

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3911246号  
(P3911246)

(45) 発行日 平成19年5月9日(2007.5.9)

(24) 登録日 平成19年2月2日(2007.2.2)

(51) Int. Cl.	F I
<b>G 1 0 L 15/28 (2006.01)</b>	G 1 0 L 15/28 3 6 0 A
<b>G 1 0 L 15/18 (2006.01)</b>	G 1 0 L 15/18 2 0 0 D
	G 1 0 L 15/18 2 0 0 Z
	G 1 0 L 15/18 4 0 0
	G 1 0 L 15/28 3 7 0 E

請求項の数 6 (全 14 頁)

<p>(21) 出願番号 特願2003-56694 (P2003-56694)</p> <p>(22) 出願日 平成15年3月4日(2003.3.4)</p> <p>(65) 公開番号 特開2004-264719 (P2004-264719A)</p> <p>(43) 公開日 平成16年9月24日(2004.9.24)</p> <p>審査請求日 平成16年6月10日(2004.6.10)</p> <p>(出願人による申告) 国等の委託研究の成果に係る特許出願(平成15年度通信・放送機構、研究テーマ「大規模コーパスベース音声対話翻訳技術の研究開発」)に関する委託研究、産業活力再生特別措置法第30条の適用を受けるもの)</p>	<p>(73) 特許権者 393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2</p> <p>(74) 代理人 100099933 弁理士 清水 敏</p> <p>(72) 発明者 大西 茂彦 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内</p> <p>(72) 発明者 菊井 玄一郎 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内</p> <p>(72) 発明者 山本 博史 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内</p> <p style="text-align: right;">最終頁に続く</p>
---	--

(54) 【発明の名称】 音声認識装置、及びコンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

所定のコーパスから生成した N - グラム (ただし N は正の整数) 言語モデルに基づき、入力音声の音声認識を行なって音声認識結果を出力するための第1の音声認識手段と、前記第1の音声認識手段の音声認識結果から、前記入力音声中に含まれる所定の音声単位の数を推定するための推定手段と、

前記所定のコーパス内の文を、各々に含まれる前記所定の音声単位の数に従って分類した複数個の部分集合からそれぞれ生成された、各々が有限状態機械からなる、音声認識のための複数個の言語モデルを記憶するための手段と、

前記複数個の言語モデルのうち、前記推定手段により推定された音声単位の数に応じて、前記推定された音声単位の数に対応するものと、前記推定された音声単位の数の前若しくは後又はその双方の予め定められた範囲の音声単位の数にそれぞれ対応するものを含む 予め定められた複数個を選択するための第1の言語モデル選択手段と、

前記第1の言語モデル選択手段により選択された複数個の言語モデルにそれぞれ基づいて入力音声の音声認識を行ない、複数個の音声認識結果を出力するための第2の音声認識手段と、

前記第1の音声認識手段の前記音声認識結果、及び前記第2の音声認識手段の前記音声認識結果のうち、所定の選択基準に従って一つを選択して出力するための認識結果選択手段とを含む、音声認識装置。

【請求項2】

10

20

前記推定手段は、前記第 1 の音声認識手段の音声認識結果から、前記入力音声に含まれるシラブルの数を推定するための手段を含み、

前記第 1 の言語モデル選択手段は、前記複数個の言語モデルのうち、前記推定するための手段により推定されたシラブルの数に応じて、前記推定されたシラブルの数に対応するものと、前記推定されたシラブルの数の前若しくは後又はその双方の予め定められた範囲のシラブルの数にそれぞれ対応するものとを含む予め定められた複数個を選択するための第 2 の言語モデル選択手段を含む、請求項 1 に記載の音声認識装置。

【請求項 3】

前記第 1 の言語モデル選択手段は、前記複数個の言語モデルのうち、前記推定手段により推定された音声単位の数に応じて、前記推定された音声単位の数に対応するものと、前記推定された音声単位の数の前後の互いに等しい範囲の音声単位の数にそれぞれ対応するものを選択するための第 2 の言語モデル選択手段を含む、請求項 1 に記載の音声認識装置。

10

【請求項 4】

前記認識結果選択手段は、

前記第 2 の音声認識手段の前記音声認識結果の各々と、前記第 1 の音声認識手段の前記音声認識結果との間の DP マッチング距離を算出するための距離算出手段と、

前記第 2 の音声認識手段の前記音声認識結果のうち、前記距離算出手段により算出された DP マッチング距離の小さなものから順番に予め定める数だけ選択するための手段と、

前記選択するための手段により選択された前記予め定められた数の前記第 2 の音声認識手段の前記音声認識結果の音響スコアの各々を、予め定められた調整式によって調整するための音響スコア調整手段と、

20

前記音響スコア調整手段により出力される複数個の音響スコアと、前記第 1 の音声認識手段の認識結果の音響スコアとを比較して、最も大きな音響スコアを選択し、選択された音響スコアに対応する音声認識結果を前記音声認識装置による音声認識結果として出力するための手段とを含む、請求項 1 ~ 請求項 3 のいずれかに記載の音声認識装置。

【請求項 5】

前記音響スコア調整手段は、前記選択するための手段により選択された前記第 2 の音声認識手段の前記音声認識結果の対数音響スコアを、 $1 +$  (ただし は予め定められた正の定数) 倍する事により調整するための手段を含む、請求項 4 に記載の音声認識装置。

30

【請求項 6】

コンピュータにより実行されると、当該コンピュータを請求項 1 ~ 請求項 5 のいずれかに記載の音声認識装置として動作させる、コンピュータプログラム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

この発明は音声認識技術に関し、特に、大規模コーパスから得た言語モデルを有効に用いる事により文認識率を向上させるための技術に関する。

【0002】

【従来の技術】

40

音声翻訳では、音声認識におけるわずか 1 単語の誤りが翻訳結果に致命的な影響を及ぼす事がしばしばある。そうした問題を回避するためには、音声認識における文認識率を向上させる事が必要である。文認識率を向上させるための一つの手法は、有限状態機械 (Finite State Automaton: FSA) などを用いて文全体の大域的なモデル化によって強い制約を与える事である。このためには、以下の二つが必要である。

【0003】

(1) 当該タスクの言語表現を広く覆う様なモデルの構築

(2) その様なモデルを使った効率の良い探索手法の開発

前者に関しては、N-グラムのような制約の緩いモデルと併用する事によって、大域的なモデルでカバーできなかった部分を救う事が試みられている (非特許文献 1 を参照されたい

50

。)。しかし、それは小規模な実験にとどまっている。

【0004】

【非特許文献1】

鶴見 他、「単語N - g r a mとネットワーク文法を併用した音声認識アルゴリズムの検討」、日本音響学会講演論文集3 - 9 - 8、pp. 145 - 146、日本音響学会2002年秋季発表会

【発明が解決しようとする課題】

文認識率を向上させるために、大量（たとえば16万文規模）の文を含むコーパスとN - グラムとを併用する場合、その様な大量の文を単純にFSAでモデル化すると、モデルサイズと探索空間の巨大化を招く事は明らかである。従って、大規模コーパスとN - グラムの様な制約の緩いモデルとを併用しながら、モデルサイズと探索空間の巨大化を抑える様な手法が求められている。

10

【0005】

それゆえに本発明の目的は、大規模コーパスとN - グラムの様な制約の緩いモデルとを併用して文認識率を向上させながら、モデルサイズと探索空間の巨大化を抑える事が可能な音声認識装置を提供する事である。

【0006】

この発明の他の目的は、大規模コーパスとN - グラムの様な制約の緩いモデルとを併用し、かつ小さなモデルを用いて実質的に大規模コーパスを用いた場合と同様の文認識率の向上を図る事が可能な音声認識装置を提供する事である。

20

【0007】

この発明のさらに他の目的は、大規模コーパスを所定の基準に従って分割した小さなコーパスから生成した言語モデルと、N - グラムの様な制約の緩いモデルとを併用する事で、実質的に大規模コーパスを用いた場合と同様の文認識率の向上を図る事が可能な音声認識装置を提供する事である。

【0008】

本発明の他の目的は、大規模コーパスを、所定の基準に従って分割した小さなコーパスから生成した言語モデルと、N - グラムの様な制約の緩いモデルとを併用し、かつ小さなコーパスから生成した言語モデルのうち、適切なものを選択して再度音声認識する事で、実質的に大規模コーパスを用いた場合と同様の文認識率の向上を図る事が可能な音声認識装置を提供する事である。

30

【0009】

【課題を解決するための手段】

本発明の第1の局面に係る音声認識装置は、所定のコーパスから生成した言語モデルに基づき、入力音声の音声認識を行なって音声認識結果を出力するための第1の音声認識手段と、第1の音声認識手段の音声認識結果から、入力音声中に含まれる所定の音声単位の数を推定するための推定手段と、所定のコーパス内の文を、各々に含まれる所定の音声単位の数に従って分類した複数個の部分集合からそれぞれ生成された、音声認識のための複数個の言語モデルを記憶するための手段と、複数個の言語モデルのうち、推定手段により推定された音声単位の数に応じて、予め定められた複数個を選択するための第1の言語モデル選択手段と、第1の言語モデル選択手段により選択された複数個の言語モデルにそれぞれ基づいて入力音声の音声認識を行ない、複数個の音声認識結果を出力するための第2の音声認識手段と、第1の音声認識手段の音声認識結果、及び第2の音声認識手段の音声認識結果のうち、所定の選択基準に従って一つを選択して出力するための認識結果選択手段とを含む。

40

【0010】

好ましくは、所定のコーパスから生成した言語モデルは、N - グラム（ただしNは1以上の整数）言語モデルである。

【0011】

より好ましくは、音声認識装置の推定手段は、第1の音声認識手段の音声認識結果から、

50

入力音声に含まれるシラブルの数を推定するための手段を含み、第1の言語モデル選択手段は、複数個の言語モデルのうち、推定するための手段により推定されたシラブルの数に応じて、予め定められた複数個を選択するための第2の言語モデル選択手段を含む。

【0012】

好ましくは、第2の言語モデル選択手段は、複数個の言語モデルのうち、推定するための手段により推定されたシラブルの数に応じて、推定されたシラブルの数に対応するものを含む予め定められた複数個を選択するための第3の言語モデル選択手段を含む。

【0013】

さらに好ましくは、第3の言語モデル選択手段は、複数個の言語モデルのうち、推定するための手段により推定されたシラブルの数に応じて、推定されたシラブルの数に対応するものと、推定されたシラブルの数の前若しくは後又はその双方の予め定められた範囲のシラブル数にそれぞれ対応するものを選択するための第4の言語モデル選択手段を含む。

10

【0014】

より好ましくは、第4の言語モデル選択手段は、複数個の言語モデルのうち、推定するための手段により推定されたシラブルの数に応じて、推定されたシラブルの数に対応するものと、推定されたシラブルの数の前後の互いに等しい範囲のシラブル数にそれぞれ対応するものを選択するための第5の言語モデル選択手段を含む。

【0015】

さらに好ましくは、認識結果選択手段は、第2の音声認識手段の音声認識結果の各々と、第1の音声認識手段の音声認識結果との間のDP(Dynamic Programming)マッチング距離を算出するための距離算出手段と、第2の音声認識手段の音声認識結果のうち、距離算出手段により算出されたDPマッチング距離の小さなものから順番に予め定める数だけ選択するための手段と、選択するための手段により選択された予め定められた数の第2の音声認識手段の音声認識結果の音響スコアの各々を、予め定められた調整式によって調整するための音響スコア調整手段と、音響スコア調整手段により出力される複数個の音響スコアと、第1の音声認識手段の認識結果の音響スコアとを比較して、最も大きな音響スコアを選択し、選択された音響スコアに対応する音声認識結果を音声認識装置による音声認識結果として出力するための手段とを含む。

20

【0016】

より好ましくは、音響スコア調整手段は、選択するための手段により選択された第2の音声認識手段の音声認識結果の対数音響スコアを、 $1 +$  (ただし は予め定められた正の定数) 倍する事により調整するための手段を含む。

30

【0017】

本発明の第2の局面に係るコンピュータプログラムは、コンピュータにより実行されると、当該コンピュータを上記したいずれかの音声認識装置として動作させるものである。

【0018】

【発明の実施の形態】

[第1の実施の形態]

<概略>

予め、本実施の形態に係る音声認識の手法の概要を述べる。本実施の形態では、まず、FSAの訓練コーパス内の全ての文を、各文が含むシラブル数によって部分集合に分け、それら部分集合にそれぞれ対応した複数のFSAを生成する。その上で、次の様にして音声認識処理を行なう。すなわち、

40

(1) まずN-グラムを用いて音声認識を行なう。

【0019】

(2) 上記(1)の結果の文が含むシラブル数を求める。

【0020】

(3) その値に応じてFSAの部分集合のうちの適切なもの(多くは複数個)を選択する。

【0021】

50

(4) 選択された F S A を用いて、(1)と同じ入力に対して音声認識を行なう。

【0022】

(5) 上記(1)と(4)とで得られた音声認識結果のうち、ある基準によって最適と思われる1文を選択して最終的な音声認識結果とする。

【0023】

<構成>

図1に、本実施の形態に係る音声認識装置30の機能的ブロック図を示す。図1を参照し、この音声認識装置30は、音声信号を受け、音声認識に適したデジタルの音声データ列に変換するための音声入力部40と、予め準備された大規模な訓練コーパス中に出現する全文のみが受理可能な(すなわち、この訓練コーパスに含まれる文のみを認識結果とする様な)マルチクラス複合2-グラム言語モデル42と、音声入力部40から与えられる音声データ列に対して、言語モデル42を用いた音声認識を行ない、最尤の認識結果文を示す単語列を出力するための音声認識部44と、音声認識部44の出力する単語列に含まれるシラブル数NSを算出するためのシラブル数算出部46とを含む。この場合、シラブル数算出部46が算出するシラブル数NSは、実際には入力音声が表示文が含むシラブル数の推定値である。

10

【0024】

上記した訓練コーパスは旅行会話基本表現集(BTEC)であり、その内容は以下の通りである。ただし下の表で「クローズド文」とは、テストセット文のうちで訓練セットにも現れた文の事をいう。

20

【0025】

【表1】

	訓練セット	テストセット
単語数	約150万	3475
異なり文数	約10万	501文
延べ文数	約16万	510文(249文はクローズド文)

30

音声認識装置30はさらに、上記した訓練コーパスを各文が含むシラブル数により部分集合に分け、それら部分集合にそれぞれ対応して生成された複数のF S A 60を記憶するF S A 記憶部48と、F S A 記憶部48の複数のF S A 60のうち、シラブル数算出部46が算出したシラブル数NSを中心とした所定の幅のシラブル数に対応するF S A 52を選択するための選択部50と、選択部50が選択した複数のF S A 52をそれぞれ用いて、音声認識部44が行なったものと同じ音声データ列に対して音声認識を行ない、結果の単語列を出力するための複数の音声認識部54と、音声認識部44の出力、及び複数の音声認識部54の出力のうち、後述する基準に従って最も適切と思われる認識結果を選択し最終認識結果として出力するための選択部56とを含む。複数の音声認識部54の各々は、音声認識部44と同じ認識エンジンによって音声認識を行なう。逆にいえば、同じ音声認識エンジンを用いるのであれば、音声認識部44と音声認識部54としてどの様なものを用いても良い。

40

【0026】

図2を参照して、選択部56は、複数の音声認識部54から与えられる複数の認識結果の各々と、音声認識部44から与えられる認識結果との間の音素並びとの間のDPマッチング距離を算出するためのDP距離算出部70と、DP距離算出部70が算出したDPマッチング距離に基づき、複数の音声認識部54からの複数の認識結果のうち、音声認識部44の認識結果の音素並びと近いもの上位3位までの候補を選択するためのDP距離比較・選択部72と、DP距離比較・選択部72により選択された上位三位までの認識結果に付随する音響スコア(認識時に認識結果と同時に得られる)の値を(1+ )倍して調整す

50

るためのスコア調整部 7 4 と、調整部 7 4 が出力する上位 3 位までの認識結果の調整済み音響スコアと、音声認識部 4 4 の認識結果に付随する音響スコアとを比較し、最も高い値を持つ音響スコアに対応する認識結果を最終認識結果として出力するためのスコア比較・選択部 7 6 とを含む。

#### 【 0 0 2 7 】

上記した様にコーパスを各文の含むシラブル数に基づいて部分集合に分割し、それらに対応した F S A を使用する様にした理由は以下の通りである。前述した通り、訓練コーパスには大量の文が含まれている。そのまま 1 つの F S A でこのコーパスをモデル化すると、モデルサイズ、探索空間ともに巨大となってしまう。そこで、F S A を分割する必要がある。

10

#### 【 0 0 2 8 】

F S A を分割する際には、言語モデル 4 2 を用いた音声認識部 4 4 の認識結果から正確に推定できるパラメータに沿って、F S A を選択できる様にすることがある。本実施の形態では、上記したパラメータとして入力音声のシラブル数 N S を採用した。入力音声のシラブル数 N S を採用したのは、N - グラムのみを言語モデルに用いた場合、認識結果の正誤はともかくとして、認識結果に含まれるシラブル数 N S が正解のそれと大差がない事が判明したためである。図 3 はそれを示す実験結果である。

#### 【 0 0 2 9 】

図 3 は、本実施の形態で用いたコーパスを用いて N - グラム ( マルチクラス複合 2 - グラム ) の言語モデルを作成し、その言語モデルを用いて音声認識を行なった結果の上位 1 0 位までの文について、そのシラブル数が正解文と最も異なるものを選び、そのシラブル数の差をヒストグラムにしたものである。図 3 からわかる様に、マルチクラス複合 2 - グラムの認識結果のシラブル数と、正解文のそれとの差が 5 以内に収まるテスト文が全体の 9 9 . 5 パーセントを占めている。したがって、マルチクラス複合 2 - グラムの認識結果のシラブル数から、正解文のシラブル数を同定 ( 推定 ) する事が十分に可能である事が分かる。

20

#### 【 0 0 3 0 】

そこで、コーパスを各文のシラブル数によって部分集合に分割した各部分集合から F S A 6 0 を得る様にした。図 4 に、異なり文数約 1 0 万の訓練コーパスを各文のシラブル数により部分集合に分けた際に、各集合に含まれる訓練文数の分布を示している。横軸はシラブル数、縦軸はそのシラブル数ごとの集合に含まれる学習文の数を示す。図 4 を参照して、この分布は、1 6 シラブルにピークを持ったガンマ分布で近似可能な分布である。集合中の文数は最大で 1 2 0 0 0 以下であり、元の訓練コーパスの 8 分の 1 以下である。個々の集合に対応した F S A モデルの個数は総計 6 5 個となった。なお、各文 ( 単語列 ) に対する言語尤度は同一とした。

30

#### 【 0 0 3 1 】

言語モデル 4 2 は上記した訓練コーパスを用いて訓練した言語モデルである。この言語モデルを用いて、この訓練コーパスと同ドメインのテストセット ( 訓練コーパスと同じく旅行という状況で使用される文の集合 ) に対して音声認識部 4 4 単独により行なった音声認識実験の結果、単語認識率は 8 8 . 7 %、文認識率は 6 5 . 7 % であった。単語認識率は高い値であるのに対して文認識率は低い値である。B T E C は定型文が多いという性質上、クローズド文が多いのだが、それらについても 8 5 . 1 % しか完全には正解となっていない。それゆえ、残りの 1 4 . 9 % について、本実施の形態の様に、訓練セットの文をモデル化したクローズドな F S A を併用して音声認識を行なう事で正解出来る可能性がある。仮にそれらについて正解できたものとする、全体として文認識率は 7 . 2 5 ポイント改善されるはずである。

40

#### 【 0 0 3 2 】

なお、本実施の形態では後述する様に音声認識部 5 4 はソフトウェアで実現する。そのため、一つの C P U で複数回の音声認識を繰返し行なう。しかし、ハードウェアに余裕があれば、音声認識処理を並列に行なってもよい。たとえば複数個の C P U があれば音声認識

50

処理を同時に並列に実行することができる。

【 0 0 3 3 】

< 動作 >

以上の構成により、ベースラインである 1 個のマルチクラス複合 2 - グラムと、65 個の F S A 群とを言語モデルとして利用できる。以下、図 1 及び図 2 に記載した音声認識装置 30 の動作について述べる。

【 0 0 3 4 】

図 1 を参照して、音声入力部 40 が音声入力を音声認識に適したデジタルデータに変換する。このデジタルデータは音声認識部 44 及び音声認識部 54 に与えられる。最初に音声認識部 44 が音声認識を行なう。この際、音声認識部 44 は、マルチクラス複合 2 - グラムの言語モデル 42 を用いて認識処理を行ない、最尤の認識結果文を得る。この認識結果文はシラブル数算出部 46 及び選択部 56 に与えられる。なおこの音声認識に伴い、音声認識の際の音響スコアも得られ、選択部 56 に与えられる。ここで、音響スコアとは、認識の際の入力波形に対する認識結果の間の尤度を音素（又はシラブル）ごとに求め、それらを乗算した値である。通常はその対数を用いる。本明細書では、音響スコアといえは対数音響スコアの事をいうものとする。

10

【 0 0 3 5 】

シラブル数算出部 46 は、与えられた認識結果文のシラブル数  $NS$  を求める。シラブル数算出部 46 は、シラブル数  $NS$  を選択部 50 に与える。

【 0 0 3 6 】

選択部 50 は、与えられたシラブル数  $NS$  を中心として  $NS \pm 5$  の範囲のシラブル数の文集合に対応した F S A 52 を選択する。すなわち、シラブル数が  $NS$  を中心として  $NS - 5$  から  $NS + 5$  までの 11 個の F S A 52 が選択される。

20

【 0 0 3 7 】

音声認識部 54 は、この 11 個の F S A 52 を用いて、音声入力部 40 から与えられたデジタルデータに対し個別に音声認識を行なう。その結果、 $NS - 5$  から  $NS + 5$  までのシラブル数の文集合に対応した F S A 52 による 11 個の認識結果文が得られる。この出力は選択部 56 に与えられる。また同時にこれらの音声認識の際の音響スコアも得られ、選択部 56 に与えられる。

【 0 0 3 8 】

図 2 を参照して、以上の処理の結果、D P 距離算出部 70 には音声認識部 44 からの一つの認識結果、及び音声認識部 54 からの 11 個の認識結果が与えられる。D P 距離算出部 70 は、11 個の F S A 52 による認識結果の各々と、音声認識部 44 の認識結果の音素並びとの間の D P マッチング距離を算出し、D P 距離比較・選択部 72 に与える。

30

【 0 0 3 9 】

D P 距離比較・選択部 72 は、これら 11 個の D P マッチング距離のうち、距離の短いものを上位 3 位まで選択し、調整部 74 に与える。

【 0 0 4 0 】

調整部 74 は、D P 距離比較・選択部 72 により選択された上位 3 位までの認識結果に付随する音響スコアを  $(1 + \alpha)$  倍する。この  $\alpha$  は、音響スコアの調整のためのパラメータである。 $\alpha$  は正の値であり、0.05 から 0.25 までの値が好ましい。特に、 $\alpha = 0.1 \sim 0.25$  が好ましく、さらに 0.15  $\sim$  0.25 が好ましい。実験によれば、最も好ましい結果が得られたのは  $\alpha = 0.2$  のときであった。ただし、 $\alpha$  の値は認識の対象により変化し得るので、この値に限定されるわけではない。調整部 74 は、この様に  $(1 + \alpha)$  を乗算する事により調整済みの音響スコアをスコア比較・選択部 76 に与える。

40

【 0 0 4 1 】

スコア比較・選択部 76 は、音声認識部 44 の認識結果の音響スコアと、調整部 74 から与えられる 3 つの調整済みの音響スコアとを比較し、最も音響スコアが高いものを最終的な認識結果として選択し出力する。

【 0 0 4 2 】

50

本実施の形態の装置で、 $\alpha$  を 0.05 から 0.25 まで 0.05 刻みで変化させ、上記した態様に従って認識結果を選択した実験結果を図 5 及び図 6 に示す。図 5 に示すのは、 $\alpha$  の各値に対する文認識率の推移である。図 6 には、参考のためにこのときの単語認識率の推移を示す。

#### 【0043】

図 5 に示す様に、本実施の形態では  $\alpha = 0.20$  のときに最高の文認識率が得られた。図 5 に示した  $\alpha = 0$  の値は、マルチクラス複合 2 - グラム言語モデル 42 のみを用いたときの文認識率である。図 5 からわかる様に、 $\alpha = 0.20$  のときには文認識率は  $\alpha = 0$  のときと比較して約 4.3 ポイント高い。この値をクロズド文のみについて計算すると、約 8.8 ポイントの文認識率の向上に値する。

10

#### 【0044】

以上の様に本実施の形態の音声認識装置 30 によれば、コーパスをシラブル数により複数の文集合に分割し、各々について FSA を言語モデルとして生成する。そして、マルチクラス複合 2 - グラムの言語モデルを用いたベースラインの音声認識結果のシラブル数 NS に応じて、FSA のうちこのシラブル数 NS と特定の関係にある文集合に対応するものを選択して再度音声認識を行なう。その認識結果のうち、ベースラインの音声認識結果の音素の並びとの DP マッチング距離が小さなものを 3 つ選択し、それらの音響スコアを  $1 + \alpha$  で調整した上で、ベースラインの音声認識結果と比較し、最も高い音響スコアを示したものを最終認識結果として選択する。その結果、ベースラインの音声認識結果と比較して上記した通り文認識率にかなりの改善が見られた。

20

#### 【0045】

< コンピュータによる実現 >

上記した本実施の形態の音声認識装置は、音声処理機能を備えたコンピュータにより実現できる。図 7 にコンピュータにより実現された音声認識装置 30 の外観を示す。図 8 はこの音声認識装置 30 のハードウェアブロック図である。

#### 【0046】

図 7 を参照して、音声認識装置 30 は、CD-ROM (Compact Disc Read-Only Memory) 駆動装置 90、FD (Flexible Disk) 駆動装置 92 を備えたコンピュータ 80 と、いずれもコンピュータ 80 に接続されたモニタ 82、マイク 84、キーボード 86、及びマウス 88 とを含む。

30

#### 【0047】

図 8 を参照して、コンピュータ 80 は、前述した CD-ROM 駆動装置 90 及び FD 駆動装置 92 に加えて、CPU (Central Processing Unit) 96 と、ROM (Read-Only Memory) 98 と、RAM (Random Access Memory) 100 と、ハードディスク 94 と、マイク 84 に接続されたサウンドボード 108 とを含む。これらはいずれもバス 106 により相互に接続されている。CD-ROM 駆動装置 90 には CD-ROM 102 が装着され、FD 駆動装置 92 には FD 104 が装着される。

#### 【0048】

以下に述べる制御構造を有するコンピュータプログラムは、たとえば CD-ROM 102 又は FD 104 の様なコンピュータ読取可能な記録媒体上に記録されて流通し、当該 CD-ROM 102 を CD-ROM 駆動装置 90 に装着したのち CD-ROM 102 からハードディスク 94 に複製される。実行時にはこのプログラムはハードディスク 94 から読出されて RAM 100 に読出され、図示しないプログラムカウンタにより指定されるアドレスから CPU 96 が読出して実行し、実行結果を RAM 100 又はハードディスク 94 に書込む。CPU 96 はさらにプログラムカウンタの値をプログラムの実行結果により書換え、さらにそのプログラムカウンタの値に基づいて次の命令を RAM 100 から読出して実行する。CPU 96 はこの様な動作原理に従って、コンピュータプログラムを実行する。なお、FSA 記憶部 48 は、ハードディスク 94、ROM 98、又は RAM 100 などにより実現される。

40

50

## 【0049】

図9に、このコンピュータプログラムの全体の制御構造を示す。図9を参照して、このコンピュータプログラムは、音声入力に対してマルチクラス複合2 - グラムの言語モデル42を用いて音声認識を行なうステップ110と、ステップ110の音声認識の結果得られた単語列のシラブル数NSを算出するステップ112と、ステップ112で算出されたシラブル数NSに従ってシラブル数NS - 5からNS + 5までの文集合から得られた11個のFSAを選択し、選択されたFSAを用いて音声認識を行なう事により、11個の音声認識結果をその音響スコアとともに出力するステップ114と、ステップ110の音声認識結果と、ステップ114で得られた11個の音声認識結果とのうち、前述した方法に従って一つの音声認識結果を選択するステップ116と、選択された音声認識結果を出力するステップ118とを含む。ステップ118が完了すると、一回の音声入力に対する処理が完了した事になる。

10

## 【0050】

図10に、図9のFSAによる認識ステップ114の詳細を示す。図10を参照して、ステップ114は、使用する文集合のシラブル数の最小値MIN (= NS - 5)を算出するステップ140と、最大値MAX (= NS + 5)を算出するステップ142と、以下のループ処理の繰返し変数Iを最小値MINに設定するステップ144と、繰返し変数Iが最大値MAXを超えたか否かを判定し、 $I > MAX$ の場合とそれ以外の場合とで制御を分岐させるステップ146とを含む。ステップ146で $I > MAX$ と判定された場合にはこのルーチンは終了する。さもなければ制御はステップ148に進む。

20

## 【0051】

FSAによる認識ステップ114はさらに、シラブル数Iの文集合から得られたFSAを用いて認識を行なうステップ148と、ステップ148での認識結果をその音響スコアと共に出力するステップ150と、繰返し変数Iに1を加算するステップ152とを含む。ステップ152の後、制御はステップ146に戻る。

## 【0052】

図11に、図9の認識結果の選択ステップ116の詳細を示す。図11を参照して、ステップ116は、FSAによる11個の認識結果の各々に対し、その音素並びと、マルチクラス複合2 - グラムの言語モデルによる認識結果の音素並びとの間のDPマッチング距離を算出するステップ180と、ステップ180で算出されたDPマッチング距離の小さなもの上位3位までを選択するステップ182と、ステップ182で選択された3個の認識結果の音響スコアを $(1 + \quad)$ 倍して調整するステップ184と、ステップ184で算出された調整後の音響スコアと、図9のステップ110で行なわれた音声認識結果に付随して得られた音響スコアとのうち、最大の音響スコアを持つ音声認識結果を最終認識結果として選択するステップ186とを含む。

30

## 【0053】

上記した制御構造を有するコンピュータプログラムをコンピュータ80に実行させる事により、既に説明した音声認識装置30の各機能が実現できる事は当業者には容易に理解できるであろう。

## 【0054】

<変形例>

上に述べた実施の形態では、言語モデル42としてマルチクラス複合2 - グラム言語モデルを用いている。しかし本発明はマルチクラス複合2 - グラムを用いたものだけに限定されるわけではない。一般的にN - グラム(ただしNは1以上の整数)の言語モデルを言語モデル42に用いる事ができる。

40

## 【0055】

また、上に述べた実施の形態では、音声認識部44による認識結果に含まれるシラブル数によってFSAを選択している。しかし本発明はFSAの選択にシラブル数を用いるものには限定されない。使用できる要素の大きさは、N - グラムによる認識結果から元の要素数を比較的正確に推定できるものであればどのようなものでもよい。たとえば、シラブル

50

に代えて認識結果に含まれる音素数を用いても良いし、日本語の場合におけるモーラの数の様に、シラブル数に類似した概念の要素を用いても良い。認識システムの性能によりたとえば単語単位で数えたときに正解に近くなる様なものがあれば、それを用いても良い。

【0056】

また、上記した実施の形態では、選択部50は、FSA記憶部48に含まれるFSA60のうち、シラブル数算出部46が算出したシラブル数NSを中心としてその前後5つずつ、合計11個のFSAを選択している。しかし本発明はその様な実施の形態に限定されない。たとえば、シラブル数算出部46が算出したシラブル数NSとその前だけ、又はその後だけを使用する場合もあり得る。また、選択するFSAの数も11個には限定されない。本実施の形態ではFSAの選択にシラブル数を用いており、実験の結果図3に示す様に推定されたシラブル数NSの前後5つずつ、合計11個を用いればほぼ正解文のシラブル数に対応するFSAを選択できる事が分かった。しかし、たとえばハードウェア資源に制限があったり処理速度に制約があったりした場合には、より少ない数のFSAを選択する様にしてもよい。また、推定する音声単位をシラブルではなく音素などとした場合には、当然に図3とは異なる実験結果が予想され、その結果、選択すべきFSAの個数も変わる事があり得る。

10

【0057】

上に述べた実施の形態では、FSAを用いた音声認識は、同じCPUを用いて繰返し処理を行なっている。しかし、前述した通りCPUを複数個利用できるのであれば、これらを並列に動作させて処理させると処理を高速に行なう事ができる。ただしその場合には、図9のステップ116の処理を開始するにあたり、各処理の間の同期をとる必要がある。この場合、FSAによる認識処理の全てが終了した時点で図9のステップ116の処理を開始する様にしてもよいし、FSAによる認識処理が終了したものについて、随時図11に示すステップ180の処理を行ない、全ての認識処理が終了した時点でステップ182以降の処理を開始する様にしてもよい。

20

【0058】

さらに、上記した実施の形態では、選択部56による選択にDPマッチング距離と音響スコアとを用いた。しかし本発明はそうした実施の形態には限定されない。たとえば、SVM(Support Vector Machine)などの機械学習システムに予め選択方法を学習させたものを用いても良い。

30

【0059】

なお、上記した説明では、コンピュータプログラムは記録媒体上に記録されて流通するものとした。しかしコンピュータプログラムの流通形態は記録媒体上に記録されるものに限定されない。たとえば、有線又は無線のネットワークを介した通信という形で流通する事もあり得る。また、コンピュータで直ちに実行可能な形式のみに限らず、ソースプログラムの形で流通し、コンピュータでコンパイルする事により実行形式とする場合もあり得る。さらに、当該コンピュータプログラムを実行するときに、一時にはプログラムのうち実行に関係する一部分のみを遠隔地からネットワークを介してRAM100に記憶して実行し、一度にプログラムの全体をRAM100又はハードディスク94に記憶しない様な運用形態もあり得る。しかしそのいずれの場合も、本発明の実施に該当する事はいうまでもない。

40

【0060】

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味及び範囲内でのすべての変更を含む。

【図面の簡単な説明】

【図1】 本発明の一実施の形態に係る音声認識装置30のブロック図である。

【図2】 図1に示す音声認識装置30のうち、選択部56のより詳細なブロック図である。

50

【図3】 N - グラムによる音声認識結果に含まれるシラブル数と、正解文のシラブル数との差の分布例を示すグラフである。

【図4】 訓練コーパス内の文のシラブル数の分布を示すグラフである。

【図5】 本発明の一実施の形態による文認識率の改善を示すグラフである。

【図6】 本発明の一実施の形態による単語認識率を示すグラフである。

【図7】 音声認識装置30を実現するコンピュータ80及びその周辺装置の外観を示す図である。

【図8】 コンピュータ80のハードウェアブロック図である。

【図9】 コンピュータ80で実行される音声認識プログラムの制御構造を示すフローチャートである。

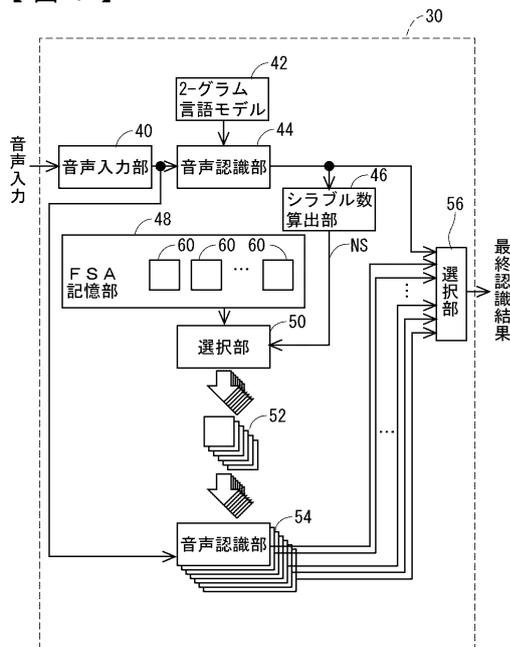
【図10】 図9に示すFSAによる認識ステップ114の詳細な制御構造を示すフローチャートである。

【図11】 図9に示す認識結果の選択ステップ116の詳細な制御構造を示すフローチャートである。

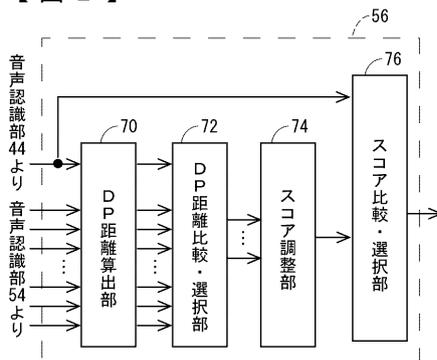
【符号の説明】

30 音声認識装置、40 音声入力部、42 マルチクラス複合2 - グラムによる言語モデル、44 音声認識部、46 シラブル数算出部、48 FSA記憶部、50 選択部、52、60 FSA、54 音声認識部、56 選択部、70 DP距離算出部、72 DP距離比較・選択部、74 スコア調整部、76 スコア比較・選択部

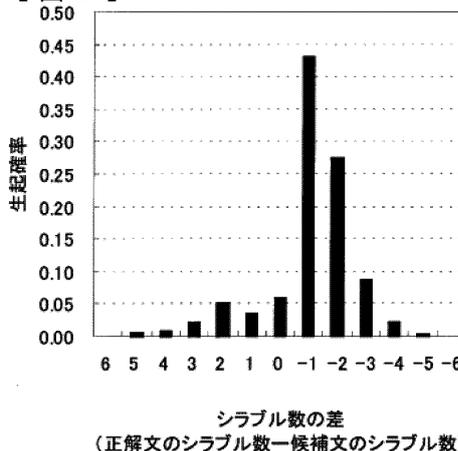
【図1】



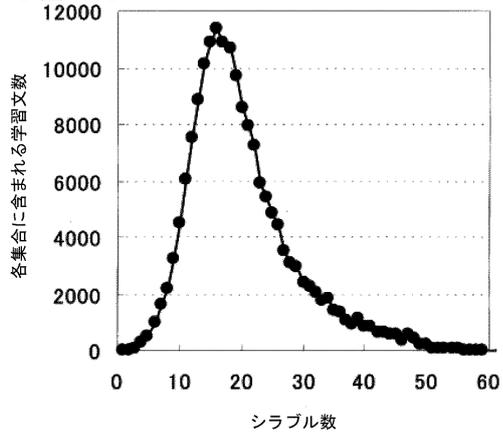
【図2】



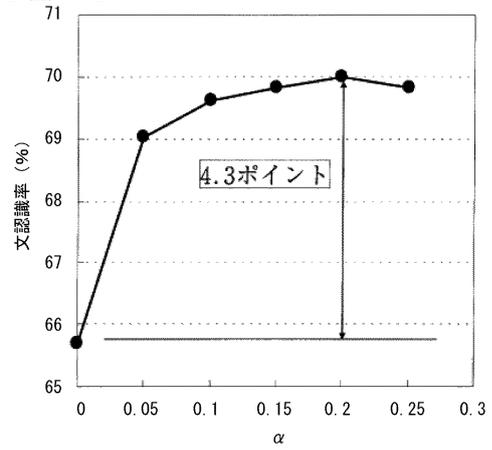
【図3】



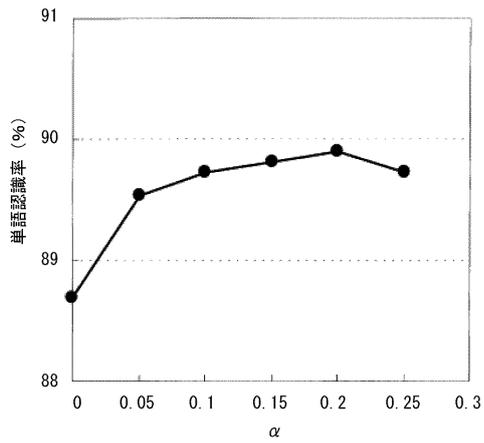
【図4】



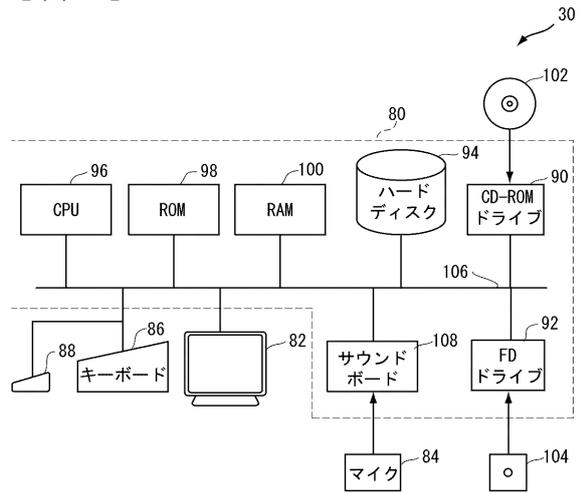
【図5】



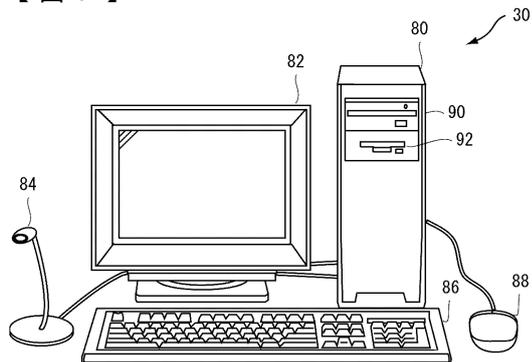
【図6】



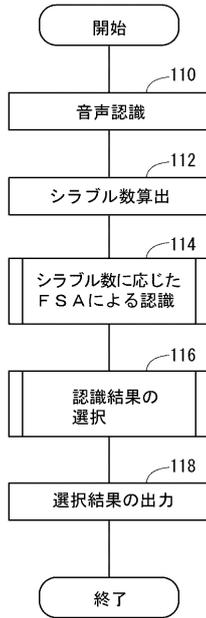
【図8】



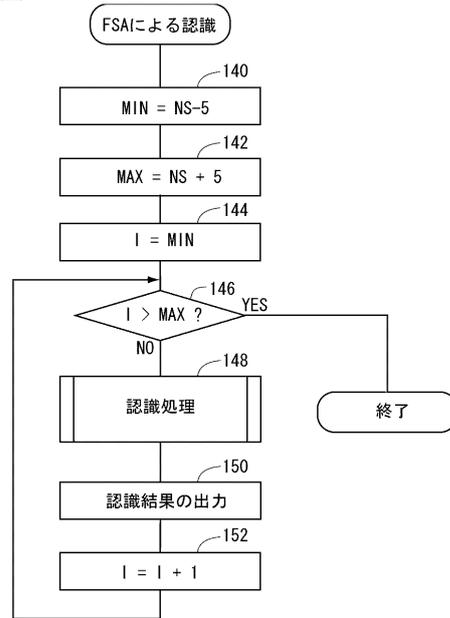
【図7】



【 図 9 】



【 図 10 】



【 図 11 】



---

フロントページの続き

審査官 荏原 雄一

(56)参考文献 特開昭57-086899(JP,A)

特開昭58-195895(JP,A)

鶴身玲典 他, "単語N-gramとネットワーク文法を併用した音声認識アルゴリズムの検討", 日本音響学会2002年度秋季研究発表会講演論文集 - I -, 2002年 9月26日, 3-9-8, p.145-146

大西茂彦 他, "文認識率の向上に向けたFSAとNGRAMの併用モデルによる大語彙連続音声認識", 日本音響学会2003年度春季研究発表会講演論文集 - I -, 2003年 3月18日, 3-Q-26, p.203-204

(58)調査した分野(Int.Cl., DB名)

G10L 15/00-15/28

JSTPlus(JDream2)