

(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2020-3536

(P2020-3536A)

(43) 公開日 令和2年1月9日(2020.1.9)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 0 G 3/04 (2006.01)	G 1 0 G 3/04	5 D 1 8 2
G 1 0 H 1/00 (2006.01)	G 1 0 H 1/00 B	5 D 4 7 8
G 1 0 L 25/51 (2013.01)	G 1 0 H 1/00 Z	
G 1 0 L 25/30 (2013.01)	G 1 0 L 25/51 3 0 0	
G 0 6 N 20/00 (2019.01)	G 1 0 L 25/30	
審査請求 未請求 請求項の数 13 O L (全 17 頁) 最終頁に続く		

(21) 出願番号 特願2018-120235 (P2018-120235)
 (22) 出願日 平成30年6月25日 (2018. 6. 25)

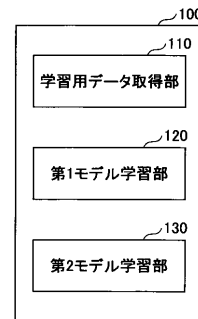
(71) 出願人 000001443
 カシオ計算機株式会社
 東京都渋谷区本町1丁目6番2号
 (74) 代理人 100107766
 弁理士 伊東 忠重
 (74) 代理人 100070150
 弁理士 伊東 忠彦
 (72) 発明者 日暮 大輝
 東京都羽村市栄町3丁目2番1号 カシオ
 計算機株式会社 羽村技術センター内
 Fターム(参考) 5D182 AC01 AD01
 5D478 DF02

(54) 【発明の名称】 学習装置、自動採譜装置、学習方法、自動採譜方法及びプログラム

(57) 【要約】 (修正有)

【課題】 各音の音高や区間が明確でないオーディオデータから楽譜を自動生成するための音響処理技術を提供する。

【解決手段】 学習装置100は、単音音源と音高情報とを第1の機械学習モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを第2の機械学習モデルの学習用データとして取得し、単音音源と前記採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得する学習用データ取得部と、単音音源のスペクトログラムを学習用入力データとして入力し、単音音源の音高の予測確率を出力するよう音高情報によって学習する第1モデル学習部と、採譜対象の音源のスペクトログラムを学習済みの第1の機械学習モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、特徴マップの固定長の区間に音符が存在する予測確率を出力するよう楽譜情報によって学習する第2モデル学習部と、を有する。



【選択図】 図2

【特許請求の範囲】

【請求項 1】

単音音源と音高情報とを第 1 の機械学習モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを第 2 の機械学習モデルの学習用データとして取得し、前記単音音源と前記採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得する学習用データ取得部と、

前記単音音源のスペクトログラムを学習用入力データとして入力し、前記単音音源の音高の予測確率を出力するよう前記音高情報によって第 1 の機械学習モデルを学習する第 1 モデル学習部と、

前記採譜対象の音源のスペクトログラムを学習済みの前記第 1 の機械学習モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するよう前記楽譜情報によって第 2 の機械学習モデルを学習する第 2 モデル学習部と、
を有する学習装置。

10

【請求項 2】

前記第 1 の機械学習モデルと前記第 2 の機械学習モデルとは、畳み込みニューラルネットワークにより構成される、請求項 1 記載の学習装置。

【請求項 3】

前記第 2 モデル学習部は、前記第 1 の機械学習モデルにより生成される異なる時間解像度を有する複数の特徴マップを前記第 2 の機械学習モデルに入力する、請求項 2 記載の学習装置。

20

【請求項 4】

前記第 2 モデル学習部は、前記第 1 の機械学習モデルと前記第 2 の機械学習モデルとを SSD (Single Shot Detection) として実現する、請求項 1 乃至 3 何れか一項記載の学習装置。

【請求項 5】

前記第 1 モデル学習部は、複数種別のオーディオ成分のそれぞれに対して前記第 1 の機械学習モデルを学習し、

前記第 2 モデル学習部は、複数種別のオーディオ成分を含む採譜対象の音源に対して各オーディオ成分種別毎に音符が存在する予測確率を出力するよう前記第 2 の機械学習モデルを学習する、請求項 1 乃至 4 何れか一項記載の学習装置。

30

【請求項 6】

単音音源から音高の予測確率を出力する第 1 の学習済み機械学習モデルと、特徴マップから前記特徴マップの固定長の区間に音符が存在する予測確率を出力する第 2 の学習済み機械学習モデルとを利用し、採譜対象の音源を前記第 1 の学習済み機械学習モデルに入力し、前記第 1 の学習済み機械学習モデルによって生成された特徴マップを前記第 2 の学習済み機械学習モデルに入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するモデル処理部と、

前記音符が存在する予測確率に基づき楽譜情報を生成する楽譜生成部と、
を有する自動採譜装置。

40

【請求項 7】

前記モデル処理部は、前記採譜対象の音源に対して前処理を実行することによってスペクトログラムを取得し、前記スペクトログラムを前記第 1 の学習済み機械学習モデルに入力する、請求項 6 記載の自動採譜装置。

【請求項 8】

前記モデル処理部は、前記特徴マップ上の各点について前記第 2 の学習済み機械学習モデルから出力された最大の予測確率を有する音符を予測音符として決定する、請求項 6 又は 7 記載の自動採譜装置。

【請求項 9】

前記楽譜生成部は、NMS (Non-Maximum Suppression) に従

50

って抽出された予測音符に基づき楽譜情報を生成する、請求項 8 記載の自動採譜装置。

【請求項 10】

プロセッサが、単音音源と音高情報とを第 1 の機械学習モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを第 2 の機械学習モデルの学習用データとして取得し、前記単音音源と前記採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得すステップと、

前記プロセッサが、前記単音音源のスペクトログラムを学習用入力データとして入力し、前記単音音源の音高の予測確率を出力するよう前記音高情報によって第 1 の機械学習モデルを学習するステップと、

前記プロセッサが、前記採譜対象の音源のスペクトログラムを学習済みの前記第 1 の機械学習モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するよう前記楽譜情報によって第 2 の機械学習モデルを学習するステップと、

を有する学習方法。

【請求項 11】

プロセッサが、単音音源から音高の予測確率を出力する第 1 の学習済み機械学習モデルに採譜対象の音源を入力するステップと、

前記プロセッサが、特徴マップから前記特徴マップの固定長の区間に音符が存在する予測確率を出力する第 2 の学習済み機械学習モデルに前記第 1 の学習済み機械学習モデルによって生成された特徴マップを入力するステップと、

前記プロセッサが、前記第 2 の学習済み機械学習モデルから出力された前記音符が存在する予測確率に基づき楽譜情報を生成するステップと、

を有する自動採譜方法。

【請求項 12】

単音音源と音高情報とを第 1 の機械学習モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを第 2 の機械学習モデルの学習用データとして取得し、前記単音音源と前記採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得すステップと、

前記単音音源のスペクトログラムを学習用入力データとして入力し、前記単音音源の音高の予測確率を出力するよう前記音高情報によって第 1 の機械学習モデルを学習するステップと、

前記採譜対象の音源のスペクトログラムを学習済みの前記第 1 の機械学習モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するよう前記楽譜情報によって第 2 の機械学習モデルを学習するステップと、

をプロセッサに実行させるプログラム。

【請求項 13】

単音音源から音高の予測確率を出力する第 1 の学習済み機械学習モデルに採譜対象の音源を入力するステップと、

特徴マップから前記特徴マップの固定長の区間に音符が存在する予測確率を出力する第 2 の学習済み機械学習モデルに前記第 1 の学習済み機械学習モデルによって生成された特徴マップを入力するステップと、

前記第 2 の学習済み機械学習モデルから出力された前記音符が存在する予測確率に基づき楽譜情報を生成するステップと、

をプロセッサに実行させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、音響処理技術に関する。

【背景技術】

10

20

30

40

50

【0002】

オーディオデータから楽譜を自動生成する自動採譜技術が従来から知られている。例えば、特開2007-033479には、同時に複数の音が演奏される場合でも単一楽器により演奏された音響信号から楽譜を自動採譜する技術が記載されている。

【先行技術文献】

【特許文献】

【0003】

【特許文献1】特開2007-033479

【発明の概要】

【発明が解決しようとする課題】

10

【0004】

しかしながら、従来の自動採譜では、楽譜に対して正確に演奏又は歌唱され、各音の音高や区間が明確なオーディオデータの場合には比較的高精度な採譜が可能であるが、例えば、各音の音高や区間が明確でないオーディオデータの場合には期待するような自動採譜が困難であった。

【0005】

上記問題点を鑑み、本開示の課題は、様々なオーディオデータからより効果的に楽譜を自動生成するための音響処理技術を提供することである。

【課題を解決するための手段】

【0006】

20

上記課題を解決するため、本開示の一態様は、単音音源と音高情報とを第1の機械学習モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを第2の機械学習モデルの学習用データとして取得し、前記単音音源と前記採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得する学習用データ取得部と、前記単音音源のスペクトログラムを学習用入力データとして入力し、前記単音音源の音高の予測確率を出力するよう前記音高情報によって第1の機械学習モデルを学習する第1モデル学習部と、前記採譜対象の音源のスペクトログラムを学習済みの前記第1の機械学習モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するよう前記楽譜情報によって第2の機械学習モデルを学習する第2モデル学習部と、を有する学習装置に関する。

30

【発明の効果】

【0007】

本開示によると、各音の音高や区間が明確でないオーディオデータから楽譜を自動生成するための音響処理技術を提供することができる。

【図面の簡単な説明】

【0008】

【図1】本開示の一実施例による学習済み機械学習モデルを有する自動採譜装置を示す概略図である。

【図2】本開示の一実施例による学習装置の機能構成を示すブロック図である。

【図3】本開示の一実施例による特徴マップ生成モデルの構成を示す概略図である。

40

【図4】本開示の一実施例による音符存在確率予測モデルの構成を示す概略図である。

【図5】本開示の一実施例による特徴マップとデフォルトボックスとの関係を示す概念図である。

【図6】本開示の一実施例による特徴マップ生成モデルの学習処理を示すフローチャートである。

【図7】本開示の一実施例による音符存在確率予測モデルの学習処理を示すフローチャートである。

【図8】本開示の一実施例による自動採譜装置の機能構成を示すブロック図である。

【図9】本開示の一実施例による自動採譜処理を示すフローチャートである。

【図10】本開示の一実施例による学習装置及び自動採譜装置のハードウェア構成を示す

50

ブロック図である。

【発明を実施するための形態】

【0009】

以下の実施例では、機械学習モデルによって音源（音の波形データであるオーディオデータ）から楽譜情報を生成する自動採譜装置が開示される。

【0010】

従来の自動採譜技術では、音高の予測に注力され、音符の切れ目を示すオンセットとオフセットとの予測は自動採譜における課題の1つであった。本開示による自動採譜装置は、音源におけるオンセットとオフセットとを機械学習モデルの1つであるSSD（Single Shot Detection）によって予測する。

10

【0011】

SSDは、1つのニューラルネットワークを用いて入力画像における物体を検出手法である。すなわち、当該ニューラルネットワークへの入力は画像であり、その出力は複数の矩形領域（SSDでは、デフォルトボックスと呼ばれる）の中心座標、高さ、幅及び物体の種類（予測確率）である。デフォルトボックスは入力画像のサイズによって予め設定された個数の候補として用意され、後処理（NMS：Non-Maximum Suppressionなど）によって大部分のデフォルトボックスを候補から外し、残ったデフォルトボックスを検出結果とするというものである。

【0012】

本開示による自動採譜装置におけるニューラルネットワークへの入力は、採譜対象の楽音の波形データ又はスペクトログラムであり、その出力は楽音のオンセット、オフセット及び音高であり、自動採譜装置は、SSDにおける中心座標及び幅に対応してオンセット及びオフセット（すなわち、楽音の形状又は長さ）を特定し、SSDにおける物体の種類に対応して音高を特定する。

20

【0013】

後述される実施例を概略すると、自動採譜装置は2つの学習済み機械学習モデル（畳み込みニューラルネットワークなど）を利用し、一方のモデルは単音音源から音高の予測確率を出力するものであり、他方のモデルは特徴マップから当該特徴マップの固定長の区間に音符が存在する予測確率を出力するものである。自動採譜装置は、採譜対象の音源を前者の学習済み機械学習モデル（特徴マップ生成モデル）に入力し、当該学習済み特徴マップ生成モデルの畳み込み層から生成された各特徴マップを後者の学習済み機械学習モデル（音符存在確率予測モデル）に入力し、各特徴マップの各点に対して当該学習済み音符存在確率予測モデルから出力された固定長の区間又はデフォルトボックスにおける各音高の音符の予測存在確率に基づき楽譜情報を生成する。

30

【0014】

学習済み特徴マップ生成モデルによって生成される特徴マップは、畳み込みの結果として異なる時間解像度を有し、固定長の区間又はデフォルトボックスは異なる時間的長さとなる。このため、音符存在確率予測モデルにより各特徴マップに対して固定長の区間と同じ長さの音符を検出することによって、異なる長さの音符のオンセット及びオフセットを特定することが可能になる。

40

【0015】

まず、図1を参照して、本開示の一実施例による自動採譜装置を説明する。図1は、本開示の一実施例による学習済み機械学習モデルを有する自動採譜装置を示す概略図である。

【0016】

図1に示されるように、本開示の一実施例による自動採譜装置200は、限定することなく、畳み込みニューラルネットワークなどの何れかのタイプのニューラルネットワークとして実現される2種類の学習済みモデルを有し、学習用データストレージ50を用いて学習装置100によって学習された機械学習モデルを利用して、採譜対象の音源から楽譜情報を生成する。

50

【0017】

次に、図2～7を参照して、本開示の一実施例による学習装置を説明する。学習装置100は、学習用データストレージ50における学習用データを利用して、自動採譜装置200に利用される特徴マップ生成モデルと音符存在確率予測モデルとを学習する。図2は、本開示の一実施例による学習装置の機能構成を示すブロック図である。

【0018】

図2に示されるように、学習装置100は、学習用データ取得部110、第1モデル学習部120及び第2モデル学習部130を有する。

【0019】

学習用データ取得部110は、単音音源と音高情報とを特徴マップ生成モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを音符存在確率予測モデルの学習用データとして取得し、単音音源と採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得する。

10

【0020】

具体的には、学習用データ取得部110は、学習用データストレージ50から、特徴マップ生成モデルを学習するための単音又はシングルノート音源（例えば、「ド」から「シ」までの12種類の音源など）の波形データと音高情報（「ド」から「シ」までの音高など）とのペアを取得し、取得した単音音源の波形データに対して前処理（例えば、短時間フーリエ変換など）を実行することによって、各単音音源のスペクトログラムと音高情報との学習用データセットを生成する。

20

【0021】

また、学習用データ取得部110は、学習用データストレージ50から、音符存在確率予測モデルを学習するための単旋律音源（歌唱音源など）の波形データと楽譜情報（音高の時系列変化など）とのペアを取得し、取得したモノフォニック音源の波形データに対して前処理（例えば、短時間フーリエ変換など）を実行することによって、モノフォニック音源のスペクトログラムと楽譜情報との学習用データセットを生成する。ここで、楽譜情報は、例えば、MIDI (Musical Instrument Digital Interface) 規格に従うものであってもよい。

【0022】

典型的には、スペクトログラムは、時間軸及び周波数軸における信号成分の強度を表し、波形データを短時間フーリエ変換することによって生成される。短時間フーリエ変換には各種パラメータが設定される必要があるが、例えば、FFT窓幅：1024、サンプリング周波数：16kHz、オーバーラップ幅：768、窓関数：ハニング窓、及びフィルタバンク：メルフィルタバンク（128バンド）などに従って、短時間フーリエ変換が実行されてもよい。スペクトログラムに変換した後、時間軸方向に一定のサンプル数（例えば、1024サンプル）だけ抽出されてもよい。また、本実施例によるスペクトログラムは、低周波数成分を精細にするよう周波数軸が対数変換されたものであってもよい。

30

【0023】

第1モデル学習部120は、単音音源のスペクトログラムを学習用入力データとして入力し、単音音源の音高の予測確率を出力するよう音高情報によって特徴マップ生成モデルを学習する。

40

【0024】

例えば、特徴マップ生成モデルは、図3に示されるように、複数の畳み込み層を含む畳み込みニューラルネットワークにより構成され、入力された単音音源のスペクトログラムを音高の予測確率に変換するSSDとして実現される。ここで、音高は連続値でなく離散値として表現され、one-hotベクトルとして表現されてもよい。なお、打楽器などの噪音音源も学習対象とする場合、噪音音源の単音又はシングルノートの音声をデータセットに含めてもよい。その場合、音高クラスとして噪音を表現するクラスを設定し、それを教師ラベルとしてもよい。

【0025】

50

第1モデル学習部120は、学習用入力データの単音音源のスペクトログラムを特徴マップ生成モデルに入力し、特徴マップ生成モデルからの出力と学習用出力データの音高情報との誤差が小さくなるように、バックプロパゲーションによって特徴マップ生成モデルのパラメータを更新する。ここで、誤差を示す損失関数として、限定することなく、特徴マップ生成モデルの出力と学習用出力データの音高との交差エントロピーが利用されてもよい。

【0026】

例えば、所定数の学習用データに対して更新処理が終了した、誤差が所定の閾値以下に収束した、誤差の改善が所定の閾値以下に収束したなどの所定の学習終了条件が充足されると、第1モデル学習部120は、更新された特徴マップ生成モデルを学習済み機械学習モデルとして設定する。

10

【0027】

第2モデル学習部130は、採譜対象の音源のスペクトログラムを学習済みの特徴マップ生成モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、特徴マップの固定長の区間に音符が存在する予測確率を出力するよう楽譜情報によって音符存在確率予測モデルを学習する。

【0028】

例えば、音符存在確率予測モデルは、図4に示されるように、複数の畳み込み層を含む畳み込みニューラルネットワークにより構成され、モノフォニック音源のスペクトログラムを学習済み特徴マップ生成モデルに入力することによって生成された特徴マップを当該特徴マップの各点を始点とする固定長の区間と同じ長さの音符が存在する予測確率に変換するSSDとして実現される。例えば、ドからシの12音で採譜する場合、特徴マップ上の各点は、ドからシの各音高及び休符（無音）の13通りの音符又は音高クラスが存在する予測確率を有する。

20

【0029】

上述したように、学習済み特徴マップ生成モデルは複数の畳み込み層を含み、各畳み込み層からモノフォニック音源のスペクトログラムの特徴マップが生成される。生成される特徴マップは、図3に示されるような畳み込み層のレベルに応じて時間解像度が異なる特徴マップとなる。典型的には、図5に示されるように、入力層に相対的に近い畳み込み層では、時間解像度が相対的に高い（図示された例では、32Hz）特徴マップが生成され、出力層に相対的に近い畳み込み層では、時間解像度が相対的に低い（図示された例では、16Hz）特徴マップが生成される。図示されるような固定長の区間又はデフォルトボックスが設定されると、時間解像度が相対的に高い特徴マップにおける区間は、時間解像度が相対的に低い特徴マップにおける区間より短い時間を占有する。このため、異なる時間的長さを有する音符の存在予測確率を導出することができ、音符の時間的長さを特定することが可能になる。

30

【0030】

第2モデル学習部130は、学習用入力データの音源のスペクトログラムを学習済み特徴マップ生成モデルに入力し、学習済み特徴マップ生成モデルによって生成された各特徴マップを音符存在確率予測モデルに入力し、音符存在確率予測モデルからの出力と学習用出力データの楽譜情報との誤差が小さくなるように、バックプロパゲーションによって音符存在確率予測モデルのパラメータを更新する。

40

【0031】

ここで、誤差を示す損失関数として、限定することなく、音符存在確率予測モデルの出力と音高の時系列変化とから算出されるタイミング誤差と信頼誤差との加重和が利用されてもよい。音高の時系列変化は、楽曲のスタートタイミング、エンドタイミング及び音高のセットが複数集まることによって表現され、楽譜情報から導出される。当該セットは発音と呼ばれてもよく、例えば、音高の時系列変化は、発音#1："0：00～0：02，A（ラ）3"、発音#2："0：03～0：05，B（シ）3"、発音#3："0：05～0：08，C（ド）4"・・・などにより表現されてもよい。図5に示されるデフォルトボ

50

ックスは、1つの発音を表現しており、複数のチャンネルを有する。デフォルトボックスの各チャンネルの最初のサンプルはそれぞれ、当該デフォルトボックスの発音のオンセットの予測値、オフセットの予測値及び音高クラスの予測確率を有する。すなわち、トータルで $2 + (\text{音高のクラス数})$ のチャンネルがある。

【0032】

第2モデル学習部130は、各発音について、オンセットとオフセットとの和が最小となるデフォルトボックスを探索し、検出されたデフォルトボックスと発音とに対してタイミング誤差と信頼誤差を求める。ここで、タイミング誤差とは、予測したオンセットを考慮したスタートタイミングのずれと、予測したオフセットを考慮したエンドタイミングのずれとの和としてもよい。ただし、差分の表現として、デフォルトボックスの長さを基準にした相対値が利用されてもよい。また、信頼誤差は、発音の音高と予測した音高とから算出される交差エントロピーであってもよい。なお、無音を表すクラスも教師ラベルとして用意されてもよく、この場合、発音のない区間を予測することができる。

10

【0033】

第2モデル学習部130は、NMS (Non - Maximum Suppression) に従って各特徴マップの各点について設定されたデフォルトボックスを減らしていき、残ったデフォルトボックスを予測発音としてもよい。具体的には、第2モデル学習部130はまず、各デフォルトボックスについて音高クラス毎の音符存在予測確率を求める。その後、第2モデル学習部130は、予測確率が所定の閾値(例えば、0.9など)以下であるデフォルトボックスを削除してもよい。第2モデル学習部は、残ったデフォルトボックスのうち積集合/和集合に閾値を設けて、閾値以上のデフォルトボックスの一方を削除し、重複したデフォルトボックスを排除する。第2モデル学習部130は、最終的に残ったデフォルトボックスを予測発音とする。

20

【0034】

例えば、所定数の学習用データに対して更新処理が終了した、誤差が所定の閾値以下に収束した、誤差の改善が所定の閾値以下に収束したなどの所定の学習終了条件が充足されると、第2モデル学習部130は、更新された音符存在確率予測モデルを学習済みモデルとして設定する。

【0035】

一実施例では、第1モデル学習部120は、複数種別のオーディオ成分のそれぞれに対して特徴マップ生成モデルを学習し、第2モデル学習部130は、複数種別のオーディオ成分を含む採譜対象の音源に対して各オーディオ成分種別毎に音符が存在する予測確率を出力するよう音符存在確率予測モデルを学習してもよい。

30

【0036】

例えば、特徴マップ生成モデルと音符存在確率予測モデルとは、モノフォニックボーカルと伴奏とを含む楽曲に対して適用されてもよい。この場合、ボーカル用特徴マップ生成モデルと伴奏用特徴マップ生成モデルとが、ボーカルの単音音源と音高情報とのペアから構成されるボーカル用学習データと、伴奏の単音音源と音高情報とのペアから構成される伴奏用学習データとを利用して、上述した学習処理と同様に学習される。一方、ボーカル用音符存在確率予測モデルと伴奏用音符存在確率予測モデルとが、学習用の音源と楽譜情報とを利用して、音源を学習済みボーカル用特徴マップ生成モデルと学習済み伴奏用特徴マップ生成モデルとに入力することによって生成された特徴マップを入力とし、上述した学習処理と同様に学習される。

40

【0037】

あるいは、特徴マップ生成モデルと音符存在確率予測モデルとは、楽器毎などの複数のパートを含む楽曲に対して適用されてもよい。上述したボーカルと伴奏とを含む楽曲に対する学習処理と同様であるが、この場合、音符存在確率予測モデルの出力は、特徴マップの固定長の区間に特定パートの特定音符が存在する予測確率であってもよい。例えば、"男声のA3の音高"、"女声のA3の音高"などの特定パートの特定音符の存在の予測確率を出力するようにしてもよい。

50

【0038】

あるいは、本開示は拍子を有する楽曲に対して適用されてもよい。この場合、音符存在確率予測モデルの出力は、拍子のオンセット及びオフセットに関するものであってもよく、例えば、デフォルトボックスが一拍である予測確率が出力されてもよい。

【0039】

図6は、本開示の一実施例による特徴マップ生成モデルの学習処理を示すフローチャートである。当該学習処理は、上述した学習装置100又は学習装置100のプロセッサによって実現される。

【0040】

図6に示されるように、ステップS101において、学習用データ取得部110は、学習用データストレージ50から単音音源と音高情報とのペアを取得する。例えば、音高は、「ド」から「シ」の12音と無音との13通りであり、当該13通りの音高に対応する単音音源が取得されてもよい。

10

【0041】

ステップS102において、学習用データ取得部110は、取得した単音音源を前処理する。具体的には、学習用データ取得部110は、単音音源の波形データに対して前処理（例えば、短時間フーリエ変換など）を実行し、単音音源のスペクトログラムを取得する。

【0042】

ステップS103において、第1モデル学習部120は、前処理された単音音源と音高情報とのペアによって特徴マップ生成モデルを学習する。例えば、特徴マップ生成モデルは、畳み込みニューラルネットワークにより構成され、入力音源を音高の予測確率に変換する。具体的には、第1モデル学習部120は、単音音源のスペクトログラムを特徴マップ生成モデルに入力し、特徴マップ生成モデルの出力と音高情報との誤差が小さくなるように、バックプロパゲーションによって特徴マップ生成モデルのパラメータを更新する。

20

【0043】

ステップS104において、第1モデル学習部120は、学習終了条件が充足されたか判断する。所定の学習終了条件は、例えば、所定数の学習用データに対して更新処理が終了した、誤差が所定の閾値以下に収束した、誤差の改善が所定の閾値以下に収束したなどであってもよい。所定の学習終了条件が充足されている場合（S104：YES）、第1モデル学習部120は、更新された特徴マップ生成モデルを学習済みモデルとして設定してもよい。他方、所定の学習終了条件が充足されていない場合（S104：NO）、当該処理はステップS101に移行し、上述した各ステップを繰り返す。

30

【0044】

図7は、本開示の一実施例による音符存在確率予測モデルの学習処理を示すフローチャートである。当該学習処理は、上述した学習装置100又は学習装置100のプロセッサによって実現される。

【0045】

図7に示されるように、ステップS201において、学習用データ取得部110は、学習用データストレージ50からモノフォニック音源と楽譜情報とのペアを取得する。例えば、モノフォニック音源は歌唱音源の波形データであってもよく、楽譜情報は当該モノフォニック音源の楽譜を示す。

40

【0046】

ステップS202において、学習用データ取得部110は、取得したモノフォニック音源を前処理する。具体的には、学習用データ取得部110は、モノフォニック音源の波形データに対して前処理（例えば、短時間フーリエ変換など）を実行し、モノフォニック音源のスペクトログラムを取得する。

【0047】

ステップS203において、第2モデル学習部130は、前処理されたモノフォニック音源を学習済み特徴マップ生成モデルに入力し、学習済み特徴マップ生成モデルによって

50

生成された特徴マップを取得する。具体的には、第2モデル学習部130は、学習済み特徴マップ生成モデルの各畳み込み層から生成された特徴マップを取得する。生成された特徴マップは、各畳み込み層の畳み込みの程度に応じて異なる時間解像度の特徴マップとなる。

【0048】

ステップS204において、第2モデル学習部130は、取得した特徴マップと楽譜情報とのペアによって音符存在確率予測モデルを学習する。例えば、音符存在確率予測モデルは、畳み込みニューラルネットワークにより構成され、入力された特徴マップを当該特徴マップの固定長の区間に音符が存在する音符存在予測確率に変換する。具体的には、第2モデル学習部130は、各特徴マップを音符存在確率予測モデルに入力し、音符存在確率予測モデルの出力と楽譜情報との誤差が小さくなるように、バックプロパゲーションによって音符存在確率予測モデルのパラメータを更新する。

10

【0049】

ステップS205において、第2モデル学習部130は、学習終了条件が充足されたか判断する。所定の学習終了条件は、例えば、所定数の学習用データに対して更新処理が終了した、誤差が所定の閾値以下に収束した、誤差の改善が所定の閾値以下に収束したなどであってもよい。所定の学習終了条件が充足されている場合(S205: YES)、第2モデル学習部130は、更新された音符存在確率予測モデルを学習済みモデルとして設定してもよい。他方、所定の学習終了条件が充足されていない場合(S205: NO)、当該処理はステップS201に移行し、上述した各ステップを繰り返す。

20

【0050】

次に、図8及び9を参照して、本開示の一実施例による自動採譜装置を説明する。図8は、本開示の一実施例による自動採譜装置の機能構成を示すブロック図である。

【0051】

図8に示されるように、自動採譜装置200は、モデル処理部210及び楽譜生成部220を有する。

【0052】

モデル処理部210は、単音音源から音高の予測確率を出力する学習済み特徴マップ生成モデルと、特徴マップから当該特徴マップの固定長の区間に音符が存在する予測確率を出力する学習済み音符存在確率予測モデルとを利用し、採譜対象の音源を学習済み特徴マップ生成モデルに入力し、当該学習済み特徴マップ生成モデルによって生成された特徴マップを学習済み音符存在確率予測モデルに入力し、特徴マップの固定長の区間に音符が存在する予測確率を出力する。

30

【0053】

具体的には、モデル処理部210は、採譜対象の音源に対して短時間フーリエ変換などの前処理を実行して当該音源のスペクトログラムを取得し、取得したスペクトログラムを学習装置100による学習済み特徴マップ生成モデルに入力して当該学習済み特徴マップ生成モデルの各畳み込み層からの特徴マップを取得する。そして、モデル処理部210は、取得した各特徴マップを学習装置100による学習済み音符存在確率予測モデルに入力し、入力した特徴マップの各点を始点とする固定長の区間又はデフォルトボックスと同じ長さの音符が存在する予測確率を取得し、取得した各特徴マップの音符存在予測確率を楽譜生成部220にわたす。例えば、音符存在予測確率は、特徴マップのデフォルトボックスに存在する各音高(例えば、「ド」、「レ」、・・・「シ」、無音など)の確率の予測値であり、高い予測確率を有する音高が当該時間的長さに対応する音符に相当すると判断できる。

40

【0054】

楽譜生成部220は、音符が存在する予測確率に基づき楽譜情報を生成する。具体的には、楽譜生成部220は、SSDに用いられるNMS(Non-Maximum Suppression)に従って学習済み音符存在確率予測モデルの出力を後処理する。典型的には、学習済み音符存在確率予測モデルから多数の予測音符候補が出力される。これら

50

の予測音符候補から予測音符を特定する必要があり、SSDではNMSを利用して予測音符候補をしばしば絞っている。

【0055】

例えば、楽譜生成部220はまず、学習済み音符存在確率予測モデルに入力された特徴マップ上の各点に対して出力された音符存在予測確率のうち最大となる音符を当該時間における予測音符とする。そして、楽譜生成部220は、特徴マップ上の各点について予測音符を決定し、各点、予測音符及び対応する音符存在予測確率のデータセットをリスト化し、音符存在予測確率に関して降順にリスト内のデータセットをソートする。そして、楽譜生成部220は、所定の抽出条件を適用し、リストから予測音符候補を絞る。例えば、楽譜生成部220は、音符存在予測確率が所定の閾値（例えば、0.9など）以下であるデータセットをリストから削除してもよい。また、楽譜生成部220は、重複して検出された音符の重複を排除するため、予測音符が同じであって、かつ、予測音符の重複度が所定の閾値（例えば、80%など）以上のデータセットがリストの上位にある場合、当該上位のリストのみを残すようにしてもよい。楽譜生成部220は、最終的なリストにおけるデータセットに基づき楽譜を生成する。

10

【0056】

図9は、本開示の一実施例による自動採譜処理を示すフローチャートである。当該自動採譜処理は、上述した自動採譜装置200又は自動採譜装置200のプロセッサによって実現される。

【0057】

図9に示されるように、ステップS301において、モデル処理部210は、採譜対象の音源を取得する。例えば、当該音源はモノフォニック音源であってもよいし、複数種類のオーディオ成分を含んでもよい。

20

【0058】

ステップS302において、モデル処理部210は、取得した音源を前処理する。具体的には、モデル処理部210は、取得した音源に対して短時間フーリエ変換などの前処理を実行し、当該音源のスペクトログラムを取得する。

【0059】

ステップS303において、モデル処理部210は、前処理した音源を学習済み特徴マップ生成モデルに入力して特徴マップを取得し、取得した特徴マップを学習済み音符存在確率予測モデルに入力して入力した特徴マップの各点を始点とする固定長の区間又はデフォルトボックスと同じ長さの音符が存在する予測確率を取得する。

30

【0060】

ステップS304において、楽譜生成部220は、特徴マップ上の各点に対して取得した音符存在予測確率に基づき予測音符を決定する。具体的には、楽譜生成部220は、各点について取得した音符存在予測確率のうち最大となる音符存在予測確率に対応する音符を当該点に対する予測音符として決定する。

【0061】

ステップS305において、楽譜生成部220は、決定された特徴マップの各点の予測音符に対して後処理を実行する。具体的には、楽譜生成部220は、SSDにおけるNMSに従って特徴マップの各点の予測音符を絞る。例えば、楽譜生成部220は、特徴マップ上の各点について決定された予測音符に基づき、各点、予測音符及び対応する音符存在予測確率のデータセットをリスト化し、音符存在予測確率に関して降順にリスト内のデータセットをソートし、音符存在予測確率が所定の閾値（例えば、0.9など）以下であるデータセットをリストから削除すると共に、予測音符が同じであって、かつ、予測音符の重複度が所定の閾値（例えば、80%など）以上のデータセットがリストの上位にある場合、当該上位のリストのみを残すようにしてもよい。

40

【0062】

ステップS306において、楽譜生成部220は、最終的なリストにおけるデータセットに基づき楽譜を生成する。

50

【0063】

上述した学習装置100及び自動採譜装置200はそれぞれ、例えば、図10に示されるように、CPU(Central Processing Unit)101、GPU(Graphics Processing Unit)102、RAM(Random Access Memory)103、通信インタフェース(IF)104、ハードディスク105、入力装置106及び出力装置107によるハードウェア構成を有してもよい。CPU101及びGPU102は、プロセッサ又は処理回路として参照されてもよく、学習装置100及び自動採譜装置200の各種処理を実行し、特に、CPU101は学習装置100及び自動採譜装置200における各種処理の実行を制御し、GPU102は機械学習モデルを学習及び実行するための各種処理を実行する。RAM103及びハードディスク105は、学習装置100及び自動採譜装置200における各種データ及びプログラムを格納するメモリとして機能し、特に、RAM103は、CPU101及びGPU102における作業データを格納するワーキングメモリとして機能し、ハードディスク105は、CPU101及びGPU102の制御プログラム及び/又は学習用データを格納する。通信IF104は、学習用データストレージ50から学習用データを取得するための通信インタフェースである。入力装置106は、情報及びデータを入力するための各種デバイス(例えば、ディスプレイ、スピーカ、キーボード、タッチ画面など)であり、出力装置107は、処理の内容、経過、結果等の各種情報を表示する各種デバイス(例えば、ディスプレイ、プリンタ、スピーカなど)である。しかしながら、本開示による学習装置100及び自動採譜装置200は、上述したハードウェア構成に限定されず、他の何れか適切なハードウェア構成を有してもよい。

10

20

【0064】

本開示の一態様では、

単音音源と音高情報とを第1の機械学習モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを第2の機械学習モデルの学習用データとして取得し、前記単音音源と前記採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得する学習用データ取得部と、

前記単音音源のスペクトログラムを学習用入力データとして入力し、前記単音音源の音高の予測確率を出力するよう前記音高情報によって第1の機械学習モデルを学習する第1モデル学習部と、

30

前記採譜対象の音源のスペクトログラムを学習済みの前記第1の機械学習モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するよう前記楽譜情報によって第2の機械学習モデルを学習する第2モデル学習部と、
を有する学習装置が提供される。

【0065】

一実施例では、

前記第1の機械学習モデルと前記第2の機械学習モデルとは、畳み込みニューラルネットワークにより構成されてもよい。

【0066】

一実施例では、

前記第2モデル学習部は、前記第1の機械学習モデルにより生成される異なる時間解像度を有する複数の特徴マップを前記第2の機械学習モデルに入力してもよい。

40

【0067】

一実施例では、

前記第2モデル学習部は、前記第1の機械学習モデルと前記第2の機械学習モデルとをSSD(Single Shot Detection)として実現してもよい。

【0068】

一実施例では、

前記第1モデル学習部は、複数種別のオーディオ成分のそれぞれに対して前記第1の機

50

械学習モデルを学習し、

前記第2モデル学習部は、複数種別のオーディオ成分を含む採譜対象の音源に対して各オーディオ成分種別毎に音符が存在する予測確率を出力するよう前記第2の機械学習モデルを学習してもよい。

【0069】

本開示の一態様では、

単音音源から音高の予測確率を出力する第1の学習済み機械学習モデルと、特徴マップから前記特徴マップの固定長の区間に音符が存在する予測確率を出力する第2の学習済み機械学習モデルとを利用し、採譜対象の音源を前記第1の学習済み機械学習モデルに入力し、前記第1の学習済み機械学習モデルによって生成された特徴マップを前記第2の学習済み機械学習モデルに入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するモデル処理部と、

前記音符が存在する予測確率に基づき楽譜情報を生成する楽譜生成部と、を有する自動採譜装置が提供される。

【0070】

一実施例では、

前記モデル処理部は、前記採譜対象の音源に対して前処理を実行することによってスペクトログラムを取得し、前記スペクトログラムを前記第1の学習済み機械学習モデルに入力してもよい。

【0071】

一実施例では、

前記モデル処理部は、前記特徴マップ上の各点について前記第2の学習済み機械学習モデルから出力された最大の予測確率を有する音符を予測音符として決定してもよい。

【0072】

一実施例では、

前記楽譜生成部は、NMS (Non - Maximum Suppression) に従って抽出された予測音符に基づき楽譜情報を生成してもよい。

【0073】

本開示の一態様では、

プロセッサが、単音音源と音高情報とを第1の機械学習モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを第2の機械学習モデルの学習用データとして取得し、前記単音音源と前記採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得すステップと、

前記プロセッサが、前記単音音源のスペクトログラムを学習用入力データとして入力し、前記単音音源の音高の予測確率を出力するよう前記音高情報によって第1の機械学習モデルを学習するステップと、

前記プロセッサが、前記採譜対象の音源のスペクトログラムを学習済みの前記第1の機械学習モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するよう前記楽譜情報によって第2の機械学習モデルを学習するステップと、を有する学習方法が提供される。

【0074】

本開示の一態様では、

プロセッサが、単音音源から音高の予測確率を出力する第1の学習済み機械学習モデルに採譜対象の音源を入力するステップと、

前記プロセッサが、特徴マップから前記特徴マップの固定長の区間に音符が存在する予測確率を出力する第2の学習済み機械学習モデルに前記第1の学習済み機械学習モデルによって生成された特徴マップを入力するステップと、

前記プロセッサが、前記第2の学習済み機械学習モデルから出力された前記音符が存在する予測確率に基づき楽譜情報を生成するステップと、

10

20

30

40

50

を有する自動採譜方法が提供される。

【0075】

本開示の一態様では、

単音音源と音高情報とを第1の機械学習モデルの学習用データとして取得し、採譜対象の音源と楽譜情報とを第2の機械学習モデルの学習用データとして取得し、前記単音音源と前記採譜対象の音源とに対して前処理を実行し、それぞれのスペクトログラムを取得するステップと、

前記単音音源のスペクトログラムを学習用入力データとして入力し、前記単音音源の音高の予測確率を出力するよう前記音高情報によって第1の機械学習モデルを学習するステップと、

前記採譜対象の音源のスペクトログラムを学習済みの前記第1の機械学習モデルに入力することによって生成される特徴マップを学習用入力データとして入力し、前記特徴マップの固定長の区間に音符が存在する予測確率を出力するよう前記楽譜情報によって第2の機械学習モデルを学習するステップと、

をプロセッサに実行させるプログラムが提供される。

10

【0076】

本開示の一態様では、

単音音源から音高の予測確率を出力する第1の学習済み機械学習モデルに採譜対象の音源を入力するステップと、

特徴マップから前記特徴マップの固定長の区間に音符が存在する予測確率を出力する第2の学習済み機械学習モデルに前記第1の学習済み機械学習モデルによって生成された特徴マップを入力するステップと、

前記第2の学習済み機械学習モデルから出力された前記音符が存在する予測確率に基づき楽譜情報を生成するステップと、

をプロセッサに実行させるプログラムが提供される。

20

【0077】

本開示の一態様では、

上述したプログラムを記憶するコンピュータ可読記憶媒体が提供される。

【0078】

以上、本開示の実施例について詳述したが、本開示は上述した特定の実施形態に限定されるものではなく、特許請求の範囲に記載された本開示の要旨の範囲内において、種々の変形・変更が可能である。

30

【符号の説明】

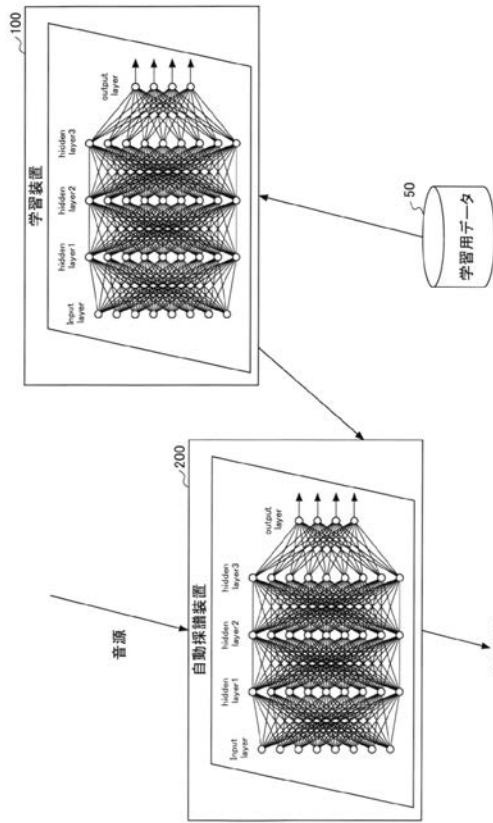
【0079】

50 学習用データストレージ

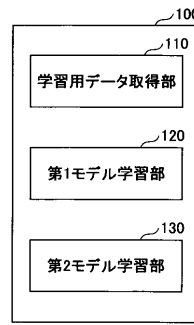
100 学習装置

200 自動採譜装置

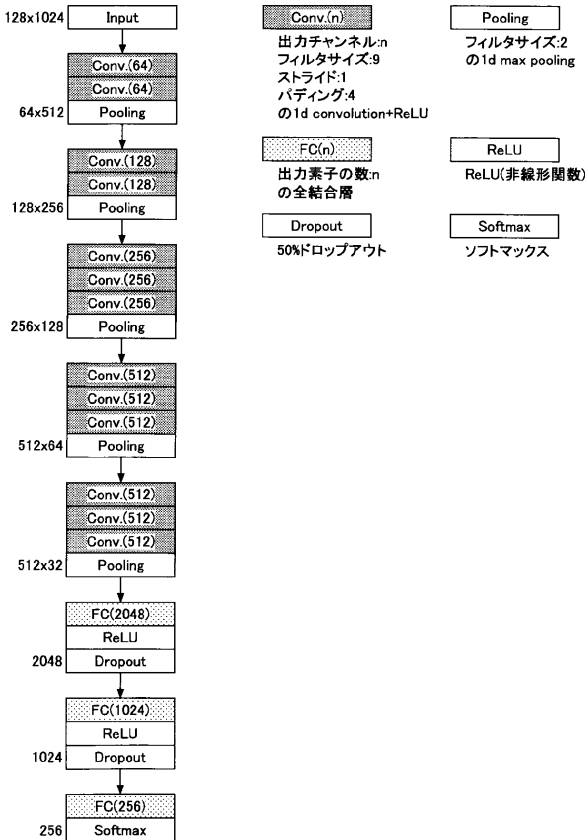
【 図 1 】



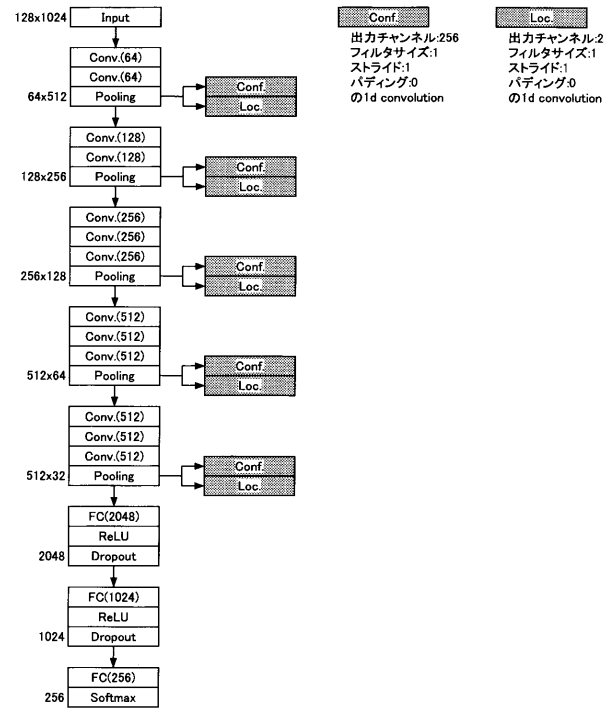
【 図 2 】



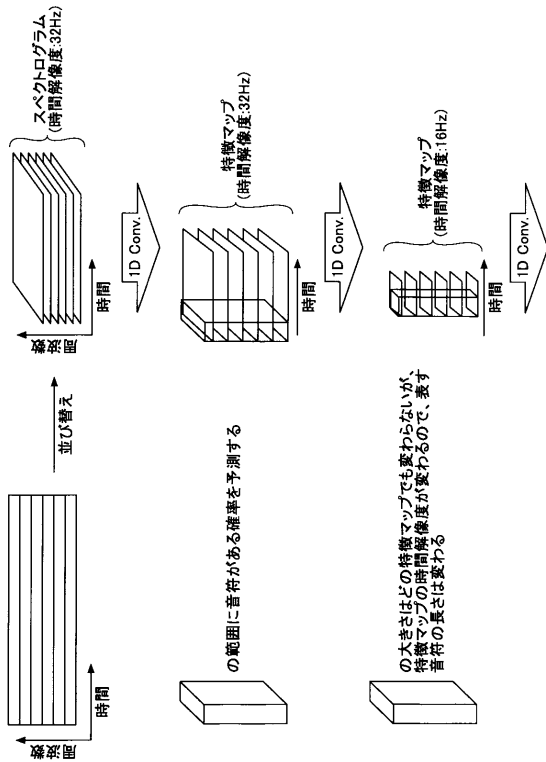
【 図 3 】



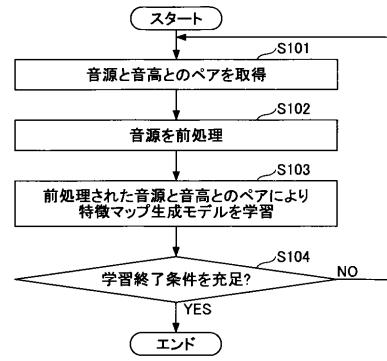
【 図 4 】



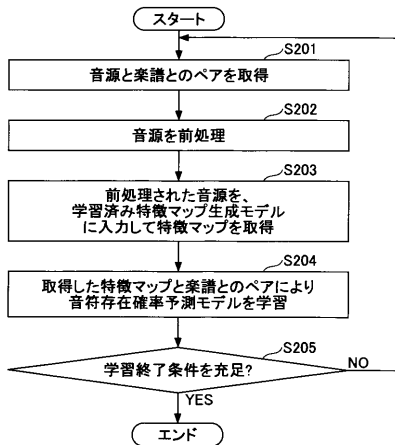
【 図 5 】



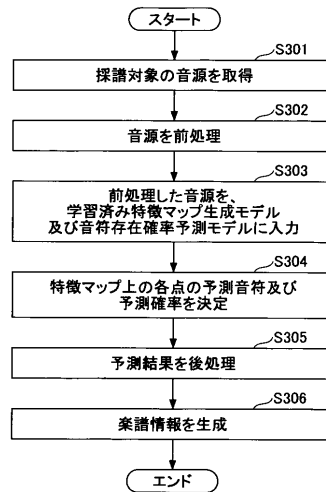
【 図 6 】



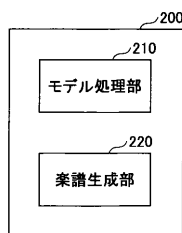
【 図 7 】



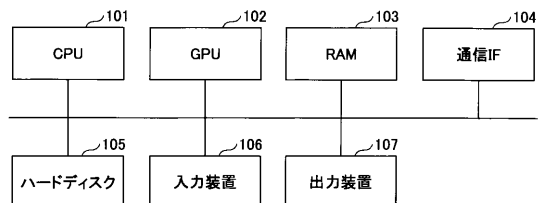
【 図 9 】



【 図 8 】



【 図 10 】



フロントページの続き

(51)Int.Cl.

F I

テーマコード(参考)

G 0 6 N 99/00 1 5 0