

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5050175号
(P5050175)

(45) 発行日 平成24年10月17日(2012.10.17)

(24) 登録日 平成24年8月3日(2012.8.3)

(51) Int. Cl.	F I			
G 1 0 L 15/28 (2006.01)	G 1 0 L	15/28	2 1 0 A	
G 1 0 L 15/06 (2006.01)	G 1 0 L	15/28	3 7 0 E	
G 1 0 L 15/22 (2006.01)	G 1 0 L	15/06	4 0 0 V	
	G 1 0 L	15/22	4 7 0 Z	

請求項の数 8 (全 18 頁)

(21) 出願番号	特願2008-173551 (P2008-173551)	(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
(22) 出願日	平成20年7月2日(2008.7.2)	(74) 代理人	100099933 弁理士 清水 敏
(65) 公開番号	特開2010-14885 (P2010-14885A)	(72) 発明者	松田 繁樹 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
(43) 公開日	平成22年1月21日(2010.1.21)	(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
審査請求日	平成23年5月26日(2011.5.26)	(72) 発明者	葦苳 豊 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

最終頁に続く

(54) 【発明の名称】 音声認識機能付情報処理端末

(57) 【特許請求の範囲】

【請求項1】

音声信号から所定の音響特徴量を抽出して記憶するための特徴量記憶手段と、
前記所定の音響特徴量を予め定められた音声認識サーバに送信するための特徴量送信手段と、

前記サーバから前記所定の音響特徴量に対する音声認識の結果のテキストを受信するための受信手段と、

音声認識のための統計的音響モデルと、音声認識のための、カテゴリ別に編成された複数個のカテゴリ別言語モデルとを記憶するためのカテゴリ別モデル記憶手段と、

前記受信手段が受信した前記テキスト中の、未知語のタグ付けがされた区間に対応する音響特徴量を前記特徴量記憶手段から読み出し、前記モデル記憶手段に記憶された前記統計的音響モデル、及び前記カテゴリ別言語モデルの内で前記未知語のカテゴリに対応する言語モデル、を使用して音声認識を行なうための音声認識手段と、

前記受信手段が受信した前記テキスト中の前記未知語のタグ付けがされた区間を、前記音声認識手段の出力で置換するための置換手段とを含む、音声認識機能付情報処理端末。

【請求項2】

前記受信手段が受信した前記テキスト中に未知語のタグ付けがされた区間があるか否かを判定するための判定手段と、

前記判定手段の判定結果に応答して、前記受信手段が受信した前記テキストと、前記置換手段によって前記未知語が置換された前記テキストとを選択的に出力するための選択手

段とをさらに含む、請求項 1 に記載の音声認識機能付情報処理端末。

【請求項 3】

前記特徴量記憶手段は、

前記音声信号を所定時間ごとに所定長のフレームにフレーム化するためのフレーム化手段と、

前記フレーム化手段によりフレーム化されたフレームの各々の音声信号から、所定の複数個の音響特徴量を抽出するための特徴量抽出手段と、

前記フレーム化手段によりフレーム化されたフレームの各々に対して前記特徴量抽出手段により抽出された前記所定の複数個の音響特徴量を所定の圧縮アルゴリズムにより圧縮するための圧縮手段と、

前記フレーム化手段によりフレーム化されたフレームの各々に対して、前記圧縮手段により圧縮された音響特徴量を記憶するための記憶手段とを含み、

前記特徴量送信手段は、前記圧縮手段により圧縮された前記音響特徴量を送信するための手段を含む、請求項 1 又は請求項 2 に記載の音声認識機能付情報処理端末。

【請求項 4】

前記音声認識手段は、

前記受信手段が受信した前記テキスト中の、未知語のタグ付けがされた区間に対応するフレーム列の各々の音響特徴量を前記特徴量記憶手段から読み出し、前記所定の圧縮アルゴリズムに対応する伸長アルゴリズムを用いて伸長するための伸長手段と、

前記伸長手段により伸長されたフレーム列の前記複数個の音響特徴量を入力として、前記モデル記憶手段に記憶された前記統計的音響モデル、及び前記カテゴリ別言語モデルの内で前記未知語のカテゴリに対応する言語モデルを使用して音声認識を行なうための手段とを含む、請求項 3 に記載の音声認識機能付情報処理端末。

【請求項 5】

前記圧縮手段は、

前記複数個の所定の音響特徴量から予め組合された 2 つずつの音響特徴量の組合せの各々に対して予め準備されたコードブックを記憶するためのコードブック記憶手段と、

前記フレーム化手段によりフレーム化されたフレームの各々について、前記複数個の所定の音響特徴量から予め組合された 2 つずつの音響特徴量の組合せの各々を、前記コードブックのうちで対応するものを用いて符号化するための符号化手段とを含み、

前記送信するための手段は、前記フレーム化手段によりフレーム化されたフレームの各々について、前記符号化手段により得られた符号からなる符号列を送信するための手段を含む、請求項 3 又は請求項 4 に記載の音声認識機能付情報処理端末。

【請求項 6】

前記複数個の所定の音響特徴量は、各フレームの音声信号の第 0 次から第 1 2 次の MFCC パラメータと、パワーとを含む、請求項 1 ~ 請求項 5 のいずれかに記載の音声認識機能付情報処理端末。

【請求項 7】

前記音声認識機能付情報処理端末で実行可能なアプリケーションプログラムにより、前記音声認識機能付情報処理端末の使用者に関連して集積された情報を記憶するための関連情報記憶手段と、

前記関連情報記憶手段に記憶された前記情報を、カテゴリに分類するための分類手段と、

前記分類手段により分類されたカテゴリごとに統計的言語モデルを作成することにより、前記複数個のカテゴリ別言語モデルを作成するための言語モデル作成手段とをさらに含む、請求項 1 ~ 請求項 6 のいずれかに記載の音声認識機能付情報処理端末。

【請求項 8】

前記未知語のタグ付けがされた区間は、未知語のタグ付けがされた音節列を含む、請求項 1 ~ 請求項 7 のいずれかに記載の音声認識機能付情報処理端末。

【発明の詳細な説明】

10

20

30

40

50

【技術分野】

【0001】

この発明は通信機能を有する情報処理端末に関し、特に、携帯に便利な小さな筐体で、しかも音声認識による入力可能な情報処理端末に関する。

【背景技術】

【0002】

計算機の性能が向上し小型化するに伴い、携帯型情報端末が普及している。現代の携帯型情報端末は、例えば携帯電話のように、主たる機能の他にインターネットに接続する機能も持っており、電子メールによるコミュニケーションの有力なツールとなっている。

【0003】

携帯型情報端末を電子メールなどのテキストベースのコミュニケーションにおいて使用する場合の最大の問題は、入力インターフェイスである。大きな情報端末であればフルキーボードを装備することも可能であるが、携帯電話のような小型の装置ではそれは難しい。そのため、テンキーパッドを複数回押すことによって一文字を入力するようなインターフェイスが一般的である。その結果、通常の情報端末と比較して入力に時間がかかるという問題がある。

【0004】

こうした問題に対処すべく、あえてフルキーボードを備えた小型の情報端末もあるが、キートップが非常に小さくなってしまいうために、結局は入力がしづらいという欠点がある。

【0005】

一方、携帯型情報端末の高能力化に伴い、いわゆる音声認識技術を携帯型情報端末における入力に使用することも考えられている。CPU(Central Processing Unit)の処理能力の向上に伴い、そのようなことも不可能ではないと考えられる。

【0006】

しかし、現代の音声認識技術の場合、音響モデル、単語辞書、及び統計的言語モデルなどを装備する必要がある。音声認識の性能を高めるためには、これらモデルは大容量にせざるを得ない。その結果、現在のところは携帯型情報端末装置単体で十分な性能の音声認識を行なうことは難しいが、コストが非常に高くなってしまふ。

【0007】

そこで、携帯型情報端末では音声認識を行わず、携帯型情報端末から送られた音声をサーバ側で音声認識する音声認識システムが特許文献1に記載されている。特許文献1に記載された音声認識システムでは、予め、サーバの検索データベースに、氏名等と、住所等と、関連情報等とを関連づけて蓄積しておく。携帯型情報端末から音声を受取ると、住所等の一部若しくは全部、又は関連情報等を取得し、この取得された住所等の一部若しくは全部、又は関連情報等に基づいて検索データベースを検索し、この検索結果を用いて、氏名等の音声データを認識する。

【特許文献1】特開2008-015439号公報

【特許文献2】特開2008-129318号公報

【非特許文献1】山本博史他4名、「複数のマルコフモデルを用いた階層化言語モデルによる未登録語認識」、電子情報通信学会論文誌、D-II、Vol. J87-D-II, No. 12, pp. 2104-2111、2004年12月

【発明の開示】

【発明が解決しようとする課題】

【0008】

特許文献1に記載の技術によれば、音声認識は携帯型情報端末装置ではなくサーバ側で行なわれる。そのため、携帯型情報端末装置に音声認識のシステムを搭載する必要はない。音声認識に必要なリソースはサーバ側に十分確保できるため、音声認識の精度は確保できる。

10

20

30

40

50

【 0 0 0 9 】

これと同様の技術として、音声認識の前段である特徴量の抽出までを携帯型情報端末で行ない、特徴量のみをサーバに送信するという、分散型音声認識という考え方もある。送信されるデータ量は音声そのものよりも特徴量の方が少なくなるため、通信容量が少なくなるという効果がある。また、サーバ側の負荷が軽くなるという効果もある。情報処理装置が高性能化していることに鑑みると、分散型音声認識システムがこれからの音声認識システムとしては有力である。

【 0 0 1 0 】

しかし、音声そのものを送信するにせよ、特徴量を送信するにせよ、音声認識の精度を十分に高めるためには、サーバ側の辞書に非常にたくさんの固有名詞を登録する必要がある。例えばある個人にとって、友人の名前又は愛称（ニックネーム）、行きつけのお店、通学する学校、よく利用する施設、駅名などの固有名詞はコミュニケーションを行なう上で非常に重要な意味を持つ。これらが十分な精度で認識されるためには、サーバ側にそれらの固有名詞を正しく登録しなければならない。万が一、ある固有名詞が登録されていない場合には、その固有名詞については正しい音声認識結果が得られない。固有名詞は無数に存在し、しかも新しく生ずるものがある。したがって、それら無数の固有名詞について最新の状態にサーバのリソースを維持する作業は大変なものになる。

10

【 0 0 1 1 】

しかも、登録される固有名詞の数が多ければ音声認識の精度が高まるかということ、必ずしもそうではない。例えば同じようでも少し異なるような固有名詞が辞書又は言語モデルに複数個登録されている場合には、音声認識の精度が低くなる可能性がある。したがって仮に固有名詞を100パーセント登録できたとしても、音声認識の精度が高まるとは限らないという問題がある。

20

【 0 0 1 2 】

このように、辞書に登録されていない単語をどのように扱うかは、未知語の問題として知られている。特許文献2には、未知語をカタカナ文字列として出力できるような言語モデルを作成するシステムが開示されている。

【 0 0 1 3 】

しかし、未知語がカタカナ文字列で出力されても、音声認識が正しく行なわれているわけではない。カタカナ文字列自体に誤りがあるかも知れず、仮にカタカナ文字列が正しくとも、日本語の場合には固有名詞としての文字列に変換されなければ正しい認識が行なわれたとはいえない。このように未知語を未知語として出力するだけでは、音声認識の精度を高めたことにならず、結局、サーバ側に十分なリソースを準備する必要があり、サーバ側のリソースの肥大化を招くことになる。

30

【 0 0 1 4 】

それゆえに本発明の目的は、分散型の音声認識を利用する情報処理端末であって、使用者にとって音声認識の精度が十分に高く、かつ音声認識を行なうサーバ側のリソースの極端な肥大化を防止できる音声認識機能付情報処理端末を提供することである。

【課題を解決するための手段】

【 0 0 1 5 】

本発明の第1の局面に係る音声認識機能付情報処理端末は、音声信号から所定の音響特徴量を抽出して記憶するための特徴量記憶手段と、所定の音響特徴量を予め定められた音声認識サーバに送信するための特徴量送信手段と、サーバから所定の音響特徴量に対する音声認識の結果のテキストを受信するための受信手段と、音声認識のための統計的音響モデルと、音声認識のための、カテゴリ別に編成された複数個のカテゴリ別言語モデルとを記憶するためのカテゴリ別モデル記憶手段と、受信手段が受信したテキスト中の、未知語のタグ付けがされた区間に対応する音響特徴量を特徴量記憶手段から読み出し、モデル記憶手段に記憶された統計的音響モデル、及びカテゴリ別言語モデルの内未知語のカテゴリに対応する言語モデル、を使用して音声認識を行なうための音声認識手段と、受信手段が受信したテキスト中の未知語のタグ付けがされた区間を、音声認識手段の出力で置換する

40

50

ための置換手段とを含む。

【0016】

この情報処理端末では、特徴量記憶手段が、音声信号から所定の音響特徴量を抽出し、記憶する。この音響特徴量は、特徴量送信手段により音声認識サーバに送信される。音声認識サーバでの音声認識結果であるテキストは受信手段により受信される。このテキスト中の未知語部分には、未知語であることを示すタグと、その未知語が属するカテゴリを示すタグとが付されている。音声認識手段は、特徴量記憶手段に記憶されている音響特徴量のうち、未知語に対応する部分を読み出し、未知語に付されていたタグのカテゴリに対応するカテゴリ別言語モデルを使用して音声認識を行なう。置換手段は、音声認識の結果で未知語部分を置換する。

10

【0017】

カテゴリ別言語モデルはこの情報処理端末に固有のものである。したがってこれらカテゴリ別言語モデルは、利用者にとって特に関連ある固有名詞などから生成される。サーバで未知語として認識された音響特徴量の部分を、情報処理端末でこのカテゴリ別言語モデルを用いて音声認識し直すことにより、サーバでは未知語であった固有名詞が正しく認識される可能性が大きくなる。そのために情報処理端末に多くのリソースを準備する必要はない。また、サーバでも未知語の音声認識を行なうためにリソースを肥大化させる必要がない。その結果、分散型の音声認識を利用する情報処理端末であって、使用者にとって音声認識の精度が十分に高く、かつ音声認識を行なうサーバ側のリソースの極端な肥大化を防止できる音声認識機能付情報処理端末を提供できる。

20

【0018】

好ましくは、音声認識機能付情報処理端末は、受信手段が受信したテキスト中に未知語のタグ付けがされた区間があるか否かを判定するための判定手段と、判定手段の判定結果に応答して、受信手段が受信したテキストと、置換手段によって未知語が置換されたテキストとを選択的に出力するための選択手段とをさらに含む。

【0019】

サーバから受けた音声認識結果に未知語がなければそれを選択し、未知語がある場合だけ情報処理端末での音声認識を行なう。情報処理端末において余分な処理をする必要がなく、音声認識の結果をより早く提示することが可能になる。

【0020】

より好ましくは、特徴量記憶手段は、音声信号を所定時間ごとに所定長のフレームにフレーム化するためのフレーム化手段と、フレーム化手段によりフレーム化されたフレームの各々の音声信号から、所定の複数個の音響特徴量を抽出するための特徴量抽出手段と、フレーム化手段によりフレーム化されたフレームの各々に対して特徴量抽出手段により抽出された所定の複数個の音響特徴量を所定の圧縮アルゴリズムにより圧縮するための圧縮手段と、フレーム化手段によりフレーム化されたフレームの各々に対して、圧縮手段により圧縮された音響特徴量を記憶するための記憶手段とを含み、特徴量送信手段は、圧縮手段により圧縮された音響特徴量を送信するための手段を含む。

30

【0021】

サーバには、圧縮された音響特徴量が送信される。その結果、情報処理端末から音声認識のためのサーバへの送信データ量を少なく抑えることができる。

40

【0022】

さらに好ましくは、音声認識手段は、受信手段が受信したテキスト中の、未知語のタグ付けがされた区間に対応するフレーム列の各々の音響特徴量を特徴量記憶手段から読み出し、所定の圧縮アルゴリズムに対応する伸長アルゴリズムを用いて伸長するための伸長手段と、伸長手段により伸長されたフレーム列の複数個の音響特徴量を入力として、モデル記憶手段に記憶された統計的音響モデル、及びカテゴリ別言語モデルの中で未知語のカテゴリに対応する言語モデルを使用して音声認識を行なうための手段とを含む。

【0023】

カテゴリ別言語モデルの中で、未知語に付されていた、カテゴリを表すタグに対応する

50

ものが選択され、それを使用して音声認識が行なわれる。情報処理端末の利用者に特に関連する情報であって、かつサーバであるカテゴリに属すると推定された単語を、そのカテゴリの単語の言語モデルを使用して音声認識するので、音声認識の結果の精度がより高くなる。

【 0 0 2 4 】

圧縮手段は、複数個の所定の音響特徴量から予め組合された２つずつの音響特徴量の組合せの各々に対して予め準備されたコードブックを記憶するためのコードブック記憶手段と、フレーム化手段によりフレーム化されたフレームの各々について、複数個の所定の音響特徴量から予め組合された２つずつの音響特徴量の組合せの各々を、コードブックのうちで対応するものを用いて符号化するための符号化手段とを含んでもよい。送信するための手段は、フレーム化手段によりフレーム化されたフレームの各々について、符号化手段により得られた符号からなる符号列を送信するための手段を含んでもよい。

10

【 0 0 2 5 】

一実施の形態では、複数個の所定の音響特徴量は、各フレームの音声信号の第 0 次から第 1 2 次の M F C C パラメータと、パワーとを含む。

【 0 0 2 6 】

好ましくは、音声認識機能付情報処理端末は、音声認識機能付情報処理端末で実行可能なアプリケーションプログラムにより、音声認識機能付情報処理端末の利用者に関連して集積された情報を記憶するための関連情報記憶手段と、関連情報記憶手段に記憶された情報を、カテゴリに分類するための分類手段と、分類手段により分類されたカテゴリごとに統計的言語モデルを作成することにより、複数個のカテゴリ別言語モデルを作成するための言語モデル作成手段とをさらに含む。

20

【 0 0 2 7 】

未知語のタグ付けがされた区間は、未知語のタグ付けがされた音節列であってもよい。

【 発明の効果 】

【 0 0 2 8 】

以上のようにこの発明によれば、分散型の音声認識を利用するシステムにおいて、サーバでは未知語であった固有名詞を情報処理端末で正しく認識できる可能性が大きくなる。そのために情報処理端末に多くのリソースを準備する必要はない。また、サーバでも未知語の音声認識を行なうためにリソースを肥大化させる必要がない。さらに、サーバからの音声認識結果に、未知語のカテゴリを示すタグを挿入することで、そのタグに対応した言語モデルを用いて情報処理端末で未知語に対して音声認識をし直すことができる。その結果、使用者にとって音声認識の精度が十分に高く、かつ情報処理端末側でも、音声認識を行なうサーバ側でも、リソースの極端な肥大化を防止できる音声認識機能付情報処理端末を提供できる。

30

【 発明を実施するための最良の形態 】

【 0 0 2 9 】

以下の説明において、全図を通じ、同一の部品には同一の参照番号を付してある。それらの名称及び機能も同一である。したがってそれらについての詳細な説明は繰返さない。

【 0 0 3 0 】

< 構成 >

図 1 に、本発明の第 1 の実施の形態に係る音声認識システム 1 0 の概略構成を示す。図 1 を参照して、音声認識システム 1 0 は、携帯型情報処理装置の一例であり、利用者の音声 3 0 から音響特徴量 3 2 を抽出する機能を持つ携帯電話機 2 0 と、携帯電話機 2 0 が抽出した音響特徴量 3 2 を受けると、この音響特徴量 3 2 に対して音声認識を行ない、認識結果のテキスト 3 4 を携帯電話機 2 0 に返信する機能を持つ音声認識サーバ 2 2 とを含む。

40

【 0 0 3 1 】

音声認識サーバ 2 2 は、音声認識結果に未知語が存在する場合には、その未知語を認識結果のテキスト内に音節列として挿入し、かつその音節列が未知語であることを示すタグ

50

と、その未知語が、予め分類されたいくつかのカテゴリの中のどのカテゴリに属するかを示すタグとをその音節列に付与する機能を持つ。音声認識サーバ 22 は、例えば周知の音声認識技術と、特許文献 2 に記載されているような未知語の認識技術及び非特許文献 1 に記載されているような、階層化言語モデルによるクラス推定とを組合せることにより実現できる。

【0032】

再び図 1 を参照して、携帯電話機 20 は、音声認識サーバ 22 から送信されてくる認識結果のテキスト 34 を受けると、この中に未知語が含まれている場合には、元の音声信号から得た音響特徴量の、その未知語部分に対して音声認識を行なって、その結果で未知語を置換する処理をして最終結果のテキスト 36 を出力する。携帯電話機 20 で行なわれるこの未知語の音声認識には、この携帯電話機 20 の使用者に関連して各種アプリケーションプログラムによって集積された情報から作成された、カテゴリ別言語モデルのうち、未知語に付されたカテゴリのタグに対応したものが使用される。このカテゴリ別言語モデルは、この携帯電話機 20 の利用者に特に関連した情報から作成されたものである。音声認識の結果として得られる固有名詞としては、この携帯電話機 20 の利用者の友人、知人、よく利用する施設、学校などに関するものが大部分であるから、携帯電話機 20 におけるこの音声認識での認識精度は高くなる。音声認識サーバ 22 のように多数の利用者による音声処理する必要はないので、携帯電話機 20 の言語モデルに登録すべき単語は少なく済む。

【0033】

図 2 に、携帯電話機 20 のうち、本発明に関連する部分の機能的構成を示す。図 2 を参照して、携帯電話機 20 は、マイクロフォン 50 と、マイクロフォン 50 からの音声信号に対して所定の音響処理を行なって音声信号の特徴量を抽出し、さらにコードブックを用いて符号化して符号列を時系列で出力する音響信号処理部 54 と、音響信号処理部 54 が符号化時に使用するコードブックを記憶したコードブックメモリ 52 と、音響信号処理部 54 が出力する符号列を一時記憶するための送信バッファ 56 と、送信バッファ 56 に記憶された符号列をパケット化して音声認識サーバ 22 に送信するための送信処理部 58 とを含む。

【0034】

携帯電話機 20 はさらに、音響信号処理部 54 が出力する符号列をフレームごとに順次記憶するための符号記憶部 60 と、音声認識サーバ 22 から音声認識結果のテキスト 34 のパケットを受信するための受信処理部 62 と、受信処理部 62 により受信された音声認識結果のテキスト 34 を一時記憶するための受信バッファ 64 と、受信バッファ 64 に記憶された音声認識結果のテキストに未知語が含まれていれば、その部分をコードブックを用いて復号し、改めて音声認識を行なって、未知語をその音声認識結果の単語で置換する未知語処理部 70 と、未知語処理部 70 が音声認識の際に利用する音響モデルを記憶する音響モデル記憶部 68 及び複数のカテゴリ別言語モデルを記憶する言語モデル記憶部 66 と、未知語処理部 70 が出力するテキストを携帯電話機 20 上で稼動している他のアプリケーションに渡す処理を行なうための出力部 72 とを含む。

【0035】

音響信号処理部 54 は、マイクロフォン 50 からの音声信号を、所定時間おきに所定時間長でフレーム化するためのフレーム化モジュール 80 と、フレーム化モジュール 80 から出力されるフレーム列の各々のフレームに対し、雑音抑圧及び特徴量抽出処理を行なって特徴量ベクトルを出力するための雑音抑圧・特徴量抽出部 82 と、雑音抑圧・特徴量抽出部 82 から出力される特徴量ベクトル列の各ベクトルに対し、コードブックメモリ 52 に記憶されたコードブックを用いた符号化を行ない、符号列を送信バッファ 56 及び符号記憶部 60 に格納するための符号化処理部 84 とを含む。

【0036】

本実施の形態では、雑音抑圧・特徴量抽出部 82 が抽出する音響特徴量は、MFCC (Mel Frequency Cepstrum Coefficient) の第 1 次 ~

10

20

30

40

50

第12次の係数、C0（第0次のMFCC係数）、及び音声信号のパワーを含む。すなわち、特徴量ベクトルは14次元である。

【0037】

未知語処理部70は、受信バッファ64に記憶された、音声認識結果のテキスト列の中で未知語のタグが付された音節列（カタカナ列）を抽出し、符号記憶部60に記憶された符号列の中から、この未知語に対応する符号列部分を切出す処理を行なう未知語切出処理部90と、未知語切出処理部90によって切出された符号列をコードブックメモリ52に記憶されたコードブックを用いて音響特徴量列に戻し、言語モデル記憶部66に記憶された複数の言語モデルの中で、未知語に付されたカテゴリタグに対応するものと、音響モデル記憶部68に記憶された音響モデルとを用いて音声認識処理を行ない、音声認識結果の単語を出力する未知語認識処理部92と、受信バッファ64に記憶されたテキストを読み込み、未知語のタグが付された音節列を、未知語認識処理部92により出力される音声認識後の単語で置換したテキストを出力するための未知語入替処理部94とを含む。

10

【0038】

未知語処理部70はさらに、受信バッファ64に記憶された音声認識後のテキストに、未知語のタグが付された音節列があるか否かを判定し、ある場合にはTRUEを、ない場合にはFALSEをとる判定結果信号を出力するための判定部96と、受信バッファ64に記憶されたテキストを受ける第1の入力と、未知語入替処理部94の出力するテキストを受ける第2の入力とを有し、判定部96から出力される判定信号がTRUEのときには未知語入替処理部94からのテキストを、FALSEのときには受信バッファ64に格納されたテキストを、それぞれ選択して出力部72に与えるための選択部98とを含む。なお、判定部96からの判定結果信号は、未知語切出処理部90、未知語認識処理部92及び未知語入替処理部94にも与えられており、これら回路は判定結果信号がTRUEのときには動作し、FALSEであるときには停止する。

20

【0039】

図3は、図2に示す言語モデル記憶部66に記憶されたカテゴリ別言語モデルを作成するためのカテゴリ別言語モデル作成部100のブロック図である。図3を参照して、図2に示す携帯電話機20には、住所録プログラムにより集積された住所録102と、メールプログラムにより集積されたメールアドレスDB104と、GPS（Global Positioning System）などの地図ソフトで使用される地図データ106とが含まれる（いずれも図2では図示していない。）。カテゴリ別言語モデル作成部100は、これらからカテゴリ別言語モデルを作成する。図3に示すように、本実施の形態では、カテゴリ別言語モデルとしては、施設名言語モデル（LM）と、日本人の姓に関する姓言語モデルと、日本人の名前に関する名前言語モデルと、日本人のニックネームに関するニックネーム言語モデルと、場所名に関する場所言語モデルとを有する。

30

【0040】

図3を参照して、カテゴリ別言語モデル作成部100は、住所録102、メールアドレスDB104、及び地図データ106から言語モデル作成のためのデータを抽出し分類して、施設名データファイル112、姓データファイル114、名データファイル116、ニックネームデータファイル118、及び場所データファイル120等、カテゴリ別のファイルに出力するための抽出部110と、抽出部110により作成されたデータファイル112～120をそれぞれ用いて、施設名言語モデル、姓言語モデル、名言語モデル、ニックネーム言語モデル、場所言語モデルなど、カテゴリ別言語モデルを言語モデル記憶部66に作成するための言語モデル作成部122とを含む。

40

【0041】

住所録102などでは、予め所定の見出しとそれに対するデータという形でデータが集積されている。内部的には、これらデータは例えばXML（eXtended Markup Language）などで保持されていることが多く、各タグをキーワードにして対応するデータを集めることにより、カテゴリ別のデータファイル112～120を集めることができる。

50

【 0 0 4 2 】

本実施の形態では、抽出部 1 1 0 を 1 本のコンピュータプログラムで実現し、住所録 1 0 2、メールアドレス DB 1 0 4 及び地図データ 1 0 6 から一度に言語モデル作成用のデータファイルを作成するが、アプリケーション別に抽出用のコンピュータプログラムを作成するようにしてもよい。

【 0 0 4 3 】

言語モデル記憶部 6 6 に記憶されるカテゴリ言語モデルはいずれも同一のフォーマットである。データファイル 1 1 2 ~ 1 2 0 も同一フォーマットである。したがってここでも言語モデル作成部 1 2 2 は 1 本のコンピュータプログラムで実現できる。言語モデルの作成時に、入力ファイル名及び言語モデル名を引数として与えれば、言語モデル作成部 1 2 2 はそれら引数にしたがって別々のデータファイルからデータを読み、指定された言語モデルを作成する。

10

【 0 0 4 4 】

図 4 は、音声認識サーバ 2 2 の機能ブロック図である。音声認識サーバ 2 2 のハードウェア構成は公知であるため、その詳細については述べない。音声認識サーバ 2 2 は、概略的には、任意の情報処理端末から音声認識の要求とともに音声認識の対象データである符号列をパケット形式で受信するための受信処理部 1 3 0 と、受信処理部 1 3 0 により受信されたパケットを一時的に記憶するための受信バッファ 1 3 2 と、図 2 に示すコードブックメモリ 5 2 に記憶されたコードブックと同一のコードブックを記憶したコードブックメモリ 1 3 4 と、受信バッファ 1 3 2 に記憶されたパケット列から、音声認識の対象となる符号列を抽出し、コードブックメモリ 1 3 4 に記憶されたコードブックを用いて音響特徴量に戻す処理を行なうためのデコーダ 1 3 6 とを含む。

20

【 0 0 4 5 】

音声認識サーバ 2 2 はさらに、音声認識に使用される、隠れマルコフモデル (H M M) からなる音響モデルを記憶した音響モデル記憶部 1 3 8 と、予め所定のコーパスから作成された、クラス (品詞) 別のバイグラムの統計的言語モデルを記憶するためのクラス言語モデル記憶部 1 4 0 と、予め所定のコーパスから作成された、単語トライグラムからなる統計的言語モデルを記憶するための単語言語モデル記憶部 1 4 4 と、携帯電話機 2 0 に記憶されているカテゴリ別の言語モデルと同様、カテゴリ別に予め作成された複数個のカテゴリ別音節モデルを記憶するためのカテゴリ別音節モデル記憶部 1 4 6 とを含む。音節モデルとは、音節単位で前後の音節との文脈を考慮して作成された言語モデルである。同一の言語では、姓、名、地名、施設名など、単語が属するカテゴリによって音韻列の生起確率は異なっている。したがって、音声認識の過程で未知語に遭遇した場合、これら音節モデルを参照してその未知語の音節列が生ずる尤度を各モデルを使用して算出し、最も高い尤度を示す音節モデルのカテゴリをその未知語のカテゴリとすることができる。

30

【 0 0 4 6 】

クラス言語モデル記憶部 1 4 0 に記憶されたクラス言語モデル (バイグラム) とは、二つの連続する単語の品詞について、どのような順序付組合せがどの程度の確率で生ずるかを表す言語モデルである。

【 0 0 4 7 】

音声認識サーバ 2 2 はさらに、音響モデル記憶部 1 3 8 に記憶された音響モデル、クラス言語モデル記憶部 1 4 0 に記憶されたクラスバイグラム、単語言語モデル記憶部 1 4 4 に記憶された単語トライグラムを用いて音声認識を行なってテキストに変換し、未知語はカタカナ列で出力するための音声認識処理部 1 4 2 を含む。音声認識処理部 1 4 2 は、未知語部分については、クラスバイグラムから算出される尤度と、音節モデルから算出される音節列の尤度とを乗算することにより、各音節列の候補の尤度を算出し、最も尤度が高い音節列を、未知語のタグを付して出力するとともに、その音節列を与える音節モデルのカテゴリを示すタグをその音節列に付与する。

40

【 0 0 4 8 】

なお、通常の音声認識処理と同様、音声認識処理部 1 4 2 が出力するテキストの各単語

50

、及び未知語を構成するカタカナ列を構成するカタカナ（音節）の各々には、元の音声信号における開始時間と終了時間とを示す情報が付加されている。

【 0 0 4 9 】

音声認識サーバ 2 2 はさらに、音声認識処理部 1 4 2 の出力する時間情報付のテキストを一時記憶するための出力バッファ 1 4 8 と、出力バッファ 1 4 8 に記憶されたテキスト列を、音声認識要求を送信してきた情報処理端末に送信するための送信処理部 1 5 0 とを含む。図 2 に示す受信処理部 6 2 が受信するのは、この送信処理部 1 5 0 により送信された、時間情報付のテキストである。

【 0 0 5 0 】

次に、図 2 に示す携帯電話機 2 0 の音声認識機能のうち、未知語処理部 7 0 の機能を実現するためのコンピュータプログラムのフローチャートを図 5 に示す。携帯電話機 2 0 の音声認識機能のうち、音響信号処理部 5 4 の部分については公知で、通常の分散処理型音声認識システムで採用されているものであるため、ここではその詳細については述べない。

10

【 0 0 5 1 】

図 5 を参照して、このプログラムは、音声認識結果の時間情報付のテキストを音声認識サーバ 2 2 から受信するステップ 1 6 0 と、受信した時間情報付のテキストを受信バッファ 6 4 に一時保存するステップ 1 6 2 と、受信したテキスト内に未知語のタグが付された部分があるか否かを判定し、判定結果に応じて制御の流れを分岐させるステップ 1 6 4 と、ステップ 1 6 4 において未知語タグが付された部分がないと判定されたことに応答して、音声認識サーバ 2 2 から受信したテキストをそのままアプリケーションに渡して処理を終了するステップ 1 8 0 とを含む。

20

【 0 0 5 2 】

このプログラムはさらに、ステップ 1 6 4 において、テキスト内に未知語のタグが付された部分があると判定されたときに実行され、その未知語のタグが付された部分の時間情報に基づいて、符号記憶部 6 0 に記憶された符号列の中で、その時間に対応する部分を切出す、すなわち読出す処理を実行するステップ 1 6 6 と、ステップ 1 6 6 に続き、その符号列をコードブックを用いて音響特徴量に伸長する処理を行なうステップ 1 6 8 と、ステップ 1 6 8 に続き、未知語部分に付されている、その未知語が属するカテゴリを示すタグに対応した言語モデルを言語モデル記憶部 6 6（図 2 参照）から選択するステップ 1 7 0 と、ステップ 1 7 0 で選択された言語モデルと、音響モデル記憶部 6 8（図 2 参照）に記憶された音響モデルとを使用して音声認識し、最尤の単語を出力するステップ 1 7 2 と、ステップ 1 7 2 で音声認識により得られた単語で、音声認識サーバ 2 2 から受信したテキスト列の内の未知語タグが付された部分を置換するステップ 1 7 4 と、ステップ 1 7 4 で未知語部分が音声認識の結果で置換されたテキストをアプリケーションに渡して処理を終了するステップ 1 7 6 とを含む。

30

【 0 0 5 3 】

< 動作 >

以上、図 1 ~ 図 5 に示した構成を有する音声認識システム 1 0 は以下のように動作する。最初に、利用者が例えばメールプログラムを起動し、メールテキストを音声で入力する場合を想定する。利用者の音声はマイクロフォン 5 0 により音声信号に変換され、フレーム化モジュール 8 0 によって所定時間おきに所定長でフレーム化される。フレーム化モジュール 8 0 が出力するフレーム列は雑音抑圧・特徴量抽出部 8 2 に与えられる。

40

【 0 0 5 4 】

雑音抑圧・特徴量抽出部 8 2 は、入力されるフレーム列の各々に対し、雑音抑圧処理を行なった後、先に述べたとおり、第 1 ~ 第 1 2 次の M F C C 係数、C 0（第 0 次の M F C C 係数）、及びエネルギーを算出して 1 4 次の音響特徴量ベクトルを生成し、符号化処理部 8 4 に与える。

【 0 0 5 5 】

符号化処理部 8 4 は、雑音抑圧・特徴量抽出部 8 2 から与えられる音響特徴量ベクトル

50

の各々に対し、特徴量を示す要素を2つずつ組合せてコードブックメモリ52に記憶されたコードブックのうちでその組合せに対応するものを用いて符号化し出力する。一つの音響特徴量ベクトルの要素は14個であり、2つずつの組合せで符号化が行なわれるので、14個の音響特徴量が全部で7個の符号からなる符号列に変換される。例えば1特徴量について8ビットが使用され、コードブックにより既定される符号が16個であれば、全部で16ビットの情報が4ビットに圧縮されることになる。これが7組あるので、全体では 7×16 ビット = 112ビットの情報が $4 \times 7 = 28$ ビットに削減されることになる。

【0056】

符号化処理部84は、このように圧縮された符号列を送信バッファ56及び符号記憶部60に格納する。

【0057】

送信処理部58は、送信バッファ56に20フレーム分の符号列が格納されると、それらから1つのパケットを組立てて音声認識サーバ22に送信する。

【0058】

音声認識サーバ22の受信処理部130は、受信したパケットを受信バッファ132に格納する。デコーダ136は、受信バッファ132に格納されたパケットから各フレーム毎の符号列を順次読出して、コードブックメモリ134に記憶されたコードブックを用いて音響特徴量に戻す。この場合、元の音響特徴量を完全に復元することはできないが、符号列をある程度の長さにしておけば、十分な精度で音声認識を行なうことができる。

【0059】

音声認識処理部142は、デコーダ136が出力する各フレームの音響特徴量に基づいて、さらにMFCC係数の差分(「 Δ 」と呼ぶ。)を算出して、12次のMFCCとそれらの差分、C0、及びパワーからなる26次元の音響特徴量ベクトルを生成する。音声認識処理部142は、このようにして生成された音響特徴量ベクトルの列に対し、音響モデル記憶部138に記憶された音響モデル、クラス言語モデル記憶部140に記憶されたクラス言語モデル、及び単語言語モデル記憶部144に記憶された単語言語モデルを用いて音声認識処理を実行する。音声認識処理部142はこの際、未知語部分については、クラス言語モデル記憶部140によって算出された尤度と、候補の音節列についてカテゴリ別音節モデル記憶部146によって算出された尤度とを乗算することによって候補の音節列の尤度を算出し、最尤の音節列を表すカタカナ列を未知語に対応する音声認識結果として出力する。音声認識処理部142は、この未知語部分には、未知語を示すタグと、さらに、最大尤度を与えた音節モデルのカテゴリを示すタグとを付して出力する。なおこのとき、音声認識処理部142は、各単語及び未知語部分の各音節について、その開始時間と終了時間とからなる時間情報を付す。

【0060】

音声認識処理部142の音声認識結果は、未知語部分を含む場合も未知語部分を含まない場合も出力バッファ148(図4)に一旦格納される。

【0061】

送信処理部150は、出力バッファ148に格納されたテキストを携帯電話機20に送信する。

【0062】

再び図2を参照して、受信処理部62は、音声認識サーバ22から音声認識結果のテキストを受信すると、時間情報とともに受信バッファ64に格納する。未知語処理部70は、受信バッファ64に格納されたテキスト中に未知語を示すタグが存在するか否かを判定する。判定部96は、もしも未知語を示すタグがなければ、FALSEの判定結果信号を出力する。その結果、未知語切出処理部90、未知語認識処理部92、及び未知語入替処理部94は動作せず、選択部98は判定結果信号がFALSEであるため、第1の入力に与えられているテキスト、すなわち受信バッファ64に記憶されている音声認識結果のテキストを選択して出力部72に与える。出力部72は、携帯電話機20上で動作しているアプリケーション(現在の説明ではメールアプリケーション)にこのテキストを渡す。ア

10

20

30

40

50

アプリケーションはこのテキストを、キーボードから入力されたものと同様の入力として取り扱う。

【 0 0 6 3 】

もしも受信バッファ 6 4 に記憶された音声認識結果のテキスト中に、未知語を示すタグが付された部分があれば、判定部 9 6 は判定結果信号を TRUE とする。未知語切出処理部 9 0 はこの判定結果信号にตอบสนองして、受信バッファ 6 4 に記憶されたテキストの中の、未知語部分の開始時間及び終了時間を参照して、対応する符号列を符号記憶部 6 0 から読出し、未知語認識処理部 9 2 に与える。

【 0 0 6 4 】

未知語認識処理部 9 2 は、未知語切出処理部 9 0 から与えられた符号列の各々の符号を、コードブックメモリ 5 2 に記憶されたコードブックを使用して音響特徴量ベクトルに伸長し、符号列に戻す。すなわち、未知語認識処理部 9 2 は、圧縮時（符号化時）に対応する伸長アルゴリズムを用いて音響特徴量ベクトルを復元する。また、未知語認識処理部 9 2 もサーバと同様に M F C C のデルタを算出する。ただし、コードブックを用いているため、ここでの復元は完全な復元ではない。

【 0 0 6 5 】

未知語認識処理部 9 2 はさらに、受信バッファ 6 4 中の未知語部分に付されている、カテゴリを現すタグを読出し、言語モデル記憶部 6 6 に記憶されているカテゴリ別言語モデルのうちで、タグに対応するものを選択する。未知語認識処理部 9 2 は、このようにして選択されたカテゴリ別言語モデルと、音響モデル記憶部 6 8 に記憶された音響モデルとを使用して未知語の音声認識を行ない、認識結果の単語を未知語入替処理部 9 4 に与える。未知語認識処理部 9 2 での音声認識では、この携帯電話機 2 0 の利用者に特に関連して、各種アプリケーションから抽出された固有名詞が音声認識結果の単語の候補となる。その結果、利用者が発話した確率の高い固有名詞が未知語の音声認識結果として得られる可能性が大きくなる。

【 0 0 6 6 】

未知語入替処理部 9 4 は、受信バッファ 6 4 に記憶されたテキストのうち、未知語のタグが付された音節列を、未知語認識処理部 9 2 による音声認識の結果得られた単語で置換し、選択部 9 8 の第 2 の入力に与える。選択部 9 8 は、判定部 9 6 からの判定結果信号が TRUE であるため、未知語入替処理部 9 4 から与えられたテキストを選択し、出力部 7 2 に与える。出力部 7 2 にテキストが与えられた後の携帯電話機 2 0 の動作は、音声認識サーバ 2 2 からの音声認識結果のテキストに未知語が含まれていない場合と同様である。

【 0 0 6 7 】

< 例 >

図 6 に、この実施の形態に係る音声認識システム 1 0 による音声認識の例を模式的に示す。図 6 を参照して、「私の名前は松田です」という音声に対する音声認識処理が携帯電話機 2 0 で実行されるものとする。この携帯電話機 2 0 がこの音声の符号列を音声認識サーバ 2 2 に送信した後、音声認識サーバ 2 2 から受信したテキスト 2 0 0 が「私の名前はマツウダです」であったものとする。このテキストでは、本来は「松田」であった部分が、サーバでの音声認識では未知語として認識されている。すると、音声認識サーバ 2 2 から送信されてきたテキスト 2 0 0 のうち、「マツウダ」という音節列 2 0 4 の部分には、未知語を示すタグ 2 0 6 と、そのカテゴリとして日本人の「姓」を示すタグ 2 0 8 とが付されている。

【 0 0 6 8 】

携帯電話機 2 0 では、符号記憶部 6 0 に記憶されている符号列 2 0 2 のうち、未知語を示すタグ 2 0 6 が付されている音節列「マツウダ」に対応する部分符号列 2 1 0 を切出し、部分符号列 2 1 0 をコードブックを参照して伸長することで音響特徴量に戻し、未知語認識処理部 9 2 で行なわれる音声認識の入力とする。

【 0 0 6 9 】

一方、「姓」を示すタグ 2 0 8 に対応するカテゴリ言語モデル、具体的には姓言語モデ

10

20

30

40

50

ル 2 1 4 が音声認識における言語モデルとして選択される。この姓言語モデル 2 1 4 には、「マツウダ」という姓はなく、例えば「松井」、「松田」、「松山」等という姓が存在しているものとする。音声認識の結果、「マツウダ」ではなく正しい「松田」という単語 2 1 2 が選択される可能性が高い。

【 0 0 7 0 】

このように携帯電話機 2 0 での音声認識処理で正しい固有名詞が選択される可能性が高いのは、この携帯電話機 2 0 の使用者に特に関連した固有名詞のみを主に集め、それらをさらにカテゴリに分類してカテゴリ別言語モデルを作成しているためである。すなわち、使用者に関連のない固有名詞などが言語モデル中に含まれないため、使用者の発話に含まれる固有名詞に関する音声認識率が高くなる。また、音声認識を行なうために必要なリソースの量も少なく済むという効果がある。

10

【 0 0 7 1 】

図 7 は、上記実施の形態に係る携帯電話機 2 0 のハードウェア構成をブロック図形式で示す。図 7 を参照して、携帯電話機 2 0 は、スピーカ 2 3 6 と、図 2 にも示したマイクロフォン 5 0 と、液晶表示装置 (L C D) 2 3 8 と、テンキー及び特殊キーなどを含むキーパッド 2 4 0 と、アンテナ 2 3 2 と、着信及びアラームなどを振動により利用者に報知するための振動部 2 4 2 と、着信及びアラームなどを音声により利用者に報知するためのリング 2 4 6 と、携帯電話機 2 0 の初期設定値、カテゴリ別言語モデル、音響モデル、及び種々のアプリケーションプログラム等を記憶するための不揮発性で書換可能なメモリ 2 4 4 と、スピーカ 2 3 6、マイクロフォン 5 0、L C D 2 3 8、アンテナ 2 3 2、振動部 2 4 2、リング 2 4 6 及びメモリ 2 4 4 を用い、携帯電話機としての機能と、複数のアプリケーションを起動し、それらの出力を L C D 2 3 8 の表示面上に表示したり、キーパッド 2 4 0 からのユーザ入力を受けたりする機能とを実現するための制御回路 2 3 0 とを含む。

20

【 0 0 7 2 】

制御回路 2 3 0 は、アンテナ 2 3 2 を介して基地局から受信した信号に基づき、他の携帯通信端末からの着信を検出して着信検出信号を出力するための着信信号検出部 2 7 0 と、回線制御信号に応答して、アンテナ 2 3 2 を介した通信回線のオン/オフを制御するための回線閉結部 2 6 8 と、回線閉結部 2 6 8 及びアンテナ 2 3 2 を介して基地局との間で授受する信号の強度を制御するための R F (R a d i o F r e q u e n c y) 処理部 2 6 4 と、基地局との信号の授受を安全に行なうために、R F 処理部 2 6 4 に与える信号及び R F 処理部 2 6 4 を介して受ける信号に所定の信号処理を施すためのベースバンド処理部 2 6 2 と、D A コンバータ及び A D コンバータを有し、マイクロフォン 5 0 及びスピーカ 2 3 6 を介した音声の入出力を行なうためのオーディオインタフェース (オーディオ I / F) 2 6 0 と、オーディオ I / F 2 6 0、ベースバンド処理部 2 6 2、R F 処理部 2 6 4、回線閉結部 2 6 8、L C D 2 3 8、振動部 2 4 2、及びリング 2 4 6 を制御することにより、ユーザからの要求に応じて発呼したり、着呼を処理したりして、ユーザと他の携帯通信端末との間の音声通信を行なったり、文字通信を行なったり、ユーザの入力する文字列に対する処理を行なったりするための通信制御部 2 7 2 とを含む。

30

【 0 0 7 3 】

通信制御部 2 7 2 の機能は、実質的にはプロセッサとソフトウェアとにより実現される。ソフトウェアは本実施の形態ではメモリ 2 4 4 に記憶されていて、適宜通信制御部 2 7 2 内の図示しないメモリに読出され、実行される。本実施の形態では、詳細は説明しないが、メモリ 2 4 4 の内容を書き換えることが可能であり、それによって携帯電話機 2 0 による種々の機能のアップグレード及び追加を行なうことができる。通信制御部 2 7 2 はまた、本実施の形態に係る携帯電話機 2 0 の音声認識のためのプログラムを実行する。

40

【 0 0 7 4 】

以上のように本実施の形態に係る音声認識システム 1 0 によれば、携帯電話機 2 0 では音声認識の前処理に相当する特徴量の抽出が行なわれる。得られた特徴量ベクトルはコードブックを用いて符号化されて記憶されるとともに、サーバ 2 2 に送信される。音声認識

50

サーバ 2 2 は、この符号列を同じコードブックを用いて特徴量に戻した上で、音声認識サーバ 2 2 に準備された豊富なリソースを使用して音声認識を行なう。音声認識の処理中に未知語に遭遇すると、音声認識サーバ 2 2 は、その未知語を構成する音節列中の音節の遷移と予め準備されたカテゴリ別の音節モデルとに基づき、その未知語がどのカテゴリに属するかを推定し、未知語部分に未知語を示すタグとカテゴリを示すタグとを付して音声認識結果のテキスト中に挿入する。音声認識サーバ 2 2 は、音声認識結果のテキストを携帯電話機 2 0 に送信する。

【 0 0 7 5 】

携帯電話機 2 0 では、このテキスト中に未知語があった場合、記憶されていた符号列の内、対応する部分を読み出して特徴量に戻して音声認識を行なう。この音声認識では、言語モデルとして未知語に付されていたカテゴリに対応するカテゴリ別言語モデルが使用される。

10

【 0 0 7 6 】

携帯電話機 2 0 において作成されたカテゴリ別言語モデルは、特にこの携帯電話機 2 0 の使用者に関連する固有名詞から作成されている。その結果、音声認識サーバ 2 2 では未知語として認識された単語でも、携帯電話機 2 0 では利用者に特に関連する固有名詞として正しく認識される可能性が大きい。また、こうして言語モデルは、利用者に関連してアプリケーションによって集積された情報から作成されるものであり、その量が際限なく大きくなる可能性は極めて低い。そのため、携帯電話機 2 0 に準備すべきリソースの量が際限なく大きくなるという心配もない。

20

【 0 0 7 7 】

その結果、分散型の音声認識を利用する情報処理端末であって、利用者にとって音声認識の精度が十分に高く、かつ音声認識を行なうサーバ側のリソースの極端な肥大化を防止できる音声認識機能付情報処理端末を提供できる。

【 0 0 7 8 】

なお、図示していないが音声認識サーバ 2 2 側では、音声認識サービスを携帯電話機 2 0 に対して提供するにあたって、利用者ごと（または携帯電話機 2 0 ごと）に課金処理を行なうことが可能であることはいうまでもない。

【 0 0 7 9 】

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のすべての変更を含む。

30

【 図面の簡単な説明 】

【 0 0 8 0 】

【 図 1 】 本発明の一実施の形態に係る音声認識システム 1 0 におけるデータの流れの概略を示す図である。

【 図 2 】 音声認識システム 1 0 で使用される携帯電話機 2 0 の機能ブロック図である。

【 図 3 】 カテゴリ別言語モデル作成部 1 0 0 の機能ブロック図である。

【 図 4 】 音声認識サーバ 2 2 の機能ブロック図である。

40

【 図 5 】 携帯電話機 2 0 において、音声認識サーバ 2 2 から音声認識結果を受けた後の未知語の音声認識及び入替処理を実現するプログラムのフローチャートである。

【 図 6 】 音声認識システム 1 0 による音声認識の過程の一例を示す図である。

【 図 7 】 携帯電話機 2 0 のハードウェアブロック図である。

【 符号の説明 】

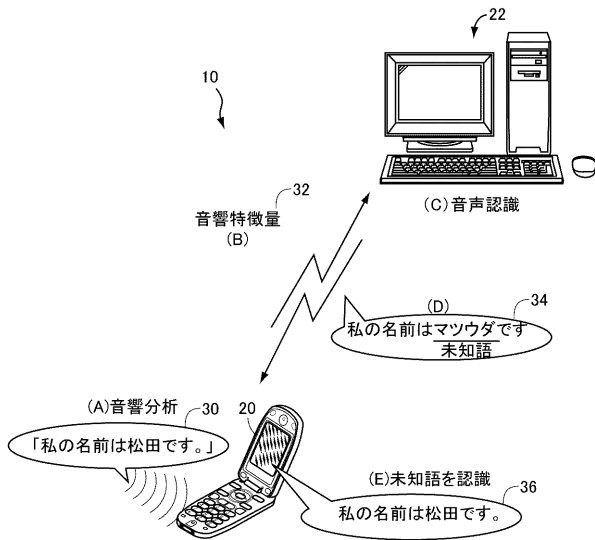
【 0 0 8 1 】

- 1 0 音声認識システム
- 2 0 携帯電話機
- 2 2 音声認識サーバ
- 5 2 コードブックメモリ

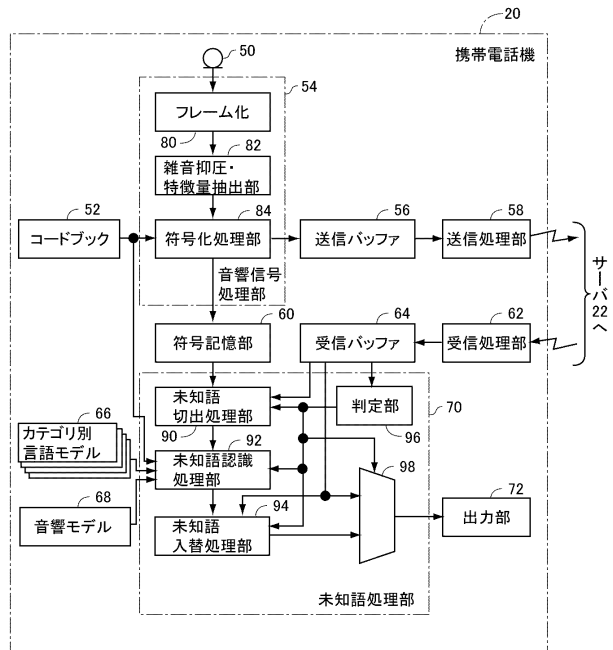
50

- 5 4 音響信号処理部
- 6 0 符号記憶部
- 6 6 言語モデル記憶部
- 6 8 音響モデル記憶部
- 7 0 未知語処理部
- 9 0 未知語切出処理部
- 9 2 未知語認識処理部
- 9 4 未知語入替処理部
- 9 6 判定部

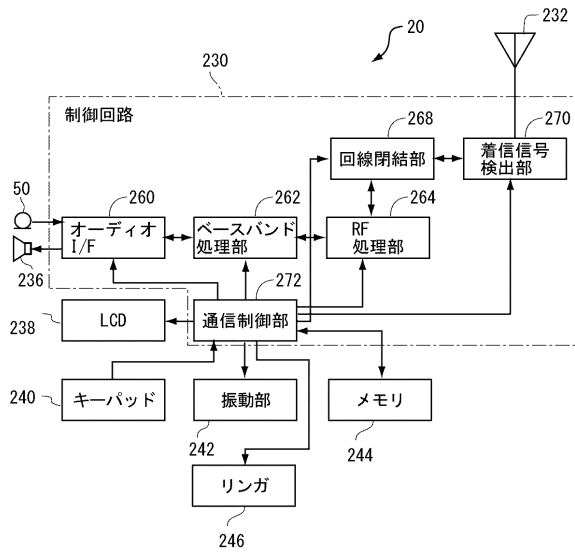
【図1】



【図2】



【図7】



フロントページの続き

- (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
- (72)発明者 林 輝昭
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

審査官 間宮 嘉誉

- (56)参考文献 特開2008-9153(JP,A)
特開2003-186494(JP,A)
特開平4-188200(JP,A)
特開2004-309523(JP,A)
特開2007-213109(JP,A)
国際公開第2005/122144(WO,A1)
特開平4-49719(JP,A)
特開2001-175286(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G10L 15/00 - 15/28