分は、2 つのモデルに共通している。

【 0 0 7 9 】

ビート時間とビート間隔を推定するため、本実施形態では、C e m g i l 他により提案された S K F モデル (手法 4) を使用する。

20

【 0 0 8 0 】

文献 4 (手法 4) A. T. Cemgil, B. Kappen, P. Desain, and H. Honing. On tempo tracking: Tempogram representation and kalman filtering. Journal of New Music Research, 28:4:259-273, 2001.

k 番目のビート時間を b_k 、その時間のビート間隔を τ_k 、そのテンポを一定とする。次のビート時間は $b_{k+1} = b_k + \tau_k$ 、次のビート間隔は $\tau_{k+1} = \tau_k$ として表される。ここで状態ベクトルを $x_k = [b_k \ \tau_k]^T$ とすると、状態遷移は次式 (1 5) のように表される。

【 0 0 8 1 】

【 数 1 5 】

30

$$x_{k+1} = F_k x_k + v_k = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} x_k + v_k \quad \dots (15)$$

【 0 0 8 2 】

式 (1 5) において、 F_k は状態遷移行列であり、 v_k は平均 0 の正規分散と共分散行列 Q から導かれた遷移誤差ベクトルである。最新の状態を x_k と仮定すると、テンポ推定部 4 5 0 は、次のビート時間 b_{k+1} を x_{k+1} の最初の成分として次式 (1 6) を用いて推定する。

40

【 0 0 8 3 】

【 数 1 6 】

$$x_{k+1} = F_k x_k \quad \dots (16)$$

【 0 0 8 4 】

ここで、観測ベクトルを $z_k = [b_k' \ \tau_k']^T$ とする。 b_k' は、マッチング部

50

4 4 0 がマッチング結果により算出したビート時間であり、 b_k' は、ビート間隔(テンポ)算出 4 3 0 がビートトラッキングにより算出したビート間隔である。テンポ推定部 4 5 0 は、観測ベクトルを、次式(17)を用いて算出する。

【0085】

【数17】

$$z_k = H_k x_k + w_k = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} x_k + w_k \quad \dots (17)$$

10

【0086】

式(17)において、 H_k は観測行列、 w_k は平均0の正規分散と共分散行列 R から導かれた観測誤差ベクトルである。本実施形態において、テンポ推定部 4 5 0 により、SKF は、観測誤差共分散行列 R^i ($i = 1, 2$) を切り換える。ここで、 i はモデル数である。予備実験から、本実施形態では、 R^i を以下の通りとした。小さな誤差モデルを $R^1 = \text{diag}(0.02, 0.005)$ 、外れ値モデルを $R^2 = \text{diag}(1, 0.125)$ とし、ここで $\text{diag}(a_1, \dots, a_n)$ は対角要素が左上から右下に a_1, \dots, a_n である $n \times n$ の対角行列である。

20

【0087】

図12は、カルマン・フィルタを適用したビートトラッキングを説明する図である。縦軸はテンポ、横軸は時間を表している。図12(a)は、ビートトラッキングにおける誤差を説明する図であり、図12(b)は、ビートトラッキング手法のみの解析結果とカルマン・フィルタ適用後の解析結果を示す図である。図12(a)において、符号501の部分是小ノイズであり、符号502の部分がビートトラッキング手法で推定したテンポにおける外れ値の例である。

図12(b)において、実線511は、ビートトラッキング手法のみによるテンポの解析結果であり、点線512は、ビートトラッキング手法による解析結果に対し、さらに本実施形態の方法でカルマン・フィルタを適用した解析結果である。図12(b)のように、本実施形態の方法を適用した結果、ビートトラッキング手法のみと比較して、テンポの外れの影響を大幅に改善できる。

30

【0088】

[ビート時間の観測]

図9(b)で説明したように、楽譜が48分音符に対応した長さのフレームに分割されているので、ビートは楽譜中の12フレーム毎にある。テンポ推定部 4 5 0 は、算出したビート時間 b_k' が、 k 番目のビートフレームに音符が存在しないとき、マッチング部 4 4 0 が行ったマッチング結果により補間する。

テンポ推定部 4 5 0 は、算出したビート時間 b_k' とビート間隔情報をマッチング部 4 4 0 に出力する。

40

【0089】

[楽譜位置推定処理の手順]

次に、楽譜位置推定装置 1 0 0 が行う楽譜位置推定処理の手順を、図13を用いて説明する。図13は、楽譜位置推定処理のフローチャートである。

まず、楽譜からの特徴量抽出 4 2 0 は、楽譜データベース 1 2 1 から楽譜データを読み出す。楽譜からの特徴量抽出 4 2 0 は、読み出した楽譜データから、式(5)~式(7)を用いて楽譜クロマベクトルとレアネスを算出し、算出した楽譜クロマベクトルとレアネスをマッチング部 4 4 0 に出力する(ステップS1)。

【0090】

次に、楽曲位置推定部 1 2 2 は、マイクロホン 3 0 が集音した音響信号に基づき、演奏

50

が継続しているか否かを判別する（ステップＳ２）。なお、この判別は、例えば、楽曲位置推定部１２２が、音響信号が継続している場合に曲が継続していると判別し、または、演奏されている曲の演奏位置が楽譜の終端ではない場合に演奏が継続していると判別する。

ステップＳ２の判別の結果、演奏が継続していないと判別された場合（ステップＳ２；Ｎｏ）、楽譜位置推定の処理を終了する。

【００９１】

ステップＳ２の判別の結果、演奏が継続していると判別された場合（ステップＳ２；Ｙｅｓ）、音響信号分離部１１０は、マイクロホン３０が集音した音響信号を、例えば１秒間分、音響信号分離部１１０が備えるバッファに記憶させる（ステップＳ３）。 10

次に、音響信号分離部１１０は、入力された音響信号と歌声生成部１３０が生成した音声信号を用いて、独立成分分析を行って残響音の抑圧と自身の歌声の抑圧を行うことで音響信号を抽出し、抽出した音響信号を楽譜位置推定部１２０に出力する。

次に、ビート間隔(テンポ)算出４３０は、入力された音楽信号に基づき、式（８）～式（１０）を用いてビートトラッキング手法によりビート間隔(テンポ)を推定し、推定したビート間隔(テンポ)をマッチング部４４０に出力する（ステップＳ４）。

【００９２】

次に、音響信号からの特徴量抽出部４１０は、入力された音響信号から式（４）を用いて、オンセット時刻情報を検出し、検出したオンセット時刻情報をマッチング部４４０に出力する（ステップＳ５）。 20

次に、音響信号からの特徴量抽出部４１０は、入力された音響信号に基づき、式（８）～式（３）を用いて音響クロマベクトルを抽出し、抽出した音響クロマベクトルをマッチング部４４０に出力する（ステップＳ６）。

【００９３】

次に、マッチング部４４０には、音響信号からの特徴量抽出部４１０から音響クロマベクトルとオンセット時刻情報と、楽譜からの特徴量抽出４２０から楽譜クロマベクトルとレアネスと、テンポ推定部４５０から推定された安定したテンポ情報とが入力される。マッチング部４４０は、式（１１）～式（１４）を用いて、入力された音響クロマベクトルと楽譜クロマベクトルとを、逐次、マッチング処理を行い、最終マッチング対（ t_n, f_m ）を推定する。マッチング部４４０は、推定した楽譜位置に対応する最終マッチング対（ t_n, f_m ）をテンポ推定部４５０と歌声生成部１３０に出力する（ステップＳ７）。 30

【００９４】

次に、テンポ推定部４５０には、ビート間隔(テンポ)算出４３０から入力されたビート間隔(テンポ)情報に基づき、式（１５）～式（３）を用いてビート時間 b_k' とビート間隔情報を算出し、算出したビート時間 b_k' とビート間隔情報をマッチング部４４０に出力する（ステップＳ８）。

また、テンポ推定部４５０には、マッチング部４４０から最終マッチング対（ t_n, f_m ）が入力される。テンポ推定部４５０は、算出したビート時間 b_k' が、 k 番目のビートフレームに音符が存在しないとき、マッチング部４４０が行ったマッチング結果により補間する。 40

なお、マッチング部４４０とテンポ推定部４５０とは、マッチング処理とテンポ推定を逐次的に行い、マッチング部４４０が、最終マッチング対（ t_n, f_m ）を推定する。

【００９５】

歌声生成部１３０の音声生成部１３２は、入力された最終マッチング対（ t_n, f_m ）に基づき、歌詞とメロディーのデータベース１３１を参照し、楽譜位置に合致する歌詞をメロディーに合わせて歌声を生成する。なお、ここで、「歌声」とは、楽譜位置推定装置１００からスピーカ２０を介して出力される音声データである。すなわち、楽譜位置推定装置１００を備えるロボット１のスピーカ２０を介して出力されるものである。便宜的に「歌声」という。また、本実施形態において、音声生成部１３２には、ＶＯＣＡＬＯＩＤ２（ＶＯＣＡＬＯＩＤ（登録商標））を用いて、歌声を生成した。ＶＯＣＡＬＯＩＤ 50

2 (VOCALOID (登録商標)) は、メロディーと歌詞を入力することでサンプリングされた人の声を元にした歌声を合成することができるエンジンのため、本実施形態では、さらに楽譜位置を情報として加え、実際の演奏から歌声が外れないようにしている。

音声生成部 132 は、生成した音声信号をスピーカ 20 から出力する。

【0096】

また、最終マッチング対 (t_n, f_m) 推定後、ステップ S2 ~ ステップ S8 を曲の演奏が終了するまで逐次的に行う。

以上の処理により、楽譜位置を推定し、推定した楽譜位置に合致する音声 (歌声) を生成し、生成した音声はスピーカ 20 から出力することで、ロボット 1 が演奏に合わせた歌唱を行うことが可能になる。また、本実施形態によれば、演奏されている音響信号に基づいて、楽譜の位置を推定するようにしたので、曲が途中から開始された場合においても、正確に楽譜の位置を推定することができる。

【0097】

[評価結果]

本実施形態における楽譜位置推定装置 100 を用いて行った評価結果について説明する。まず、実験条件について説明する。評価に用いた楽曲は、後藤らにより作成された RWC 研究用音楽データベース (RWC-MDB-P-2001; <http://staff.aist.go.jp/m.goto/RWC-MDB/index-j.html>) からのポピュラー音楽 100 曲を使用した。また、使用した楽曲については、歌唱部分や演奏部分を含むこれらの楽曲のフルバージョンを使用した。

【0098】

楽譜同期の正解データは、評価者が各楽曲の MIDI ファイルから生成した。これら MIDI ファイルは、実際の演奏に厳密に同期される。誤差は、秒単位で、本実施形態により抽出されたビート時間と、正解データとの相違の絶対値として定義される。誤差は楽曲毎に平均化される。

【0099】

評価は、以下の 4 種類について行い、評価結果を比較した。

(i) 本実施形態の方法; SKF およびレアネス使用

(ii) SKF 無使用; テンポ推定への修正なし

(iii) レアネス無使用; 全音符は同等のレアネスを有する状態

(iv) ビートトラッキング手法; この手法は、音楽の最初からビートを数えることにより楽譜位置を判断する。

【0100】

さらに、楽譜位置推定装置 100 のマイクロホン 30 が集音する音は、室内環境における残響に影響を与えるのかをも評価するため、以下の 2 種類の音楽信号を使用して評価を行った。

(v) クリーン音楽信号: 残響なしの音楽信号

(vi) 残響あり音楽信号: 残響付きの音楽信号

残響は、インパルス応答畳み込みによりシミュレートしたものを使用した。図 14 は、楽譜位置推定装置 100 を備えるロボット 1 と音源の設置関係を説明する図である。図 14 のように、評価用の音源は、ロボット 1 の正面から 100 [cm] 離れた位置に設置したスピーカ 601 から出力した音源を用いた。この生成されたインパルス応答は、実験室において測定した。実験室での残響時間 (RT20) は、156 [msec] である。講堂又は音楽ホールであれば、より長い残響時間となると考えられる。

【0101】

図 15 は、2 種類の音楽信号 ((v) と (vi)) と 4 つの手法 ((i) ~ (iv)) の結果を示している。各値は 100 曲についての累積絶対値誤差の平均値と標準偏差である。クリーン信号及び残響あり信号の双方において、本実施形態による方法 (i) の誤差は、ビートトラッキング手法 (iv) の誤差より少ない。本実施形態による方法 (i) は、誤差をクリーン信号で 29%、残響あり信号で 14% 改善している。本実施形態による方法 (i) は、SKF を無使用の手法 (ii) より誤差が少ないことから、SKF を用いるこ

10

20

30

40

50

とで誤差低減されていることがわかる。同様に、本実施形態の方法(i)とレアネス無使用の手法(i i i)の結果を比較すると、レアネスが誤差を減少させている。

さらに、SKF無使用の手法(i i)は、レアネス無使用の手法(i i i)より誤差が大きいので、SKFはレアネスより一層効果的であると言える。これは、しばしばレアネスが、楽譜中のフレームとドラム音のような誤ったオンセット時刻との間で高い類似性を誘引するからである。仮にドラム音が、高いレアネスを伴い、クロマベクトル成分中に大きなパワーを持つとすると、これが誤ったマッチングとなる。この問題を避けるため、楽譜位置推定装置100では、単一音名でなく組み合わせた音名へのレアネスの考慮をすることが可能である。

【0102】

図16は、クリーン信号時の各手法の累積絶対値誤差平均値で分類された楽曲数を示している。図17は、残響あり信号時の各手法の累積絶対値誤差平均値で分類された楽曲数を示している。図16と図17において、より少ない平均誤差を有する楽曲数が多いほど、よりよい演奏を示している。クリーン信号では、本実施形態の方法(i)では、2秒以下の誤差を有する楽曲が31曲あるのに対し、ビートトラッキングのみの手法(i v)では9曲であった。

残響あり信号では、同様に、本実施形態の方法(i)では、2秒以下の誤差を有する楽曲が36曲あるのに対し、ビートトラッキングのみの手法(i v)では12曲であった。このように、より少ない誤差で楽譜位置を推定できる点から、本実施形態の方法は、ビートトラッキング手法より優れている。これは音楽に合わせて自然な歌声を発生させること

【0103】

本実施形態の方法による分類において、クリーン信号と残響あり信号との間に大差ないが、本実施形態の方法は、図15のように、残響あり信号においてより多くの誤差を有する。したがって、実験室の残響が多くの誤差を含む曲に主に影響している。残響は少ない誤差を含む曲にはあまり影響しない。音楽ホール内のような、より長い残響を有する環境においては、楽譜同期の精度に悪影響を与えることも考えられる。

このため、本実施形態では、楽譜位置を推定するために、音響信号分離部110が独立成分分析を行って残響音の抑圧を行った後の音響信号を用いているので、この場合においても、残響の影響を軽減して、精度の高い楽譜同期を行うことができる。

【0104】

このため、ドラム音ありとドラム音なしの楽曲の誤差を比較することにより、本実施形態の方法の精度が、楽曲中でドラムが演奏されるか否かに依存することを評価した。ドラム音ありとドラム音なしの楽曲数は各89と11である。ドラムあり楽曲の累積絶対値誤差平均値は7.37[秒]であり標準偏差は9.4[秒]である。一方、ドラムなし楽曲の平均累積誤差は22.1[秒]であり標準偏差は14.5[秒]である。ビートトラッキングによるテンポ推定は、ドラム音がない時、非常に大きな変動を生じやすい。これは、高い累積誤差を引き起こす、不正確なマッチングの原因となる。

本実施形態では、ドラムなどの低音領域の影響を軽減するため、図10で説明したように、高周波成分に重み付けを行い、重み付けしたパワーからオンセット時刻を検出する

【0105】

本実施形態では、楽譜位置推定装置100をロボット1に適用し、ロボット1が演奏に合わせて歌う(歌声を、スピーカ20を介して出力する)例について説明したが、推定した楽譜位置情報に基づき、さらにロボット1が演奏に合わせて自身の可動部を動かし、あたかもロボット1が演奏に合わせて、リズムに合わせて体を動かしているようにロボット1が備える制御部により制御するようにしてもよい。

【0106】

また、本実施形態では、楽譜位置推定装置100をロボット1に適用する例を説明したが、他の装置に適用してもよく、例えば携帯電話等に適用するようにしてもよく、あるい

10

20

30

40

50

は演奏に合わせて歌う歌唱装置に適用してもよい。

【 0 1 0 7 】

また、本実施形態では、マッチング部 4 4 0 において、レアネスを用いて重み付けを行う例を説明したが、重み付けは、他の要素により行うようにしてもよい。また、音符の出現頻度が低いと判定された場合においても、特定の前後するフレームにおいて、出現頻度の低い音符と判定された音符の出現頻度が高い場合などは、出現頻度が高い音符や、出現頻度が平均的なものを用いてもよい。

【 0 1 0 8 】

また、本実施形態では、ビート間隔(テンポ)算出 4 3 0 で、楽譜が 4 8 分音符に対応した長さのフレームに分割する例を説明したが、他の分割値でも良い。また、バッファを 1 秒間行う例を説明したが、バッファする時間は 1 秒でなくてもよく、処理に用いる時間以上分のデータを含むようにしてもよい。

【 0 1 0 9 】

なお、実施形態の図 2 と図 7 の各部の機能を実現するためのプログラムをコンピュータ読み取り可能な記録媒体に記録して、この記録媒体に記録されたプログラムをコンピュータシステムに読み込ませ、実行することにより各部の処理を行ってもよい。なお、ここでいう「コンピュータシステム」とは、OS や周辺機器等のハードウェアを含むものとする。

また、「コンピュータシステム」は、WWWシステムを利用している場合であれば、ホームページ提供環境（あるいは表示環境）も含むものとする。

また、「コンピュータ読み取り可能な記録媒体」とは、フレキシブルディスク、光磁気ディスク、ROM (Read Only Memory)、CD-ROM等の可搬媒体、USB (Universal Serial Bus) I/F (インタフェース) を介して接続されるUSBメモリー、コンピュータシステムに内蔵されるハードディスク等の記憶装置のことをいう。さらに「コンピュータ読み取り可能な記録媒体」とは、インターネット等のネットワークや電話回線等の通信回線を介してプログラムを送信する場合の通信線のように、短時間の間、動的にプログラムを保持するもの、その場合のサーバやクライアントとなるコンピュータシステム内部の揮発性メモリーのように、一定時間プログラムを保持しているものも含むものとする。また上記プログラムは、前述した機能の一部を実現するためのものであっても良く、さらに前述した機能をコンピュータシステムにすでに記録されているプログラムとの組み合わせで実現できるものであっても良い。

【符号の説明】

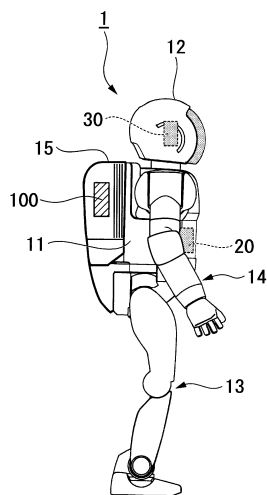
【 0 1 1 0 】

- 1・・・ロボット
- 1 1・・・基体部
- 1 2・・・頭部（可動部）
- 1 3・・・脚部（可動部）
- 1 4・・・腕部（可動部）
- 1 5・・・収納部
- 2 0・・・スピーカ
- 3 0・・・マイクロホン
- 1 0 0・・・楽譜位置推定装置
- 1 1 0・・・音響信号分離部
- 1 1 1・・・自己生成音抑制フィルタ部
- 1 2 0・・・楽譜位置推定部（楽譜情報取得部、音響信号の特徴量抽出部、楽譜情報の特徴量抽出部、ビート位置推定部、マッチング部）
- 1 2 1・・・楽譜データベース
- 1 2 2・・・楽曲位置推定部（音響信号の特徴量抽出部、楽譜情報の特徴量抽出部、ビート位置推定部、マッチング部）
- 1 3 0・・・歌声生成部

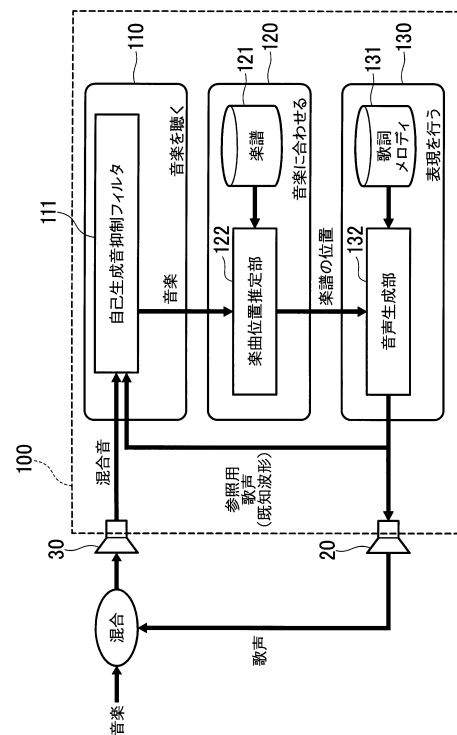
- 1 3 1 . . . 歌詞とメロディーのデータベース
- 1 3 2 . . . 音声生成部
- 4 1 0 . . . 音響信号からの特徴量抽出部（音響信号の特徴量抽出部）
- 4 2 0 . . . 楽譜からの特徴量抽出（楽譜情報の特徴量抽出部）
- 4 3 0 . . . ビート間隔(テンポ)算出
- 4 4 0 . . . マッチング部
- 4 4 1 . . . 類似度計算部
- 4 4 2 . . . 重み付け計算部
- 4 5 0 . . . テンポ推定部（ビート位置推定部）
- 4 5 1 . . . 小さな観測誤差モデル
- 4 5 2 . . . 大きな観測誤差モデル

10

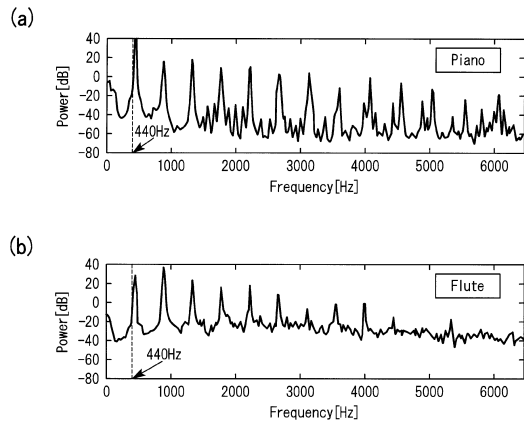
【図 1】



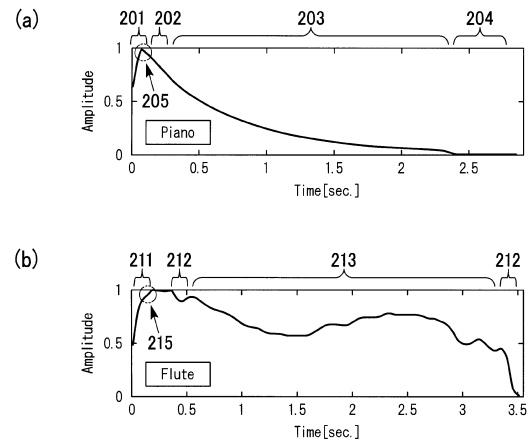
【図 2】



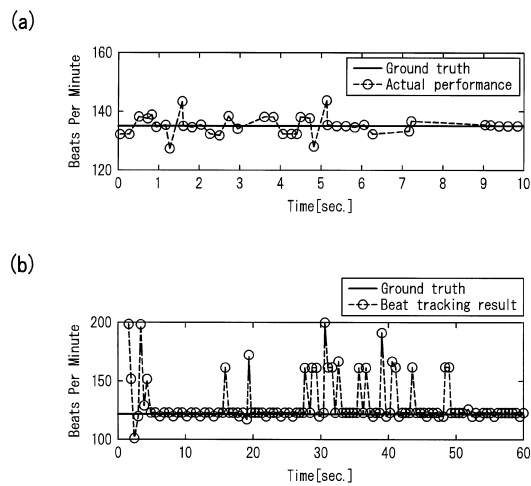
【図 3】



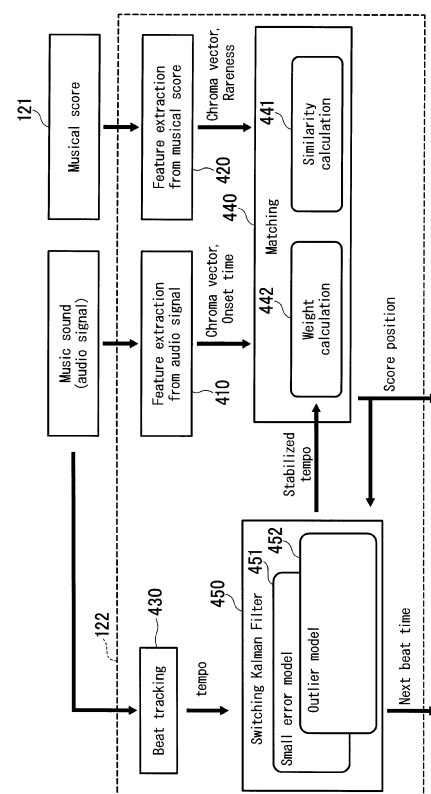
【図 4】



【図 6】



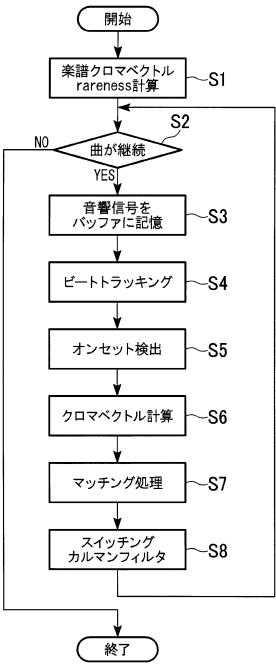
【図 7】



【図 8】

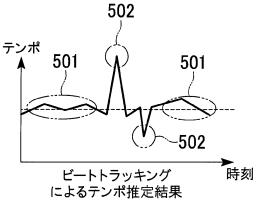
Symbols	Definitions
i	index for 12 pitch names(C,C#,...B)
t	time frame of audio signal
n	index for onsets in audio signal
t_n	n -th onset time in audio signal
f	frame index of musical score
m	index for onsets in musical score
f_m	m -th onset time frame in musical score

【図 13】

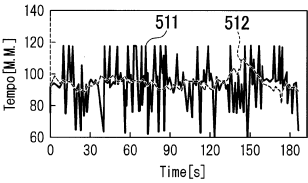


【図 12】

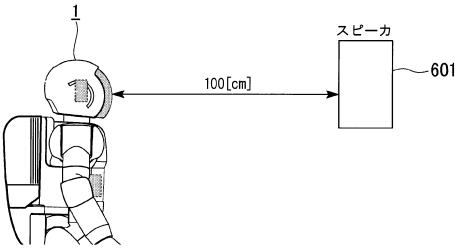
(a)



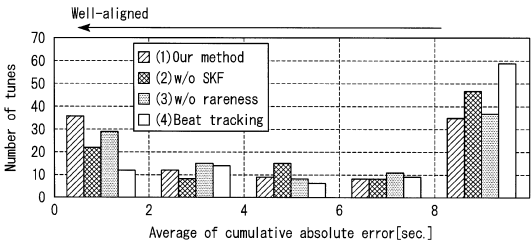
(b)



【図 14】



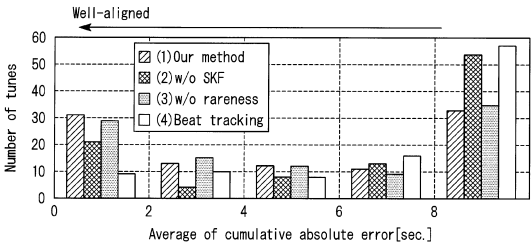
【図 17】



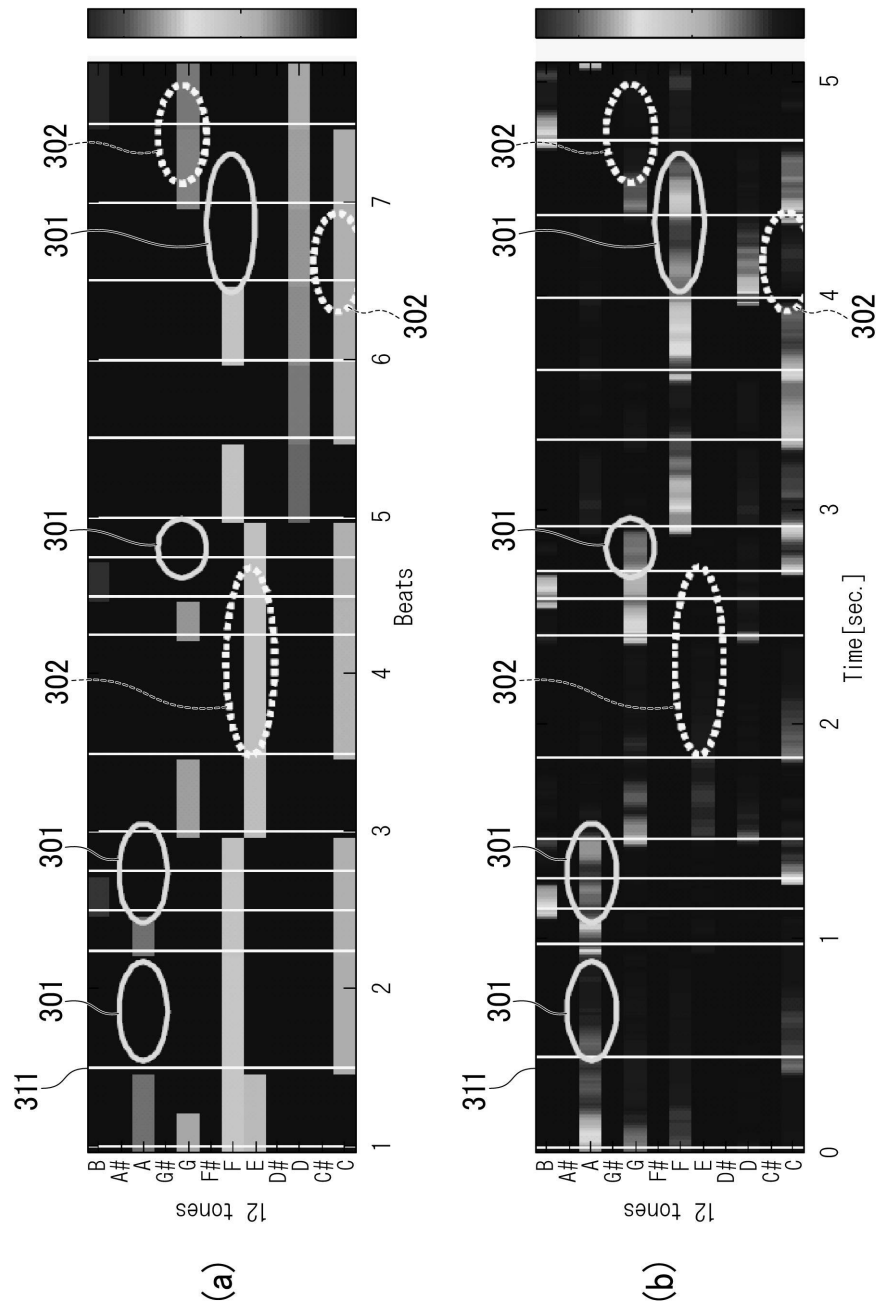
【図 15】

	Clean signal		Reverberated signal	
	Ave.	Std. dev.	Ave.	Std. dev.
(1)Our method	8.9	11.0	9.9	12.6
(2)w/o SKF	11.6	12.8	10.5	11.2
(3)w/o rareness	9.7	13.5	10.3	13.3
(4)Beat tracking	12.5	9.7	11.5	9.1

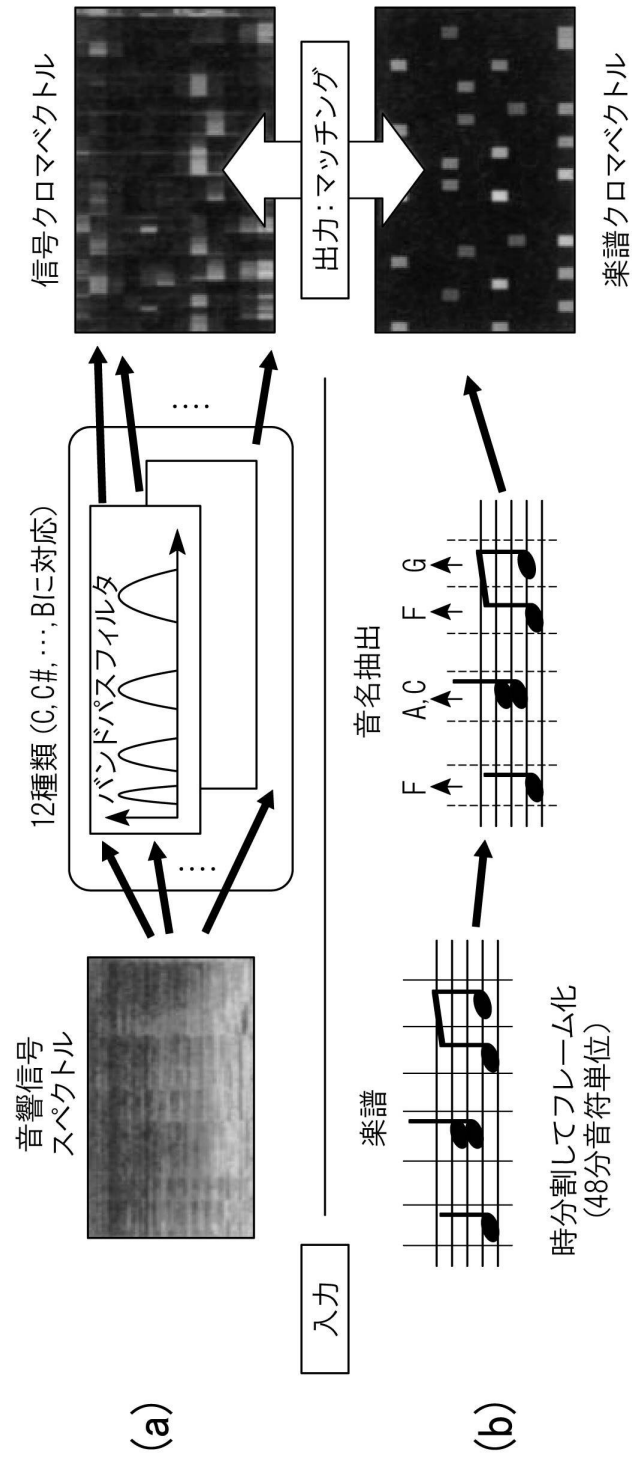
【図 16】



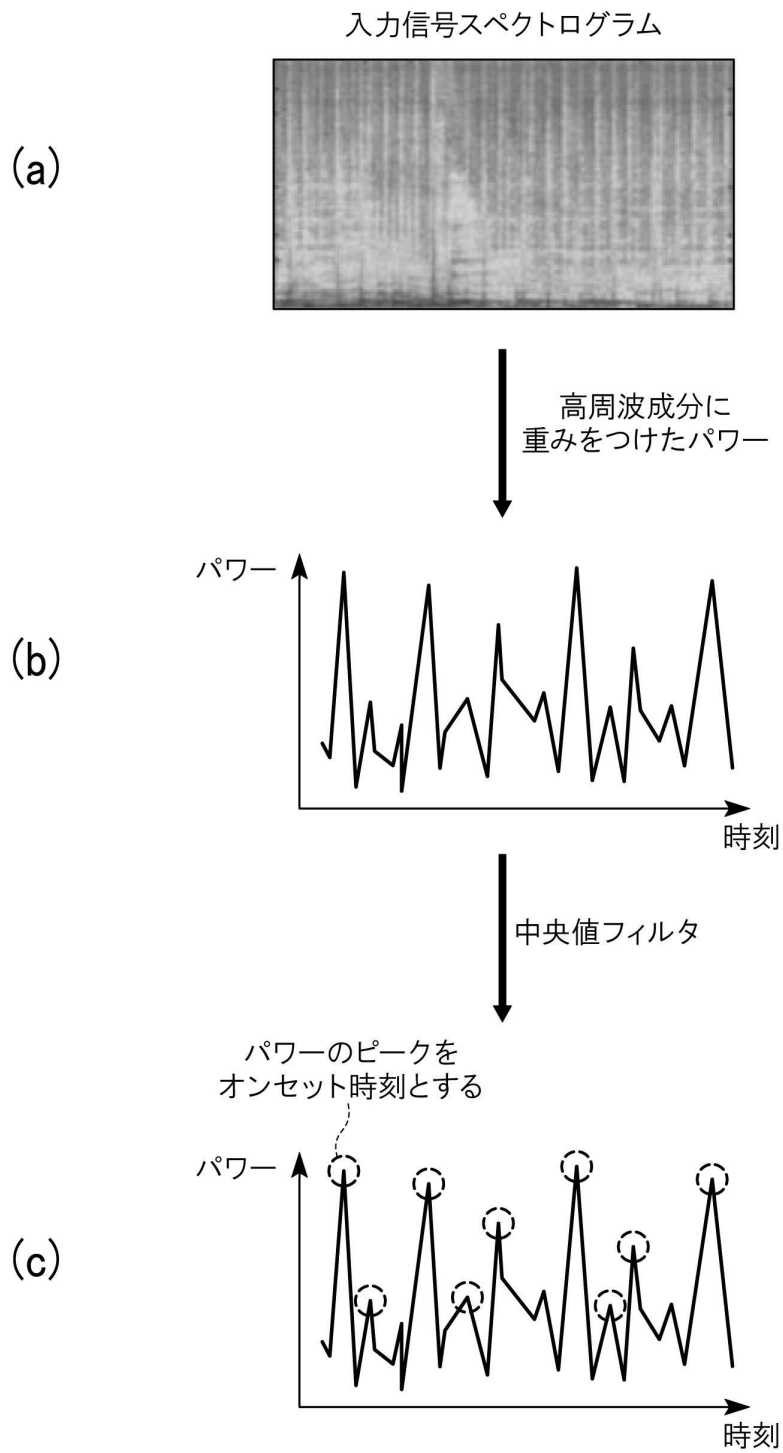
【図 5】



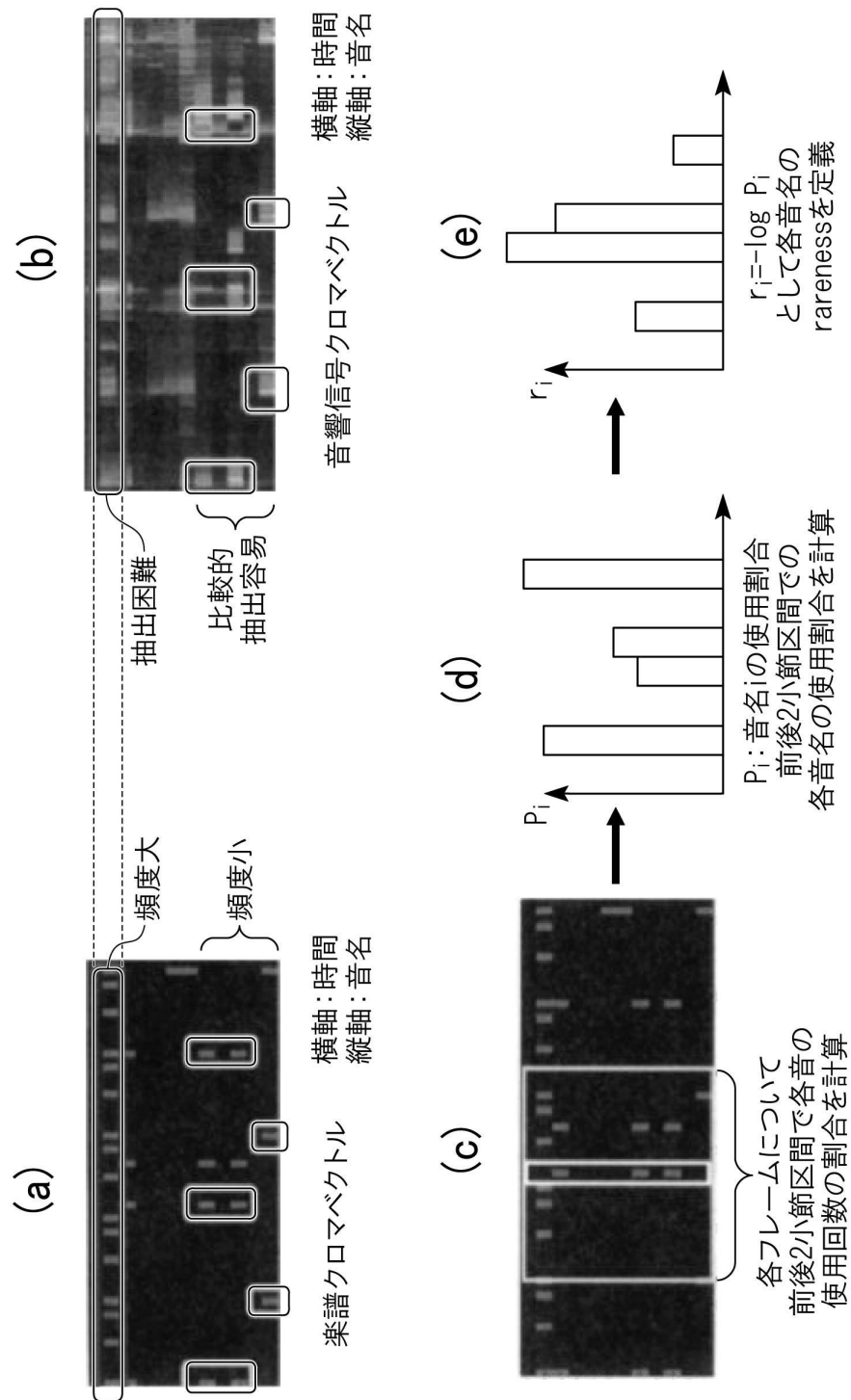
【図9】



【図 10】



【図 11】



フロントページの続き

- (72)発明者 中臺 一博
埼玉県和光市本町 8 - 1 株式会社ホンダ・リサーチ・インスティテュート・ジャパン内
- (72)発明者 大塚 琢馬
京都府京都市左京区吉田本町 国立大学法人京都大学 大学院情報学研究科内
- (72)発明者 奥乃 博
京都府京都市左京区吉田本町 国立大学法人京都大学 大学院情報学研究科内

審査官 上田 雄

- (56)参考文献 米国特許第 0 6 1 0 7 5 5 9 (U S , A)
大塚琢馬、村田和馬、武田龍、中臺一博、高橋徹、尾形哲也、奥乃博、歌唱ロボットのためのビート情報とメロディ・ハーモニー情報の統合による音楽音響信号と楽譜の実時間同期手法の開発、情報処理学会第 7 1 回全国大会講演論文集、日本、社団法人情報処理学会、2 0 0 9 年 3 月 1 0 日、pp. 2-243 - 2-244
Juan Pablo Bello, Guiliano Monti and Mark Sandler, Techniques for automatic music transcription, Proc. ISMIR 2000, 米国、2 0 0 0 年 1 0 月 2 3 日
A. T. Cemgil, B. Kappen, P. Desain and H. Honing, On tempo tracking: Tempogram Representation and Kalman filtering, Journal of New Music Research, 2 0 0 1 年 1 2 月, Volume 28, Issue 4, pp. 259-273

(58)調査した分野(Int.Cl., D B 名)

G 1 0 G 1 / 0 0 - 3 / 0 4
G 1 0 L 2 5 / 5 1