

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4631076号  
(P4631076)

(45) 発行日 平成23年2月16日(2011.2.16)

(24) 登録日 平成22年11月26日(2010.11.26)

(51) Int.Cl. F I  
**G 1 O L 15/06 (2006.01)** G 1 O L 15/06 3 O O D  
 G 1 O L 15/06 4 O O V

請求項の数 5 外国語出願 (全 10 頁)

<p>(21) 出願番号 特願2004-318208 (P2004-318208)</p> <p>(22) 出願日 平成16年11月1日(2004.11.1)</p> <p>(65) 公開番号 特開2006-126730 (P2006-126730A)</p> <p>(43) 公開日 平成18年5月18日(2006.5.18)</p> <p>審査請求日 平成19年10月12日(2007.10.12)</p> <p>(出願人による申告)平成16年度独立行政法人情報通信研究機構、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条の適用を受ける特許出願</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 張 勁松 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 フランク・スーン 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 中村 哲 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】音素単位セットを最適化する方法及びシステム

(57) 【特許請求の範囲】

【請求項1】

予め定められた言語の音素単位セットを最適化する方法であって、コンピュータに、  
 コンピュータ読出可能なフォーマットで基本単位セットを準備するステップと、  
 前記基本単位セットにリーブ・ワン・アウト法を適用することによって複数の基本単位サブセットを生成するステップと、  
 前記基本単位サブセットの各々について言語的識別力の所定の尺度を計算するステップと、

前記基本単位セットを、前記基本単位サブセットのうち最も高い言語的識別力を備えたもので置換えるステップと、

前記生成するステップ、計算するステップ、及び置換えるステップを、所定の基準が満たされるまで繰返すステップとを実行させる、予め定められた言語の音素単位セットを最適化する方法。

【請求項2】

前記計算するステップが、前記基本単位セットと、前記基本単位サブセットの各々との間の相互情報量を計算するステップを含む、請求項1に記載の方法。

【請求項3】

前記置換えるステップが、前記基本単位セットを、前記基本単位サブセットのうち前記計算するステップで計算された相互情報量の最も高い値を有するもので置換えるステップを含む、請求項2に記載の方法。

## 【請求項 4】

前記基本単位セットは前記予め定められた言語のための基本音素セットである、請求項 1 ~ 請求項 3 のいずれかに記載の方法。

## 【請求項 5】

予め定められた言語の単位セットを最適化するシステムであって、

基本単位セットをコンピュータ読出可能なフォーマットで記憶するための記憶手段と、  
前記基本単位セットにリーブ・ワン・アウト法を適用することによって複数の基本単位サブセットを生成するための生成手段と、

前記基本単位サブセットの各々について言語的識別力の所定の尺度を計算するための計算手段と、

前記記憶手段に記憶された前記基本単位セットを、最も高い言語的識別力を有する基本単位サブセットで置換えるための置換手段と、

前記記憶手段、生成手段、計算手段及び置換手段を、所定の基準が満たされるまで繰返し動作するよう制御するための制御手段とを含む、システム。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

この発明は自動音声認識 (Automatic Speech Recognition: ASR) に関し、特に、ASR で用いられる音素セット等の音素単位セットの最適化に関する。

## 【背景技術】

## 【0002】

ASR はマン - マシン - インタラクションにおける必須のツールである。ASR によって、コンピュータは自然言語によるオペレータの指令を理解することができ、オペレータはコンピュータのための複雑なコマンドシステムを学ぶ必要がなくなる。

## 【0003】

図 6 は基本的な ASR の機構を示す。図 6 を参照して、ASR システム 162 は、入力音声 X 160 をデコードし、認識された (デコードされた) 単語  $\hat{W}$  164 (文中「 $\hat{\quad}$ 」の記号は本来文字 W の上に付されるものである。) を、以下の式 166 を用いて出力する。

## 【0004】

## 【数 1】

$$\begin{aligned}\hat{W} &= \arg \max_W P(X, W) \\ &= \arg \max_W P(X | W) \cdot P(W)\end{aligned}$$

ここで  $P(X | W)$  は音響モデル確率を示し、 $P(W)$  は言語モデル確率を示す。これらのモデルは対象となる言語の単語を、それぞれの音素のシーケンスと共に記載するレキシコンを用いて構築される。音素は予め定められた基本音素セットのうちから選択される。

## 【0005】

大語彙連続音声認識 (Large Vocabulary Continuous Speech Recognition: LVCSR) システムでは、広く受け入れられた音素セットが用いられる。

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【0006】

簡単な LVCSR タスクと、より複雑な LVCSR タスクとで同じ音素セットを用いるべきか、という問題がある。数字認識タスク等の小さな語彙のタスクでは、数字等の単語が基本単位として用いられる。同様に、簡単な LVCSR タスクでは、簡単な音素セット

10

20

30

40

50

を用いることが有利かもしれない。

【0007】

A S Rに関する多くの研究では、いくつかの発見的手法により決定された音素セットが試され、A S R認識性能に基づいて、1セットが選択される。

【0008】

音素セットにより多くの単位が含まれれば、音素学的により識別性のある情報を提供するであろう。しかしこれは、より詳細な音響的差異を使用するという意味でもある。音声認識の場合、より詳細な、またはより小さい音響差異をモデル化する必要が生じると、A M (Acoustic Model: 音響モデル)の頑健性が低下する傾向がある。

【0009】

音素セットに含まれる単位数が少なければ、より大きな音素セットに比べて、各音素A Mは、より多くのトレーニングデータを有することが通常である。さらに、音素の数が少ない場合、音素間での差異は多くの音素間での差異より大きくなる傾向がある。この結果、音素セットが小さくなればA Mはより頑健になり得る。しかし、音素セットサイズを小さくすると別の問題が生じる。すなわち、言語空間内における識別力が失われることである。例えば、日本語の長母音「A」と短母音「a」とが一つの母音にマージされるので、単語間の混同が増加するであろう。

【0010】

この点に関して、最新のA S R最適化は以下の考え方により行なわれる。上述の式を以下の形に書くことができる。

【0011】

【数2】

$$\begin{aligned} \arg \max_W P(X, W) &= \arg \max_W P(X | W) P(W) \\ &= \arg \max_W \sum_F P(X | F, W) P(F, W) \\ &\approx \arg \max_W \sum_F P(X | F) P(F | W) P(W) \\ &\approx \arg \max_W \max_F P(X | F) P(F | W) P(W), \end{aligned}$$

ここでFは基本単位シーケンスを示し、P(X | F)は頑健な音響モデル化の優勢なトピックを示し、P(F | W)は発音モデル化の注目のトピックを示し、P(W)は顕著な言語モデル化を示す。多くの場合、Fは音素セットである。

【0012】

しかし、先行技術では、種々の基本単位のセットを用いた場合に関する比較についてはヒューリスティックな試みがいくつかあったものの、特に確率を用いたA S Rの枠組み全体を考慮して基本単位セットの最適化を行なうことはほとんど全くされていないといえる。

【0013】

従って、この発明の目的の一つは、特定のA S Rタスクのための基本単位セットを最適化する方法と装置とを提供することである。

【0014】

この発明の別の目的は特定のA S Rタスクのための音素セットを最適化する方法と装置とを提供することである。

【課題を解決するための手段】

【0015】

この発明の一面によれば、予め定められた言語の音素単位セットを最適化する方法は、コンピュータに、コンピュータ読出可能なフォーマットで基本単位セットを準備するス

10

20

30

40

50

テップと、基本単位セットにリーブ・ワン・アウト法を適用することによって複数個の基本単位サブセットを生成するステップと、基本単位サブセットの各々について言語的識別力の所定の尺度を計算するステップと、基本単位セットを、基本単位サブセットのうち最も高い言語的識別力を備えたもので置換えるステップと、生成するステップ、計算するステップ、及び置換えるステップを、所定の基準が満たされるまで繰返すステップとを実行させる。

【0016】

好ましくは、計算するステップは、基本単位セットと、基本単位サブセットの各々との間の相互情報量を計算するステップを含む。

【0017】

より好ましくは、置換えるステップは、基本単位セットを、基本単位サブセットのうち計算するステップで計算された相互情報量の最も高い値を有するもので置換えるステップを含む。

【0018】

さらに好ましくは、基本単位セットは予め定められた言語のための基本音素セットである。

【0019】

この発明の別の局面によれば、予め定められた言語の単位セットを最適化するシステムは、基本単位セットをコンピュータ読出可能なフォーマットで記憶するための記憶手段と、基本単位セットにリーブ・ワン・アウト法を適用することによって複数個の基本単位サブセットを生成するための生成手段と、基本単位サブセットの各々について言語的識別力の所定の尺度を計算するための計算手段と、記憶手段に記憶された基本単位セットを、最も高い言語的識別力を有する基本単位サブセットで置換えるための置換手段と、記憶手段、生成手段、計算手段及び置換手段を、所定の基準が満たされるまで繰返し動作するよう制御するための制御手段とを含む。

【発明を実施するための最良の形態】

【0020】

A S Rの場合、二つの単語を識別するのに2種類の識別のための手段がある。一つは発音であり、他方は単語の文脈、すなわち言語モデル(Language Model: LM)である。一対の単語を音響スコアで識別することが困難な場合、例えば、同音語や類音語の場合、文脈的な単語情報があれば識別が容易になるであろう。例えば、「橋」と「箸」とは明らかに異なる文脈の単語である。

【0021】

上述の議論に基づき、この実施例は特定のA S Rタスクのための音素セットの最適な設計、すなわちタスクに基づく音素設計を提案する。基本的な考え方は、ある大きな音素セットから1音素を削除しても言語的識別力が大きく減じられることがなければ、音素セットサイズを減じるためにその音素を削除してもよい、というものである。

【0022】

この実施例では、最大相互情報量(Mutual Information: MI)基準に基づく音素セット設計を採用する。すなわち、MIを基本単位サブセットの言語的識別力の尺度として用いる。この実施例は中国語の最適化された音素セットを設計することに関するものである。

【0023】

基本単位セットは二つの具体的な局面で重要となる。すなわち、これは音響空間全体の主たる分類を規定し、さらに、言語空間の分類の重要な手がかりを提供する。

【0024】

図1は異なる基本単位セット、 $\mathcal{F}_1 = \{f_1, f_2, \dots, f_N\}$ 及び $\mathcal{F}_2 = \{p_1, p_2, \dots, p_M\}$ による直観的な影響力を示す。図1を参照して、 $\mathcal{F}_2$ の音素数Mは $\mathcal{F}_1$ の音素数Nよりはるかに大きいと仮定する(すなわち、 $N \ll M$ )。  $\mathcal{F}_1$ は音響空間 $\mathcal{F}$ をN個のサブスペース $f_1, f_2, \dots, f_N$ に分割し、 $\mathcal{F}_2$ は同じ音響空間 $\mathcal{F}$ をより小さいサブ

10

20

30

40

50

スペース  $p_1, p_2, \dots, p_M$  に分割する。従って、 $p_1$  は頑健な音響モデルを提供することができるが、その一方で、識別力は  $p_2$  のそれに比して弱い。

【0025】

図2はこの実施例の単位セットのトレーニングのための構成全体を示す。図2を参照して、トレーニングシステムは、トレーニング用の最新のASRシステム40と、言語モデルのための記憶部42と、レキシコンベースのデコードシステム44とを含む。

【0026】

トレーニング用ASRシステム40は、入力されたテキストWを音素シーケンスFに変換するための音声生成及びASRモジュール50と、音素シーケンスFによって形成される単語ラティス内のデコードされた単語テキストのうちで最も確率の高い単語テキスト $\hat{W}$ を、言語モデル42を参照しつつラティスの各経路をスコアリングすることによって選択するための単語ラティススコアリングモジュール52とを含む。

【0027】

レキシコンベースのデコードシステム44は、見出し語の各々を、それぞれの音素セット $p_1$ 及び $p_2$ を用いて記述する辞書62及び64と、辞書62及び64をそれぞれ用いて、入力テキストWを音素シーケンス $F_1$ 及び $F_2$ に変換するためのレキシコンベースの変換モジュール60と、音素シーケンス $F_1$ 及び $F_2$ によって形成される単語ラティス内の単語テキストのうちで最も確率の高い単語テキスト $W_1$ 及び $W_2$ を、言語モデル42を参照しつつラティスの各経路をスコアリングすることによって選択するための単語ラティススコアリングモジュール66とを含む。図2では説明を簡潔にするため、二つの辞書のみを示す。この実施例は中国語のASRシステムに関し、音素セット $p_1$ は声調情報を含み、一方音素セット $p_2$ はこれを含まない。

【0028】

トレーニング用ASRシステム40はトレーニングテキストWのコーパスを受け、以下の最大化式に従って、デコードされた単語 $\hat{W}$ を出力する。

【0029】

【数3】

$$\hat{W} = \arg \max_W \max_F P_A(X|F)P_P(F|W)P_L(W).$$

確率 $P(W|F)$ を最大にする音素セットが最適な音素セット $\hat{\Phi}$ として選択される。すなわち、

【0030】

【数4】

$$\begin{aligned} \hat{\Phi} &= \arg \max_{\Phi} P(W|F) \\ &= \arg \max_{\Phi} \frac{P(F|W)P(W)}{P(F)} \\ &= \arg \max_{\Phi} \frac{P(F|W)P(W)}{\sum_{W_i} p(F|W_i)p(W_i)} \end{aligned}$$

トレーニング用ASRシステム40とレキシコンベースのデコードシステム44との動作により、上述の式に従って、 $P(W|F)$ の要素を計算し、最適な音素セット $\hat{\Phi}$ を選択することができる。

【0031】

図3はこの実施例の音素セット最適化システム80の全体構造を示す図である。図3を

10

20

30

40

50

参照して、音素セット最適化システム 80 は、基本単位セット 90 の記憶装置と、トレーニングテキスト 92 の記憶装置と、基本単位セット 90 及びトレーニングテキスト 92 を用いて音素セットを最適化し、最適化された音素セット 94 を出力するための音素セット最適化モジュール 96 とを含む。

【0032】

音素セット最適化モジュール 96 は、コンピュータ上で実行されるソフトウェアで実現可能である。ソフトウェアの制御の流れを図 4 のフロー図で示す。図 4 を参照して、音素セット最適化モジュール 96 は以下のステップを実行する。初期音素セット  $S_0$  (すなわち基本単位セット 90) で作業中の音素セット  $S$  を置換える (ステップ 100)。音素サブセット  $S_i$  ( $i = 1$  から  $N$  の要素数まで;  $S_i = S - \{e_i\}$ ;  $e_i$  は  $S$  中の  $i$  番目の音素) を生成する (ステップ 102)。作業中のセット  $S$  とサブセット  $S_i$  の各々との間の相互情報量  $MI_i$  を計算する (ステップ 104)。以下の式を満たす指数  $M$  を特定する (ステップ 106)。

【0033】

【数 5】

$$M = \arg \max_i MI_i;$$

その後  $M$  番目の音素サブセット  $S_M$  を選択し、選択されたサブセット  $S_M$  中の音素を用いてレキシコン及びテキストコーパスを作り変える (ステップ 108)。作り変える過程において、レキシコンとテキストコーパスとは、レキシコンとテキストコーパス中で用いられている削除された音素を、それぞれ最も近い音素とマージするように更新される。

【0034】

音素セット最適化モジュール 96 はさらに、予め定められた停止条件が満たされたか否かを判断するステップを実行する (ステップ 110)。もし条件が満たされれば、音素セット最適化モジュール 96 は動作を停止する。さもなければ、制御はステップ 112 に進み、ここで選択されたサブセット  $S_M$  で作業中のセット  $S$  を置換え、その後制御はステップ 102 に戻る。

【0035】

予め定められた数だけ繰返したあと、動作は停止する。これに代えて、相互情報量の減少が予め定められたしきい値を超えた場合に動作を停止することもできる。

【0036】

音素セット最適化モジュール 96 は以下のように動作する。始めに、ステップ 100 で、基本単位セット 90 が作業用セット  $S$  として選択される。ステップ 102 で音素サブセット  $S_1$  から  $S_N$  までが生成される。サブセット  $S_i$  は作業中のセット  $S$  から音素  $e_i$  を除くことで生成される。言換えれば、 $S_i$  は作業中のセット  $S$  にリーブ・ワン・アウト法を適用することによって生成される。

【0037】

ステップ 104 で、作業中のセット  $S$  とサブセット  $S_1$  から  $S_N$  の各々との間の相互情報量  $MI_i$  が計算される。ステップ 106 で、相互情報量  $MI_i$  中で対応の相互情報量  $MI_M$  を最大にする指数  $M$  が選択される。

【0038】

ステップ 108 で、 $M$  番目の音素サブセット (サブセット  $S_M$ ) が選択され、選択された音素サブセット  $S_M$  を用いてレキシコンとテキストコーパスとが作り変えられる。

【0039】

ステップ 110 で、停止条件が満たされたか否かが判断される。もし条件が満たされていないならば、制御はステップ 112 に進み、ここで  $S$  が  $S_M$  と置換される。その後、制御はステップ 102 に戻り、ステップ 102 から 108 までが繰返される。停止条件が満たされると、動作は停止する。

【0040】

こうして、詳細な音素分類に基づいたものであつてかつサイズの大きい初期単位セット

10

20

30

40

50

90 から始めて、音素セット最適化モジュール96は何らかの基準に従って繰返しなが  
音素セットを減じることができる。

【0041】

図5はこの実施例の検証実験の結果を示す。この実験では、声調情報を含む元の203  
単位からなるセットを減少させる。声調情報を含まない59単位のセットを比較のために  
用いた。これら二つのセットは最新の中国語ASRシステムで広く用いられているもので  
ある。検証用テキストコーパスは1,614個の短文を含み、単語数は合計で9,484  
個である。

【0042】

図5を参照して、59の声調なしの単位セットC(ボックス132で示す)と比較して  
、元の203の声調付きセットは、ビット表現でより高い相互情報量を有する。線130  
で示す削減過程で、同じ59単位の数を備えて生成された単位セットは、図5の点Aで示  
すように、声調なしの単位セットに比べてより高い相互情報量を維持した。言換えれば、  
生成されたセットAは、数が同じであるにもかかわらず、伝統的な59の声調なし単位セ  
ットよりも良好な言語的識別力を有する。図5の点Bの単位セットは、声調なし単位セ  
ットCとほぼ同じ量の相互情報量を維持しているが、単位数は遥かに少ない。数は47であ  
り、従ってこれはCセットより効率が良い。

【0043】

上述の通り、この実施例のシステムと方法とは、相互情報量を減じることなく、音素セ  
ット中の音素の数をうまく減じることができる。タスクを特定したテキストをトレーニ  
ングに用いれば、音素セットはそのタスク用に最適化でき、その音素セットを用い  
れば、そのタスクについて十分な識別力を有する頑健な音響モデルを得ることができ  
る。また、十分に詳細な識別力を有する言語モデルを提供できる。

【0044】

上述の実施例では音素セットを最適化した。この発明は音素セットの最適化に限定さ  
れない。この発明は、ASRにおいて音素セットに置換可能ないずれの基本音素単位セ  
ットの最適化にも適用可能である。例えば、語彙が比較的小さい場合には、単位セ  
ットは語彙中の単語(単語発音)のセットであり得る。

【0045】

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および範囲内  
でのすべての変更を含む。

【図面の簡単な説明】

【0046】

【図1】異なる基本単位セットからの直観的な影響力を示す図である。

【図2】この実施例の単位セットのトレーニングの全体構成を示す図である。

【図3】この実施例の音素セット最適化システム80の全体構造を例示する図である。

【図4】この実施例の音素セット最適化モジュール96を実現するソフトウェアの制御フ  
ローを示す図である。

【図5】この実施例の検証実験結果をグラフの形で示す図である。

【図6】先行技術による基本ASRスキームを示す図である。

【符号の説明】

【0047】

40 トレーニング用ASRシステム

42 言語モデル

44 レキシコンベースのデコードシステム

50 ASRモジュール

52 単語ラティススコアリングモジュール

60 レキシコンベースの変換モジュール

10

20

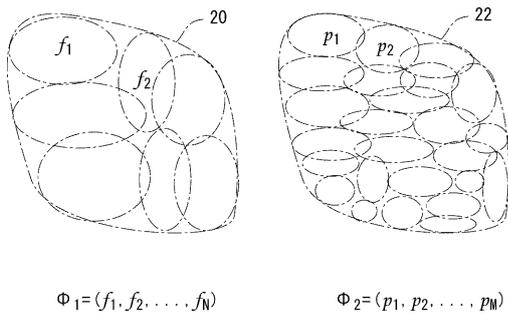
30

40

50

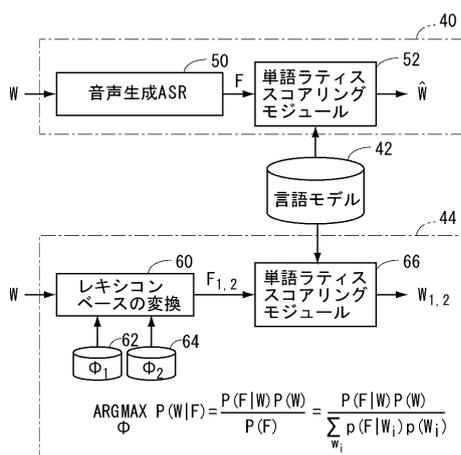
- 6 2、6 4 辞書
- 6 6 単語ラティススコアリングモジュール
- 8 0 音素セット最適化システム
- 9 0 基本単位セット
- 9 2 トレーニングテキスト
- 9 4 最適化音素セット
- 9 6 音素セット最適化モジュール

【図1】

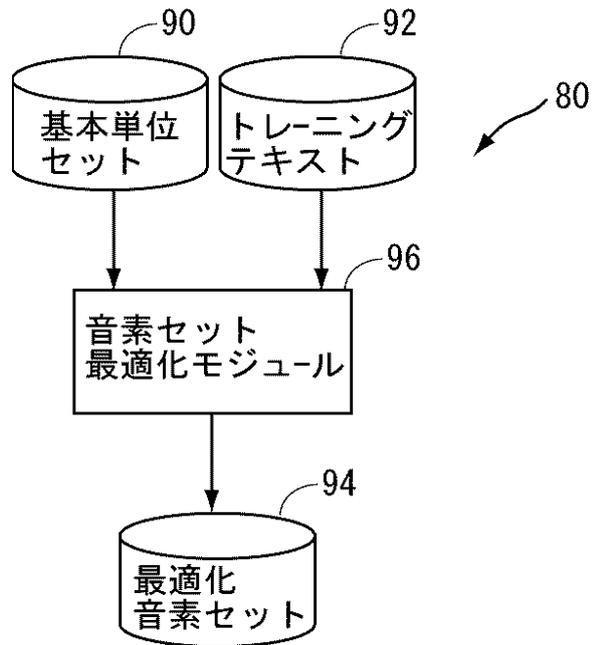


$N \ll M$

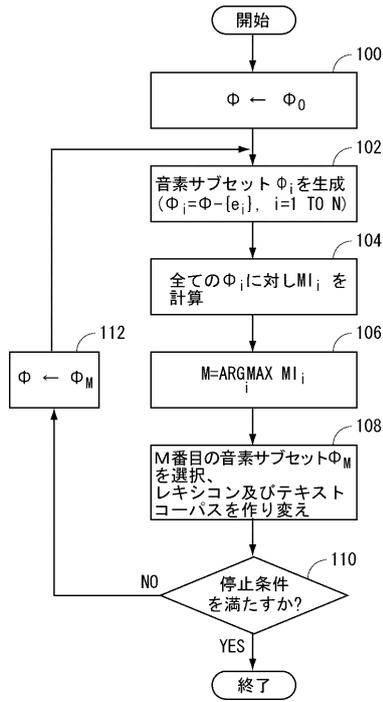
【図2】



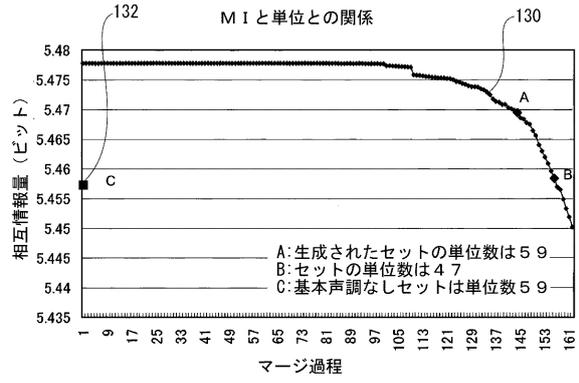
【図3】



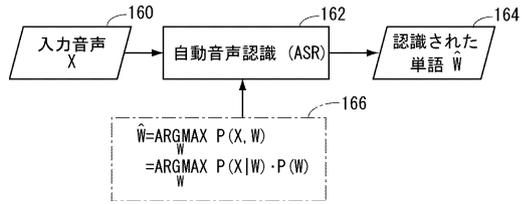
【図4】



【図5】



【図6】



---

フロントページの続き

審査官 山下 剛史

- (56)参考文献 特表平10-501078(JP,A)  
特開平09-288492(JP,A)  
特開平04-295893(JP,A)  
特開平11-352994(JP,A)  
特開平11-272291(JP,A)  
上田 修功, 最小分類誤り基準に基づくニューラルネットワーク識別機の最適線形統合法, 電子情報通信学会論文誌, 日本, 電子情報通信学会, 1999年 3月25日, Vol. J82-D-III No. 3, p. 522-530

(58)調査した分野(Int.Cl., DB名)

G10L 15/00-17/00

JSTPlus/JMEDPlus/JST7580(JDreamII)