

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4351385号
(P4351385)

(45) 発行日 平成21年10月28日(2009.10.28)

(24) 登録日 平成21年7月31日(2009.7.31)

(51) Int.Cl.		F I			
G 1 0 L	15/10	(2006.01)	G 1 0 L	15/10	3 0 0 G
G 1 0 L	15/06	(2006.01)	G 1 0 L	15/06	3 0 0 C
G 1 0 L	15/14	(2006.01)	G 1 0 L	15/14	2 0 0 A

請求項の数 21 (全 21 頁)

(21) 出願番号	特願2000-513270 (P2000-513270)	(73) 特許権者	500046438
(86) (22) 出願日	平成10年9月16日 (1998.9.16)		マイクロソフト コーポレーション
(65) 公表番号	特表2001-517816 (P2001-517816A)		アメリカ合衆国 ワシントン州 9805
(43) 公表日	平成13年10月9日 (2001.10.9)		2-6399 レッドモンド ワン マイ
(86) 国際出願番号	PCT/US1998/019346		クロソフト ウェイ
(87) 国際公開番号	W01999/016052	(74) 代理人	100089705
(87) 国際公開日	平成11年4月1日 (1999.4.1)		弁理士 社本 一夫
審査請求日	平成17年9月13日 (2005.9.13)	(74) 代理人	100071124
(31) 優先権主張番号	08/934,622		弁理士 今井 庄亮
(32) 優先日	平成9年9月19日 (1997.9.19)	(74) 代理人	100076691
(33) 優先権主張国	米国 (US)		弁理士 増井 忠武
		(74) 代理人	100075236
			弁理士 栗田 忠彦
		(74) 代理人	100075270
			弁理士 小林 泰

最終頁に続く

(54) 【発明の名称】 連続および分離音声を認識するための音声認識システム

(57) 【特許請求の範囲】

【請求項 1】

音声認識システムを実現する方法であって、

複数の離散的に発話された訓練ワードを示す分離音声訓練データを受け取るステップと

、
複数の連続的に発話された訓練ワードを示す連続音声訓練データを受け取るステップと

、
前記連続音声訓練データおよび前記分離音声訓練データにおける音声単位を表す出力確率分布を含む複数の音声単位モデルを与えるステップであって、前記モデルが、前記分離音声訓練データおよび前記連続音声訓練データの両方に基づいて訓練される、ステップと

10

、
前記分離音声訓練データおよび前記連続音声訓練データに基づいて、様々な長さのワード・フレーズに含まれるワードの近似ワード持続時間を示す、複数のワード持続時間モデルを与えるステップであって、前記長さは、前記ワード・フレーズにおけるワード・カウントにより定まる、ステップと、

前記訓練した音声単位モデル及び前記複数のワード持続時間モデルに基づいて音声認識する認識装置を設けるステップと、を含む方法。

【請求項 2】

請求項 1 記載の方法において、前記分離音声訓練データを受け取るステップが、第 1 の複数の音響信号を受け取るステップを含み、連続音声訓練データを受け取るステップが、

20

第2の複数の音響信号を受け取るステップを含み、複数の音声単位モデルを与えるステップが、

前記第1および第2の複数の音響信号に基づいて、複数の音響モデルを生成するステップ、
を含む、方法。

【請求項3】

請求項2記載の方法において、複数の音響モデルを生成するステップが、

前記第1および第2の複数の音響信号に基づいて、前記連続音声訓練データおよび分離音声訓練データにおける音素を表わす複数の出力確率分布を生成するステップ、
を含む、方法。

10

【請求項4】

請求項1記載の方法において、分離音声訓練データを受け取るステップが、

複数の離散的に発話された訓練ワードに関連する無声前後関係情報を含む分離音声データを受け取るステップ、
を含む、方法。

【請求項5】

請求項4記載の方法において、離散音声データを受け取るステップが、

ユーザが、前記複数の訓練ワードの各々の間にポーズを入れて前記複数の訓練データを離散的に発話したことを示す、前記離散音声データを受け取るステップ、
を含む、方法。

20

【請求項6】

請求項1記載の方法において、連続音声訓練データを受け取るステップが、

ユーザが複数の訓練ワードを流暢に発話したことを示す連続音声データを受け取るステップ、
を含む、方法。

【請求項7】

請求項1記載の方法であって、更に、

前記音声単位モデルを訓練する前に、認識する予想音声に基づいて、前記連続音声訓練データおよび前記分離音声訓練データに重み付けするステップ、
を含む方法。

30

【請求項8】

請求項1記載の方法であって、更に、

ユーザが複数の訓練ワードを異なる様式で発話したことを示す追加音声訓練データを受け取るステップ、
を含む方法。

【請求項9】

請求項8記載の方法において、追加音声訓練データを受け取るステップが、

前記ユーザが前記複数の訓練ワードを第1振幅および第2振幅で発話したことを示す、前記追加音声訓練データを受け取るステップ、
を含み、前記第2振幅が前記第1振幅よりも大きい、方法。

40

【請求項10】

請求項8記載の方法において、追加音声訓練データを受け取るステップが、

前記ユーザが前記複数の訓練ワードを流暢に、第1ペースおよび第2ペースで発話したことを示す、前記追加音声訓練データを受け取るステップ、
を含み、前記第2ペースが前記第1ペースよりも速い、方法。

【請求項11】

請求項3記載の方法において、複数の音声単位モデルを与えるステップが、更に、

前記出力分布の各々を、前記訓練ワードの1つのものの少なくとも一部を形成する音素における所定数の状態の1つと関連付けるステップ、
を含む方法。

50

【請求項 1 2】

請求項 1 1 記載の方法であって、更に、

音素毎に、選択した音素を含む前記訓練ワードの全てからの前記選択した音素に関連する出力分布をグループ化して、出力分布グループを形成するステップと、

各音素における状態毎に、前記出力分布グループにおいて選択した状態に関連する出力分布を、前記選択した音素に関連する言語学的前後関係情報に基づいてセノンに分離することによって、前記選択した音素において選択した状態について、セノン・ツリーを作成するステップと、

を含む、方法。

【請求項 1 3】

請求項 1 2 記載の方法において、音声を認識する認識装置を設けるステップが、

認識すべき目標ワードにおける各連続目標音素の各連続状態に対する出力分布を受け取るステップと、

目標音素毎に、該目標音素を表わす可能性が最も高い、ある数の音素候補を特定するステップと、

前記音素候補の状態に関連するセノンを、前記目標音素の対応する状態と関連する前記出力分布と比較するステップと、

前記目標音素の前記出力分布と最も密接に一致するセノンを有する最尤音素を特定するステップと、

を実行するように、前記音声認識装置を構成するステップを含む、方法。

【請求項 1 4】

請求項 1 3 記載の方法において、比較するステップが、

各音素候補における各状態に関連する前記セノン・ツリーを、前記目標音素の言語学的前後関係情報に基づいて通り抜けて、前記目標音素における各状態毎にセノンを特定するステップと、

前記目標音素における状態に関連する前記出力分布を、前記音素候補において特定したセノンに関連する前記出力分布と比較するステップと、

を含む、方法。

【請求項 1 5】

請求項 1 3 記載の方法において、ある数の音素候補を特定するステップが、

前記分離音声訓練データおよび前記連続音声訓練データに基づいて、前記訓練ワードにおける音素を示す複数の単音モデルを形成するステップと、

前記目標音素に関連する前記出力分布を、前記単音モデルと比較するステップと、

前記目標音素に関連する前記出力分布に密接に一致する単音モデルを有する、ある数の音素候補を特定するステップと、

を含む、方法。

【請求項 1 6】

請求項 1 記載の方法において、音声を認識する音声認識装置を設けるステップが、

認識すべき複数の目標ワードを受け取るステップと、

前記目標ワードにおける目標ワード・フレーズを示すフレーズ境界を検出するステップと、

前記目標ワード・フレーズの近似持続時間を判定するステップと、

前記目標ワード・フレーズによって表わされるワード・フレーズ候補を示す、複数のワード・フレーズ仮説を得るステップと、

前記ワード・フレーズ仮説におけるワードの近似ワード・カウントおよび持続時間を判定するステップと、

前記ワード・フレーズ仮説における前記ワードのワード持続時間を、前記ワード・フレーズ仮説におけるワード数に等しいワード・カウントを有するワード持続時間モデルと比較して、前記ワード・フレーズ仮説における前記ワード持続時間がどれ位緊密に前記ワード持続時間モデルと一致するかに基づいて、最尤ワード・フレーズ仮説を得るステップと

10

20

30

40

50

、
 を実行するように前記認識装置を構成するステップを含む、方法。

【請求項 17】

請求項 1 記載の方法において、複数のワード持続時間モデルを与えるステップが、
 前記分離音声訓練データおよび前記連続音声訓練データにおいて訓練ワード・フレーズ
 を検出するステップと、

前記訓練ワード・フレーズにおけるワード数を判定するステップと、
 複数の前記検出した訓練ワード・フレーズにおける前記ワードの近似ワード持続時間を
 判定するステップと、

前記訓練ワード・フレーズにおけるワード数、および前記訓練ワード・フレーズにおけ
 る前記ワードの持続時間についてパラメータ化した、複数のワード持続時間分布を判定す
 るステップと、
 を含む、方法。

10

【請求項 18】

音声認識方法であって、
 認識すべきワードを示す訓練データを受け取るステップと、
前記訓練データにおいてポーズを検出して、複数の訓練ワード・フレーズを識別するス
 テップと、

前記訓練ワード・フレーズの各々におけるワード数を判定するステップと、
前記訓練ワード・フレーズの各々におけるワード数に基づいて、前記訓練ワード・フレ
 ーズに対応する複数のワード持続時間の分布を生成するステップと、

20

認識すべき音声を示す入力データを受け取るステップと、
 前記入力データに基づいて、前記音声においてポーズを検出して、フレーズの持続時間
 を特定するステップと、

前記検出したポーズの間にある前記入力データによって表わされるワード・フレーズ候
 補を表わす複数のフレーズ仮説を生成するステップと、

各フレーズ仮説における各ワードに関連するワード持続時間を、前記フレーズ仮説にお
 けるワード数に基づいておよび前記フレーズ持続時間に基づいて、前記複数のワード持続
 時間の分布から選択された、前記フレーズ仮説におけるワード数に等しいワード数を有す
 るフレーズについての予想ワード持続時間と比較するステップと、

30

前記ワード持続時間の前記予想ワード持続時間との比較に基づいて、各フレーズ仮説に
 スコアを割り当てて、前記入力データにより表わされる最尤フレーズ仮説を得るステッ
 プと、

を含む方法。

【請求項 19】

請求項 18 記載の方法において、各フレーズ仮説におけるワード持続時間を予想ワード
 持続時間と比較するステップが、

仮説毎に当該仮説におけるワード数および前記フレーズの持続時間に基づいて、前記仮
 説におけるワードのワード持続時間を判定するステップと、

前記仮説におけるワード数に等しい、フレーズ毎のワード数に関連する前記複数の分布
 から、選択された 1 つを選ぶステップと、

40

前記仮説に対して判定した前記ワード持続時間を、前記選択した分布と比較するステッ
 プと、

を含む、方法。

【請求項 20】

請求項 19 記載の方法において、各フレーズ仮説にスコアを割り当てるステップが、

前記仮説に対して判定したワード持続時間がどの位緊密に前記選択した分布と一致する
 かを示すスコアを、各ワード仮説に割り当てるステップ、

を含む、方法。

【請求項 21】

50

音声認識を行う方法であって、

複数の離散的に発話された訓練ワードを示す分離音声訓練データを受け取るステップであって、前記分離音声訓練データが第1の複数の出力分布を含み、各出力分布が、前記離散的に発話された訓練ワードの1つのものの少なくとも一部を形成する音素における所定数の状態の1つに関連する、ステップと、

複数の連続的に発話された訓練ワードを示す連続音声訓練データを受け取るステップであって、前記連続音声訓練データが第2の複数の出力分布を含み、該第2複数の出力分布の各々が、前記連続的に発話された訓練ワードの1つのものの少なくとも一部を形成する音素における所定数の状態の1つと関連する、ステップと、

選択した音素を含む前記訓練ワードの全てから、前記選択した音素に関連する出力分布をグループ化して、出力分布グループを形成するステップと、

前記選択した音素における選択した状態について、セノン・ツリーを作成するステップであって、前記選択した音素に関連する言語学的前後関係情報に基づいて、前記出力分布グループにおける前記選択した状態に関連する前記出力分布を分離することによって作成する、ステップと、

前記分離音声訓練データおよび前記連続音声訓練データに基づいて、様々な長さのワード・フレーズに含まれるワードの近似ワード持続時間を示す、複数のワード持続時間モデルを与えるステップであって、前記長さは、前記ワード・フレーズにおけるワード・カウントにより定まる、ステップと、

前記出力分布グループ、前記セノン・ツリー及び前記複数のワード持続時間モデルに基づいて音声認識する認識装置を設けるステップと、

を含む方法。

【発明の詳細な説明】

【0001】

(発明の背景)

本発明は、コンピュータ音声認識に関する。更に特定すれば、本発明は、連続音声および分離音声双方を認識する方法に関するものである。

【0002】

現在最も成功している音声認識システムは、隠れマルコフ・モデル(HMM: hidden Markov model)として知られる確率モデルを採用するものである。隠れマルコフ・モデルは、複数の状態を含み、同一状態への遷移を含む、各遷移から他のあらゆる遷移への遷移毎に、遷移確率を定義する。各一意の状態には、確率的に観察(observation)が関連付けられる。状態間の遷移確率(観察が1つの状態から次の状態に遷移する確率)は、全てが同一ではない。したがって、状態および観察確率間の遷移確率が与えられた際に、ビタビ・アルゴリズムのような探索技法を用いて、確率全体が最大となる最尤状態シーケンス(most likely state sequence)を判定する。

【0003】

現行の音声認識システムでは、音声は、隠れマルコフ・プロセスによって発生されるものと見られている。その結果、音声スペクトルの観察シーケンスをモデル化するためにHMMが採用され、特定のスペクトルにHMMにおける1つの状態を確率的に関連付けてきた。言い換えると、所与の音声スペクトルの観察シーケンスについて、対応するHMMには最尤状態シーケンスがある。

【0004】

この対応するHMMは、したがって、観察シーケンスに関連付けられる。この技法は、HMMにおける別個の各状態シーケンスを音素のようなサブワード単位に関連付ければ、最尤サブワード単位シーケンスを求めることができるように、拡張することができる。更に、サブワード単位をどのように組み合わせるワードを形成するかというモデルを用い、次いでワードをどのように組み合わせる文章を形成するかという言語モデルを用いることによって、完全な音声認識を達成することができる。

【0005】

10

20

30

40

50

実際に音響信号を処理する場合、信号は、通常、フレームと呼ばれる連続時間間隔でサンプリングする。フレームは、通常、複数のサンプルを含み、重複したり、あるいは連続する場合もある。各フレームには、音声信号の一意の部分が関連付けられている。各フレームによって表わされる音声信号の部分を分析し、対応する音響ベクトルを得る。音声認識の間、音声単位モデルの探索を行い、音響ベクトル・シーケンスに関連する可能性が最も高い状態シーケンスを判定する。

【 0 0 0 6 】

音響ベクトル・シーケンスに対応する可能性が最も高い状態シーケンスを見出すために、ビタビ・アルゴリズムを用いることができる。ビタビ・アルゴリズムは、最初のフレームから開始し、一度に1フレームずつ時間に同期して進める計算を実行する。考慮する対象の状態シーケンスにおける（即ち、HMMにおける）状態毎に、確率スコアを計算する。したがって、ビタビ・アルゴリズムが音響信号をフレーム毎に分析するに連れて、可能な状態シーケンスの各々について、蓄積確率スコア（cumulative probability score）を連続的に計算する。発声の終了時までには、ビタビ・アルゴリズムが計算した最も高い確率スコアを有する状態シーケンス（またはHMMあるいは一連のHMM）が、発声全体に対する最尤状態シーケンスを与える。次に、この最尤状態シーケンスを、対応する発話サブワード単位（spoken subword unit）、ワード、またはワード・シーケンスに変換する。

【 0 0 0 7 】

ビタビ・アルゴリズムは、指数計算を、モデルにおける状態および遷移の数、ならびに発声の長さに比例する計算に減少させる。しかしながら、大きな語彙では、状態および遷移の数が大きくなり、全ての可能な状態シーケンスに対し各フレーム内の各状態における確率スコアを更新するために必要な計算は、通常10ミリ秒の持続時間である、1フレームの持続時間よりも何倍も長くなる。

【 0 0 0 8 】

したがって、最尤状態シーケンスを判定するために必要な計算を大幅に削減するために、ブルーニング（pruning）またはビーム探索（beam searching）と呼ばれる技法が開発された。この種の技法は、非常に可能性が低い状態シーケンスに対する確率スコアの計算を不要にする。これは、通常、各フレームにおいて、考慮対象の各残留状態シーケンス（または潜在的シーケンス）に対する確率スコアを、当該フレームに関連する最高スコアと比較することによって行われる。特定の潜在的シーケンスに対する状態の確率スコアが十分に低い場合（当該時点において他の前栽駅シーケンスに対して計算した最大空く率と比較して）、ブルーニング・アルゴリズムは、このようにスコアが低い状態シーケンスは、完全な最尤状態シーケンスの一部である可能性は低いと見なす。通常、この比較を行うには、最小スレシホールド値を用いる。最小スレシホールド値未満のスコアを有する潜在的状態シーケンスは、探索プロセスから除外する。スレシホールド値は、いずれの所望のレベルにも設定することができ、主に所望のメモリおよび計算削減、ならびにメモリおよび計算削減によって得られる所望の誤り率上昇に基づいて設定する。保持する状態シーケンスを能動的仮説（active-hypothesis）と呼ぶ。

【 0 0 0 9 】

音声認識に求められる計算量（magnitude）を更に削減するための別の従来からの技法に、プレフィクス・ツリー（prefix tree）の使用を含むものがある。プレフィクス・ツリーは、音声認識システムの辞書（lexicon）を、ツリー構造として表わし、システムが遭遇する可能性のあるワード全てを、このツリー構造で表わす。

【 0 0 1 0 】

このようなプレフィクス・ツリーでは、（音素のような）各サブワード単位は、通常、特定の（HMMのような）音響モデルに関連付けられたブランチによって表わされる。音素ブランチを、ノードにおいて、後続の音素ブランチに接続する。同じ最初の音素を共有する辞書における全てのワードは、同じ最初のブランチを共有する。同じ最初の音素および二番目の音素を有する全てのワードは、同じ最初のブランチおよび2番目のブランチを共有する。対象的に、共通の第1音素を有するが、異なる第2音素を有するワードは、プレ

10

20

30

40

50

フィクス・ツリーにおいて同じ第1ブランチを共有するが、プレフィクス・ツリーの最初のノードにおいて分岐 (diverge) する等となる。ツリー構造はこのように続き、システムが遭遇する可能性のあるワード全てを、ツリーの終端ノード (即ち、ツリー上のリーフ (leaf)) によって表わすようにしている。

【0011】

プレフィクス・ツリー構造を採用することによって、初期ブランチ数は、システムの辞書または語彙における典型的なワード数よりは遥かに少なくなることは明白である。実際、初期ブランチ数は、探索する語彙または辞書のサイズには無関係に、音素の総数 (約40ないし50) を超過する可能性はない。しかしながら、異音変動 (allophonic variation) を用いた場合、用いる異音によっては、ブランチの初期数は大きくなる可能性はある。

10

【0012】

前述の技法を採用する音声認識システムは、通常、2つの種類に分類することができる。第1の種類は、流暢音声 (fluent speech) を認識可能な連続音声認識 (CSR) システムである。CSRシステムは、連続音声データに基づいて訓練され (即ち、音響モデルを生成する)、一人以上の読み手が訓練データを連続的に即ち流暢にシステムに読み込んでいく。訓練中に生成した音響モデルを用いて音声認識する。

【0013】

第2の種類システムは、分離音声認識 (ISR) システムであり、通常、分離した音声 (即ち、離散音声) のみを認識するために採用する。ISRシステムは、離散即ち分離音声データに基づいて訓練され (即ち、音響モデルを生成し)、この場合一人以上の読み手には、各ワードの間にポーズを入れて、離散的即ち分離して訓練データをシステムに読み込むように要求する。また、ISRシステムは、通常、連続音声認識システムよりも精度が高くかつ効率的である。何故なら、ワードの境界が一層明確であり、したがって探索空間が一層厳しい制約を受けるからである。また、分離音声認識システムは、連続音声認識の特殊な場合と考えられてきた。何故なら、連続音声認識システムは、一般に、分離音声も同様に受け入れることができるからである。これらは、単に、分離音声を認識しようとするときに、同様に動作しないだけである。

20

【0014】

CSRシステムのユーザは、通常、システムが誤りを犯し始めるまで、または、ユーザが文書の組み立てを思案するまで、流暢に発話しがちであることが観察されている。その時点で、ユーザは、ワード間にポーズを入れると言ってもよい程に、速度を落とすことが多い。双方の場合において、ユーザは、ワード間にポーズを入れて、よりゆっくりと区別して発話することにより、ユーザは認識システムを助けていると信じているが、実際には、ユーザは、システムの能力を超えてシステムにストレスを与えているのである。

30

【0015】

しかしながら、単に分離音声認識システムを用いて連続音声を認識しようとするのは、適当ではない。ISRシステムは、通常、連続音声を認識しようとする場合には、CSRシステムよりも遥かに性能が劣る。その理由は、ISR訓練データには、交差ワード同時調音 (crossword coarticulation) がないからである。

(発明の概要)

40

音声認識は、複数の離散的に発話した訓練ワードを示す分離音声訓練データを受け取り、複数の連続的に発話した訓練ワードを示す連続音声訓練データを受け取ることによって行われる。分離音声訓練データおよび連続音声訓練データに基づいて、複数の音声単位モデルを訓練する。訓練した音声単位モデルに基づいて、音声を認識する。

【0016】

好適な実施形態の1つでは、認識対象音声におけるポーズを識別し、フレーズの持続時間を判定する。ポーズの入力データによって表わされる、フレーズ候補を示す複数のフレーズ仮説 (phrase hypothesis) を生成する。各フレーズ仮説における各ワードに関連するワード持続時間を、フレーズ仮説内のワード数に等しいワード数を有するフレーズに対する予想ワード持続時間と比較する。ワード持続時間の予測ワード持続時間との比較に

50

基づいて、各フレーズ仮説にスコアを割り当てる。

(好適な実施形態の詳細な説明)

図 1 および関連する論述は、本発明を実現可能な、適切な計算機環境の端的で概略的な説明を行うことを意図するものである。必須ではないが、本発明は、少なくとも部分的に、パーソナル・コンピュータによって実行するプログラム・モデルのような、コンピュータ実行可能命令に全体的に関連して説明する。一般的に、プログラム・モジュールとは、ルーチン・プログラム、オブジェクト、コンポーネント、データ構造等を含み、特定のタスクを実行したり、あるいは特定の抽象データ型を実装するものである。更に、本発明は、ハンド・ヘルド機器、マルチプロセッサ・システム、マイクロプロセッサを用いた、即ち、プログラム可能な民生用電子機器、ネットワーク PC、ミニコンピュータ、メインフレーム・コンピュータ等を含む、その他のコンピュータ・システム構成を用いても実施可能であることを当業者は認めよう。また、本発明は、分散型計算機環境においても実施可能であり、その場合、通信ネットワークを通じてリンクしてあるリモート処理用機器によってタスクを実行する。分散型計算機環境では、プログラム・モジュールは、ローカル記憶装置およびリモート・メモリ記憶装置双方に配置することもできる。

10

【 0 0 1 7 】

図 1 を参照すると、本発明を実現するためのシステム例は、従来のパーソナル・コンピュータ 2 0 の形態の汎用計算機を含み、演算装置 2 1、システム・メモリ 2 2、およびシステム・メモリを含む種々のシステム・コンポーネントを演算装置 2 1 に結合するシステム・バス 2 3 を含む。システム・バス 2 3 は、メモリ・バスまたはメモリ・コントローラ、周辺バス、および種々のバス・アーキテクチャのいずれかを用いたローカル・バスを含む数種類のバス構造のいずれかとすればよい。システム・メモリは、リード・オンリ・メモリ (ROM) 2 4 およびランダム・アクセス・メモリ (RAM) 2 5 を含む。起動中におけるように、パーソナル・コンピュータ 2 0 内部の要素間で情報を転送する際に役立つ基本ルーチンを含む基本入出力システム 2 6 (BIOS) を ROM 2 4 に格納してある。更に、パーソナル・コンピュータ 2 0 は、図示しないハード・ディスクに対して読み出しおよび書き込みを行うハード・ディスク・ドライブ 2 7、ならびに CD ROM またはその他の光媒体のようなリムーバブル光ディスク 3 1 に対して読み出しおよび書き込みを行う光ディスク・ドライブ 3 0 を含む。ハード・ディスク・ドライブ 2 7、磁気ディスク・ドライブ 2 8、および光ディスク・ドライブ 3 0 は、それぞれ、ハード・ディスク・ドライブ・インターフェース 3 2、磁気ディスク・ドライブ・インターフェース 3 3、および光ドライブ・インターフェース 3 4 を介してシステム・バス 2 3 に接続してある。これらのドライブおよびそれらと関連するコンピュータ読取可能媒体は、コンピュータ読取可能命令、データ構造、プログラム・モジュール、およびパーソナル・コンピュータ 2 0 のためのその他のデータの揮発性格納を行う。

20

30

【 0 0 1 8 】

ここに記載する環境の一例では、ハード・ディスク、リムーバブル磁気ディスク 2 9 およびリムーバブル光ディスク 3 1 を採用するが、磁気カセット、フラッシュ・メモリ・カード、デジタル・ビデオ・ディスク、ベルヌーイ・カートリッジ、ランダム・アクセス・メモリ (RAM)、リード・オンリ・メモリ (ROM) 等のような、コンピュータによるアクセスが可能なデータを格納することができる、その他の種類のコンピュータ読取可能媒体も、動作環境の一例では使用可能であることは、当業者には認められよう。

40

【 0 0 1 9 】

オペレーティング・システム 3 5、1 つ以上のアプリケーション・プログラム 3 6、その他のプログラム・モジュール 3 7、およびプログラム・データ 3 8 を含む、多数のプログラム・モジュールは、ハード・ディスク、磁気ディスク 2 9、光ディスク 3 1、ROM 2 4 または RAM 2 5 に格納することができる。ユーザは、キーボード 4 0 およびポインティング・デバイス 4 2 のような入力デバイスを介して、コマンドおよび情報をパーソナル・コンピュータ 2 0 に入力することができる。その他の入力デバイス (図示せず) には、マイクロフォン、ジョイスティック、ゲーム・パッド、衛星ディッシュ、スキャナ等が含

50

まれる場合もある。これらおよびその他の入力デバイスは、多くの場合システム・バスに結合してあるシリアル・ポート・インターフェース46を介して演算装置21に接続するが、パラレル・ポート、ゲーム・ポートまたはユニバーサル・シリアル・バス(USB)のような他のインターフェースを介して接続することも可能である。モニター47またはその他の種類の表示装置も、ビデオ・アダプタ48のようなインターフェースを介して、システム・バス23に接続してある。モニター47に加えて、パーソナル・コンピュータは、通常、スピーカおよびプリンタのような、その他の周辺出力装置(図示せず)も含む。

【0020】

パーソナル・コンピュータ20は、リモート・コンピュータ49のような、1つ以上のリモート・コンピュータへの論理接続を用いて、ネットワーク化環境においても動作することができる。リモート・コンピュータ49は、別のパーソナル・コンピュータ、サーバ、ルータ、ネットワークPC、ピア・デバイス、またはその他の共通ネットワーク・ノードとすることができ、通常、パーソナル・コンピュータ20に関して先に述べた要素の多くまたは全てを含むが、図1にはメモリ記憶装置50のみを示してある。図1に示す論理接続は、ローカル・エリア・ネットワーク(LAN)51およびワイド・エリア・ネットワーク(WAN)52を含む。このようなネットワーク化環境は、オフィスの企業規模のコンピュータ・ネットワーク、イントラネットおよびインターネットでは一般的である。

【0021】

LANネットワーク化環境において用いる場合、パーソナル・コンピュータ20はネットワーク・インターフェース即ちアダプタ53を経由してローカル・エリア・ネットワーク51に接続する。WANネットワーク化環境において用いる場合、パーソナル・コンピュータ20は通常モデム54、またはインターネットのような広域ネットワーク52を通じて通信を確立するその他の手段を含む。モデム54は、内蔵型でも外付けでもよく、シリアル・ポート・インターフェース46を介してシステム・バス23に接続してある。ネットワーク化環境では、パーソナル・コンピュータ20に関して図示したプログラム・モジュール、またはその部分をリモート・メモリ記憶装置に格納することも可能である。図示のネットワーク接続は一例であり、コンピュータ間に通信リンクを確立するその他の手段も使用可能であることは認められよう。

【0022】

更に、図1における環境を音声認識システムとして実現する場合、他のコンポーネントも望ましいこともある。このようなコンポーネントは、マイクロフォン、サウンド・カードおよびスピーカを含み、その一部については以下で更に詳しく説明する。

【0023】

図2は、本発明の一形態による、音声認識システム60のブロック図を示す。音声認識システム60は、マイクロフォン62、アナログ/デジタル(A/D)変換器64、訓練モジュール65、特徴抽出モジュール66、無声(silence)検出モジュール68、セノン・ツリー(senone tree)記憶モジュール70、単音モデル記憶モジュール72、三音マッピング記憶モジュール74、プレフィクス・ツリー記憶モジュール76、ワード持続時間モデル記憶モデル78、サーチ・エンジン80、および出力装置82を含む。システム60全体またはシステム60の一部は、図1に示す環境において実現可能であることを注記しておく。例えば、マイクロフォン62は、好ましくは、適切なインターフェースを介して、更にA/D変換器64を介して、パーソナル・コンピュータ20への入力デバイスとして設けるとよい。訓練モジュール65、特徴抽出モジュール66および無声検出モジュール68は、コンピュータ20内のハードウェア・モジュール(CPUとは別個のプロセッサまたはCPU21内に実装したプロセッサ)、または図1に開示した情報記憶装置のいずれかに格納し、CPU21またはその他の適切なプロセッサによるアクセスが可能なソフトウェア・モジュールのいずれとしてもよい。加えて、セノン・ツリー(senone tree)記憶モジュール70、単音モデル記憶モジュール72、三音マッピング記憶モジュール74、プレフィクス・ツリー記憶モジュール76、およびワード持続時間モデル記憶モジュール78も、図1に示すいずれかの適切なメモリ素子に格納することが好ましい

10

20

30

40

50

。更に、サーチ・エンジン 80 は、CPU 21 (1つ以上のプロセッサを含むことができる) 内に実装することが好ましく、またはパーソナル・コンピュータ 20 が採用する専用音声認識プロセッサによって実行することも可能である。加えて、出力装置 82 は、好適な実施形態の 1 つでは、モニタ 47 として、またはプリンタとして、あるいはその他のいずれかの適切な出力装置として実現することができる。

【0024】

多くの場合、システム 60 は、最初に訓練データを用いて訓練を受ける。図 3 および図 4 は、本発明の好適な実施形態の 1 つにしたがって利用する、訓練データ収集およびシステム訓練手順を示すフロー図である。システム 60 を訓練するためには、最初に、図 3 に関して説明するように訓練データを収集する。好適な実施形態では、訓練データは、訓練ワードを話者が連続的に即ち流暢にシステム 60 に読み込む連続 (即ち、流暢) 訓練データ、および話者がワード間にポーズを入れながら離散的即ち分離的に訓練ワードをシステム 60 に読み込む分離 (即ち、離散) 訓練データを含む。

10

【0025】

このため、第 1 の話者を選択する。これをブロック 84 で示す。次いで、話者に、システム 60 のマイクロフォン 62 に向かって流暢に訓練文章を読むように要求する。これをブロック 86 で示す。ブロック 88 で示すように、訓練文章を記録する。システム 60 が受け取った各訓練ワードの音素的転写を、キーボード 40 のようなユーザ入力デバイスによって、訓練器 65 およびシステム 60 に入力する。これをブロック 90 で示す。次に、訓練文章を同様に読むように追加の話者に要求するか否かについて判定を行う。話者独立システムでは、多数の話者を用いることが好ましい。しかしながら、話者依存システムでは、多数の話者は任意であり、訓練文章は単一の話者によってのみ発話すればよい。

20

【0026】

いずれの場合でも、別の話者が流暢に訓練文章をシステム 60 に読み込む場合、新たな話者を選択し、ブロック 86、88 および 90 を通るプロセスを繰り返す。これをブロック 92 および 94 で示す。

【0027】

一旦連続訓練データをシステム 60 内に読み込んだならば、ブロック 96 で示すように、再度第 1 の話者を選択する。次いで、選択した話者は訓練ワードの一群をシステム 60 に読み込む。これらは、各ワード間にポーズを入れて、離散的即ち分離して読み込む。これをブロック 98 で示す。分離訓練データは、好適な実施形態の 1 つでは、連続訓練データにおいて見られるものと同ーワードを構成する。しかしながら、分離訓練データは、連続訓練データと同一である必要はなく、全て異なるワードの集合で形成することも可能である。いずれの場合でも、各ワードをシステムに読み込みながら、これをシステム 60 によって記録する。これをブロック 100 で示す。

30

【0028】

再び、システム 60 は、システム 60 に読み込まれた各訓練ワードの音素的転記を、キーボード 40 のようなユーザ入力デバイスから受け取る。これをブロック 102 で示す。

【0029】

次に、追加の話者が分離音声訓練データをシステム 60 に供給するか否かについて判定を行う。そうする場合、新たな話者を選択し、その話者が最初の話者と同様に、分離音声訓練データを入力する。追加の話者が分離訓練データをシステム 60 に入力しないと判定した場合、データ収集プロセスは完了したことになる。これを図 3 のブロック 104 および 106 で示す。

40

【0030】

また、訓練データは、話者によってマイクロフォンを通じて入力する代わりに、出力分布の形態で、フロッピ・ディスク・ドライブのような入力デバイスを通じて、システム 60 に直接ロード可能であることも注記しておく。

【0031】

訓練ワードをマイクロフォン 62 を介してシステム 60 に入力しながら、A/D変換器 6

50

4によってデジタル・サンプルに変換し、次いで特徴抽出モジュール66によって特徴ベクトルに変換する(または、ベクトル量子化および訓練データから得られるコードブックを用いて、コードワードに量子化する)。特徴ベクトル(またはコードワード)は、訓練モジュール65に供給する。また、訓練モジュール65は、ユーザ入力デバイスから音素的転写も受け取る。次に、訓練モジュール65は、訓練データにおける特徴ベクトル(またはコードワード)および音素的転写を用いて、1組の単音モデル、セノン・ツリー、三音マッピング・メモリ、プレフィクス・ツリー、およびワード持続時間モデルを、訓練データに基づいて構築する。これらの品目は、全て、認識を実行する際にサーチ・エンジン80が使用する。

【0032】

図4は、訓練モジュール65が、単音モデル、セノン・ツリーおよび三音マッピング・メモリを計算する際のプロセス全体を示すフロー図である。最初に、訓練モジュール65は、共同(pooled)訓練データを受け取る。共同とは、連続および分離音声訓練データ両方を意味するものとする。これを、図4のブロック108で示す。訓練データは、特徴抽出モジュール66によって、前述のように出力分布に変換される。したがって、訓練モジュール65は、特徴ベクトル(またはコードワード)およびそれに与えられた音素的転写を用いて、共同訓練データにおける各ワードについて、1つ以上の隠れマルコフ・モデルを算出する。隠れマルコフ・モデルは、共同訓練データにおいて見出される音素に関連付けられ、各音素について算出する出力および発生頻度に基づいて算出する。

【0033】

本発明の好適な実施形態の1つでは、訓練モジュール65は、訓練データ・セットにおいて見出される各音素を、単音モデルとしてモデル化する。単音モデルは、モデル内の各状態に対する出力確率分布を含む。これを図4のブロック110および112で示す。単音モデルは、認識方式において用い、音素のセノン評価が始まる前に、入力音声の発声に対する最尤一致音素を判定する。次いで、ブロック113で示すように、単音モデルをメモリ72に格納する。

【0034】

次に、各音素における各状態について、訓練モジュール65はセノン・ツリーを作成する。セノン・ツリーを作成する技法については、図5に関して更に詳細に説明する。セノン・ツリーの作成は、図4ではブロック114で表わす。次に、ブロック116で示すように、セノン・ツリーをメモリ70に格納する。

【0035】

一旦セノン・ツリーを作成したなら、次に訓練器65は、所望の三音全て(訓練データ内で見られるものおよび見られないもの双方)を、メモリ70に格納してあるセノン・ツリーによって表わされるセノン・シーケンスにマッピングする。これを行うために、訓練器65は所望の三音(対応する右および左の関係を有する音素)を選択し、メモリ70に格納してあるセノン・ツリーを通り抜ける(traverse)。セノン・ツリーを通り抜けた結果として、訓練モジュール65は、モデル化した三音における各状態に対応するセノンを獲得し、したがって各三音を表わすセノンのシーケンスを獲得する。このセノンのシーケンスを、三音マッピング・メモリ74における対応する三音にマッピングする。これをブロック118によって示す。三音マッピング・シーケンスについても、図6に関して、更に詳しく説明する。

【0036】

次に、訓練モジュール65は、プレフィクス・ツリーを組み立て、このプレフィクス・ツリーをメモリ76に格納する。これをブロック120で示す。最後に、訓練モジュール65は、ワード持続時間モデルを算出し、このワード持続時間モデルをメモリ78に格納する。これを図4のブロック122で示す。ワード持続時間モデルの算出については、図7および図8に関して更に詳しく説明する。

【0037】

単音モデル、セノン・ツリー、三音マッピング、プレフィクス・ツリーおよびワード持続

10

20

30

40

50

時間モデルを算出した後、音声認識を実行するようにシステム60を構成する。音声認識タスクについては、図9および図10において更に詳しく説明する。

【0038】

図5は、訓練モジュール65が、共同訓練データに含まれる各音素における各状態についてセノン・ツリーを作成する際のプロセスを、更に詳細に示すフロー図である。英語には約50の音素があることが、一般的に認められている。好適な実施形態では、各音素に3つの状態を関連付ける。したがって、訓練モジュール65は150のセノン・ツリーを作成しなければならない。また、好適な実施形態では、50個の音素の各々は、数個の異なる前後関係において、共同訓練データ（即ち、連続訓練データおよび分離訓練データ）内に現れる。したがって、三状態隠れマルコフ・モデルに基づいて音素をモデル化する場合、隠れマルコフ・モデルの各々における各状態に関連する出力分布は、訓練データに現れる際の音素の前後関係によっては異なる場合がある。この情報に基づいて、図5に関して説明するようにセノン・ツリーを構築する。

10

【0039】

最初に、訓練データに表わされている50個の音素から1つを選択する。これをブロック124で示す。次に、ブロック126で示すように、選択した音素の最初の状態を選択する。

【0040】

選択した音素において選択した状態に関連する出力分布は、共同訓練データ内の音素の全ての発生に対して、検索し、共に集合化する。これをブロック28で示す。次いで、選択した状態に対して集合化した出力分布を、言語学的な前後関係に関する質問に基づいて互いに分離する。この質問は、セノン・ツリーを生成しようとしている特定の音素の前後関係に関する言語学的情報を求める質問である。個々の出力分散の各々に対する質問の回答に基づいて、これらの出力分布を第1（親）グループから2つの（子）グループに分離する。

20

【0041】

適正な言語学的質問を選択する方法についてこれより説明する。端的に言えば、言語学的質問は、専門の言語学者が作り、前後関係的効果（contextual effect）の言語学的部類を捕獲するように設計することが望ましい。例えば、Hon（ホン）およびLee（リー）のCMU ROBUST VOCABULARY-INDEPENDENT SPEECH RECOGNITION SYSTEM（CMUロバスタな語彙独立音声認識システム）と題する論文（IEEE Int'l Conf. On Acoustics, Speech and Signal Processing, Toronto, Canada, 1991, pps889-892）に見ることができる。親グループを子グループに分割するために、訓練モジュール65は、多数の言語学的質問の内どれが、親グループにとって最良の質問であるかについて判定を行う。好適な実施形態では、最良の質問は、親グループおよび子グループ間で最大のエントロピ減少をもたらす質問であると決定する。言語学的質問は全て「はい」または「いいえ」で答える質問であるので、親ノードの分割から2つの子ノードが得られる。

30

【0042】

グループの分割は、所定の分岐スレシホールドにしたがって停止する。このようなスレシホールドは、例えば、グループ内の出力分布数が所定値未満となった場合、またはグループ分割によって生じたエントロピ減少が他のスレシホールド未満となった場合を含むことができる。所定の分岐スレシホールドに達した場合、得られる最終グループは、全て、クラスタ化した出力分布即ちセノンを表わすリーフ・グループとなる。クラスタ化出力分布に基づいて、セノンを表わす単一の出力分布を選択する。これをブロック130および132で示す。また、セノン・ツリーにおける質問を組み合わせたりあるいは結合して、複合質問の形成も可能であることを注記しておく。更に、複合質問は、親グループから子グループへのエントロピ減少に基づいて、より良い複合質問に分離することも可能である。

40

【0043】

選択した音素の選択した状態に対してセノン・ツリーを作成した後、このセノン・ツリーをメモリ70に格納する。これをブロック134で示す。このプロセスは、語彙における

50

各音素の状態毎に繰り返し、各音素の状態毎にセノン・ツリーを作成する。これを図5のブロック136および138で示す。

【0044】

語彙における各音素の状態毎にセノン・ツリーを作成した後、システム60によって認識すべき各三音を、特定のセノン・シーケンスにマッピングしなければならない。言い換えると、認識すべき三音毎に、当該三音における各状態について、メモリ70内に格納してある適切なセノン・ツリーを注意深く考察することによって、適切なセノンを特定しなければならない。

【0045】

最初に、システム60は、認識すべき各三音の音素的転写を、キーボード40のような転写入力デバイスを介して、ユーザから受け取る。次いで、この三音素の中央の音素の各状態に対応するセノン・ツリーを通り抜ける。単にセノン・ツリーのノードに関連する言語学的質問に答えることによって、セノン・ツリーを通り抜ける。三音の各連続状態に適切なセノン・ツリーを特定した後、特定したセノンを組み合わせてセノン・シーケンスを形成し、メモリ74内の当該三音にマッピングする。

【0046】

図6は、どのようにしてセノン・ツリーを作成し、通り抜けるかを理解するのに役立つ一例を示す。図6は、ワード「welcome」の一部として、文字「c」の発話音に対する音素/k/に対するセノン・ツリーを示す。図6は、/k/音素の最初の状態に対するセノン・ツリーを示す。図6に示すセノン・ツリーにおける質問の多くは、前述の技法にしたがって形成した複合質問であることは認められよう。

【0047】

ワード「welcome」の文字「lco」によって形成される三音/L, K, UH/に対して適切なセノン・シーケンスを決定するために、/k/音素の各セノン・ツリーを通り抜けなければならない。図6に示すセノン・ツリーは、/K/音素の最初の状態に関連する。ルート・ノード140に関連する言語学的質問は、三音の左側の音が自鳴音かまたは鼻音かである。/L/は自鳴音であるので、ツリーの通り抜けは子ノード142に移動する。

【0048】

子ノード142は、ノード140において出された質問に対する肯定の回答に対応する。ノード142において出される質問は、左側の音素(/L/)は後音素(back phoneme) (即ち、左側の音素は、舌の位置を口の後方に向けて発話する音素である)であるか否かについて尋ねる。/L/は後音素であるので、通り抜けはノード144に進む。これは、ノード142において出された質問に対する肯定の回答に対応する。右側の音(三音の/UH/音素)がLまたはWでないとすると、/L/音素は、ノード142によって出される質問において指定される音素のいずれでもないので、ノード142における質問に対する回答は否定となる。これによって、セノン2として示すセノンに至る。これを、/L, K, UH/三音の最初の状態に対する適切なセノンとして特定する。同様のツリー通り抜けは、/K/音素の他の状態の各々についても進められる。システム60に入力した三音モデルの全てのマルコフ・モデル全てについて、リーフ(即ち、セノン)に到達するまで、対応するセノン・ツリーを通り抜ける。各三音について定義したセノン・シーケンスをメモリ70に格納する。

【0049】

認識装置が発音プレフィクス・ツリー・デコーダに基づく好適な実施形態では、次に、システム60が認識する語彙または辞書を表わすために、プレフィクス・ツリーを組み立てる。プレフィクス・ツリーの組み立ては、好ましくは、ルート・ノードからリーフに進み、入力データを示す可能性が最も高いワードに到達することができるようにする。好適な実施形態では、プレフィクス・ツリーは、複数の文脈依存無声音(silence phone)を含み、辞書内のワードの一部として無声が埋め込まれるようにモデル化する(メモリ72に格納した単音モデルと同様)。プレフィクス・ツリー60を通り抜けた後、システム60

10

20

30

40

50

は、認識対象のいずれかの所与のフレーズに対して認識した最尤ワードまたはワード・シーケンスを構成する、能動的仮説を維持することが好ましい。

【0050】

次に、システム60は、好適な実施形態の1つでは、複数のワード持続時間モデルを組み立てる。これは、プレフィクス・ツリー・デコーダから現れる能動的仮説間で選択を行うために用いることができる。ワード持続時間モデルをメモリ78に格納する。図7は、ワード持続時間モデルの組み立てを更に詳細に示すフロー図である。

【0051】

システム60に入力した訓練データは、好ましくは、異なる持続時間の分離ワード、およびポーズによって分離したワード・シーケンス（即ち、フレーズ）を含むことが好ましく、ワード・シーケンスは、シーケンス毎に種々の異なるワード・カウントを有する。訓練モジュール65は、ワード・カウント n を有する各離散フレーズにおけるワードの平均持続時間をモデル化する。したがって、訓練モジュール65は、最初に、共同訓練データにおける異なる長さのフレーズ（これは、1ワードの長さを有するフレーズを含む）について、ワード毎の平均持続時間を算出する。これを図7のブロック144で示す。次に、訓練モジュール65は、フレーズ当たりのワード数によってパラメータ化した、ワード持続時間の分布族（distribution family）を生成する。これをブロック146で示す。次に、訓練モジュール65は、分布族をワード持続時間モデル・メモリ78に格納する。これをブロック148で示す。

【0052】

図8は、訓練モジュール65が算出する分布族をより明確に示すグラフである。図8は、 x 軸上にワード持続時間を有し、 y 軸上に n -ワード・フレーズの発生回数を有するグラフ上にプロットした3つの分布150、152および154を示す。分布150、152および154は、概略的にガンマ分布の形態となっており、分布150は一ワード・フレーズの平均持続時間に関連し、分布152は二ワード・フレーズにおける各ワードの平均持続時間に関連し、分布154は n ワード・フレーズ（ n は2よりも大きい整数である）における各ワードの平均持続時間に関連する。このように、図8は、一ワード・フレーズにおける各ワードの平均持続時間は、二ワード・フレーズにおける各ワードの平均持続時間よりも多少長いことをグラフで示す。また、フレーズ内のワード数が2を超過する場合、このようなフレーズにおける各ワードの平均持続時間は、一ワード・フレーズまたは二ワード・フレーズのいずれかにおけるワードの平均持続時間よりも多少短くなる。

【0053】

認識の間、プレフィクス・ツリーを通り抜けた後に保持してある能動的仮説における平均ワード持続時間を、訓練モデル65によって計算したワード持続時間モデルと比較する。次に、当該特定の仮説におけるワード毎の平均持続時間が、適切なワード持続時間モデルと密接に一致するか（または密接に一致しないか）否かに基づいて、各仮説にスコアを割り当てる（または、減点を適用する）。これについては、本明細書の後ろの方で更に詳しく説明する。

【0054】

一旦訓練モジュール65が単音モデル、セノン・ツリー、三音マッピング、プレフィクス・ツリー、およびワード持続時間モデルを生成したなら、音声を認識するためにシステム60を適切に構成する。

【0055】

図9は、システム60を用いて音声を認識する好適な技法の1つを示すフロー図である。最初に、ユーザがマイクロフォン62に供給した可聴ボイス信号の形態で、音声をシステム60に入力する。マイクロフォン62は、可聴音声信号をアナログ電子信号に変換し、A/D変換器64に供給する。A/D変換器64は、アナログ信号をデジタル信号シーケンスに変換し、特徴抽出モジュール66に供給する。好適な実施形態では、特徴抽出モジュール66は、従来からのアレイ・プロセッサであり、デジタル信号に対してスペクトル分析を行い、周波数スペクトルの各周波数帯域毎に絶対値（magnitude value）を計

10

20

30

40

50

算する。好適な実施形態の1つでは、約16キロヘルツのサンプル・レートで、A/D変換器64によって特徴抽出モジュール66に供給する。A/D変換器64は、商業的に入手可能な周知のA/D変換器として実施する。

【0056】

特徴抽出モジュール66は、A/D変換器64から受け取ったデジタル信号を、複数のデジタル・サンプルを含むフレームに分割する。各フレームの持続時間は、約10ミリ秒である。次に、特徴抽出モジュール66によって、各フレームを、複数の周波数帯域についてスペクトル特性を反映する特徴ベクトルに符号化することが好ましい。特徴抽出モジュール66は、更に、ベクトル量子化技法および訓練データから得られるコードブック（個々には示さない）に基づいて、特徴ベクトルをコードワードに符号化することも可能である。分析した特定のフレームの特徴ベクトル（またはコードワード）を用いて、出力分布を隠れマルコフ・モデルと比較することができる。特徴抽出モジュール66は、約10ミリ秒毎に1つの割合で、特徴ベクトルを供給することが好ましい。

10

【0057】

特徴抽出モジュール66がA/D変換器64からのデジタル・サンプルを処理している際、無声（または境界）検出モジュール68もサンプルを処理している。無声検出モジュール68は、特徴抽出モジュール66を実現するために用いたプロセッサと同一または異なるプロセッサ上で実現することができる。無声検出モジュール68は、周知の方法で動作する。端的に言うと、無声検出モジュール68は、A/D変換器が供給するデジタル・サンプルを処理して無声（即ち、ポーズ）を検出し、ユーザが発声したワードまたはフレーズ間の境界を判定する。次に、無声検出モジュール68は、ワードまたはフレーズの境界検出を示す境界検出信号をサーチ・エンジン80に供給する。このように、サーチ・エンジン80は、認識すべき目標ワードに関連する出力分布の形態で、音声データを受け取る。これを図9のブロック156で示す。

20

【0058】

次に、サーチ・エンジン80は、受け取った出力分布を、単音メモリ72に格納してある単音モデルと比較する。発話した目標ワードの連続する目標音素毎に、そして目標音素の連続する目標状態毎に、サーチ・エンジン80は、目標状態に対する出力分布を、メモリ72に格納してある各音素の単音モデルの対応する状態と比較する。次に、サーチ・エンジン80は、目標状態の出力分布に最も密接に一致する状態を有する、所定数の音素単音モデルを選択し、目標音素が表わす音素候補（likely phoneme）を得る。これを図9にブロック158で示す。

30

【0059】

次に、サーチ・エンジン80は、音素候補の1つを選択し、当該音素における最初の状態を選択する。これをブロック160および162で示す。次に、サーチ・エンジン80は、選択した状態に対してセノン・ツリーによって生成したセノンを検索する。

【0060】

次に、サーチ・エンジン80は、最初の目標状態の目標出力分布を、選択した音素モデルの最初の状態に対応するセノン・ツリーの各セノンと比較する。次に、サーチ・エンジン80は、目標状態の出力分布と最も密接に一致するセノンであればどれであっても、最良の一致セノンとして選択し、この最良の一致セノンについて、一致確率スコアを計算し格納する。これをブロック164および166で示す。

40

【0061】

選択した音素が1つよりも多い状態を有する場合、サーチ・エンジン80は、選択した音素に残っている状態毎に同じステップを実行する。こうして、サーチ・エンジン80は、選択した音素における状態毎に、最も密接に一致するセノンを選択し、最良の一致セノンに対する一致確率スコアを計算し格納する。これをブロック168で示す。選択した音素における全ての状態を比較し終わった後、サーチ・エンジン80は、判定した確率スコアに基づいて、選択した音素に対して、セノン・シーケンス候補を特定したことになる。これをブロック170で示す。次に、サーチ・エンジン80は、メモリ74に格納してある

50

情報にアクセスし、判定したセノン・シーケンス候補にマッピングされている、三音候補を検索する。これをブロック172で示す。

【0062】

次に、サーチ・エンジン80は、音素候補を全て処理し終わったか否かについて判定を行う。し終わっていない場合、サーチ・エンジン80は前述の処理を繰り返し、音素候補毎に、比較の間に判定した確率スコアに基づいて、セノン・シーケンス候補に到達する（したがって、目標音素に関連するN個の三音候補に到達する）。これをブロック174および176で示す。

【0063】

一旦N個の三音候補を特定したなら、サーチ・エンジン80はメモリ76内のプレフィクス・ツリーにアクセスする。プレフィクス・ツリーを通り抜けた後、サーチ・エンジン80は能動的仮説を特定する。好適な実施形態の1つでは、サーチ・エンジン80は次に、North American Business News Corpus（北アメリカビジネス・ニュース・コーパス）から導出し、CSR-III Text Language Model（CSR-III テキスト言語モデル）（1994年University of Penn.）と題し、Linguistic Data Consortiumが発行した刊行物に詳細に明記されている、60,000ワード三重字言語モデル（trigram language model）のような、辞書および言語モデルに単純にアクセスする。この言語モデルを用いて、入力データが表わす最尤ワードまたはワード・シーケンスを特定し、サーチ・エンジン80によってこれを出力装置82に供給する。

【0064】

しかしながら、本発明の別の形態および別の好適な実施形態によれば、サーチ・エンジン80は、メモリ78内のワード持続時間モデルも利用して、入力データによって表わされる最尤ワードまたはワード・シーケンスを、更に精度高く特定する。図10は、マルチワード・フレーズおよび単一ワード・フレーズ間の判別を行うために、持続時間モデルをどのように用いるのかを示すフロー・チャートである。この説明の目的のため、ワード・カウントXの離散フレーズが、無声で開始しかつ終了するY個の流暢に発話されたワードのシーケンスであるとする。

【0065】

持続時間モデルの適用は、好ましくは、離散フレーズの境界において行う。入力データにおいてポーズを検出することによって、フレーズを検出する。最初に、入力データ内のポーズを、無声検出モジュール68によって検出する。これをブロック180で示す。次に、サーチ・エンジン80は、検出したポーズが、スレシホールド持続時間 $d(p)$ よりも短い持続時間 $d(P)$ を有するか否かについて判定を行う。スレシホールド持続時間 $d(p)$ は、偽りのポーズ、またはフレーズ間の境界を正確に反映しないポーズの検出を回避するように、訓練データに基づいて経験的に決定する。これをブロック182で示す。 $d(P)$ が $d(p)$ 未満である場合、処理はブロック80に戻り、別のポーズの検出を待つ。

【0066】

しかしながら、 $d(P)$ が $d(p)$ 未満でない場合、サーチ・エンジン80は、現在のポーズと、スレシホールド持続時間 $d(p)$ を超過した最後のポーズとの間の期間を示す、フレーズの持続時間（セグメント持続時間） $d(S)$ を計算する。これをブロック184で示す。次に、サーチ・エンジン80は、セグメント持続時間 $d(S)$ がスレシホールド・セグメント持続時間 $d(s)$ よりも長い場合について判定を行う。 $d(p)$ の場合と同様、 $d(s)$ の決定も、セグメント持続時間が、発見的方法を適用すべきでないような長さには決してならないように、訓練データに基づいて経験的に行う。言い換えると、ワード持続時間モデルは、持続時間が短いフレーズに適用する方が、持続時間が非常に長いフレーズに適用するよりも、高い効果が得られると考えられている。セグメント持続時間 $d(S)$ がセグメント・スレシホールド $d(s)$ よりも長い場合、処理はブロック180に戻り、別のポーズの検出を待つ。

【0067】

しかしながら、 $d(S)$ がスレシホールド・セグメント持続時間 $d(s)$ 未満である場合、

10

20

30

40

50

サーチ・エンジン 80 は、入力データによって表わされる n 個の最尤ワードまたはワード・フレーズの 1 つを示す、現フレーズ仮説 H を選択する。これをブロック 188 で示す。次に、サーチ・エンジン 80 は、 H のワード・カウント ($w_c(H)$) を判定し、 H の各ワードの平均持続時間を、 $w_c(H)$ および $d(S)$ に基づいて計算し、 $w_c(H)$ に等しいワード・カウントを有するフレーズに対応する、メモリ 78 内に格納してあるワード持続時間分布とこれを比較する。これをブロック 190 で示す。

【0068】

この比較に基づいて、サーチ・エンジン 80 は次に関数 $i_p(w_c(H), d(S))$ に応じて、この仮説 H にスコア (または減点) を割り当てる。関数 $i_p(w_c(H), d(S))$ は、 H の平均ワード持続時間が、対応するワード持続時間モデルとどの程度緊密に一致するかを示す。好適な実施形態では、 $i_p(w_c(H), d(S))$ は、システム 60 に入力した訓練データに基づいて経験的に求めた傾斜減少関数である。これをブロック 192 で示す。サーチ・エンジン 80 は、ブロック 194 で示すように、能動的仮説の各々についてこのプロセスを繰り返し、最尤仮説を選択する際にこの情報を用いる。次に、サーチ・エンジン 80 は、最尤仮説を出力装置 82 に、入力データが表わす最尤フレーズとして、供給する。これをブロック 194 および 196 で示す。

【0069】

したがって、本発明は、従来のシステムに対して大きな利点をもたらすことがわかる。本発明は、分離音声データおよび連続音声データを訓練データ・セットとして収集するデータ収集方法を用いる。通常の方法を強化し、読み手にワード間にポーズを入れたり、流暢に発話するように要求することによって、連続音声に関連する有音 (non-silence) 前後関係だけでなく、離散音声に関連する無音前後関係 (silence context) も、システムにおける音響モデルを訓練する際に用いられる。連続音声訓練データおよび分離音声訓練データに対する訓練データは、同じワードまたは異なるワードのいずれを含むことも可能であることを注記しておく。この共同訓練データ・セットは、音素モデルの訓練、セノン・ツリーの生成およびセノンの訓練、ならびに三音の適切なセノン・シーケンスへのマッピングに用いられる。

【0070】

また、異なる種類の訓練データ (連続および分離) の効果は、認識の間に予想される音声の種類に応じて別々に重み付けが可能であることも注記しておく。重み付けは、重み係数を割り当てることによって、または単に訓練データ・セットにおける各種類のデータのシステムに供給した量によって行うことができる。好適な実施形態の 1 つでは、双方の種類の訓練データに等しく重み付けする。

【0071】

更に、好適な実施形態の 1 つでは、本発明はワード持続時間モデルを採用する。ワード持続時間モデルは、訓練中に生成し、フレーズの境界に適用し、認識システムの精度を更に高めるようにすることが好ましい。

【0072】

また、本発明の技法は、他の種類の訓練データをシステムに導入するためにも、同様に使用可能である。例えば、ユーザに分離または連続音声として訓練データを入力するように指図するだけでなく、ユーザに、大声で、優しく、もっとゆっくりと、またはもっと素早く、あるいは別の言い方で訓練データを入力するように指図することも可能である。この訓練データの全ては、前述と同様に使用し、システムにおいて用いる音響モデルを訓練し、更に一層ロバストな認識システムを得ることが可能となる。

【0073】

以上好適な実施形態を参照しながら本発明について説明してきたが、本発明の精神および範囲から逸脱することなく、形態および詳細において変更も可能であることを、当業者は認めよう。

【図面の簡単な説明】

【図 1】 本発明による音声認識システムを実現する環境例のブロック図である。

10

20

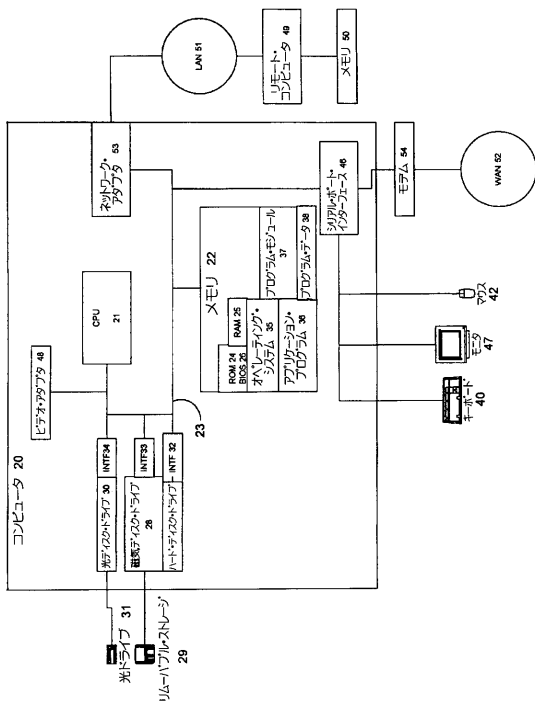
30

40

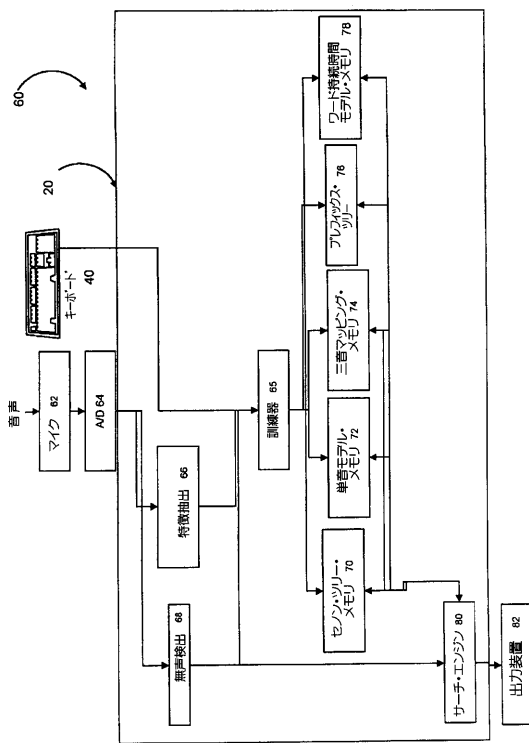
50

- 【図 2】 図 1 に示すシステムの一部の更に詳細なブロック図である。
- 【図 3】 本発明の一形態によるデータ収集手順を示すフロー図である。
- 【図 4】 本発明の一形態による、共同訓練データを用いた音響モデルの訓練およびセノンのマッピングを示すフロー図である。
- 【図 5】 本発明によるセノン・ツリーの作成を示すフロー図である。
- 【図 6】 本発明によるセノン・ツリーの図である。
- 【図 7】 本発明によるワード持続時間モデルの作成を示すフロー図である。
- 【図 8】 図 7 に示す手順にしたがって作成した、複数のワード持続時間モデルのグラフである。
- 【図 9】 本発明の一形態による音声認識手順の一部を示すフロー図である。
- 【図 10】 本発明の一形態によるワード持続時間モデルの適用を示すフロー図である。

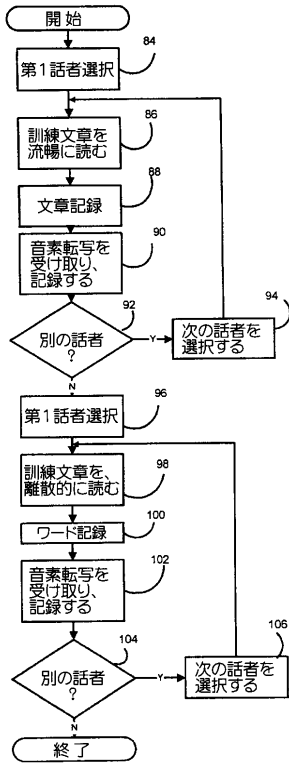
【 図 1 】



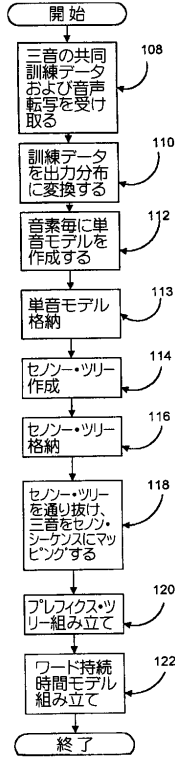
【 図 2 】



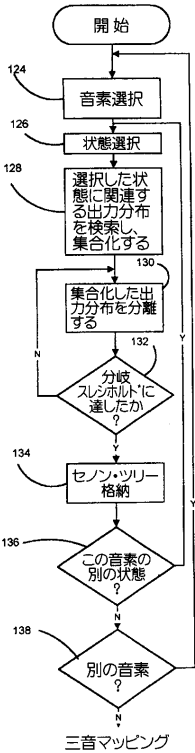
【図3】



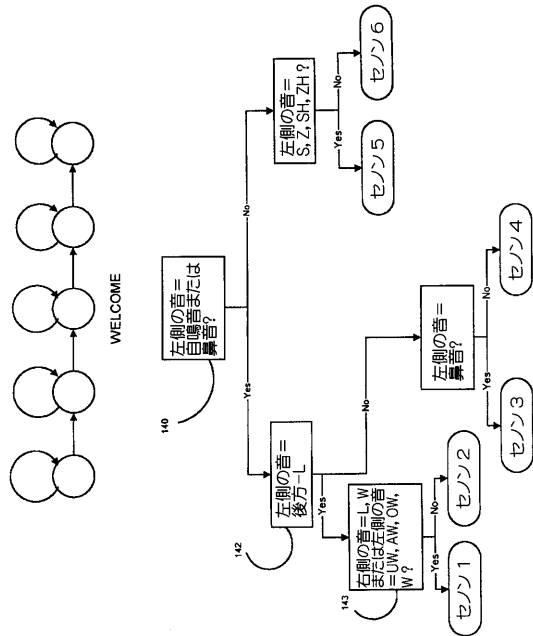
【図4】



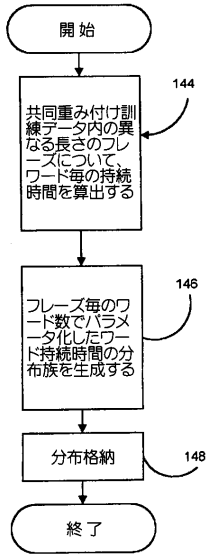
【図5】



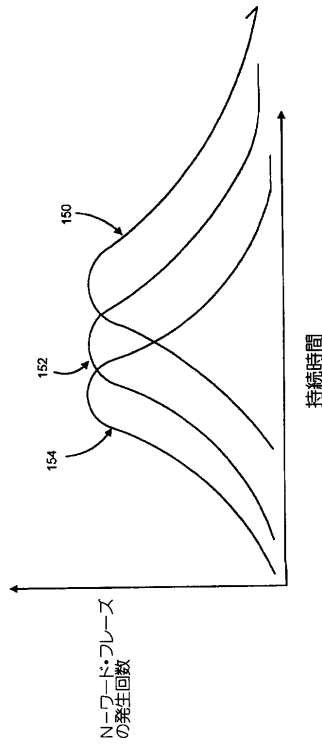
【図6】



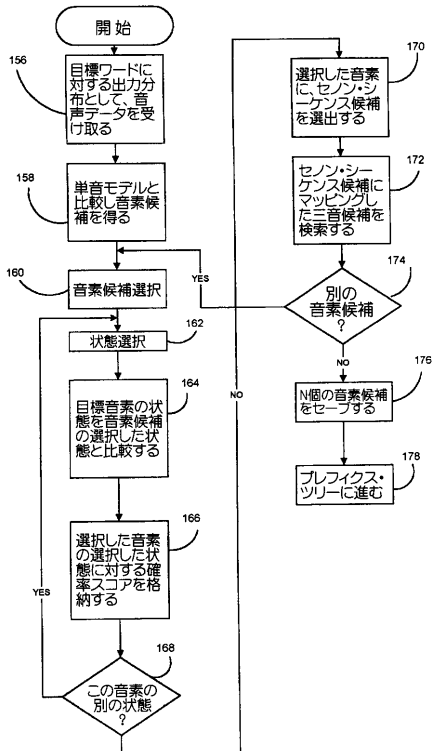
【図7】



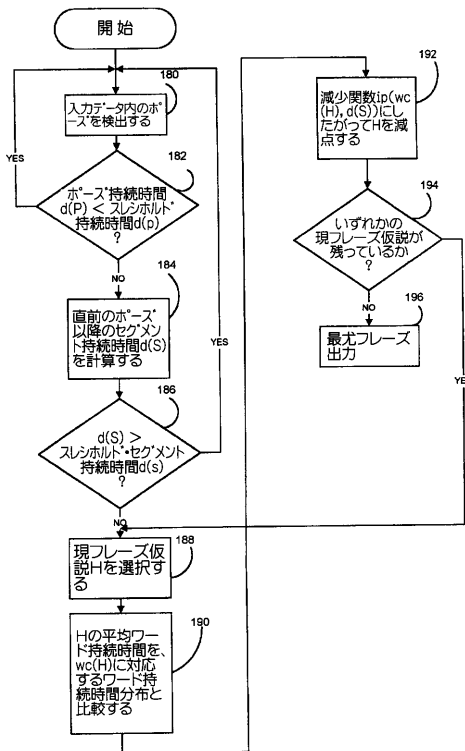
【図8】



【図9】



【図10】



フロントページの続き

- (72)発明者 ファン, シュードン
アメリカ合衆国ワシントン州98072, ウッディンヴィル, ノースイースト・ワンハンドレッド
トゥエンティファースト・ストリート 20020
- (72)発明者 アルレヴァ, フィレノ・エイ
アメリカ合衆国ワシントン州98052, レッドモンド, ノースイースト・フォーティエイス・ス
トリート 16516
- (72)発明者 ジャン, リ
アメリカ合衆国ワシントン州98052, レッドモンド, ノースイースト・シックスティシックス
ス・コート 15360
- (72)発明者 ファン, メイ・ユー
アメリカ合衆国ワシントン州98052, レッドモンド, ノースイースト・シックスティエイス・
ストリート 14802

審査官 涌井 智則

- (56)参考文献 特開平08-211893(JP, A)
特開平04-326400(JP, A)
特開平03-206500(JP, A)
特開昭61-254993(JP, A)
特開昭63-026699(JP, A)
特表平06-501319(JP, A)
特開平8-221090(JP, A)
特開平9-22297(JP, A)
Mei-Yuh HWANG, Xuedong HUANG, Fileno ALLEVA, PREDICTING UNSEEN TRIPHONES WITH SENONES
, Proceedings of the 1993 IEEE International Conference on Acoustics, Speech, and Sign
al Processing, 米国, 1993年 4月27日, Vol.2, p.311-314
三樹聡, 菅村昇, 中津良平, 発声変動への適応化手法を用いた音声認識, 電子情報通信学会技術
研究報告, 日本, 社団法人電子情報通信学会, 1990年 6月29日, Vol.90, No.
112, p.1-8
Richard P. LIPPMANN, Edward A. MARTIN, Douglas B. PAUL, MULTI-STYLE TRAINING FOR ROBUS
T ISOLATED-WORD SPEECH RECOGNITION, Proceedings of the 1987 IEEE International Confere
nce on Acoustics, Speech, and Signal Processing, 1987年 4月, Vol12, p.705-708

(58)調査した分野(Int.Cl., DB名)

G10L 15/00-15/28

CiNii