

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2015-11170

(P2015-11170A)

(43) 公開日 平成27年1月19日(2015.1.19)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 O L 15/30 (2013.01)	G 1 O L 15/28 2 1 O A	
G 1 O L 15/00 (2013.01)	G 1 O L 15/00 2 O O A	
G 1 O L 15/10 (2006.01)	G 1 O L 15/10 2 O O W	

審査請求 未請求 請求項の数 6 O L (全 16 頁)

(21) 出願番号 特願2013-136306 (P2013-136306)
 (22) 出願日 平成25年6月28日 (2013. 6. 28)

(71) 出願人 307041344
 株式会社A T R - T r e k
 神奈川県川崎市川崎区砂子2-4-10
 (74) 代理人 100099933
 弁理士 清水 敏
 (72) 発明者 古谷 利昭
 大阪府大阪市淀川区西中島6丁目1番1号

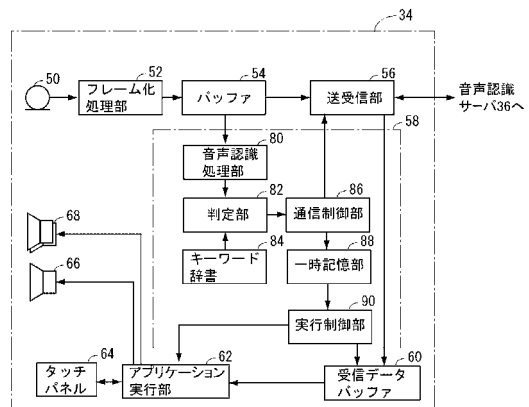
(54) 【発明の名称】 ローカルな音声認識を行なう音声認識クライアント装置

(57) 【要約】

【課題】 ローカルにも音声認識機能を持ち、音声認識サーバの音声認識機能の起動を自然に行なえ、通信回線の負荷を抑えながら精度も高く維持できるクライアントを提供する。

【解決手段】 音声認識クライアント装置 3 4 は、音声認識サーバ 3 6 との通信により、音声認識サーバ 3 6 による音声認識結果を受信するクライアントであり、音声を音声データに変換するフレーム化処理部 5 2 と、音声データに対する音声認識を行なうローカルな音声認識処理部 8 0 と、音声データを音声認識サーバに送信し、当該音声認識サーバによる音声認識結果を受信する送受信部 5 6 と、音声データに対する音声認識処理部 8 0 の認識結果により、送受信部 5 6 による音声データの送信を制御する判定部 8 2 及び通信制御部 8 6 とを含む。

【選択図】 図 2



【特許請求の範囲】**【請求項 1】**

音声認識サーバとの通信により、当該音声認識サーバによる音声認識結果を受信する音声認識クライアント装置であって、

音声データを音声データに変換する音声変換手段と、

前記音声データに対する音声認識を行なう音声認識手段と、

前記音声データを前記音声認識サーバに送信し、当該音声認識サーバによる音声認識結果を受信する送受信手段と、

前記音声データに対する前記音声認識手段の認識結果により、前記送受信手段による音声データの送信を制御する送受信制御手段とを含む、音声認識クライアント装置。

10

【請求項 2】

前記送受信制御手段は、

前記音声認識手段による音声認識結果中にキーワードが存在することを検出して、検出信号を出力するキーワード検出手段と、

前記検出信号に応答して、前記音声データのうち、前記キーワードの発話区間の先頭と所定の関係にある部分を前記音声認識サーバに送信するよう前記送受信手段を制御する送信開始制御手段とを含む、請求項 1 に記載の音声認識クライアント装置。

【請求項 3】

前記送信開始制御手段は、前記検出信号に応答して、前記音声データのうち、前記キーワードの発話終了位置を先頭とする部分を前記音声認識サーバに送信するよう前記送受信手段を制御する手段を含む、請求項 2 に記載の音声認識クライアント装置。

20

【請求項 4】

前記送信開始制御手段は、前記検出信号に応答して、前記音声データのうち、前記キーワードの発話開始位置を先頭とする部分を前記送受信手段を制御する手段を含む、請求項 2 に記載の音声認識クライアント装置。

【請求項 5】

前記送受信手段が受信した前記音声認識サーバによる音声認識結果の先頭部分が、前記キーワード検出手段が検出したキーワードと一致するか否かを判定する一致判定手段と、

前記一致判定手段による判定結果にしたがって、前記送受信手段が受信した前記音声認識サーバによる音声認識結果を利用する処理と、前記音声認識サーバによる音声認識結果を破棄する処理とを選択的に実行する手段とをさらに含む、請求項 4 に記載の音声認識クライアント装置。

30

【請求項 6】

前記送受信制御手段は、

前記音声認識手段による音声認識結果中に第 1 のキーワードが存在することを検出して第 1 の検出信号を、何らかの処理を依頼することを表す第 2 のキーワードが存在することを検出して第 2 の検出信号を、それぞれ出力するキーワード検出手段と、

前記第 1 の検出信号に応答して、前記音声データのうち、前記第 1 のキーワードの発話区間の先頭と所定の関係にある部分を前記音声認識サーバに送信するよう前記送受信手段を制御する送信開始制御手段と、

40

前記送受信手段により前記音声データの送信が開始された後に前記第 2 の検出信号が発生されたことに応答して、前記音声データの第 2 のキーワードの発話の終了位置で前記送受信手段による音声データの送信を終了させる送信終了制御手段とを含む、請求項 1 に記載の音声認識クライアント装置。

【発明の詳細な説明】**【技術分野】****【0001】**

この発明は音声認識サーバと通信することにより音声認識機能を備えた音声認識クライアント装置に関し、特に、サーバとは別にローカルな音声認識機能を備えた音声認識クライアント装置に関する。

50

【背景技術】

【0002】

ネットワークに接続される携帯電話等の携帯型端末装置の数が爆発的に増加している。携帯型端末装置は、事実上、小型のコンピュータとすることができる。特に、いわゆるスマートフォン等では、インターネット上のサイトの検索、音楽・ビデオの視聴、メールの交換、銀行取引、スケッチ、録音・録画等、デスクトップコンピュータと同等の充実した機能が利用できる。

【0003】

しかしこのように充実した機能を利用するための1つのネックが、携帯型端末装置の筐体の小ささである。携帯型端末装置はその宿命として筐体が小さい。そのため、コンピュータのキーボードのように高速に入力をするためのデバイスを搭載することができない。タッチパネルを使用した様々な入力方式が考えられており、以前と比較して素早く入力できるようにはなっているが、依然として入力はそれほど容易でない。

10

【0004】

こうした状況で入力のための手段として注目されているのが音声認識である。音声認識の現在の主流は、多数の音声データを統計的に処理して作成した音響モデルと、大量の文書から得た統計的言語モデルとを使用する統計的音声認識装置である。こうした音声認識装置は、非常に大きな計算パワーを必要とするため、大容量で計算能力が十分に高いコンピュータでのみ実現されていた。携帯型端末装置で音声認識機能を利用する場合には、音声認識サーバと呼ばれる、音声認識機能をオンラインで提供するサーバが利用され、携帯型端末装置はその結果を利用する音声認識クライアントとして動作する。音声認識クライアントが音声認識をする際には、音声データをローカルに処理して得た音声データ、符号データ、又は音声の特徴量(素性)を音声認識サーバにオンラインで送信し、音声認識結果を受け取ってそれに基づいた処理を行なっている。これは、携帯型端末装置の計算能力が比較的 low、利用できる計算資源も限られていたためである。

20

【0005】

しかし、半導体技術の進歩により、CPU(Central Processing Unit)の計算能力は非常に高くなり、また、メモリ容量も従来と比較して桁違いに大きくなってきた。しかも消費電力は少なくなっている。そのため、携帯型端末装置でも音声認識が十分に利用可能となっている。しかも、携帯型端末装置では使用するユーザが限定されるため、音声認識の話者を予め特定し、その話者に適合した音響モデルを準備したり、特定の語彙を辞書に登録したりすることで、音声認識の精度を高めることができる。

30

【0006】

もっとも、利用できる計算資源の点では音声認識サーバの方が圧倒的に有利であるため、音声認識の精度の点では、携帯型端末装置よりも音声認識サーバで行なわれる音声認識の方が優れている点は間違いない。

【0007】

このように、携帯型端末装置に搭載される音声認識の精度が比較的 low、という欠点を補うための提案が、後掲の特許文献1に開示されている。特許文献1は音声認識サーバと通信するクライアントに関する。このクライアントは、音声データを処理して音声データに変換し、音声認識サーバに送信する。音声認識サーバからその音声認識結果を受信すると、その音声認識結果には、文節の区切り位置、文節の属性(文字種)、単語の品詞、文節の時間情報等が付されている。クライアントは、サーバからの音声認識結果に付されているこのような情報を利用して、ローカルに音声認識を行なう。この際、ローカルに登録されている語彙又は音響モデルを使用できるので、語彙によっては音声認識サーバで誤って認識された語を正しく認識できる可能性がある。

40

【0008】

特許文献1に開示されたクライアントでは、音声認識サーバからの音声認識結果と、ローカルに行なった音声認識結果とを比較し、両者の認識結果が異なった箇所についてはユーザによりいずれかを選択させる。

50

【先行技術文献】

【特許文献】

【0009】

【特許文献1】特開2010-85536号公報、特に段落0045～0050、図4

【発明の概要】

【発明が解決しようとする課題】

【0010】

特許文献1に開示されたクライアントは、音声認識サーバによる認識結果をローカルな音声認識結果で補完できるという優れた効果を奏する。しかし、現在の携帯型端末装置における音声認識の利用方法を見てみると、こうした機能を持つ携帯型端末の操作に関して

10

【0011】

特許文献1には、ローカルでどのようにして音声認識を開始するかについての開示はない。現在利用可能な携帯型端末装置では、音声認識を開始するためのボタンを画面に表示させ、このボタンがタッチされたら音声認識機能を起動するものが主流である。又は、音声認識を開始させるための専用のハードウェアボタンを設けたものもある。ローカルな音声認識機能を持たない携帯電話で動作するアプリケーションの中には、ユーザが発話姿勢をとったとき、すなわち携帯電話を耳にあてたときをセンサで感知し、音声入力とサーバへの音声データの送信とを開始するものもある。

20

【0012】

しかし、これらはいずれも音声認識機能を起動するにあたって特定の動作をユーザに要求するものである。これからの携帯型端末装置では、多様な機能を利用するために、音声認識機能を従来以上に活用することが予測され、そのためには音声認識機能の起動をより自然なものにする必要がある。一方で、携帯型端末装置と音声認識サーバとの間の通信量はできるだけ抑える必要があるし、音声認識の精度は高く維持する必要もある。

【0013】

それゆえにこの発明の目的は、音声認識サーバを利用するとともに、ローカルにも音声認識機能を持つ音声認識クライアント装置であって、音声認識機能の起動を自然に行なえ、通信回線の負荷を抑えながら音声認識の精度も高く維持できる音声認識クライアント装置を提供することである。

30

【課題を解決するための手段】

【0014】

本発明の第1の局面に係る音声認識クライアント装置は、音声認識サーバとの通信により、当該音声認識サーバによる音声認識結果を受信する音声認識クライアント装置である。この音声認識クライアント装置は、音声を音声データに変換する音声変換手段と、音声データに対する音声認識を行なう音声認識手段と、音声データを音声認識サーバに送信し、当該音声認識サーバによる音声認識結果を受信する送受信手段と、音声データに対する音声認識手段の認識結果により、送受信手段による音声データの送信を制御する送受信制御手段とを含む。

40

【0015】

ローカルな音声認識手段の出力に基づいて、音声データを音声認識サーバに送信するかが制御される。音声認識サーバを利用するためには、発話することを除き特別な操作は必要ない。音声認識手段の認識結果が特定のものでなければ音声認識サーバへの音声データの送信が行なわれない。

【0016】

その結果、本発明によれば、音声認識機能の起動を自然に行なえ、通信回線の負荷を抑えながら音声認識の精度も高く維持できる音声認識クライアント装置を提供できる。

【0017】

好ましくは、送受信制御手段は、音声認識手段による音声認識結果中にキーワードが存

50

在することを検出して、検出信号を出力するキーワード検出手段と、検出信号に応答して、音声データのうち、キーワードの発話区間の先頭と所定の関係にある部分を音声認識サーバに送信するよう送受信手段を制御する送信開始制御手段とを含む。

【0018】

ローカルな音声認識手段の音声認識結果中にキーワードが検出されると、音声データの送信が開始される。音声認識サーバの音声認識を利用するために、特別なキーワードを発話するだけでよく、ボタンを押す等、音声認識を開始するための明示的な操作をする必要がない。

【0019】

より好ましくは、送信開始制御手段は、検出信号に応答して、音声データのうち、キーワードの発話終了位置を先頭とする部分を音声認識サーバに送信するよう送受信手段を制御する手段を含む。

10

【0020】

キーワードの次の部分から音声認識サーバに音声データを送信することにより、キーワード部分の音声認識を音声認識サーバでは行わずに済む。音声認識結果にキーワードが含まれないため、キーワードに続けて発話した内容に関する音声認識結果をそのまま利用できる。

【0021】

さらに好ましくは、送信開始制御手段は、検出信号に応答して、音声データのうち、キーワードの発話開始位置を先頭とする部分を送信するよう送受信手段を制御する手段を含む。

20

【0022】

キーワードの発話開始位置を先頭として音声認識サーバに送ることにより、音声認識サーバで再びキーワード部分の確認を行ったり、音声認識サーバの音声認識結果を利用して携帯型端末でローカルな音声認識の結果の正確さを検証したりできる。

【0023】

音声認識クライアント装置は、送受信手段が受信した音声認識サーバによる音声認識結果の先頭部分が、キーワード検出手段が検出したキーワードと一致するか否かを判定する一致判定手段と、一致判定手段による判定結果にしたがって、送受信手段が受信した音声認識サーバによる音声認識結果を利用する処理と、音声認識サーバによる音声認識結果を破棄する処理とを選択的に実行する手段とをさらに含む。

30

【0024】

ローカルな音声認識結果と、音声認識サーバによる音声認識結果とが異なる場合、より精度が高いと思われる音声認識サーバの結果を用いて発話者の発話を処理するか否かを判定する。ローカルな音声認識結果が誤っている場合には、音声認識サーバの音声結果は何ら利用されず、携帯型端末は何事もなかったように動作する。したがって、ローカルな音声認識による音声認識結果の誤りにより、ユーザの意図しないような処理を音声認識クライアント装置が実行することが予防できる。

【0025】

好ましくは、送受信制御手段は、音声認識手段による音声認識結果中に第1のキーワードが存在することを検出して第1の検出信号を、何らかの処理を依頼することを表す第2のキーワードが存在することを検出して第2の検出信号を、それぞれ出力するキーワード検出手段と、第1の検出信号に応答して、音声データのうち、第1のキーワードの発話区間の先頭と所定の関係にある部分を音声認識サーバに送信するよう送受信手段を制御する送信開始制御手段と、送受信手段により音声データの送信が開始された後に第2の検出信号が発生されたことに応答して、音声データの第2のキーワードの発話の終了位置で送受信手段による音声データの送信を終了させる送信終了制御手段とを含む。

40

【0026】

音声データを音声認識サーバに送信するにあたり、ローカルな音声認識手段による音声認識結果に第1のキーワードが検出されたときには、その第1のキーワードの発話開始位

50

置と所定の関係にある部分の音声データが音声認識サーバに送信される。その後、ローカルな音声認識手段による音声認識結果に、何らかの処理を依頼することを表す第2のキーワードが検出されたときには、それ以後の音声データの送信は行なわれない。音声認識サーバを利用するにあたり、第1のキーワードを発話するのみでよいだけでなく、第2のキーワードを発話することにより音声データの送信をその時点で終了できる。発話の終了を検知するために所定の無音区間を検出したりする必要はなく、音声認識のレスポンスを向上させることができる。

【図面の簡単な説明】

【0027】

【図1】本発明の第1の実施の形態に係る音声認識システムの概略構成を示すブロック図である。

10

【図2】第1の実施の形態に係る携帯端末装置である携帯電話の機能的ブロック図である。

【図3】逐次方式の音声認識の出力の仕方の概略を説明する模式図である。

【図4】第1の実施の形態において、音声認識サーバへの音声データの送信開始及び送信終了タイミングと送信内容とを説明するための模式図である。

【図5】第1の実施の形態において、音声認識サーバへの音声データの送信開始及び終了を制御するプログラムの制御構造を示すフローチャートである。

【図6】第1の実施の形態において、音声認識サーバの結果とローカルな音声認識結果とを利用して携帯型端末装置を制御するプログラムの制御構造を示すフローチャートである。

20

【図7】本発明の第2の実施の形態に係る携帯型端末装置である携帯電話の機能的ブロック図である。

【図8】第2の実施の形態において、音声認識サーバへの音声データの送信開始及び送信終了タイミングと送信内容とを説明するための模式図である。

【図9】第2の実施の形態において、音声認識サーバへの音声データの送信開始及び終了を制御する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】第1及び第2の実施の形態に係る装置の構成を示すハードウェアブロック図である。

【発明を実施するための形態】

30

【0028】

以下の説明及び図面では、同一の部品には同一の参照番号を付してある。したがって、それらについての詳細な説明は繰返さない。

【0029】

< 第1の実施の形態 >

[概略]

図1を参照して、第1の実施の形態に係る音声認識システム30は、ローカルな音声認識機能を持つ音声認識クライアント装置である携帯電話34と、音声認識サーバ36とを含む。両者はインターネット32を介して相互に通信可能である。この実施の形態では、携帯電話34はローカルな音声認識の機能を持ち、音声認識サーバ36との間の通信量を抑えながら、自然な形でユーザによる操作に対する応答を実現する。なお、以下の実施の形態では、携帯電話34から音声認識サーバ36に送信される音声データは音声信号をフレーム化したデータであるが、例えば音声信号を符号化した符号化データでもよいし、音声認識サーバ36で行なわれる音声認識処理で使用される特徴量でもよい。

40

【0030】

[構成]

図2を参照して、携帯電話34は、マイクロフォン50と、マイクロフォン50から出力される音声信号をデジタル化し、所定フレーム長及び所定シフト長でフレーム化するフレーム化処理部52と、フレーム化処理部52の出力である音声データを一時的に蓄積するバッファ54と、バッファ54に蓄積された音声データを音声認識サーバ36に送信す

50

る処理と、音声認識サーバ36からの音声認識結果等を含むネットワークからのデータを無線により受信する送受信部56とを含む。フレーム化処理部52の出力する各フレームには、各フレームの時間情報が付されている。

【0031】

携帯電話34はさらに、バッファ54に蓄積された音声データによるローカルな音声認識をバックグラウンドで行ない、音声認識結果の中に所定のキーワードが検出されたことに応答して、送受信部56による音声認識サーバ36への音声信号の送信開始及び送信終了を制御する処理と、音声認識サーバからの受信結果とローカルな音声認識の結果とを照合し、その結果にしたがって携帯電話34の動作を制御するための制御部58と、送受信部56が音声認識サーバ36から受信した音声認識結果を一時的に蓄積する受信データバッファ60と、ローカルな音声認識結果と音声認識サーバ36からの音声認識結果との照合に基づいて制御部58が実行指示信号を発生したことに応答して、受信データバッファ60の内容を用いたアプリケーションを実行するアプリケーション実行部62と、アプリケーション実行部62に接続されたタッチパネル64と、アプリケーション実行部62に接続された受話用のスピーカ66と、同じくアプリケーション実行部62に接続されたステレオスピーカ68とを含む。

【0032】

制御部58は、バッファ54に蓄積された音声データに対してローカルな音声認識処理を実行する音声認識処理部80と、音声認識処理部80の出力する音声認識結果に、音声認識サーバ36への音声データの送受信を制御するための所定のキーワード(開始キーワード及び終了キーワード)が含まれているか否かを判定し、含まれている場合には検出信号をそのキーワードとともに出力する判定部82と、判定部82が判定の対象とする開始キーワードを1又は複数個記憶するキーワード辞書84とを含む。なお、音声認識処理部80は、無音区間が所定のしきい値時間以上続くと発話が終了したとみなし、発話終了検出信号を出力する。判定部82は、発話終了検出信号を受信すると、通信制御部86に対して音声認識サーバ36へのデータの送信を終了する指示を出すものとする。

キーワード辞書84に記憶される開始キーワードは、通常の発話とできるだけ区別するために、名詞を用いるものとする。携帯電話34に何らかの処理を依頼することを考えると、この名詞としては特に固有名詞を使用することが自然であり好ましい。固有名詞でなく、特定のコマンド用語を用いるようにしてもよい。

終了キーワードとしては、日本語の場合には、開始キーワードとは異なり、より一般的に動詞の命令形、動詞の基本形+終止形、依頼表現、又は疑問表現等、通常の日本語で他人に何かを依頼する表現を採用する。すなわち、これらのいずれかを検出したときに、終了キーワードを検出したものと判定する。こうすることにより、ユーザが自然な話し方で携帯電話に処理を依頼することが可能になる。こうした処理を可能とするためには、音声認識処理部80が、認識結果の各単語にその単語の品詞、動詞の活用形、助詞の種類等を示す情報を付すようなものであればよい。

【0033】

制御部58はさらに、判定部82から検出信号と検出されたキーワードとを受信したことに応答し、検出されたキーワードが開始キーワードか終了キーワードかにしたがって、バッファ54に蓄積された音声データを音声認識サーバ36に送信する処理を開始又は終了するための通信制御部86と、判定部82が音声認識処理部80による音声認識結果内に検出したキーワードのうち、開始キーワードを記憶する一時記憶部88と、受信データバッファ60が受信した音声認識サーバ36の音声認識結果のテキストの先頭部分と、一時記憶部88に記憶された、ローカル音声認識結果の開始キーワードとを比較し、両者が一致したときには受信データバッファ60に記憶されたデータの内、開始キーワードの後に続く部分を使用して所定のアプリケーションを実行するようアプリケーション実行部62を制御するための実行制御部90とを含む。本実施の形態では、どのようなアプリケーションを実行するかはアプリケーション実行部62が受信データバッファ60に記憶された内容によって判定する。

【 0 0 3 4 】

音声認識処理部 8 0 が、バッファ 5 4 に蓄積された音声データに対する音声認識をするにあたり、音声認識結果を出力する仕方には 2 通りある。発話ごと方式と逐次方式とである。発話ごと方式は、音声データ内に所定時間を超える無音区間があったときに、それまでの音声の音声認識結果を出力し、次の発話区間から新たに音声認識を開始する。逐次方式は、随時バッファ 5 4 に蓄積されている音声データ全体に対する音声認識結果を所定時間間隔（たとえば 1 0 0 ミリ秒ごと）で出力する。したがって、発話区間が長くなると音声認識結果のテキストもそれにつれて長くなる。本実施の形態では、音声認識処理部 8 0 は逐次方式を採用している。なお、発話区間が非常に長くなると、音声認識処理部 8 0 による音声認識が困難になる。したがって音声認識処理部 8 0 は、発話区間が所定時間長以上になると、強制的に発話が終了したのものとしてそれまでの音声認識を終了し、新たな音声認識を開始するものとする。なお、音声認識処理部 8 0 による音声認識の出力が発話ごとの方式である場合でも、以下の機能は本実施の形態のものと同様に実現できる。

10

【 0 0 3 5 】

図 3 を参照して、音声認識処理部 8 0 の出力タイミングについて説明する。発話 1 0 0 が、第 1 の発話 1 1 0 と第 2 の発話 1 1 2 とを含み、両者の間に無音区間 1 1 4 があるものとする。音声認識処理部 8 0 は、バッファ 5 4 に音声データが蓄積されていくと、音声認識結果 1 2 0 で示されるように、1 0 0 ミリ秒ごとに、バッファ 5 4 に蓄積された音声全体に対する音声認識結果を出力する。この方式では、音声認識結果の一部が途中で修正される場合もある。例えば、図 3 に示す音声認識結果 1 2 0 の場合、2 0 0 ミリ秒時点で出力された「熱い」という単語が 3 0 0 ミリ秒時点では「暑い」に修正されている。この方式では、無音区間 1 1 4 の時間長が所定のしきい値より大きい場合には、発話が終了したものとみなされる。その結果、バッファ 5 4 に蓄積されていた音声データはクリアされ（読捨てられ）、次の発話に対する音声認識処理が開始される。図 3 の場合には、次の音声認識結果 1 2 2 が新たな時間情報とともに音声認識処理部 8 0 から出力される。判定部 8 2 は、音声認識結果 1 2 0 又は音声認識結果 1 2 2 等の各々について、音声認識結果が出力されるごとに、キーワード辞書 8 4 に記憶された開始キーワードのいずれかと一致しているか、又は終了キーワードの条件を充足しているか否かを判定し、開始キーワード検出信号又は終了キーワード検出信号を出力する。ただし、本実施の形態では、開始キーワードは音声認識サーバ 3 6 への音声データの送信が行なわれていないときにしか検出されず、終了キーワードは開始キーワードが検出された後でなければ検出されない。

20

30

【 0 0 3 6 】

〔動作〕

携帯電話 3 4 は以下のように動作する。マイクロフォン 5 0 は常に周囲の音声を検知して音声信号をフレーム化処理部 5 2 に与える。フレーム化処理部 5 2 は、音声信号をデジタル化及びフレーム化し、バッファ 5 4 に順次入力する。音声認識処理部 8 0 は、バッファ 5 4 に蓄積されていく音声データの全体について、1 0 0 ミリ秒ごとに音声認識を行ない、その結果を判定部 8 2 に出力する。ローカルな音声認識処理部 8 0 は、しきい値時間以上の無音区間を検知するとバッファ 5 4 をクリアし、発話の終了を検出したことを示す信号（発話終了検出信号）を判定部 8 2 に出力する。

40

【 0 0 3 7 】

判定部 8 2 は、音声認識処理部 8 0 からローカルな音声認識結果を受信すると、その中にキーワード辞書 8 4 に記憶された開始キーワードがあるか、又は終了キーワードとしての条件を充足する表現があるかを判定する。判定部 8 2 は、音声認識サーバ 3 6 に音声データを送信していない期間にローカルな音声認識結果内に開始キーワードを検出した場合、開始キーワード検出信号を通信制御部 8 6 に与える。一方、判定部 8 2 は、音声認識サーバ 3 6 に音声データを送信している間にローカルな音声認識結果内に終了キーワードを検出すると、終了キーワード検出信号を通信制御部 8 6 に与える。判定部 8 2 はまた、音声認識処理部 8 0 から発話終了検出信号を受信したときには、音声認識サーバ 3 6 への音声データの送信を終了するよう通信制御部 8 6 に対して指示を与える。

50

【 0 0 3 8 】

通信制御部 8 6 は、判定部 8 2 から開始キーワード検出信号が与えられると、送受信部 5 6 を制御してバッファ 5 4 に蓄積されているデータのうち、検出された開始キーワードの先頭位置からデータを読み出して、音声認識サーバ 3 6 に送信する処理を開始させる。このとき、通信制御部 8 6 は、判定部 8 2 から与えられた開始キーワードを一時記憶部 8 8 に保存する。通信制御部 8 6 は、判定部 8 2 から終了キーワード検出信号が与えられると、送受信部 5 6 を制御して、バッファ 5 4 に蓄積されているデータのうち、検出された終了キーワードまでの音声データを音声認識サーバ 3 6 に送信させた後に送信を終了させる。判定部 8 2 から発話終了検出信号による送信終了の指示が与えられると、通信制御部 8 6 は、送受信部 5 6 を制御して、バッファ 5 4 に記憶されている音声データのうち、発話の終了が検出された時間までの音声データを全て音声認識サーバ 3 6 に送信させた後に送信を終了させる。

10

【 0 0 3 9 】

受信データバッファ 6 0 は、通信制御部 8 6 によって音声認識サーバ 3 6 への音声データの送信が開始された後、音声認識サーバ 3 6 から送信されてくる音声認識結果のデータを蓄積する。実行制御部 9 0 は、受信データバッファ 6 0 の先頭部分が、一時記憶部 8 8 に保存されている開始キーワードと一致するか否かを判定する。両者が一致していると、実行制御部 9 0 は、アプリケーション実行部 6 2 を制御し、受信データバッファ 6 0 のうちで、開始キーワードと一致した部分の次からのデータを読み出すようにさせる。アプリケーション実行部 6 2 は、受信データバッファ 6 0 から読み出したデータに基づいてどのようなアプリケーションを実行するかを判定し、そのアプリケーションに音声認識結果を渡して処理させる。処理の結果は、例えばタッチパネル 6 4 に表示されたり、スピーカ 6 6 又はステレオスピーカ 6 8 から音声の形で出力されたりする。

20

【 0 0 4 0 】

例えば図 4 を参照して、具体的な例を説明する。ユーザが発話 1 4 0 を行なったものとする。発話 1 4 0 は、「v G a t e 君」という発話部分 1 5 0 と、「このあたりのラーメン屋さん調べて」という発話部分 1 5 2 とを含む。発話部分 1 5 2 は、「このあたりのラーメン屋さん」という発話部分 1 6 0 と、「調べて」という発話部分 1 6 2 とを含む。

【 0 0 4 1 】

ここでは、開始キーワードとして例えば「v G a t e 君」、「羊君」等が登録されているものとする。すると、発話部分 1 5 0 が開始キーワードと一致しているため、発話部分 1 5 0 が音声認識された時点で音声データ 1 7 0 を音声認識サーバ 3 6 に送信する処理が開始される。音声データ 1 7 0 は、図 4 に示すように発話 1 4 0 の音声データの全体を含み、その先頭は開始キーワードに対応する音声データ 1 7 2 である。

30

【 0 0 4 2 】

一方、発話部分 1 6 2 のうち、「調べて」という表現は依頼表現であり終了キーワードとしての条件を充足する。したがって、この表現がローカル音声認識結果中に検出された時点で、音声データ 1 7 0 を音声認識サーバ 3 6 に送信する処理は終了する。

【 0 0 4 3 】

音声データ 1 7 0 の送信が終了すると、音声データ 1 7 0 に対する音声認識結果 1 8 0 が音声認識サーバ 3 6 から携帯電話 3 4 に送信され、受信データバッファ 6 0 に蓄積される。音声認識結果 1 8 0 の先頭部分 1 8 2 は、開始キーワードに対応する音声データ 1 7 2 の音声認識結果である。この先頭部分 1 8 2 が、発話部分 1 5 0 (開始キーワード)に対するクライアント音声認識結果と一致すると、音声認識結果 1 8 0 の内、先頭部分 1 8 2 の次の部分からの音声認識結果 1 8 4 がアプリケーション実行部 6 2 (図 1 参照)に送信され、適切なアプリケーションにより処理される。先頭部分 1 8 2 が発話部分 1 5 0 (開始キーワード)に対するクライアント音声認識結果と一致していないと、受信データバッファ 6 0 はクリアされ、アプリケーション実行部 6 2 は何ら動作しない。

40

【 0 0 4 4 】

以上のようにこの実施の形態によれば、ローカル音声認識により発話中に開始キーワー

50

ドが検出されると音声データを音声認識サーバ36に送信する処理が開始される。ローカル音声認識により発話中に終了キーワードが検出されると、音声認識サーバ36への音声データの送信が終了される。音声認識サーバ36から送信されてくる音声認識結果の先頭部分と、ローカル音声認識により検出された開始キーワードとが比較され、両者が一致していれば、音声認識サーバ36の音声認識結果を用いて何らかの処理が実行される。したがって、この実施の形態では、携帯電話34に何らかの処理を実行させようとする場合、ユーザは他に何もせず、単に開始キーワードと実行内容を発話するだけでよい。開始キーワードがローカル音声認識で正しく認識されれば、携帯電話34による音声認識の結果を用いた所望の処理が実行され、結果が携帯電話34により出力される。音声入力の開始のためのボタンを押したりする必要はなく、携帯電話34をより簡単に使用できる。

10

【0045】

こうした処理で問題になるのは、開始キーワードが誤って検出された場合である。前述したように、一般的に、携帯型端末でローカルに実行される音声認識の精度は、音声認識サーバで実行される音声認識の精度よりも低い。したがってローカル音声認識で誤って開始キーワードが検出される可能性がある。そうした場合、誤って検出された開始キーワードに基づいて何らかの処理を実行し、その結果を携帯電話34が出力すると、それはユーザが意図しない動作となってしまう。そのような動作は好ましくない。

【0046】

本実施の形態では、仮にローカル音声認識で開始キーワードが誤検出されたとしても、音声認識サーバ36からの音声認識結果の先頭部分が開始キーワードと一致していなければ携帯電話34はその結果による処理は何も実行しない。携帯電話34の状態は何も変化せず、見かけ上全く何もしていないように見える。したがって、ユーザは、上に記載したような処理が実行されたことには全く気付かない。

20

【0047】

さらに、上記実施の形態では、開始キーワードがローカル音声認識で検出された場合に音声データを音声認識サーバ36に送信する処理を開始し、終了キーワードがローカル音声認識で検出された場合に送信処理を終了する。音声の送信を終了するためにユーザが特別な操作をする必要がない。所定時間以上の空白を検出したときに送信を終了する場合と比較して、終了キーワードを検出すると直ちに音声認識サーバ36への音声データの送信を終了できる。その結果、携帯電話34から音声認識サーバ36への無駄なデータ送信を防止できるし、音声認識のレスポンスも向上する。

30

【0048】

[プログラムによる実現]

上記第1の実施の形態に係る携帯電話34は、後述するような、コンピュータと同様の携帯電話ハードウェアと、その上のプロセッサにより実行されるプログラムとにより実現できる。図5に、図1の判定部82及び通信制御部86の機能を実現するプログラムの制御構造をフローチャート形式で示し、図6に、実行制御部90の機能を実現するプログラムの制御構造をフローチャート形式で示す。ここでは両者を別プログラムとして記載しているが、両者をまとめることもできるし、それぞれさらに細かい単位のプログラムに分割することもできる。

40

【0049】

図5を参照して、判定部82及び通信制御部86の機能を実現するプログラムは、携帯電話34の電源投入時に起動されると、使用するメモリアリアの初期化等を実行するステップ200と、システムからプログラムの実行を終了することを指示する終了信号を受信したか否かを判定し、終了信号を受信したときには必要な終了処理を実行してこのプログラムの実行を終わるステップ202と、終了信号が受信されていないときに、音声認識処理部80からローカル音声認識結果を受信したか否かを判定し、受信していなければ制御をステップ202に戻すステップ204とを含む。前述したとおり、音声認識処理部80は所定時間ごとに音声認識結果を逐次的に出力する。したがってステップ204の判定は、所定時間ごとにYESとなる。

50

【 0 0 5 0 】

このプログラムはさらに、ステップ 2 0 4 でローカル音声認識の結果を受信したと判定されたことに応答して、キーワード辞書 8 4 に記憶された開始キーワードのいずれかがローカル音声認識結果に含まれるか判定し、含まれていない場合には制御をステップ 2 0 2 に戻すステップ 2 0 6 と、開始キーワードのいずれかがローカル音声認識結果にあったときに、その開始キーワードを一時記憶部 8 8 に保存するステップ 2 0 8 と、バッファ 5 4 (図 2) に記憶されている音声データのうち、開始キーワードの先頭部分から音声認識サーバ 3 6 への音声データの送信を開始させるよう送受信部 5 6 に指示するステップ 2 1 0 とを含む。以後、処理は携帯電話 3 4 への音声データ送信中の処理に移る。

【 0 0 5 1 】

音声データ送信中の処理は、システムの終了信号を受信したか否かを判定し、受信したときには必要な処理を実行してこのプログラムの実行を終了するステップ 2 1 2 と、終了信号を受信されていないときに、音声認識処理部 8 0 からローカル音声認識結果を受信したか否かを判定するステップ 2 1 4 と、ローカル音声認識結果を受信したときに、その中に終了キーワードの条件を充足する表現があるか否かを判定し、なければ制御をステップ 2 1 2 に戻すステップ 2 1 6 と、ローカル音声認識結果中に終了キーワードの条件を充足する表現があったときに、バッファ 5 4 に記憶されている音声データのうち、終了キーワードが検出された部分の末尾までを音声認識サーバ 3 6 に送信して送信を終了し、制御をステップ 2 0 2 に戻すステップ 2 1 8 とを含む。

【 0 0 5 2 】

このプログラムはまた、ステップ 2 1 4 でローカル音声認識結果を音声認識処理部 8 0 から受信していないと判定されたときに、発話なしで所定時間が経過したか否かを判定し、所定時間が経過していなければ制御をステップ 2 1 2 に戻すステップ 2 2 0 と、発話なしで所定時間が経過したときに、バッファ 5 4 に記憶されている音声データの音声認識サーバ 3 6 への送信を終了し、制御をステップ 2 0 2 に戻すステップ 2 2 2 とを含む。

【 0 0 5 3 】

図 6 を参照して、図 2 の実行制御部 9 0 を実現するプログラムは、携帯電話 3 4 の電源投入時に起動され、必要な初期化処理を実行するステップ 2 4 0 と、終了信号を受信したか否かを判定し受信したときにはこのプログラムの実行を終了するステップ 2 4 2 と、終了信号を受信していないときに、音声認識サーバ 3 6 から音声認識結果のデータを受信したか否かを判定し、受信していなければ制御をステップ 2 4 2 に戻すステップ 2 4 4 とを含む。

【 0 0 5 4 】

このプログラムはさらに、音声認識サーバ 3 6 から音声認識結果のデータを受信したときに、一時記憶部 8 8 に保存されていた開始キーワードを読み出すステップ 2 4 6 と、ステップ 2 4 6 で読み出された開始キーワードが音声認識サーバ 3 6 からの音声認識結果のデータの先頭部分と一致するか否かを判定するステップ 2 4 8 と、両者が一致したときに、音声認識サーバ 3 6 による音声認識結果のうち、開始キーワードの終端部の次の位置から終了までのデータを受信データバッファ 6 0 から読み出すようアプリケーション実行部 6 2 を制御するステップ 2 5 0 と、ステップ 2 4 8 で開始キーワードが一致しないと判定されたときに、受信データバッファ 6 0 に記憶された音声認識サーバ 3 6 による音声認識結果をクリアする (又は読捨てる) ステップ 2 5 4 と、ステップ 2 5 0 又はステップ 2 5 4 の後に、一時記憶部 8 8 をクリアして制御をステップ 2 4 2 に戻すステップ 2 5 2 とを含む。

【 0 0 5 5 】

図 5 に示すプログラムによれば、ローカルな音声認識結果が開始キーワードとマッチしているとステップ 2 0 6 で判定されると、ステップ 2 0 8 でその開始キーワードが一時記憶部 8 8 に保存され、ステップ 2 1 0 以後で、バッファ 5 4 に記憶された音声データのうち、開始キーワードと一致した先頭部分からの音声データが音声認識サーバ 3 6 に送信される。音声データの送信中にローカルな音声認識結果中に終了キーワードとしての条件を充足する表現が検出されると (図 5 のステップ 2 1 6 で Y E S)、バッファ 5 4 に記憶さ

10

20

30

40

50

れた音声データのうち、終了キーワードの部分の終端まで音声認識サーバ36に送信された後、送信が終了する。

【0056】

一方、音声認識サーバ36から音声認識結果を受信したときに、図6のステップ248の判定が肯定なら、音声認識結果のうち、開始キーワードと一致した部分の末尾以後が受信データバッファ60からアプリケーション実行部62に読出され、アプリケーション実行部62が音声認識結果の内容に応じた適切な処理を実行する。

【0057】

したがって、図5及び図6に制御構造を示すプログラムを携帯電話34で実行することにより、上記した実施の形態の機能を実現できる。

10

【0058】

<第2の実施の形態>

上記実施の形態では、ローカル音声認識で開始キーワードを検出すると、その開始キーワードを一時的に一時記憶部88に保存している。そして、音声認識サーバ36から音声認識結果が返ってきたときに、音声認識結果の先頭部分と一時的に保存された開始キーワードとが一致するか否かにより、音声認識サーバ36の音声認識結果を使用した処理を実行するか否かを判定している。しかし本発明はそのような実施の形態には限定されない。そのような判定を行わず、音声認識サーバ36の音声認識結果をそのまま利用する実施の形態も考えられる。これは、特にローカル音声認識でのキーワード検出の精度が十分に高いときに有効である。

20

【0059】

図7を参照して、この第2の実施の形態に係る携帯電話260は、第1の実施の形態の携帯電話34とほぼ同様な構成である。しかし、音声認識サーバ36による音声認識結果と開始キーワードとの照合に必要な機能ブロックを含まず、より簡略となっている点で携帯電話34と異なっている。

【0060】

具体的には、携帯電話260は、図1に示す制御部58を簡略化し、音声認識サーバ36からの音声認識結果と開始キーワードとの照合を行わないようにした制御部270を制御部58に代えて持つ点と、制御部58の制御によらず、音声認識サーバ36からの音声認識結果を一時的に保持し、全て出力する受信データバッファ272を図1の受信データバッファ60に代えて持つ点と、制御部270の制御を受けず、音声認識サーバ36からの音声認識結果を全て処理するアプリケーション実行部274を図1のアプリケーション実行部62に代えて持つ点で第1の実施の形態の携帯電話34と異なっている。

30

【0061】

制御部270は、図1に示す一時記憶部88及び実行制御部90を持たない点、及び、図1の通信制御部86に代えて、ローカルな音声認識結果内に開始キーワードが検出されたときに、バッファ54に記憶されている音声データの内で、開始キーワードに対応する位置の直後からのデータを音声認識サーバ36に送信する処理を開始するよう送受信部56を制御する機能を持つ通信制御部280を持つ点で図1の制御部58と異なっている。なお、通信制御部280もまた、制御部58と同様、ローカルな音声認識結果の中に終了キーワードが検出されたときには、音声認識サーバ36への音声データの送信を終了するよう送受信部56を制御する。

40

【0062】

図8を参照して、この実施の形態に係る携帯電話260の動作の概略について説明する。発話140の構成は図4に示すものと同様であるものとする。本実施の形態に係る制御部270は、発話140中の発話部分150に開始キーワードが検出されたときに、音声データのうち、開始キーワードが検出された部分の次から終了キーワードが検出された直後(図8に示す発話部分152に相当)までの音声データ290を音声認識サーバ36に送信する。すなわち、音声データ290には開始キーワード部分の音声データは含まれない。その結果、音声認識サーバ36から返信される音声認識結果292にも開始キーワー

50

ドは含まれない。したがって、発話部分 150 の部分のローカル音声認識の結果が正しければ、サーバからの音声にも開始キーワードは含まれず、音声認識結果 292 の全体をアプリケーション実行部 274 が処理しても特に不都合は生じない。

【0063】

図 9 に、この実施の形態に係る携帯電話 260 の判定部 82 及び通信制御部 280 の機能を実現するためのプログラムの制御構造をフローチャート形式で示す。この図は、第 1 の実施の形態の図 5 に示すものに相当する。なおこの実施の形態では、第 1 の実施の形態の図 6 に制御構造を示すようなプログラムは必要ない。

【0064】

図 9 を参照して、このプログラムは、図 5 に制御構造を示すものからステップ 208 を削除し、ステップ 210 に代えて、バッファ 54 に記憶された音声データのうち、開始キーワードの終端の次の位置から音声認識サーバ 36 に音声データを送信するように送受信部 56 を制御するステップ 300 を含む。その他の点では、このプログラムは図 5 に示すものと同じ制御構造を示す。このプログラムの実行時の制御部 270 の動作も、既に説明したのから十分に明らかである。

【0065】

この第 2 の実施の形態では、音声データの送信を開始するためにユーザが何らかの操作を特に行なう必要がないという点と、音声データを音声認識サーバ 36 に送信するにあたり、データ量を少なく抑えることができるという点で第 1 の実施の形態と同じ効果を得ることができる。またこの第 2 の実施の形態では、ローカル音声認識のキーワード検出の精度が高ければ、簡単な制御でサーバを用いた音声認識結果を利用した様々な処理を利用できるという効果を奏する。

【0066】

[携帯電話のハードウェアブロック図]

図 10 に、第 1 の実施の形態に係る携帯電話 34 及び第 2 の実施の形態に係る携帯電話 260 を実現する携帯電話のハードウェアブロック図を示す。以下の説明では、携帯電話 34 及び 260 を代表して携帯電話 34 について説明する。

【0067】

図 10 を参照して、携帯電話 34 は、マイクロフォン 50 及びスピーカ 66 と、マイクロフォン 50 及びスピーカ 66 が接続されたオーディオ回路 330 と、オーディオ回路 330 が接続されたデータ転送用及び制御信号転送用のバス 320 と、GPS 用、携帯電話回線用、及びその他規格にしたがった無線通信用のアンテナを備え、様々な通信を無線により実現する無線回路 332 と、無線回路 332 と携帯電話 34 の他のモジュールとの間を仲介する処理を行なう、バス 320 に接続された通信制御回路 336 と、通信制御回路 336 に接続され、携帯電話 34 に対する利用者の指示入力を受けて入力信号を通信制御回路 336 に与える操作ボタン 334 と、バス 320 に接続され、様々なアプリケーションを実行するための CPU (図示せず)、ROM (読出専用メモリ : 図示せず) 及び RAM (Random Access Memory : 図示せず) を備えたアプリケーション実行用 IC (集積回路) 322 と、アプリケーション実行用 IC 322 に接続されたカメラ 326、メモリカード入出力部 328、タッチパネル 64 及び DRAM (Dynamic RAM) 338 と、アプリケーション実行用 IC 322 に接続され、アプリケーション実行用 IC 322 により実行される様々なアプリケーションを記憶した不揮発性メモリ 324 とを含む。

【0068】

不揮発性メモリ 324 には、図 1 に示す音声認識処理部 80 を実現するローカル音声認識処理プログラム 350 と、判定部 82、通信制御部 86 及び実行制御部 90 を実現する発話送受信制御プログラム 352 と、キーワード辞書 84 と、キーワード辞書 84 に記憶されるキーワードを保守するための辞書保守プログラム 356 とが記憶されている。これらプログラムは、いずれもアプリケーション実行用 IC 322 による実行時にはアプリケーション実行用 IC 322 内の図示しないメモリにロードされ、アプリケーション実行用

10

20

30

40

50

IC322内のCPUが持つプログラムカウンタと呼ばれるレジスタにより指定されるアドレスから読出され、CPUにより実行される。実行結果は、DRAM338、メモリカード入出力部328に装着されたメモリカード、アプリケーション実行用IC322内のメモリ、通信制御回路336内のメモリ、オーディオ回路330内のメモリのうち、プログラムにより指定されるアドレスに格納される。

【0069】

図2及び図7に示すフレーム化処理部52はオーディオ回路330により実現される。バッファ54及び受信データバッファ272は、DRAM338若しくは通信制御回路336又はアプリケーション実行用IC322内のメモリにより実現される。送受信部56は無線回路332及び通信制御回路336により実現される。図1の制御部58及びアプリケーション実行部62に、並びに図7の制御部270及びアプリケーション実行部274は、本実施の形態ではいずれもアプリケーション実行用IC322により実現される。

10

【0070】

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の全ての変更を含む。

【符号の説明】

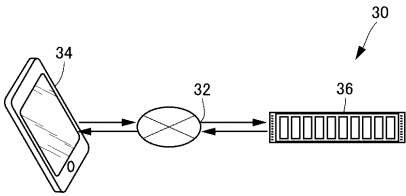
【0071】

- 30 音声認識システム
- 34 携帯電話
- 36 音声認識サーバ
- 50 マイクロフォン
- 54 バッファ
- 56 送受信部
- 58 制御部
- 60 受信データバッファ
- 62 アプリケーション実行部
- 80 音声認識処理部
- 82 判定部
- 84 キーワード辞書
- 86 通信制御部
- 88 一時記憶部
- 90 実行制御部

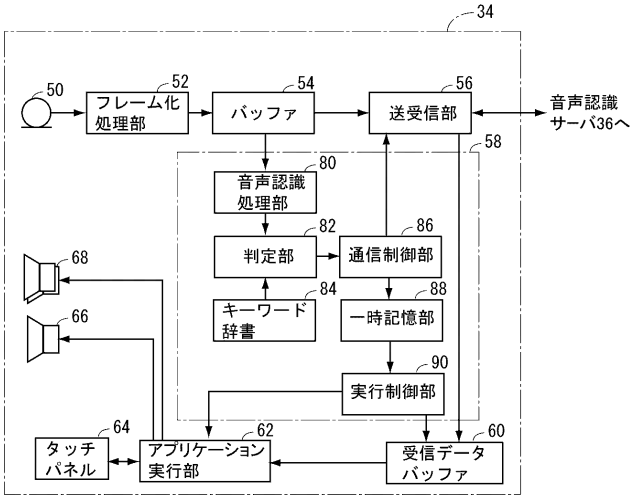
20

30

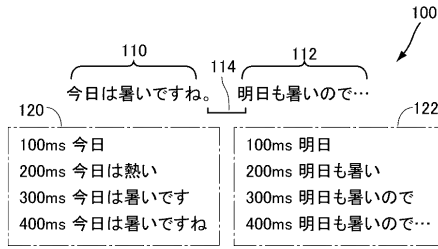
【図1】



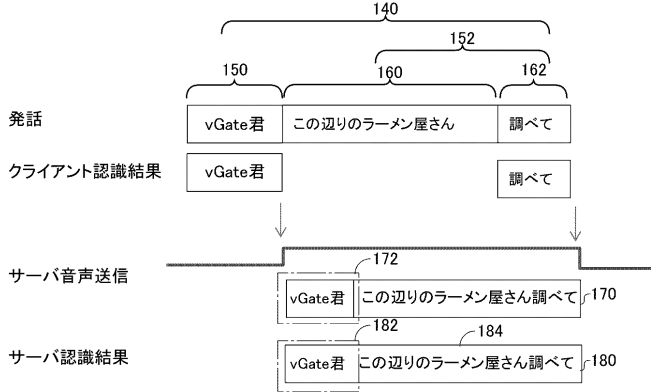
【図2】



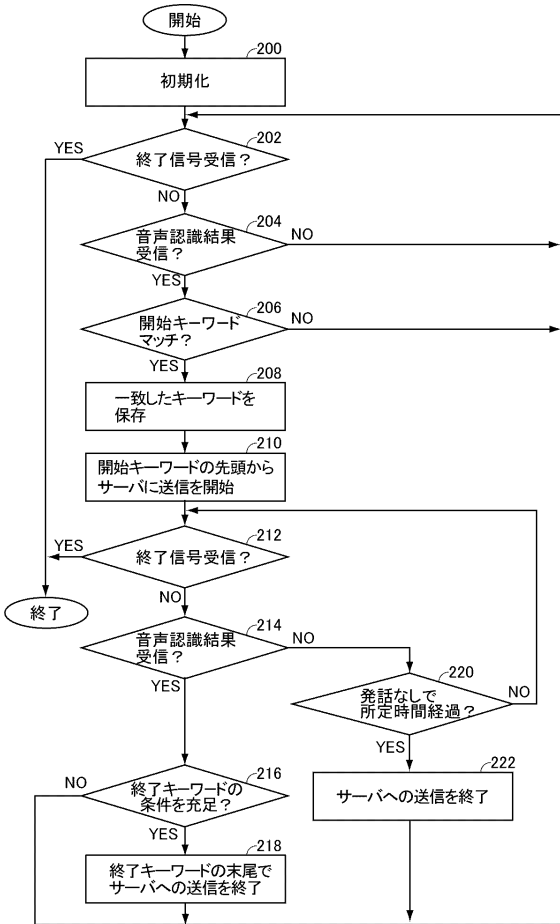
【図3】



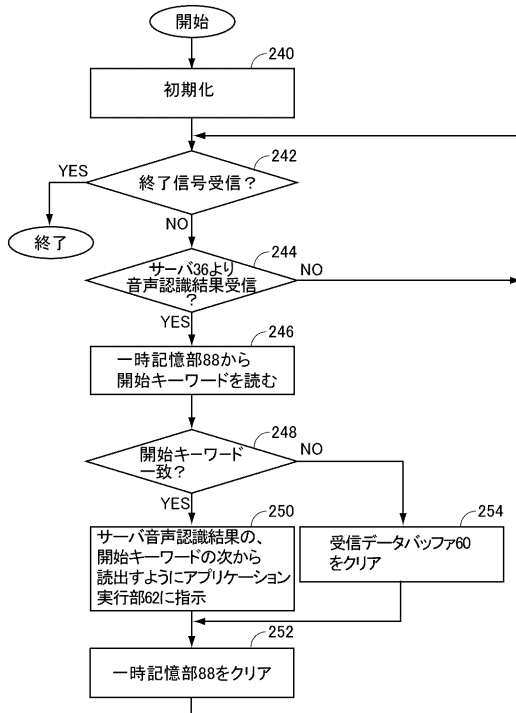
【図4】



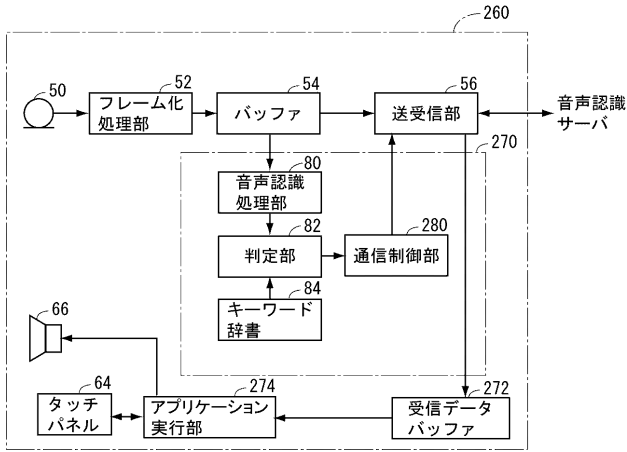
【図5】



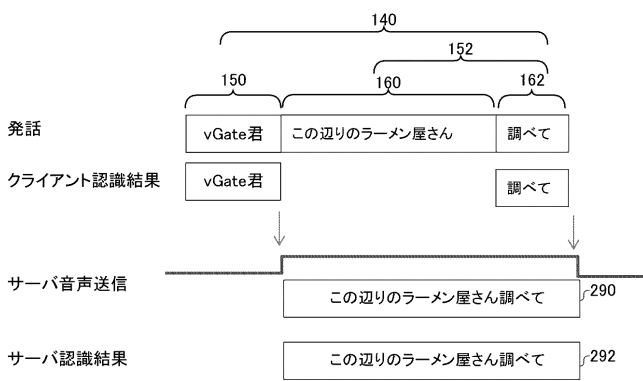
【図6】



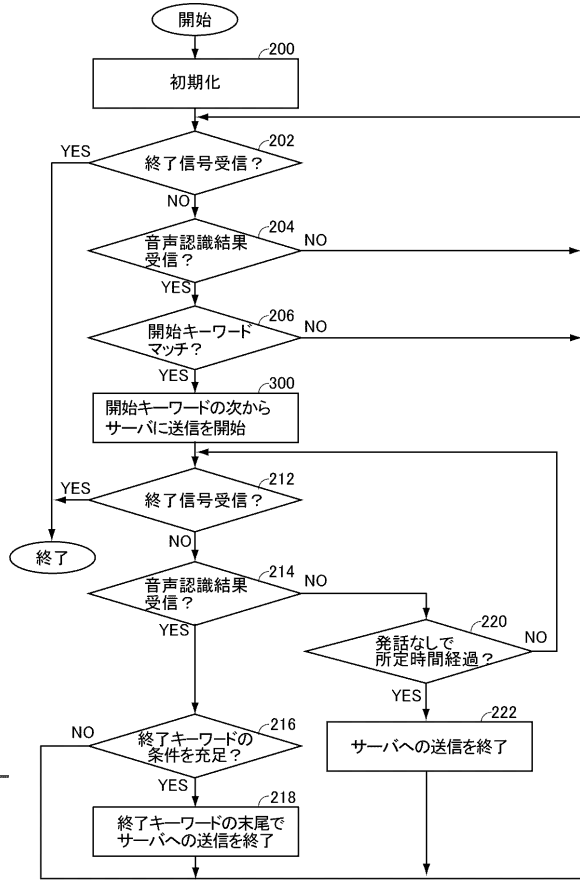
【図7】



【図8】



【図9】



【図10】

