

(19) 日本国特許庁 (J P)

(12) 特 許 公 報 (B 2)

(11) 特許番号

第2886121号

(45) 発行日 平成11年(1999) 4月26日

(24) 登録日 平成11年(1999) 2月12日

(51) Int.Cl. <sup>6</sup>	識別記号	F I	
G 1 0 L 3/00	5 3 5 5 6 1	G 1 0 L 3/00	5 3 5 5 6 1 G

請求項の数4 (全 14 頁)

(21) 出願番号	特願平7-292685	(73) 特許権者	593118597 株式会社エイ・ティ・アール音声翻訳通信研究所 京都府相楽郡精華町大字乾谷小字三平谷5番地
(22) 出願日	平成7年(1995)11月10日	(72) 発明者	政瀧 浩和 京都府相楽郡精華町大字乾谷小字三平谷5番地 株式会社エイ・ティ・アール音声翻訳通信研究所内
(65) 公開番号	特開平9-134192	(72) 発明者	匂坂 芳典 京都府相楽郡精華町大字乾谷小字三平谷5番地 株式会社エイ・ティ・アール音声翻訳通信研究所内
(43) 公開日	平成9年(1997)5月20日	(74) 代理人	弁理士 青山 葆 (外2名)
審査請求日	平成7年(1995)11月10日	審査官	樫本 剛

最終頁に続く

(54) 【発明の名称】 統計的言語モデル生成装置及び音声認識装置

1

(57) 【特許請求の範囲】

【請求項1】 所定の話者の発声音声文を書き下した学習用テキストデータに基づいて、すべての語彙を品詞毎にクラスタリングされた品詞クラスに分類し、それらの品詞クラス間のバイグラムを初期状態の統計的言語モデルとして生成する生成手段と、  
上記生成手段によって生成された初期状態の統計的言語モデルに基づいて、単語の品詞クラスからの分離することができる第1の分離クラス候補と、1つの単語と1つの単語との結合、1つの単語と複数の単語の単語列との結合、複数の単語の単語列と1つの単語との結合、複数の単語の単語列と、複数の単語の単語列との結合を含む接続単語又は接続単語列の結合によって単語の品詞クラスから分離することができる第2の分離クラス候補とを探索する探索手段と、

2

上記探索手段によって探索された第1と第2の分離クラス候補に対して、次単語の予測の難易度を表わす所定のエントロピーを用いて、クラスを分離することによる当該エントロピーの減少量を計算する計算手段と、  
上記計算手段によって計算された上記第1と第2の分離クラス候補に対するエントロピーの減少量の中で最大のクラス分離を選択して、選択されたクラスの分離を実行することにより、品詞のバイグラムと可変長Nの単語のN-グラムとを含む統計的言語モデルを生成する分離手段と、  
上記分離手段によって生成された統計的言語モデルのクラス数が所定のクラス数になるまで、上記分離手段によって生成された統計的言語モデルを処理対象モデルとして、上記探索手段の処理と、上記計算手段の処理と、上記分離手段の処理とを繰り返すことにより、所定のクラ

10

3

ス数を有する統計的言語モデルを生成する制御手段とを備えたことを特徴とする統計的言語モデル生成装置。

【請求項 2】 入力される発声音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識手段を備えた音声認識装置において、上記音声認識手段は、品詞のバイグラムと可変長 N の単語の N - グラムとを含む統計的言語モデルを用いて音声認識することを特徴とする音声認識装置。

【請求項 3】 上記統計的言語モデルは、請求項 1 記載の統計的言語モデル生成装置によって生成されたことを特徴とする音声認識装置。

【請求項 4】 入力される発聲音声文の音声信号に基づいて上記発聲音声文の単語仮説を検出し尤度を計算することにより、連続的に音声認識する音声認識手段を備えた連続音声認識装置において、上記音声認識手段は、請求項 1 記載の統計的言語モデル生成装置によって生成された統計的言語モデルを参照して、終了時刻が等しく開始時刻が異なる同一の単語の単語仮説に対して、当該単語の先頭音素環境毎に、発声開始時刻から当該単語の終了時刻に至る計算された総尤度のうちの最も高い尤度を有する 1 つの単語仮説で代表させるように単語仮説の絞り込みを行うことを特徴とする連続音声認識装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、学習用テキストデータに基づいて統計的言語モデルを生成する統計的言語モデル生成装置、及び上記統計的言語モデルを用いて、入力される発聲音声文の音声信号を音声認識する音声認識装置に関する。

【0002】

【従来の技術】近年、連続音声認識装置において、その性能を高めるために言語モデルを用いる方法が研究されている。これは、言語モデルを用いて、次単語を予測し探索空間を削減することにより、認識率の向上および計算時間の削減の効果を狙ったものである。最近盛んに用いられている言語モデルとして N - グラム ( N - gram ) がある。これは、大規模なテキストデータを学習し、直前の N - 1 個の単語から次の単語への遷移確率を統計的に与えるものである。複数 L 個の単語列  $w_1^L = w_1, w_2, \dots, w_L$  の生成確率  $P(w_1^L)$  は次式で表される。

【0003】

【数 1】

$$P(w_1^L) = \prod_{t=1}^L P(w_t | w_{1+1-N}^{t-1})$$

【0004】ここで、 $w_t$  は単語列  $w_1^L$  のうち t 番目の 1 つの単語を表し、 $w_i^j$  は i 番目から j 番目の単語列を

4

表わす。上記数 1 において、確率  $P(w_t | w_{1+1-N}^{t-1})$  は、N 個の単語からなる単語列  $w_{1+1-N}^{t-1}$  が発声された後に単語  $w_t$  が発声される確率であり、以下同様に、確率  $P(A | B)$  は単語又は単語列 B が発声された後に単語 A が発声される確率を意味する。また、数 1 における「 $\prod$ 」は  $t = 1$  から L までの確率  $P(w_t | w_{1+1-N}^{t-1})$  の積を意味し、以下同様である。

【0005】N - グラムは極めて単純なものでありながら、構築の容易さ、統計的音響モデルとの相性の良さ、認識率向上や計算時間の短縮の効果が大きい等の理由で、連続音声認識には非常に有効である(例えば、従来文献 1「L. R. Bahl ほか、“A Maximum Likelihood Approach to Continuous Speech Recognition”, IEEE Transaction on Pattern Analysis and Machine Intelligence, pp. 179 - 190, 1983 年」、従来文献 2「P. C. Woodland ほか、“THE 1994 HTK Large Vocabulary Speech Recognition System”, Proceedings of ICASSP95’, Vol. 1, pp. 73 - 76, 1995 年」、従来文献 3「村上ほか、“単語の trigram を利用した文音声認識と自由発話認識への拡張”, 電子情報通信学会技術研究報告, SP93 - 127, pp71 - 78, 平成 6 年」参照。)

【0006】一般に、N - グラムの言語モデルは、N を大きくすると長い単語連鎖を取り扱うことにより次単語の精度は高くなるが、パラメータ数が多くなり、学習データ量が少ない場合は出現頻度の低い単語に信頼できる遷移確率を与えることはできない。例えば語彙数が 5, 000 語のとき、トライグラム ( trigram ) ( N = 3 ) の全ての単語の遷移組は  $(5, 000)^3 = 1, 250$  億であるから、信頼できる遷移確率を求めるためには、数千億単語以上からなる膨大なテキストデータが必要となる。これだけの膨大なテキストデータを集めるのは事実上不可能である。逆に、N を小さくすると、遷移確率の信頼性は高くなるが、短い単語連鎖しか取り扱うことができず、次単語の予測精度は低くなる。

【0007】

【発明が解決しようとする課題】この問題を解決するため、次のような方法が提案されている。

( 1 ) 補間による未学習遷移確率の推定方法

この方法は、例えば、Deleted Interpolation ( 削除補間法 ) ( 例えば、従来文献 4「F. Jelinek ほか、“Interpolated estimation of Markov Source Parameters from Sparse Data”, Proceedings of Workshop Pattern Recognitio

n in Practice, pp. 381 - 37, 1980年」参照。)や、Back-off Smoothing法(従来文献5「S.M.Katz, "Estimation of Probabilities from Sparse Data for the Language model Component of a Speech Recognizer", IEE E Transaction on Acoustics, Speech, and Signal Processing, Vol. ASSP-35, No. 3, pp. 400 - 401, 1987年3月」参照。)等に代表される方法で、小さいNのN-グラム(N-gram)の値で遷移確率を補間することにより、学習用テキストデータには存在しない単語遷移に対しても、遷移確率を与えることができる。しかしながら、出現頻度の低い単語に関しては信頼できる遷移確率を与えられない恐れがある。

L

$$P(w_1^L) = \prod_{t=1}^L P(w_t | c_t) \cdot P(c_t | c_{1-N+1}^{t-1})$$

【0010】ここで、 $c_t$ は単語 $w_t$ の属するクラスを表し、 $c_i^j$ はi番目からj番目のクラス列を表わす。上記数2で、 $P(c_t | c_{1-N+1}^{t-1})$ は、直前の(N-1)個の単語の属するクラスから次の単語の属するクラスへの遷移確率を表す。クラス数が50のとき、トライグラムの全てのクラス間の遷移の組は $50^3 = 125,000$ であるから、数十万単語程度と単語N-グラムに比べてかなり小規模なテキストデータで遷移確率が求められると考えられる。しかしながら、単語間の特有な接続関係を表現することができないので、次単語の予測精度は悪くなると考えられる。

【0011】本発明の目的は以上の問題点を解決し、従来例に比較して遷移確率の予測精度及び信頼性を改善することができる統計的言語モデルを生成することができる統計的言語モデル生成装置、及び、当該統計的言語モデルを用いて従来例に比較して高い音声認識率で音声認識することができる音声認識装置を提供することにある。

【0012】

【課題を解決するための手段】本発明に係る請求項1記載の統計的言語モデル生成装置は、所定の話者の発声音声を書き下した学習用テキストデータに基づいて、すべての語彙を品詞毎にクラスタリングされた品詞クラスに分類し、それらの品詞クラス間のバイグラムを初期状態の統計的言語モデルとして生成する生成手段と、上記生成手段によって生成された初期状態の統計的言語モデルに基づいて、単語の品詞クラスからの分離することができる第1の分離クラス候補と、1つの単語と1つの単語との結合、1つの単語と複数の単語の単語列との結

\*【0008】(2)クラスN-グラムによるパラメータ数の削減方法

この方法は、相互情報量に基づくクラスタリング(例えば、従来文献6「P.F.Brownほか, "Class-Based n-gram models of natural language", Computational Linguistics, Vol. 18, No. 4, pp. 467 - 479, 1992年」参照。)や、品詞(従来文献7「周ほか, "確率モデルによる日本語の大語彙連続音声認識", 情報処理学会, 第51回全国大会講演論文集, pp. 119 - 120, 平成7年」参照。)等によるクラス間のN-グラムを考えたもので、L個の単語の文生成確率 $P(w_1^L)$ は一般に次式で表される。

【0009】

【数2】

合、複数の単語の単語列と1つの単語との結合、複数の単語の単語列と、複数の単語の単語列との結合とを含む接続単語又は接続単語列の結合によって単語の品詞クラスから分離することができる第2の分離クラス候補とを検索する検索手段と、上記検索手段によって検索された第1と第2の分離クラス候補に対して、次単語の予測の難易度を表わす所定のエントロピーを用いて、クラスを分離することによる当該エントロピーの減少量を計算する計算手段と、上記計算手段によって計算された上記第1と第2の分離クラス候補に対するエントロピーの減少量の中で最大のクラス分離を選択して、選択されたクラスの分離を実行することにより、品詞のバイグラムと可変長Nの単語のN-グラムとを含む統計的言語モデルを生成する分離手段と、上記分離手段によって生成された統計的言語モデルのクラス数が所定のクラス数になるまで、上記分離手段によって生成された統計的言語モデルを処理対象モデルとして、上記検索手段の処理と、上記計算手段の処理と、上記分離手段の処理とを繰り返すことにより、所定のクラス数を有する統計的言語モデルを生成する制御手段とを備えたことを特徴とする。

【0013】本発明に係る請求項2記載の音声認識装置は、入力される発聲音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識手段を備えた音声認識装置において、上記音声認識手段は、品詞のバイグラムと可変長Nの単語のN-グラムとを含む統計的言語モデルを用いて音声認識することを特徴とする。

【0014】また、請求項3記載の音声認識装置においては、上記統計的言語モデルは、請求項1記載の統計的

言語モデル生成装置によって生成されたことを特徴とする。

【0015】本発明に係る請求項4記載の連続音声認識装置は、入力される発声音声文の音声信号に基づいて上記発聲音声文の単語仮説を検出し尤度を計算することにより、連続的に音声認識する音声認識手段を備えた連続音声認識装置において、上記音声認識手段は、請求項1記載の統計的言語モデル生成装置によって生成された統計的言語モデルを参照して、終了時刻が等しく開始時刻が異なる同一の単語の単語仮説に対して、当該単語の先頭音素環境毎に、発声開始時刻から当該単語の終了時刻に至る計算された総尤度のうちの最も高い尤度を有する1つの単語仮説で代表させるように単語仮説の絞り込みを行うことを特徴とする。

【0016】

【発明の実施の形態】以下、図面を参照して本発明に係る実施形態について説明する。図1に本発明に係る一実施形態の連続音声認識装置のブロック図を示す。本実施形態の連続音声認識装置は、公知のワン・パス・ピタビ復号化法を用いて、入力される発聲音声文の音声信号の特徴パラメータに基づいて上記発聲音声文の単語仮説を検出し尤度を計算して出力する単語照合部4を備えた連続音声認識装置において、単語照合部4からバッファメモリ5を介して出力される、終了時刻が等しく開始時刻が異なる同一の単語の単語仮説に対して、統計的言語モデル22を参照して、当該単語の先頭音素環境毎に、発声開始時刻から当該単語の終了時刻に至る計算された総尤度のうちの最も高い尤度を有する1つの単語仮説で代表させるように単語仮説の絞り込みを行う単語仮説絞込部6を備えたことを特徴とする。

【0017】ここで用いる統計的言語モデル22は、学習用テキストデータに基づいて言語モデル生成部20により生成されたものであって、統計的言語モデル22は、品詞クラス間のバイグラム( $N=2$ )を基本としたものであるが、単独で信頼できる単語は品詞クラスより分離させ、単独のクラスとして取り扱い、さらに、予測精度を向上させるため、頻出単語列に関してはそれらの単語を結合して一つのクラスとして取り扱い、長い単語連鎖の表現を可能にさせ、こうして、生成されたモデルは、品詞バイグラムと可変長単語 $N$ -グラムとの特徴を併せ持つ統計的言語モデルとなり、遷移確率の精度と信頼性とのバランスをとられたものであることを特徴とする。

【0018】まず、本実施形態において用いる可変長 $N$ -グラムの概念について以下に説明する。 $N$ -グラムは、( $N-1$ )重のマルコフモデルであり、これは、過去( $N-1$ )回の状態遷移を記憶するように単純(1重)マルコフモデルの各状態が分離されたものと解釈される。例として、図3にバイグラムをマルコフモデルとして図式化した状態遷移図を示し、図4にトライグラム

をマルコフモデルとして図式化した状態遷移図を示す。

【0019】図3においては、状態 $s_1$ においてシンボル $a$ を出力されたとき状態 $s_1$ のままであるが、状態 $s_1$ でシンボル $b$ を出力した状態 $s_2$ に遷移する。状態 $s_2$ でシンボル $b$ を出力したときは状態 $s_2$ のままであるが、状態 $s_2$ でシンボル $a$ を出力したとき状態 $s_1$ に戻る。図4のトライグラムは、バイグラムの状態 $s_1$ を状態 $s_{11}$ と状態 $s_{12}$ とに分離しかつ、状態 $s_2$ を状態 $s_{21}$ と状態 $s_{22}$ とに分離したものと考えられる。さらに、全ての状態の分離を進めることにより、より高次の $N$ -グラムとなる。

【0020】図5に示す可変長 $N$ -グラムは、単純マルコフモデルの状態を部分的に分離させたものである。すなわち、図3のバイグラムにおいて、状態 $s_2$ から、シンボル $a$ が出力される際に、続けてシンボル $b$ を出力する場合(これを $a\ b$ と表わし、シンボル $a\ b$ を出力するという。)、続けて $b$ 以外のシンボルを出力する場合(これを $a(/b)$ と表し、シンボル $a(/b)$ を出力するという。ここで、 $/$ は否定の意味を表しパー(上線)である。)とに分け、前者の場合、状態 $s_1$ から状態 $s_{12}$ に遷移させる一方、後者の場合、状態 $s_2$ から状態 $s_{11}$ に遷移させる。すなわち、前者の場合において、状態 $s_1$ から状態 $s_{12}$ へと分離させ、シンボル $a$ を出力する残りの遷移( $a(/b)$ )を状態 $s_{11}$ に残したものである。なお、このモデルにおいて、状態 $s_{11}$ でシンボル $a\ b$ を出力したとき状態 $s_{12}$ に遷移する一方、状態 $s_{11}$ でシンボル $a(/b)$ を出力したとき状態 $s_{11}$ のままである。また、状態 $s_{12}$ でシンボル $a\ b$ を出力したとき状態 $s_{12}$ のままである一方、状態 $s_{12}$ でシンボル $a(/b)$ を出力したとき状態 $s_{11}$ に遷移する。

【0021】このモデルは、複数の連続したシンボルを新しいシンボルとみなすことで、単純マルコフモデルの構造のまま、長い連鎖を表すことができるという特徴がある。同様の状態分離を繰り返すことで、局所的にさらに長い連鎖を表すことができる。これが可変長 $N$ -グラムである。すなわち、シンボルを単語とみなした言語モデルとしての可変長単語 $N$ -グラムは、単語列(1単語も含む)間のバイグラムと表される。

【0022】次いで、可変長 $N$ -グラムの動作について説明する。本実施形態で用いる統計的言語モデル22は、品詞クラスと単語との可変長 $N$ -グラムであり、次の3種類のクラス間のバイグラムとして表現する。

- (1) 品詞クラス(以下、第1のクラスという。)
- (2) 品詞クラスから分離した単語のクラス(以下、第2のクラスという。)
- (3) 接続単語が結合してできたクラス(以下、第3のクラスという。)

【0023】上記第1のクラスに属する単語は、主として出現頻度の小さいもので、単語単独で取り扱うよりも遷移確率の信頼性が高められる。また、第2のクラスに属する単語は、主として出現頻度が高いもので、単独で

取り扱っても十分な信頼性があり、さらに、接続単語が結合して上記第3のクラスに分類されることにより、可変長N-グラムとして動作し、次単語の予測精度が高められる。ただし、本実施形態において、接続する品詞クラスと品詞クラス、および、品詞クラスと単語の結合は\*

\* 考えない。複数L個の単語からなる文の生成確率  $P(w_1^L)$  は、次式で与えられる。

【0024】

【数3】

$$P(w_1^L) = \prod_{t=1}^L P(ws_t | c_t) \cdot P(c_t | c_{t-1})$$

【0025】ここで、 $ws_t$  は文章を上記のクラスに分類した時の、t番目の単語列(単独の単語も含める)を意味する。従って、 $P(ws_t | c_t)$  は、t番目のクラスがわかったときに単語列 $ws_t$ が出現する確率であり、 $P(c_t | c_{t-1})$  は1つ前の(t-1)番目のクラスから当該t番目のクラスの単語が出現する確率である。また、文章のKは単語列の個数を表し、K=Lである。従って、数3の は  $t=1$  からKまでの積である。

10 ここで、例として、次の7単語からなる発声音声文の文章を考える。

【0026】

【数4】「わたくし-村山-と-言-い-ま-す」

【0027】この文章の生成確率  $P(w_1^L)$  は、数3を用いて、次の式で与えられる。

【0028】

【数5】

$$P(w_1^L) = P(\text{わたくし} | \{\text{わたくし}\}) \cdot P(\{\text{わたくし}\}) \\ \cdot P(\text{村山} | \langle \text{固有名詞} \rangle) \cdot P(\langle \text{固有名詞} \rangle | \{\text{わたくし}\}) \\ \cdot P(\text{と} | \{\text{と}\}) \cdot P(\{\text{と}\} | \langle \text{固有名詞} \rangle) \\ \cdot P(\text{言います} | [\text{言います}]) \cdot P([\text{言います}] | \{\text{と}\})$$

【0029】ただし、 $\langle \rangle$  ,  $\{ \}$  ,  $[ \ ]$  はそれぞれ、第1のクラス、第2のクラス、第3のクラスに属していることを表す。ただし、各単語および単語列は次のように属している。

(3) 「言います」は動詞と、動詞の接尾辞と、助動詞と、助動詞の接尾辞との組み合わせであり、第3のクラスに属する。ここで、第2と第3のクラスにおいて、単語とクラスの出現頻度は等しいので、 $P(\text{わたくし} | \{\text{わたくし}\}) = 1$ 、 $P(\text{と} | \{\text{と}\}) = 1$ 、 $P(\text{言います} | [\text{言います}]) = 1$  であり、従って、上記数5は次の式のようになる。

【0030】

30 【数6】

$$P(w_1^L) = P(\text{わたくし}) \\ \cdot P(\text{村山} | \langle \text{固有名詞} \rangle) \cdot P(\langle \text{固有名詞} \rangle | \text{わたくし}) \\ \cdot P(\text{と} | \langle \text{固有名詞} \rangle) \\ \cdot P(\text{言います} | \text{と})$$

【0031】次いで、本実施形態で用いる可変長N-グラムである統計的言語モデル22を生成するための言語モデル生成処理について参照して説明する。本実施形態で用いる統計的言語モデル22は、品詞クラスのバイグラムを初期状態とし、エントロピーの最小化の基準によるクラス分離という形で生成される。エントロピーの減少は正になることが保証されており、クラス分離によって、学習用テキストデータに関してエントロピーは単調に減少する。ここで用いるエントロピーは、一般には、

「あいまいさ」の尺度を表わすものであり、言語モデルにおいて、エントロピーが小さいことは、言語としてあいまいさが小さく、次の単語の予測が容易であることを意味する。すなわち、エントロピーは次単語の予測の難易度を表わす。yという条件のもとでのxの確率である条件付き確率  $P(x | y)$  のエントロピー  $H(X | Y)$  は次式で表される。

【0032】

【数7】

$$H(X | Y) = - \sum_y P(y) \sum_x P(x | y) \log_2 P(x | y)$$

【0033】従って、上記数7に基づいて、本実施形態で用いるエントロピーは次式で計算される。

【0034】

【数8】

$$\begin{aligned}
& H(\{c_i\}) \\
& = -\sum_i P(c_i) \\
& \quad \cdot \sum_k P(w_k | c_j) \cdot P(c_j | c_i) \log_2 \{P(w_k | c_j) \cdot P(c_j | c_i)\}
\end{aligned}$$

ここで、 $w_k \quad c_j$

【0035】図6は、言語モデル生成部20によって実行される言語モデル生成処理の詳細を示すフローチャートであり、以下、図6を参照して当該処理について説明する。まず、ステップS1では、所定の話者の発声音声文を書き下した学習用テキストデータに含まれる全語彙を品詞クラス(ここで、品詞クラスとは、品詞毎にクラスタリングされたクラスをいう。)に分類し、それらのクラス間のバイグラムを初期状態の統計的言語モデルとする。次いで、次のステップS2乃至S4でクラスの分離を行う。すなわち、ステップS2で、クラス分離することが可能な分離クラス候補を検索することによりリストアップを行う。ここでは、次の2種類のクラス分離を考える。

(1) 単語の品詞クラスからの分離(以下、第1のクラス分離という。)、(2) 接続単語又は接続単語列の結合によるクラス分離(以下、第2のクラス分離という。)。ここで、接続単語又は接続単語列の結合とは、接続する(時間的に隣接して入力される)1つの単語と1つの単語との結合、1つの単語と複数の単語の単語列との結合、複数の単語の単語列と1つの単語との結合、複数の単語の単語列と、複数の単語の単語列との結合を含む。

【0036】前者の単語の品詞クラスからの分離においては、当初品詞クラスに属している単語が、そのクラスから分離し、分離した単語は、その単語で単独のクラスを形成する。

【0037】

【数9】  $c = \{w_x\} + c \setminus \{w_x\}$

ここで、 $w_x \quad c$

【0038】ここで、 $c \setminus \{w_x\}$ はクラスcから単語 $w_x$ のクラスを除いたクラスであることを意味し、単語 $w_x$ はクラスcに属している。従って、数9の意\*

$$\begin{aligned}
& H \\
& = H(\{c_i\}) - H(\{c_i \setminus c\} + \{w_x\} + \{c \setminus w_x\})
\end{aligned}$$

【数12】

$$\begin{aligned}
& H \\
& = H(\{c_i\}) - H(\{c_i \setminus w_x\} + \{w_x, w_y\} + \{w_x, /w_y\})
\end{aligned}$$

【0044】ここで、数11及び数12において、 $H(\{c_i\})$ は元のすべての品詞クラス $c_i$ についてのエントロピーであり、数11において $H(\{c_i \setminus c\} + \{w_x\} + \{c \setminus w_x\})$ は元のすべての品詞クラス $c_i$ から単語 $w_x$ のクラスを分離したときのエントロピーであり、数11のHは単語 $w_x$ のクラスを分離したと

\* 味するところは、例えば、名詞のクラスcは、「机」という単語 $w_x$ のクラス $\{w_x\}$ と、「机」という単語 $w_x$ のクラス $\{w_x\}$ をクラスcから除いたクラスとに分離することを意味する。

10

【0039】後者の接続単語又は接続単語列の結合によるクラス分離においては、既に初期クラスより分離されている単語クラス及び単語列クラスについて、接続した2クラス間の結合を考える。結合した単語列は、その単語列で単独のクラスを形成する。

【0040】

【数10】  $\{w_x\} \quad \{w_x, w_y\} + \{w_x, /w_y\}$

20

【0041】ここで、 $\{w_x, w_y\}$ は接続単語列 $w_x, w_y$ のクラスを表し、 $\{w_x, /w_y\}$ は単語 $w_x$ の次に単語 $w_y$ が後続しない単語 $w_x$ のクラスを表わす。すなわち、 $/w_y$ は単語 $w_y$ 以外の単語を表わす。数10の意味するところは、例えば、「机」という単語のクラス $\{w_x\}$ は、「机の」という単語列のクラス $\{w_x, w_y\}$ と、「机の」以外の例えば「机は」、「机が」などの単語列のクラス $\{w_x, /w_y\}$ とに分離することを意味する。上記数10は、単語の結合に関する式であるが、単語列と単語の結合、および、単語列と単語列との結合も同様に表される。従って、第2のクラス分離では、これらのクラス分離を含む。

30

【0042】次いで、ステップS3で、ステップS2でリストアップされた上記第1と第2の分離クラス候補に対して次の数11及び数12を用いてエントロピー減少量を計算する。ここで、上記第1のクラス分離である初期クラスの分離に対して数11を用いる一方、上記第2のクラス分離である接続単語又は接続単語列の結合によるクラス分離に対して数12を用いる。

【0043】

【数11】

きのエントロピーの減少量である。また、数12において $H(\{c_i \setminus w_x\} + \{w_x, w_y\} + \{w_x, /w_y\})$ は、元のすべての品詞クラス $c_i$ から単語列 $\{w_x, w_y\}$ のクラスを分離したときのエントロピーであり、数12のHは単語列 $\{w_x, w_y\}$ のクラスを分離したときのエントロピーの減少量である。

50

【0045】次いで、ステップS4においては、ステップS2でリストアップされたすべての分離クラス候補の中で、ステップS3で計算したエントロピー減少量Hを最大にするクラスのみを実際にクラス分離する。そして、ステップS5で分離クラス数が所定のしきい値の所望分離クラス数(例えば、500、1000など)以上になったか否かを判断し、なっていないときは、ステップS2に戻って上記の処理を繰り返す。一方、ステップS5で所望分離クラス数以上になっているときは、ステップS6で、得られた統計的言語モデル22をメモリに格納した後、当該言語モデル生成処理を終了する。この言語モデル生成処理のアルゴリズムは、品詞間、および、品詞と単語間の結合は行なわないため、生成完了時点では、品詞のバイグラムと可変長単語のN-グラムの特徴を併せた統計的言語モデル22となる。

【0046】図1において、単語照合部4に接続され、例えばハードディスクメモリに格納される音素HMM11は、各状態を含んで表され、各状態はそれぞれ以下の情報を有する。

- (a) 状態番号
- (b) 受理可能なコンテキストクラス
- (c) 先行状態、及び後続状態のリスト
- (d) 出力確率密度分布のパラメータ
- (e) 自己遷移確率及び後続状態への遷移確率

なお、本実施形態において用いる音素HMM11は、各分布がどの話者に由来するかを特定する必要があるため、所定の話者混合HMMを変換して生成する。ここで、出力確率密度関数は3次元の対角共分散行列をもつ混合ガウス分布である。

【0047】また、単語照合部4に接続され、例えばハードディスクに格納される単語辞書12は、音素HMM11の各単語毎にシンボルで表した読みを示すシンボル列を格納する。

【0048】図1において、話者の発声音声はマイクロホン1に入力されて音声信号に変換された後、特徴抽出部2に入力される。特徴抽出部2は、入力された音声信号をA/D変換した後、例えばLPC分析を実行し、対数パワー、16次ケプストラム係数、対数パワー及び16次ケプストラム係数を含む3次元の特徴パラメータを抽出する。抽出された特徴パラメータの時系列はバッファメモリ3を介して単語照合部4に入力される。

【0049】単語照合部4は、ワン・パス・ビタビ復号化法を用いて、バッファメモリ3を介して入力される特徴パラメータのデータに基づいて、音素HMM11と単語辞書12とを用いて単語仮説を検出し尤度を計算して出力する。ここで、単語照合部4は、各時刻の各HMMの状態毎に、単語内の尤度と発声開始からの尤度を計算する。尤度は、単語の識別番号、単語の開始時刻、先行単語の違い毎に個別にもつ。また、計算処理量の削減のために、音素HMM11及び単語辞書12とに基づいて

計算される総尤度のうちの低い尤度のグリッド仮説を削減する。単語照合部4は、その結果の単語仮説と尤度の情報を発声開始時刻からの時間情報(具体的には、例えばフレーム番号)とともにバッファメモリ5を介して単語仮説絞込部6に出力する。

【0050】単語仮説絞込部6は、単語照合部4からバッファメモリ5を介して出力される単語仮説に基づいて、統計的言語モデル22を参照して、終了時刻が等しく開始時刻が異なる同一の単語の単語仮説に対して、当該単語の先頭音素環境毎に、発声開始時刻から当該単語の終了時刻に至る計算された総尤度のうちの最も高い尤度を有する1つの単語仮説で代表させるように単語仮説の絞り込みを行った後、絞り込み後のすべての単語仮説の単語列のうち、最大の総尤度を有する仮説の単語列を認識結果として出力する。本実施形態においては、好ましくは、処理すべき当該単語の先頭音素環境とは、当該単語より先行する単語仮説の最終音素と、当該単語の単語仮説の最初の2つの音素とを含む3つの音素並びをいう。

【0051】例えば、図2に示すように、(i-1)番目の単語 $W_{i-1}$ の次に、音素列 $a_1, a_2, \dots, a_n$ からなるi番目の単語 $W_i$ がくるときに、単語 $W_{i-1}$ の単語仮説として6つの仮説 $W_a, W_b, W_c, W_d, W_e, W_f$ が存在している。ここで、前者3つの単語仮説 $W_a, W_b, W_c$ の最終音素は/x/であるとし、後者3つの単語仮説 $W_d, W_e, W_f$ の最終音素は/y/であるとする。終了時刻 $t_e$ と先頭音素環境が等しい仮説(図2では先頭音素環境が"x/a<sub>1</sub>/a<sub>2</sub>"である上から3つの単語仮説)のうち総尤度が最も高い仮説(例えば、図2において1番上の仮説)以外を削除する。なお、上から4番めの仮説は先頭音素環境が異なるため、すなわち、先行する単語仮説の最終音素がxではなくyであるので、上から4番めの仮説を削除しない。すなわち、先行する単語仮説の最終音素毎に1つのみ仮説を残す。図2の例では、最終音素/x/に対して1つの仮説を残し、最終音素/y/に対して1つの仮説を残す。

【0052】以上の実施形態においては、当該単語の先頭音素環境とは、当該単語より先行する単語仮説の最終音素と、当該単語の単語仮説の最初の2つの音素とを含む3つの音素並びとして定義されているが、本発明はこれに限らず、先行する単語仮説の最終音素と、最終音素と連続する先行する単語仮説の少なくとも1つの音素とを含む先行単語仮説の音素列と、当該単語の単語仮説の最初の音素を含む音素列とを含む音素並びとしてもよい。

【0053】以上の実施形態において、特徴抽出部2と、単語照合部4と、単語仮説絞込部6と、言語モデル生成部20とは、例えば、デジタル電子計算機で構成され、バッファメモリ3, 5は例えばハードディスクメモリで構成され、音素HMM11と単語辞書12と学習用

テキストデータ 2 1 と統計的言語モデル 2 2 とは、例えばハードディスクメモリなどの記憶装置に記憶される。

【0054】以上実施形態においては、単語照合部 4 と単語仮説絞込部 6 とを用いて音声認識を行っているが、本発明はこれに限らず、例えば、音素 HMM 1 1 を参照する音素照合部と、例えば One Pass DP アルゴリズムを用いて統計的言語モデル 2 2 を参照して単語の音声認識を行う音声認識部とで構成してもよい。

【0055】

【実施例】本発明者は、本実施形態で用いる統計的言語モデル 2 2 の性能を確認するため、パープレキシティおよびパラメータ数について従来の単語 N - グラムとの比較を行った。実験に用いたデータは本出願人が所有する自然発話旅行会話データベース（従来文献 8 「Morimoto ほか，“A Speech and Language Database for Speech Translation Research”，ICSLP，pp1791 - 1794，1994 年」参照。）であって、846 対話、354，700 語から構成される。このうち、統計的言語モデル 2 2 を生成するための学習用テキストデータ（トレーニングセットともいう。）として、828 対話、347，299 語を使用し、残りのデータをテスト用テキストデータ（テストセットともいう。）とした。本実施形態に係る統計的言語モデル 2 2 は、初期クラスを活用形も含めた 80 品詞とし、1000 個まで分離を行い、100 個おきにデータを採取した。また、本実施形態に係る統計的言語モデル 2 2 と、単語 N - グラムとともに、未知単語遷移に対する対策として、クラスおよび単語の遷移確率を削除補間法（従来文献 4 参照。）によって補間し、テストセットにおいて、未知語が出現したときは、所定の固定値（例えば、 $7.0 \times 10^{-6}$ ）を与えた。ここで、本発明に係る統計的言語モデル 2 2 を評価するために、パープレキ\*

各言語モデルの性能比較

	バイグラム	トライグラム	本実施形態（分離クラス数）		
			0	500	1000
テストセット パープレキシティ	20.31	16.96	41.68	17.61	16.75
トレーニングセット パープレキシティ	13.50	5.99	48.77	18.77	15.05
パラメータ数(1)	$4.10 \times 10^7$	$2.62 \times 10^{11}$	$1.28 \times 10^4$	$3.43 \times 10^5$	$1.17 \times 10^6$
パラメータ数(2)	52, 244 7, 991	27, 830	165, 139 43, 075		

【0063】ここで、パラメータ数(1)は全クラス

\* シティを用いる。例えば、複数 n 個の単語からなる長い単語列  $w_1^n = w_1 w_2 \dots w_n$  があるときのエントロピー  $H(n)$  は次式で表される。

【0056】

【数13】

$$H(n) = - (1/n) \cdot \log_2 P(w_1^n)$$

【0057】ここで、 $P(w_1^n)$  は単語列  $w_1^n$  の生成確率であり、パープレキシティ  $PP(n)$  は次式で表される。

【0058】

【数14】  $PP(n) = 2^{H(n)}$

【0059】ここで、単語列としてテスト用テキストデータを用いたときのパープレキシティをテストセットパープレキシティといい、単語列として学習用テキストデータを用いたときのパープレキシティをトレーニングセットパープレキシティという。

【0060】当該実験結果におけるテストセットパープレキシティの値の変化の様子を図 7 に示す。図 7 から明らかのように、分離クラス数が増加するに従って、テストセットパープレキシティは減少し、分離クラス数が 200 で単語バイグラムと、分離クラス数が 600 で単語トライグラムと同程度の値となることが分かる。分離クラス数が 600 以上になると、パープレキシティの減少の割合が極端に小さくなるため、分離クラス 600 程度で、本実施形態の統計的言語モデル 2 2 が最も有効に働いていると考えられる。従って、本実施形態の統計的言語モデル 2 2 は単語バイグラム以上、単語トライグラムと同程度の予測精度の言語モデルと考えられる。

【0061】表 1 にまた、分離クラス数が 0, 500, 1000 の時のパープレキシティの値、およびパラメータ数を示す。

【0062】

【表 1】

50 (単語) の遷移の組み合わせ数を意味し、パラメータ数

(2)は、トレーニングセットにおいて実際に存在するクラス(単語)遷移の組み合わせ数を意味する。表1より、本実施形態の統計的言語モデル22は、テストセットとトレーニングセットとのパープレキシティの差が、単語バイグラム及び単語トライグラムと比較して非常に小さいことが分かる。また、パラメータ数は、1000クラス分離した時でも、単語バイグラムよりも少なく、単語トライグラムよりもはるかに少ない。したがって、本実施形態の統計的言語モデル22は、与えられたパラメータで言語特徴を効率的に表現できる優れた言語モデルであると言える。従って、当該統計的言語モデル22は従来の単語バイグラム、単語トライグラムよりも信頼性が高い言語モデルであると考えられる。

【0064】また、本実施形態の統計的言語モデル22の信頼性を確認するため、学習単語数を変化させてテストセットパープレキシティの値の変化を調べた結果を図8に示す。この図8から明らかのように、全ての学習セット(約35万語)を用いたときは、単語バイグラム \* 実験条件

分析条件 サンプリング周波数：12KHz，  
ハミング窓：20ms，  
フレーム周期：10ms

使用パラメータ 16次LPCケプストラム+16次 ケプストラム  
+logパワー+ logパワー

音響モデル HM網の男女別不特定話者モデル  
400状態，5混合

【0067】表2において、HM網の男女別不特定話者モデルについては、従来文献9「小坂ほか，“話者混合SSSによる不特定話者音声認識”，日本音響学会講演論文集，2-5-9，pp135-136，平成4年」に開示されている。この実験では、単語グラフを用いた連続音声認識法を用いて音響モデルおよび言語モデルを連続音声認識装置に適用した。また、認識の対象は、  
正解単語含有率

\* と、本実施形態の統計的言語モデル22(200クラス)(カッコ内の数字は分離クラス数を表す、以下同様である。)とは、ほぼ同じパープレキシティ値であるが、学習単語数を減少させても当該統計的言語モデル22のパープレキシティの増加は比較的小さく、単語バイグラムよりも値が低くなる事が分かる。同様に、単語トライグラムと、当該統計的言語モデル22(600クラス)とを比較しても、学習単語数が減少すると、当該統計的言語モデル22の方が低いパープレキシティを呈する。

【0065】次いで、本発明者は、本実施形態の統計的言語モデル22を図1の連続音声認識装置に適用し、統計的言語モデル22の効果を確認した。音素認識の実験条件を表2に示す。また、音響をパラメータもあわせて表2に示す。

【0066】  
【表2】

計的言語モデル22のトレーニングセット中の16対話であり、学習に用いられていないテストセットは18対話である。各言語モデルで尤度1位の文認識候補の正解単語含有率を表3に示す。

【0068】  
【表3】

		バイグラム	本実施形態(分離クラス数)	
			0	500
辞書サブセット	テストセット	71.4	67.3	72.2
	トレーニングセット	69.7	69.4	63.4
辞書フルセット	テストセット	-	57.1	58.4
	トレーニングセット	-	54.6	56.0

【0069】表3において、辞書サブセットは認識対象 50 に含まれる単語のみを辞書に登録したもの(750

語)、辞書フルセットは、統計的言語モデルの生成のための学習に用いた全単語よりなる辞書(6,400語)を表す。ただし、従来の単語バイグラムは、メモリ容量と計算時間の都合上で、辞書フルセットの辞書の認識は、今回の実験では計算を行っていない。この場合は、言い換えれば、大容量のメモリと莫大な処理時間が必要である。

【0070】テストセットに関しては、パープレキシティの低い順、すなわち本実施形態の統計的言語モデル22(0クラス)単語バイグラム本実施形態の統計的言語モデル22(500クラス)の順で正解単語含有率が良くなっており、本実施形態の統計的言語モデル22(500クラス)は、単語のバイグラムよりも若干ではあるが正解単語含有率が向上している。トレーニングセットに関しては、本実施形態の統計的言語モデル22(500クラス)は単語バイグラムよりも高いパープレキシティであるが、正解単語含有率は高くなっている。また、本実施形態の統計的言語モデル22はパラメータ数が少ないので、大語彙の認識への拡張が容易ある。したがって、本実施形態の統計的言語モデル22は連続音声認識に対しても単語バイグラムより有効な言語モデルであると考えられる。

【0071】以上説明したように、N-グラムの精度・信頼性の向上を目的とした可変長N-グラムの統計的言語モデル22の生成装置及びこれを用いた連続音声認識装置を実現することができる。当該統計的言語モデル22は、品詞バイグラムを初期状態とし、品詞クラスからの単語分離、および、接続単語の結合という、2種類の状態分離を行なうことにより生成されるもので、品詞バイグラムと可変長単語N-グラムの特徴を併せ持つモデルである。当該統計的言語モデル22の評価実験の結果、当該統計的言語モデル22は、単語バイグラム以上、単語トライグラムと同等のパープレキシティを、はるかに少ないパラメータで実現できることが分かり、目的とした性能が実現されていることが確認できた。また、連続音声認識に適用した結果、単語バイグラムと同じ程度の正解単語含有率を得ることができた。当該統計的言語モデル22は少ないパラメータで実現できるため、大語彙の音声認識にも容易に拡張可能である。

【0072】従って、遷移確率の予測精度及び信頼性を改善することができる統計的言語モデル22を生成することができる統計的言語モデル生成装置を提供することができるとともに、当該統計的言語モデル22を用いてより高い音声認識率で連続的に音声認識することができる連続音声認識装置を提供することができる。

【0073】

【発明の効果】以上詳述したように本発明に係る請求項1記載の統計的言語モデル生成装置によれば、所定の話者の発声音声を書き下した学習用テキストデータに基づいて、すべての語彙を品詞毎にクラスターリングされた

品詞クラスに分類し、それらの品詞クラス間のバイグラムを初期状態の統計的言語モデルとして生成する生成手段と、上記生成手段によって生成された初期状態の統計的言語モデルに基づいて、単語の品詞クラスからの分離することができる第1の分離クラス候補と、1つの単語と1つの単語との結合、1つの単語と複数の単語の単語列との結合、複数の単語の単語列と1つの単語との結合、複数の単語の単語列と、複数の単語の単語列との結合とを含む接続単語又は接続単語列の結合によって単語の品詞クラスから分離することができる第2の分離クラス候補とを探索する探索手段と、上記探索手段によって探索された第1と第2の分離クラス候補に対して、次単語の予測の難易度を表わす所定のエントロピーを用いて、クラスを分離することによる当該エントロピーの減少量を計算する計算手段と、上記計算手段によって計算された上記第1と第2の分離クラス候補に対するエントロピーの減少量の中で最大のクラス分離を選択して、選択されたクラスの分離を実行することにより、品詞のバイグラムと可変長Nの単語のN-グラムとを含む統計的言語モデルを生成する分離手段と、上記分離手段によって生成された統計的言語モデルのクラス数が所定のクラス数になるまで、上記分離手段によって生成された統計的言語モデルを処理対象モデルとして、上記探索手段の処理と、上記計算手段の処理と、上記分離手段の処理とを繰り返すことにより、所定のクラス数を有する統計的言語モデルを生成する制御手段とを備える。従って、遷移確率の予測精度及び信頼性を改善することができる統計的言語モデルを生成することができる。また、当該統計的言語モデルは少ないパラメータで実現できるため、大語彙の音声認識にも容易に拡張可能であるという特有の利点を有する。

【0074】本発明に係る請求項2記載の音声認識装置においては、入力される発声音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識手段を備えた音声認識装置において、上記音声認識手段は、品詞のバイグラムと可変長Nの単語のN-グラムとを含む統計的言語モデルを用いて音声認識する。従って、遷移確率の予測精度及び信頼性を改善することができる統計的言語モデルを用いて音声認識するので、より高い音声認識率で音声認識することができる音声認識装置を提供することができる。

【0075】また、請求項3記載の音声認識装置においては、上記統計的言語モデルは、請求項1記載の統計的言語モデル生成装置によって生成された。従って、遷移確率の予測精度及び信頼性を改善することができる統計的言語モデルを用いて音声認識するので、より高い音声認識率で音声認識することができる音声認識装置を提供することができる。

【0076】本発明に係る請求項4記載の連続音声認識装置は、入力される発声音声文の音声信号に基づいて上

記発声音声文の単語仮説を検出し尤度を計算することにより、連続的に音声認識する音声認識手段を備えた連続音声認識装置において、上記音声認識手段は、請求項 1 記載の統計的言語モデル生成装置によって生成された統計的言語モデルを参照して、終了時刻が等しく開始時刻が異なる同一の単語の単語仮説に対して、当該単語の先頭音素環境毎に、発声開始時刻から当該単語の終了時刻に至る計算された総尤度のうちの最も高い尤度を有する 1 つの単語仮説で代表させるように単語仮説の絞り込みを行う。すなわち、先行単語毎に 1 つの単語仮説で代表させる従来技術の単語ペア近似法に比較して、単語の先頭音素の先行音素（つまり、先行単語の最終音素）が等しいものをひとまとめに扱うために、単語仮説数を削減することができ、近似効果は大きい。特に、語彙数が増加した場合において削減効果が大きい。従って、当該連続音声認識装置を、間投詞の挿入や、言い淀み、言い直しが頻繁に生じる自然発話の認識に用いた場合であっても、単語仮説の併合又は分割に要する計算コストは従来例に比較して小さくなる。すなわち、音声認識のために必要な処理量が小さくなり、それ故、音声認識のための記憶装置において必要な記憶容量は小さくなる一方、処理量が小さくなるので音声認識のための処理時間を短縮することができる。さらに、遷移確率の予測精度及び信頼性を改善することができる統計的言語モデルを用いて音声認識するので、より高い音声認識率で連続的に音声認識することができる連続音声認識装置を提供することができる。

【図面の簡単な説明】

【図 1】 本発明に係る一実施形態である連続音声認識 \*

\* 装置のブロック図である。

【図 2】 図 1 の連続音声認識装置における単語仮説絞込部 6 の処理を示すタイミングチャートである。

【図 3】 バイグラムの統計的言語モデルを示す状態遷移図である。

【図 4】 トライグラムの統計的言語モデルを示す状態遷移図である。

【図 5】 図 1 の連続音声認識装置において用いる可変長 N - グラムの下のモデルを示す状態遷移図である。

10 【図 6】 図 1 の言語モデル生成部 2 0 によって実行される言語モデル生成処理を示すフローチャートである。

【図 7】 図 1 の言語モデル生成部 2 0 によって生成される統計的言語モデルにおける分離クラス数に対するテストセットパープレキシティを示すグラフである。

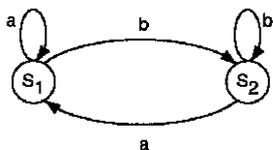
【図 8】 図 1 の言語モデル生成部 2 0 によって生成される統計的言語モデルにおける学習データの単語数に対するテストセットパープレキシティを示すグラフである。

【符号の説明】

- 20 1 ... マイクロホン、
- 2 ... 特徴抽出部、
- 3, 5 ... バッファメモリ、
- 4 ... 単語照合部、
- 6 ... 単語仮説絞込部、
- 1 1 ... 音素 HMM、
- 1 2 ... 単語辞書、
- 2 0 ... 言語モデル生成部、
- 2 1 ... 学習用テキストデータ、
- 2 2 ... 統計的言語モデル。

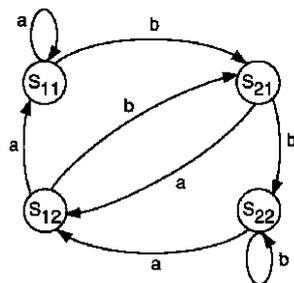
【図 3】

バイグラム



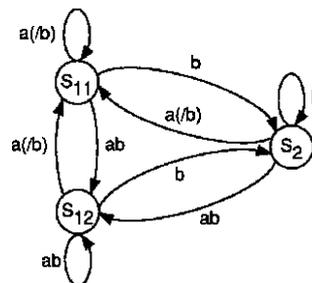
【図 4】

トライグラム

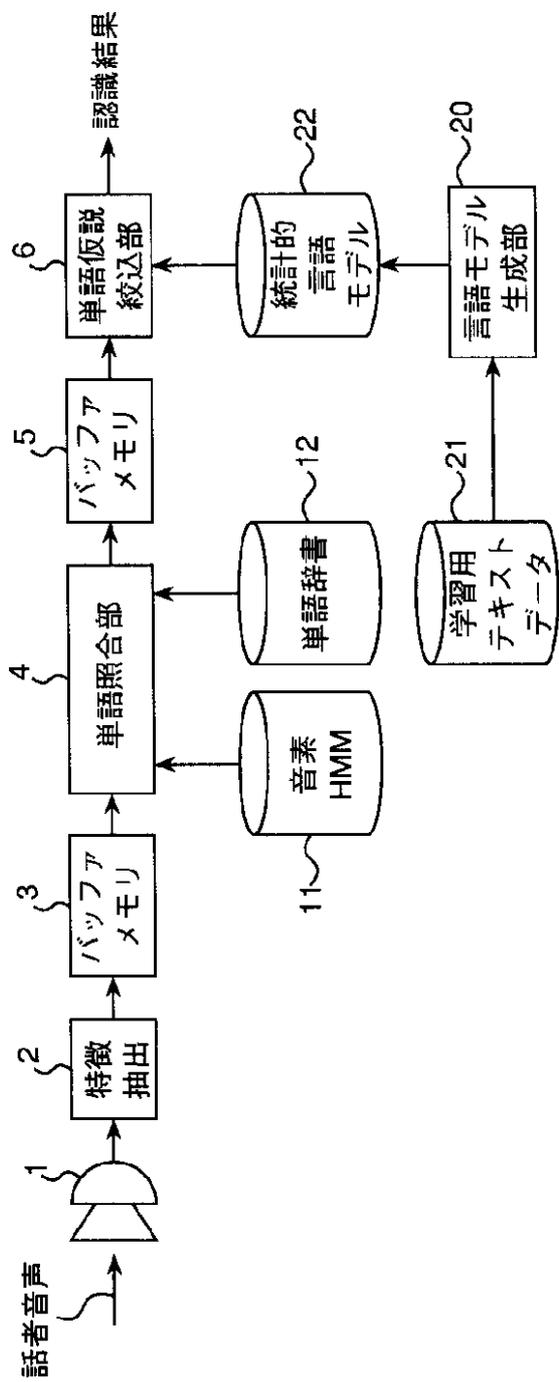


【図 5】

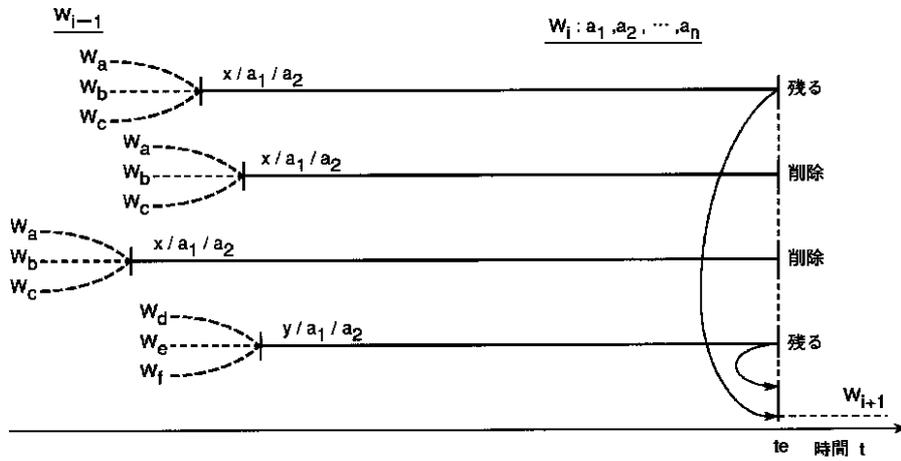
可変長 N-グラム



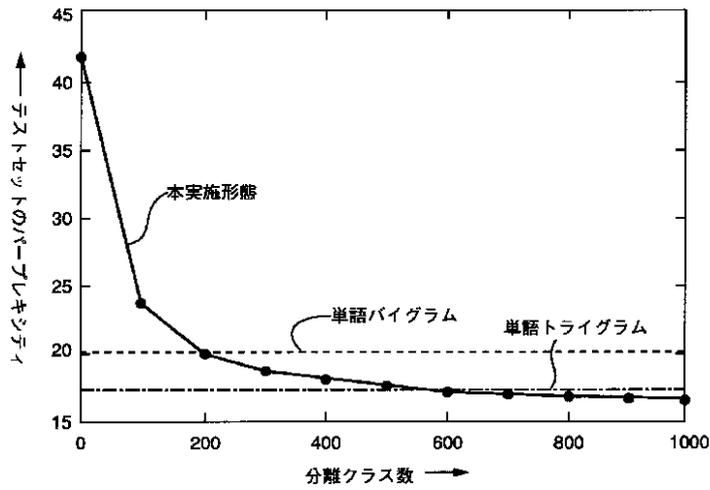
【図 1】



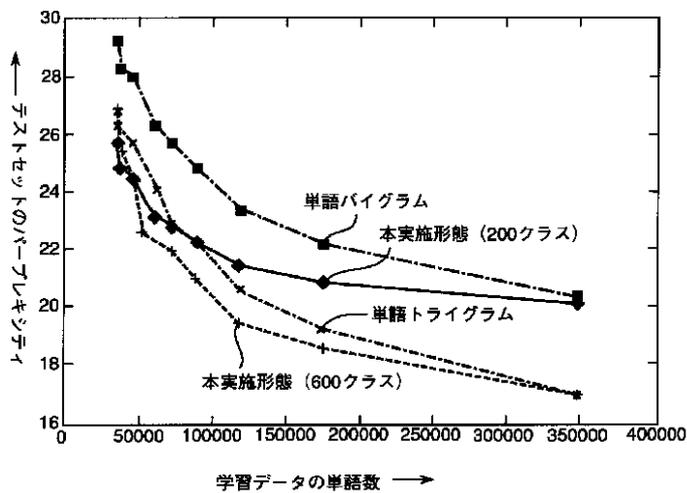
【図 2】



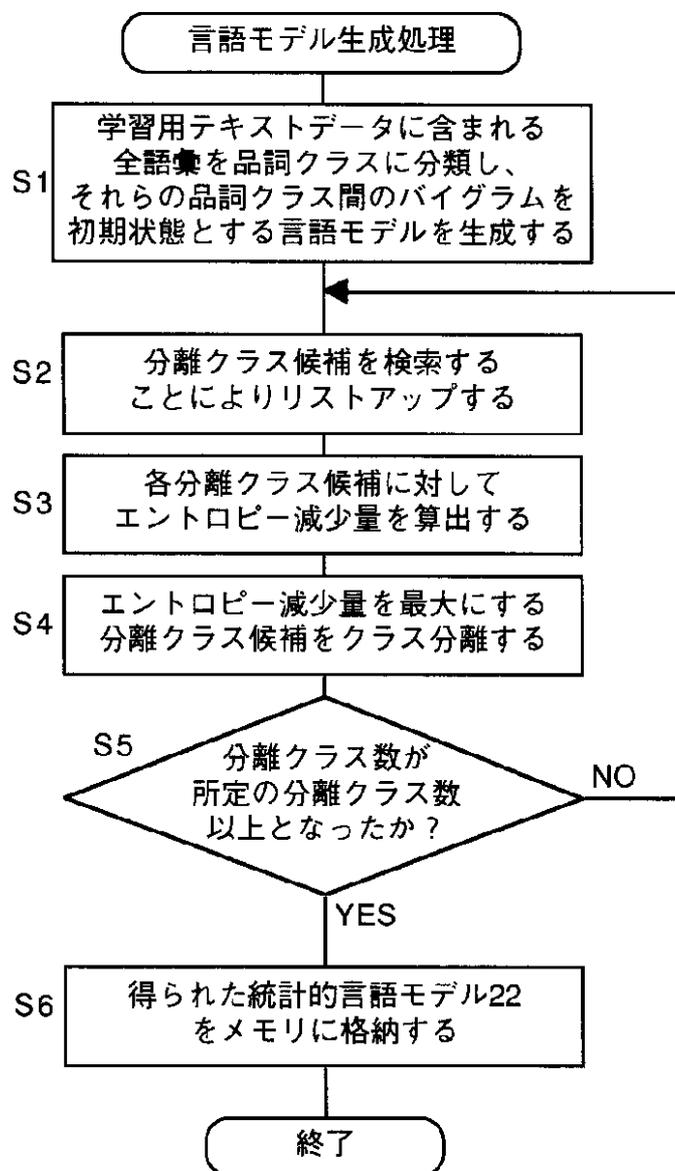
【図 7】



【図 8】



【図 6】



フロントページの続き

(72)発明者 松永 昭一  
 京都府相楽郡精華町大字乾谷小字三平谷  
 5番地 株式会社エイ・ティ・アール音  
 声翻訳通信研究所内

(56)参考文献 特開 平 5 - 108704 ( J P , A )  
 特開 平 5 - 250405 ( J P , A )  
 日本音響学会講演論文集 (平成 8 年 3  
 月) 1 - P - 17 , p . 195 ~ 196

(58)調査した分野(Int.Cl.<sup>6</sup>, D B 名)  
 G10L 3/00 535  
 G10L 3/00 561  
 J I C S T ファイル ( J O I S )