

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4416643号  
(P4416643)

(45) 発行日 平成22年2月17日(2010.2.17)

(24) 登録日 平成21年12月4日(2009.12.4)

(51) Int.Cl.		F I		
<b>G 0 6 F</b>	<b>3/16</b>	<b>(2006.01)</b>	<b>G 0 6 F</b>	<b>3/16 3 2 0 A</b>
<b>G 0 6 F</b>	<b>3/01</b>	<b>(2006.01)</b>	<b>G 0 6 F</b>	<b>3/01</b>
<b>G 1 0 L</b>	<b>15/24</b>	<b>(2006.01)</b>	<b>G 1 0 L</b>	<b>15/24 R</b>

請求項の数 5 (全 14 頁)

(21) 出願番号	特願2004-379948 (P2004-379948)	(73) 特許権者	000001007 キヤノン株式会社 東京都大田区下丸子3丁目30番2号
(22) 出願日	平成16年12月28日(2004.12.28)	(74) 代理人	100090538 弁理士 西山 恵三
(65) 公開番号	特開2006-48628 (P2006-48628A)	(74) 代理人	100096965 弁理士 内尾 裕一
(43) 公開日	平成18年2月16日(2006.2.16)	(72) 発明者	池田 裕美 東京都大田区下丸子3丁目30番2号キヤ ノン株式会社内
審査請求日	平成19年12月10日(2007.12.10)		
(31) 優先権主張番号	特願2004-191632 (P2004-191632)	審査官	円子 英紀
(32) 優先日	平成16年6月29日(2004.6.29)		
(33) 優先権主張国	日本国(JP)		

最終頁に続く

(54) 【発明の名称】 マルチモーダル入力方法

(57) 【特許請求の範囲】

【請求項 1】

音声認識された情報と G U I 入力された情報とを統合した認識結果を出力する情報処理装置の情報処理方法であって、

音声情報を受信する音声情報受信工程と、

前記音声情報を受信している時間内にユーザに操作された 1 または複数のボタンに対応する G U I 入力情報を受信する G U I 情報受信工程と、

前記音声情報を音声認識し、尤度が高い順に複数の解釈結果の候補を取得する音声認識工程と、

前記尤度が最も高い解釈結果の候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致しているか否かを判断する判断工程と、

前記尤度が最も高い解釈結果の候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致していないと判断された場合、解釈結果の候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致している別の解釈結果の候補を検索する検索工程と、

検索された解釈結果の候補に含まれる未確定語を前記 G U I 入力情報で置換えた情報を、認識結果として出力する出力工程とを有する情報処理方法。

【請求項 2】

前記検索工程は、前記尤度が高い解釈結果の候補から順に、解釈結果の候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致している別の解釈結果の候補を検索することを特徴とする請求項 1 記載の情報処理方法。

10

20

**【請求項 3】**

前記 G U I 入力情報の数とは、前記 G U I 入力手段から、前記音声情報を受信している時間内にユーザに操作されたボタンの個数または回数であることを特徴とする請求項 1 記載の情報処理方法。

**【請求項 4】**

音声入力手段から音声情報を受信する音声情報受信手段と、

G U I 入力手段から、前記音声情報を受信している時間内にユーザに操作された 1 または複数のボタンに対応する G U I 入力情報を受信する G U I 情報受信手段と、

前記音声情報を音声認識し、尤度が高い順に複数の解釈結果の候補を取得する音声認識手段と、

前記尤度が最も高い解釈結果の候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致しているか否か判断する判断手段と、

前記尤度が最も高い解釈結果の候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致していないと判断された場合、解釈結果の候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致している別の解釈結果の候補を検索する検索手段と、

検索された解釈結果の候補に含まれる未確定語を前記 G U I 入力情報で置換えた情報を、認識結果として出力する出力手段とを有する情報処理装置。

**【請求項 5】**

請求項 1 乃至請求項 3 に記載の情報処理方法をコンピュータに実行させるためのプログラム。

**【発明の詳細な説明】****【技術分野】****【0001】**

本発明は、マルチモーダル・ユーザインタフェースに係る技術に関するものである。

**【背景技術】****【0002】**

G U I 入力や音声入力といった複数種類の入力手段から、ユーザの所望の入力手段をもって情報の入力を可能にするマルチモーダル・ユーザインタフェースは、ユーザにとって利便性が高いものである。特に、複数種類の入力手段を同時に用いて入力を行った場合の利便性は高く、例えば音声で「これをここに移動」等の指示を発声しながら、G U I で「これ」に対応する対象と、「ここ」に対応する対象をクリックする等の操作を行うことにより、コマンド等の専門的な言語に不慣れなユーザであっても自由に対象を操作することができる。このような操作を可能にするためには複数種類の入力手段による入力を統合するための処理が必要になる。

**【0003】**

複数種の入力手段による入力を統合する処理の例として、マウスイベントの種類や速度に関する情報を用いる方法（特許文献 1、特許文献 2）のほか、音声認識結果に対して言語解析を行う方法（特許文献 3）や文脈情報を用いる方法（特許文献 4）、入力時刻の近いものをまとめて意味解析単位として出力する方法（特許文献 5）、入力データの認識結果遅着を考慮した方法（特許文献 6）、利用者の意図を統計的な学習により検出する方法（特許文献 7、特許文献 8）、文法解析方法を用いた方法（特許文献 9）、言語解析を行って意味構造を用いる方法（特許文献 10）等や、マウスなどのポインティングデバイスによるポインティング入力をリストに登録し、音声入力データ中の指示表現の数とリスト中の数とを比較し、ポインティング入力数が少ない場合に、次のポインティング入力を得ることで数を合わせ、音声入力とポインティング入力を統合する方法（特許文献 11）が開示されている。

**【特許文献 1】**特開平 8 - 2 8 6 8 8 7 号公報

**【特許文献 2】**特開平 9 - 8 1 3 6 4 号公報

**【特許文献 3】**特許第 2 9 9 3 8 7 2 号公報

**【特許文献 4】**特許第 3 3 7 5 4 4 9 号公報

10

20

30

40

50

【特許文献 5】特許第 3 3 6 3 2 8 3 号公報  
【特許文献 6】特開平 1 0 - 1 9 8 5 4 4 号公報  
【特許文献 7】特開平 1 1 - 2 8 8 3 4 2 号公報  
【特許文献 8】特開 2 0 0 1 - 1 0 0 8 7 8 号公報  
【特許文献 9】特開平 6 - 2 8 2 5 6 9 号公報  
【特許文献 1 0】特開 2 0 0 0 - 2 3 1 4 2 7 号公報  
【特許文献 1 1】特開平 7 - 1 1 0 7 3 4 号公報  
【発明の開示】  
【発明が解決しようとする課題】  
【0 0 0 4】

10

上記従来例では、各入力の入力時刻や入力順序を考慮しているが、1つの入力結果に対する複数の候補を解析するには複雑な処理を行わなければならない。また、音声入力を正確に認識できることを前提としているが、現在の音声認識技術では 1 0 0 % 正しく認識するのは困難である。そのため誤認識への対応が重要となるが、上記従来例には誤認識を起こした場合の対応や誤認識率を下げることにについて言及していない。

【0 0 0 5】

特許文献 1 1 では、音声入力データ中の指示入力の数に対してポインティング入力の数  
が足りない場合に次のポインティング入力を待って統合する技術が記載されているが、上  
述したように基本的に音声入力データ中の指示入力の数  
が正確に認識できることを前提としており、誤認識に関する記載はなく、また誤認識率を下げる  
ことについては記載されていない。特許文献 1 1 は、ポインティング入力  
の数  
が音声入力データ中の指示入力の数よりも多い場合には、エラー処理を行い入力をやり直す構成とな  
っているが、入力をやり直すことはユーザにとって負担となるため、このような事態を減らす技術が重要となる。

20

【0 0 0 6】

本発明は、このような事情を鑑みてなされたものであり、少なくとも 2 種類の入力手段  
からの入力が意図する指示内容の認識精度を向上することを目的とする。

【課題を解決するための手段】

【0 0 0 7】

上記課題を解決するために、本発明の情報処理方法は、音声認識された情報と G U I 入  
力された情報とを統合した認識結果を出力する情報処理装置の情報処理方法であって、音  
声情報を受信する音声情報受信工程と、前記音声情報を受信している時間内にユーザに操  
作された 1 または複数のボタンに対応する G U I 入力情報を受信する G U I 情報受信工程  
と、前記音声情報を音声認識し、尤度が高い順に複数の解釈結果の候補を取得する音声認  
識工程と、前記尤度が最も高い解釈結果の候補に含まれる未確定語の数と、前記 G U I 入  
力情報の数とが一致しているか否か判断する判断工程と、前記尤度が最も高い解釈結果の  
候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致していないと判断され  
た場合、解釈結果の候補に含まれる未確定語の数と、前記 G U I 入力情報の数とが一致し  
ている別の解釈結果の候補を検索する検索工程と、検索された解釈結果の候補に含まれる  
未確定語を前記 G U I 入力情報で置換えた情報を、認識結果として出力する出力工程とを  
有することを特徴とする。

30

40

【発明の効果】

【0 0 0 9】

以上説明したように、本発明によれば、少なくとも 2 種類の入力手段からの入力が意図  
する指示内容の認識精度を向上することができる。

【発明を実施するための最良の形態】

【0 0 1 0】

以下、図面を参照して、本発明に係るマルチモーダル入力方法の好適な実施例について説  
明する。

【実施例 1】

【0 0 1 1】

50

図1は、本発明の実施例1におけるシステムの基本構成を示す図である。ここでは、音声入力とGUI入力を受け付けるシステムを例にあげて説明する。GUI入力部101、GUI入力解釈部102、音声入力部103、音声認識・解釈部104、マルチモーダル入力統合部105、記憶部106、マークアップ解釈部107、制御部108、音声合成部109、表示部110、通信部111から構成される。

【0012】

GUI入力部101はGUI上で指示を与えるボタン群やキーボード、マウス、タッチパネル、ペン、タブレット等から構成され、各種の指示を本装置に入力するための入力インタフェースとして機能する。本システムはこれら入力部から入力情報を受信する。GUI入力解釈部102は、GUI入力部101から入力された情報に対して解釈を行う。上記解釈については、例えば手書き認識技術等、公知の技術を利用する。

10

【0013】

音声入力部103はマイクロフォンやA/D変換器等により構成されており、ユーザの音声を入力する。音声認識・解釈部104は音声入力部103より入力された音声に対して音声認識を行う。上記音声認識技術については公知の技術を利用する。マルチモーダル入力統合部105は、GUI入力解釈部102、音声認識・解釈部104で解釈された情報を統合する。

【0014】

記憶部106は、各種の情報を保存するためのハードディスクドライブ装置や、システムに各種の情報を提供するためのCD-ROMやDVD-ROM等の記憶媒体等により構成されている。またこのハードディスクドライブ装置や記憶媒体には、各種のアプリケーションプログラム、ユーザインタフェース制御プログラム、そして各プログラムを実行する際に必要な各種のデータ等が記憶されており、これらは後段の制御部108の制御により、本システムに読み込まれる。マークアップ解釈部107はマークアップで記述された文書を解釈する。制御部108はワークメモリやCPU、MPU等により構成されており、記憶部105に記憶されたプログラムやデータを読み出して各種の処理を実行する。またGUI入力解釈部、音声認識・解釈部、マルチモーダル入力統合部などの制御も実行する。

20

【0015】

音声合成部109はスピーカやヘッドフォン、D/A変換器等により構成されており、制御部108の制御により読み上げテキストから音声データを作成してD/A変換し、音として外部に出力する処理を行う。上記音声合成技術については公知の技術を利用する。表示部110は液晶ディスプレイ等の表示装置から構成され、画像や文字等により構成される各種の情報を表示する。なお、表示部110としてタッチパネル式の表示装置を用いてもよく、その場合、表示部110はGUI入力部101としての機能（各種の指示を本システムに入力する機能）をも有することになる。通信部111は、インターネットやLAN等のネットワークを介して他の装置とのデータ通信を行うためのネットワークインタフェースである。

30

【0016】

以下では、上記マルチモーダル入力統合部105における統合処理方法について図2のフローチャートを用いて説明する。GUI入力解釈部102、音声認識・解釈部104で解釈された情報、つまり認識候補がマルチモーダル入力統合部105に渡されると、まず、GUI入力のintegration=0の解釈結果を出力する（ステップS201）。上記integrationは他の入力手段の入力結果と統合する必要があるか否かを示す情報であり、必要がある場合は“1”、必要がない場合は“0”が、GUI入力解釈部102、音声認識・解釈部104において入力される。他の入力手段の入力結果と統合する必要があるか否かを判別する方法については、例えば、値を格納する場所が決まっているか否かを判別する等、公知の技術を利用する。

40

【0017】

続いて、すべての音声認識・解釈結果においてintegration=0の場合（ス

50

ステップS202にてYES)、第1位の音声認識・解釈結果を出力して終了する(ステップS203)。音声認識・解釈結果にintegration=1の結果がある場合(ステップS202にてNO)、音声入力時間T内のGUI入力の中でintegration=1の個数NUMをカウントする(ステップS204)。ここで、音声入力時間Tは、図3(縦軸:音声入力のパワー、横軸:時間)の301に示すように閾値th以上のパワーが検出されている時間とする。あるいは、図3の302のように、閾値th以上のパワーが検出されている時間の前後に任意の時間(例えば数秒)を足す等、音声入力部にて設定した値でもよい。図3の301の例では、GUI入力がすべてintegration=1とするとNUM=2となる。

#### 【0018】

NUM=0であれば(ステップS205にてYES)、integration=0の音声認識・解釈結果が存在するかどうかをチェックする(ステップS206)。存在する場合は(ステップS206にてYES)integration=0の音声認識・解釈結果の中で最も確信度の高い結果を出力して終了する(ステップS207)。存在しない場合は(ステップS206にてNO)結果を統合できなかった旨のエラーを出力して終了する(ステップS208)。

#### 【0019】

NUM=0でなければ(ステップS205にてNO)、Nに1を代入して(ステップS209)ステップS210へと進む。N位(最初は1位)の音声認識・解釈結果が存在しない場合は(ステップS210にてNO)、結果を統合できなかった旨のエラーを出力して終了する(ステップS208)。存在する場合は(ステップS210にてYES)、ステップS211へと進む。ステップS211ではN位の音声認識・解釈結果のintegrationが1か(統合する必要があるか否か)を判別する(ステップS211)。統合する必要がある場合は(ステップS211にてYES)、Nに1を足し(ステップS212)、ステップS210へと進む。統合する必要がある場合は(ステップS211にてYES)、“?”の数(図4の例では401~403のテーブルにおけるunknownの値)が前述のNUMと同じかどうかを判別する(ステップS213)。同じ値でなければ(ステップS213にてNO)、Nに1を足し(ステップS212)、ステップS210へと進む。同じ値の場合は(ステップS213にてYES)、N位の音声認識・解釈結果とGUIの解釈結果を統合して出力する(ステップS214)。統合の具体例としては、音声入力「ここから」とGUI入力「恵比寿」を統合し、「恵比寿から」という結果になる。

#### 【0020】

以下では、上記マルチモーダル入力統合の例を図を用いて示す。図4~図7では、音声入力とボタン入力がなされた場合の、解釈処理結果の流れを示している。

#### 【0021】

1つ目の例を、図4を用いて説明する。図4の401、402は音声入力に対する音声認識・解釈結果を表すテーブルであり、確信度第1位の結果を401、確信度第2位の結果を402に示している。また、403はボタン入力に対する解釈結果を表すテーブルであり、この例では確信度が第1位の結果のみの場合を示している。

#### 【0022】

以下、401~403に示す各テーブルの項目について説明する。“rank”は確信度の順位(高いものから1位,2位,...)、“unknown”は確定していない値(後述の“value=?”)の数、“value”は解釈値、“time(start)”は入力開始時刻、“time(end)”は入力終了時刻、“score”は確信度、“integration”は統合が必要か否か(必要:1,不要:0)を表す。GUI入力解釈部102、音声認識・解釈部104にて解釈された結果が上記テーブルに入力され、マルチモーダル入力統合部105へと渡される。上記テーブルがXMLで記述されている場合はマークアップ解釈部107にて解釈される。

#### 【0023】

マルチモーダル入力統合部105では、前述のフローチャート図2に示す流れで処理を

10

20

30

40

50

行う。GUI入力解釈処理結果である402は音声入力時間T(02:10:00~02:12:00)内に入力されたものであり、integration=1であるのでステップ204にてNUM=1とし、続いてN=1とする(ステップS209)。第1位の音声認識・解釈結果(401)が存在するので(ステップS210にてYES)、ステップS211へと進む。続いてintegration=1(ステップS211にてYES)、unknown=1=NUMであるので(ステップS213にてYES)、ステップS214へと進む。ステップS214では、音声入力の解釈結果「東京からここまで」とボタン入力の解釈結果「恵比寿」を統合し、「東京から恵比寿まで」を出力する。

【0024】

同様に、図5の例では、503と504よりNUM=2(ステップS204)である。第1位の音声認識・解釈結果(501)はunknown=1であり、NUM=2とは異なるので(ステップS213にてNO)、続いて第2位の音声認識・解釈結果(502)を調べる。502の結果はunknown=2=NUMであるので(ステップS213にてYES)、ステップS214にて、音声入力の解釈結果「ここからここまで」とボタン入力の解釈結果「恵比寿」「横浜」を統合し、「恵比寿から横浜まで」を出力する。

【0025】

図6の例では、604はintegration=0であるのでボタン入力の解釈結果「1」を出力する。また、603より、NUM=1(ステップS204)である。第1位の音声認識・解釈結果(601)はunknown=1=NUMであるので(ステップS213にてYES)、ステップS214にて、音声入力の解釈結果「東京からここまで」とボタン入力の解釈結果「恵比寿」を統合して「東京から恵比寿まで」を出力する。

【0026】

図7の例では、703はintegration=0であるのでボタン入力の解釈結果「1」を出力し、NUM=0とする(ステップS204)。音声認識・解釈結果701、702にintegration=0の結果が存在しないので(ステップS206にてNO)、結果を統合できなかった旨のエラーを出力して終了する(ステップS208)。

【0027】

以上のように、実施例1によれば、音声入力時間中のボタン入力の個数情報を利用して音声認識・解釈結果を選択することで、音声認識処理により生ずる候補に対して優先度をつけることができ、認識の精度を向上することができる。その結果、複数の候補から正しい認識結果が出力される可能性が高くなり、ユーザの再入力の手間を省くなどの効果が生まれる。

【実施例2】

【0028】

続いて、本発明に係る情報処理システムの実施例2について説明する。前述した実施例1では、GUI入力がボタン入力で認識率100%である場合の例を示した。しかしながら、実際のマルチモーダル・ユーザインタフェースでは、統合する入力の解釈の確信度がいずれも100%でない場合がある。このような場合は、第1位の解釈結果から順番に第1の実施形態と同様の処理を行えばよい。

【0029】

本実施例では、音声入力とペン入力となされた場合の例について図8~図10を用いて説明する。図8に示すテーブル801~804の各項目は、前述の図4~図7に示した各テーブルの項目と同様であり、801は第1位の音声認識・解釈結果、802は第2位の音声認識・解釈結果、803は第1位のGUI入力解釈結果、804は第2位のGUI入力解釈結果である。

【0030】

GUI入力解釈結果の第1位の解釈結果から順番に第1の実施形態と同様の処理を行う。まず、第1位のGUI入力解釈結果803は音声入力時間T(02:10:00~02:12:00)内に入力されたものであり、integration=1である。また803よりvalueの数は1つであるので、ステップ204にてNUM=1とし、続いて

10

20

30

40

50

N = 1とする(ステップS209)。第1位の音声認識・解釈結果(801)が存在するので(ステップS210にてYES)、ステップS211へと進む。続いてintegration = 1(ステップS211にてYES)、unknown = 1 = NUMであるので(ステップS213にてYES)、ステップS214へと進む。ステップS214では、音声入力の解釈結果「ここ」とペン入力の解釈結果「恵比寿」を統合し、「恵比寿」を出力する。

#### 【0031】

図9の例では、まず第1位のGUI入力解釈結果903より、NUM = 1(ステップS204)である。第1位の音声認識・解釈結果(901)はunknown = 1以上であり、NUMと一致するので(ステップS213にてYES)、ステップS214にて、音声入力の解釈結果「このへん」とペン入力の解釈結果「恵比寿」を統合して「恵比寿」を出力する。

#### 【0032】

また、図10では、GUI入力としてペン入力とボタン入力の両方がなされた場合の例を示している。1005はintegration = 0であるのでボタン入力の解釈結果「1」を出力する。また、第1位のGUI入力解釈結果1003より、NUM = 1(ステップS204)である。第1位の音声認識・解釈結果(1001)はunknown = 2以上であり、NUMとは異なるので(ステップS213にてNO)、続いて第2位の音声認識・解釈結果(1002)を調べる。第2位の音声認識・解釈結果(1002)は、unknown = 3であり、NUMとは異なるので(ステップS213にてNO)統合できない。次に第2位のGUI入力解釈結果1004より、NUM = 2(ステップS204)とする。第1位の音声認識・解釈結果(1001)はunknown = 2以上であり、NUMと一致するので(ステップS213にてYES)、ステップS214にて、音声入力の解釈結果「これらを」とペン入力の解釈結果「A, B」を統合して「A, Bを」を出力する。

#### 【0033】

以上のように、実施例2によれば、統合する入力の解釈の確信度がいずれも100%でない場合においても、音声入力時間中のGUI入力個数の情報を利用して音声認識・解釈結果を選択することで、音声認識結果の解釈の精度を向上することができる。

#### 【実施例3】

#### 【0034】

上記実施例では、GUI入力を受け付ける例をあげて説明したが、本発明はこれに限定されるものではなく、キーボードやテンキーなどの物理的なキー入力を受け付ける構成としてもかまわない。ここでは、テンキーと音声入力によって操作可能な複写機を例にあげて説明する。複写機における各指示コマンドが以下に示すようにテンキーの各キーに割り当てられていることを前提とする。キー1：用紙選択、キー2：枚数(部数)、キー3：倍率、キー4：濃さ、キー5：両面、キー6：ソータ、キー7：ステイブルソート、キー8：応用モード。

#### 【0035】

ユーザは、10ページからなるA5サイズの資料を左上にステイブルしたものを5部コピーしたい場合、キー1を押して「A5」、キー2を押して「5部」、キー7を押して「左上」と発声することで設定することができるが、この操作に慣れてきた場合は、1つ1つを入力するよりも「A5、5部、左上」のように連続発声できたほうが効率的に作業を進めることができる。しかしながら現在の音声認識の精度は100%ではないため、『5枚、左上』や『A5、5部、左上、濃く』等の認識誤りが発生し、誤った認識候補が発生する可能性がある。

#### 【0036】

本実施例ではこのような場面において、音声入力とキー入力をキー入力の個数を用いて統合する。ユーザは、キー1、2、7を押しながら「A5、5部、左上」と発声する。キーの押し方は3つ同時に押していても良いし、連続的に押しても構わない。ここでキー入

力の入力数は3である。音声入力のリコグニション候補が1位「5部、左上」、2位「A5、5部、左上」、3位「A5、5部、左上、濃く」、4位「A2、50部、左上」であった場合、これとキー入力数3を統合することで、数が一致しない「A5、5部」、「A5、5部、左上、濃く」が除去されるかもしくは、数が一致する「A5、5部、左上」、「A2、50部、左上」が選択されることにより、1位「A5、5部、左上」、2位「A2、50部、左上」となってリコグニション候補が絞られ、ここからリコグニション尤度の最も高いリコグニション候補がリコグニション結果として選ばれることによって、ユーザが発声した「A5、5部、左上」が正しくリコグニションされることとなる。

#### 【0037】

この他、携帯電話に表示された番号付きのメニューに対して携帯電話のボタンで各メニューを指定できるような場面を想定し、例えばボタンで、3番、5番を押しながら、「これとこれのヘルプがみたい」と発声する場合などにも本発明が適用できることは言うまでもない。

#### 【実施例4】

#### 【0038】

上記実施例では、GUIや物理的なキー入力数の情報を用いて、音声認識・解釈結果の第1位～第N位の候補の中から適切なものを選択する例を示したが、どの入力手段の入力個数情報をどの入力手段の入力情報に適用するかは上記例に限られない。例えば、音声認識・解釈結果より入力個数（前述の例でいえばunknownの値）を判別して手書き文字入力のリコグニション結果の第1位～第N位の候補の中から適切なものを選択してもよい。一般にボタン入力は音声入力（認識）に比べて確実性が高いことから、ボタン入力の個数情報を他方のモダリティに適用する等、任意に決めてもよいが、いずれの入力手段も曖昧性をもつような場合、どの入力手段の入力個数情報をどの入力手段の入力情報に適用すればいいのか、一意に決めることは適切でない。このような場合は、第1位の確信度と第2位以下の確信度の差が大きい方のモダリティを選択する等、確信度の結果から毎回決めてもよい。

#### 【0039】

図11を用いてそのような例について説明する。図11は音声入力で「ここ」と発声しながら、ペン入力で「恵比寿」に丸をつけた例を示している。ユーザは「恵比寿」にのみ丸をつけるつもりが「渋谷」にも少しかかってしまった状態である。音声入力の解釈処理により第1位が「ここ」、第2位が「ここここ」であり、それぞれSCOREが90、55である。ペン入力の解釈処理により、第1位が「渋谷、恵比寿」、第2位が「恵比寿」であり、SCOREが95、90である。第1位と第2位の確信度の差が音声入力の方が大きく、第1位が正解である確率が高いことから、入力数が正しい可能性も高いといえる。つまり入力数の確からしさの値が高いため、ここでは入力数は音声入力の方を信頼する。入力数の確からしさの値は、上述したように第1位と第2位の差から求めても良いし、例えば、確信度が上位のリコグニション候補が共通して含む入力数がより多いものが確からしさの値が高くなるよう求めてもよい。例えば、1位と2位の差が大きくても1位と2位で入力数が異なるものは確からしさの値を低くし、1位から4位までの確信度の差が少なくても全て入力数が等しい場合は入力数の確からしさの値が高くなるように求めても良い。また、上述したように、「一般にボタン入力は音声入力（認識）に比べて確実性が高い」等の情報を加味して求めても良い。音声入力の第1位の入力数は1であり、GUI入力のうち入力数が1である第2位が選ばれ、「ここ」と「恵比寿」が統合されて統合結果が「恵比寿」となる。

#### 【実施例5】

#### 【0040】

上記実施例では、入力数を取得する入力手段が1種類の場合について説明してきたが、本発明はこれに限られる物ではない。例えばGUI入力とキー入力を併せ持つ複写機においては、これら両方の入力からの入力数を考慮することも可能である。音声入力で「B5、片面から両面をこの枚数で」と入力しながら、GUI入力で用紙サイズと両面の指示を



選択し、キー入力で10と入力した場合は、音声入力に含まれる入力数は3であり、GUI入力とキー入力をあわせた入力数が3となり、これらの数が一致しないものを音声入力の認識結果から除外することで認識精度を向上することが可能となる。

【実施例6】

【0041】

上記実施例では、音声入力と他の入力手段を統合する例を挙げて説明してきたが、本発明はこれに限定されるものではない。例えば、ジェスチャ入力と視線入力でのマルチモーダル入力に適用した場合にも適用可能である。ここでは、視線入力で操作対象を指定し、ジェスチャ入力でその操作対象に指示を与えるタスクを考えてみる。視線入力でオブジェクトを指定する場合は、そのオブジェクトを長く見ていた場合は選択されたとみなすなどの処理によって指定するが、それが正しく認識されず、操作対象の認識候補が複数得られる場合がある。それに対してジェスチャ入力により2つの指示が入力された場合は、操作対象が2つである可能性が高いため、視線入力の認識候補のうち、操作対象が2つの候補以外を除外することで認識精度を向上することができる。

10

【0042】

なお、本発明の目的は、前述した実施例の機能を実現するソフトウェアのプログラムコードを記録した記憶媒体を、システムあるいは装置に供給し、そのシステムあるいは装置のコンピュータ（またはCPUやMPU）が記憶媒体に格納されたプログラムコードを読み出し実行することによっても、達成されることは言うまでもない。

20

【0043】

この場合、記憶媒体から読み出されたプログラムコード自体が前述した実施形態の機能を実現することになり、そのプログラムコードを記憶した記憶媒体は本発明を構成することになる。

【0044】

プログラムコードを供給するための記憶媒体としては、例えば、フレキシブルディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、磁気テープ、不揮発性のメモリカード、ROMなどを用いることができる。

【0045】

また、コンピュータが読み出したプログラムコードを実行することにより、前述した実施形態の機能が実現されるだけでなく、そのプログラムコードの指示に基づき、コンピュータ上で稼働しているOS（オペレーティングシステム）などが実際の処理の一部または全部を行い、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

30

【0046】

さらに、記憶媒体から読出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPUなどが実際の処理の一部または全部を行い、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

40

【図面の簡単な説明】

【0047】

【図1】本発明の実施例1における情報処理システムの基本構成を示す図である。

【図2】本発明の実施例1に係る情報処理システムにおけるマルチモーダル入力統合部の処理の流れを説明するためのフローチャートである。

【図3】本発明の実施例1に係る入力の例を示す図である。

【図4】本発明の実施例1に係る情報処理システムにおけるマルチモーダル統合入力統合の例を示す図である。

【図5】本発明の実施例1に係る情報処理システムにおけるマルチモーダル統合入力統合の例を示す図である。

【図6】本発明の実施例1に係る情報処理システムにおけるマルチモーダル統合入力統合

50

の例を示す図である。

【図 7】本発明の実施例 1 に係る情報処理システムにおけるマルチモーダル統合入力統合の例を示す図である。

【図 8】本発明の実施例 2 に係る情報処理システムにおけるマルチモーダル統合入力統合の例を示す図である。

【図 9】本発明の実施例 2 に係る情報処理システムにおけるマルチモーダル統合入力統合の例を示す図である。

【図 10】本発明の実施例 2 に係る情報処理システムにおけるマルチモーダル統合入力統合の例を示す図である。

【図 11】本発明の実施例 4 に係る情報処理システムにおけるマルチモーダル統合入力統合の例を示す図である。

10

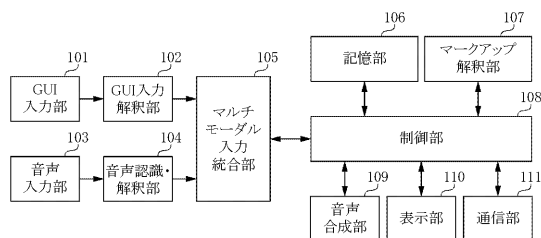
【符号の説明】

【 0 0 4 8 】

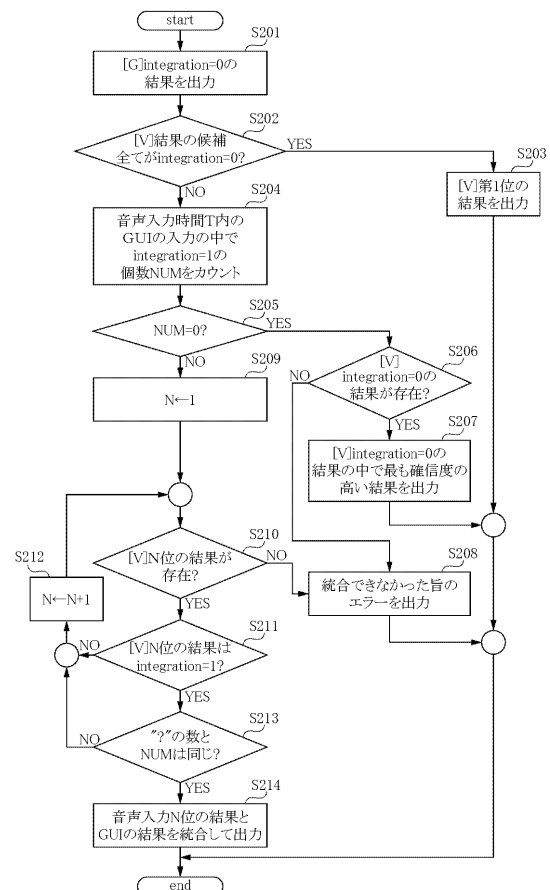
- 1 0 1    G U I 入力部
- 1 0 2    G U I 入力解釈部
- 1 0 3    音声入力部
- 1 0 4    音声認識・解釈部
- 1 0 5    マルチモーダル入力統合部
- 1 0 6    記憶部
- 1 0 7    マークアップ解釈部
- 1 0 8    制御部
- 1 0 9    音声合成部
- 1 1 0    表示部
- 1 1 1    通信部

20

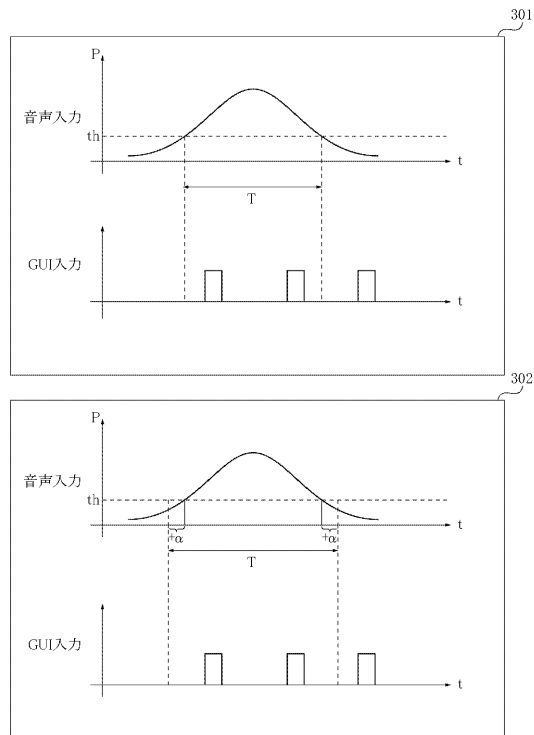
【図 1】



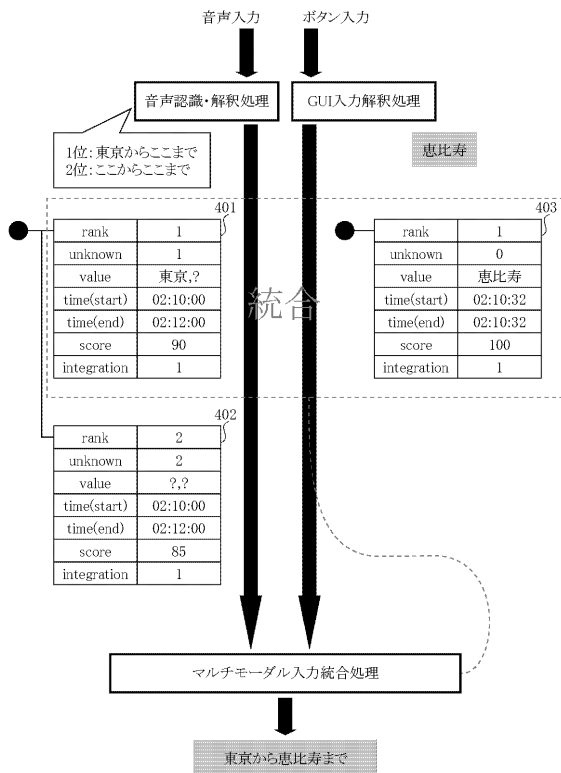
【図 2】



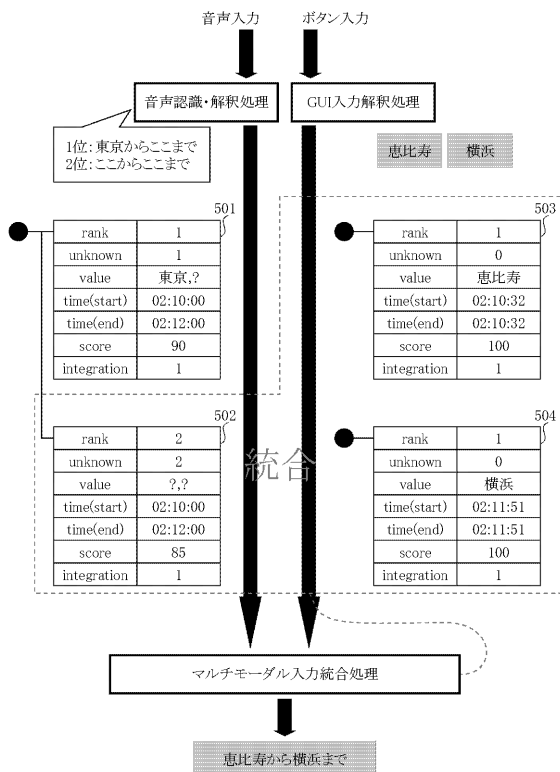
【図 3】



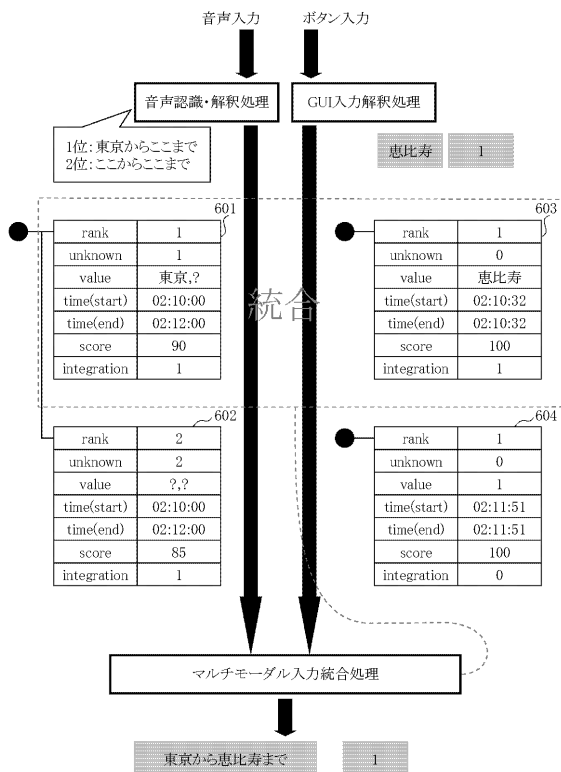
【図 4】



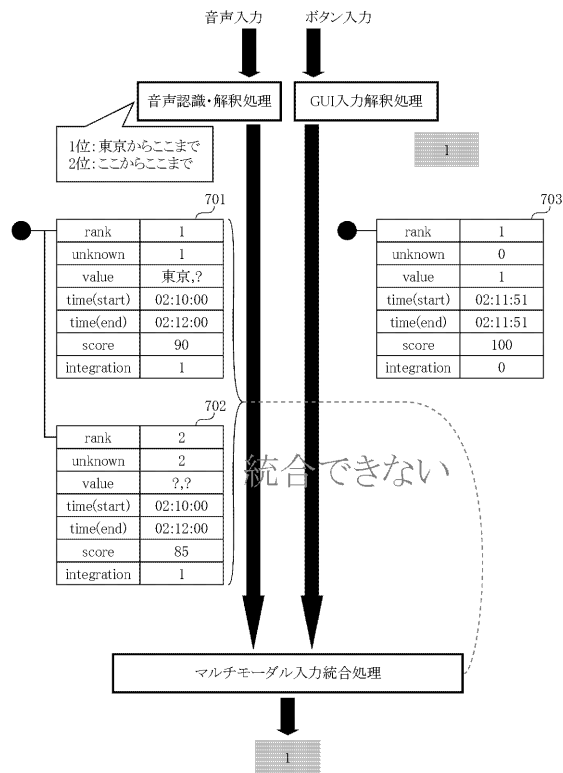
【図 5】



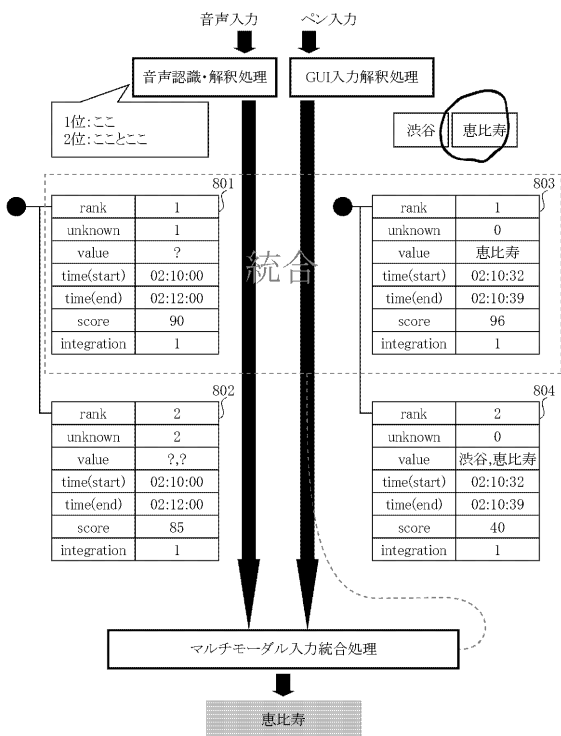
【図 6】



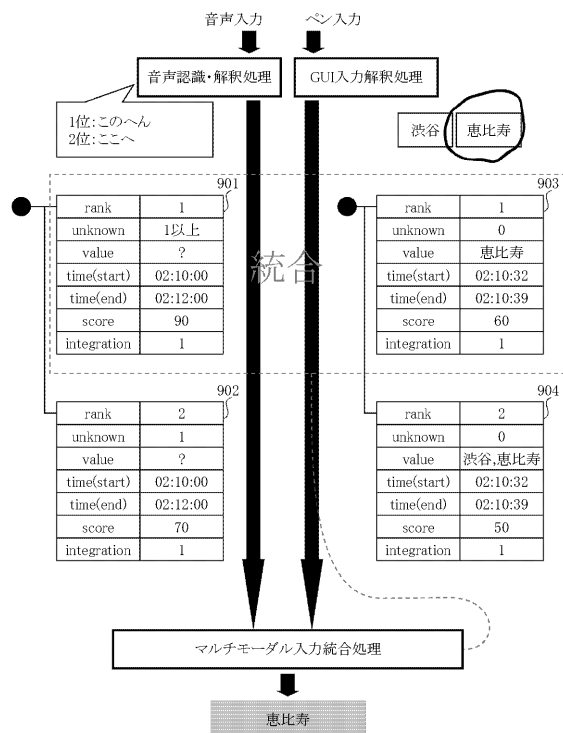
【図 7】



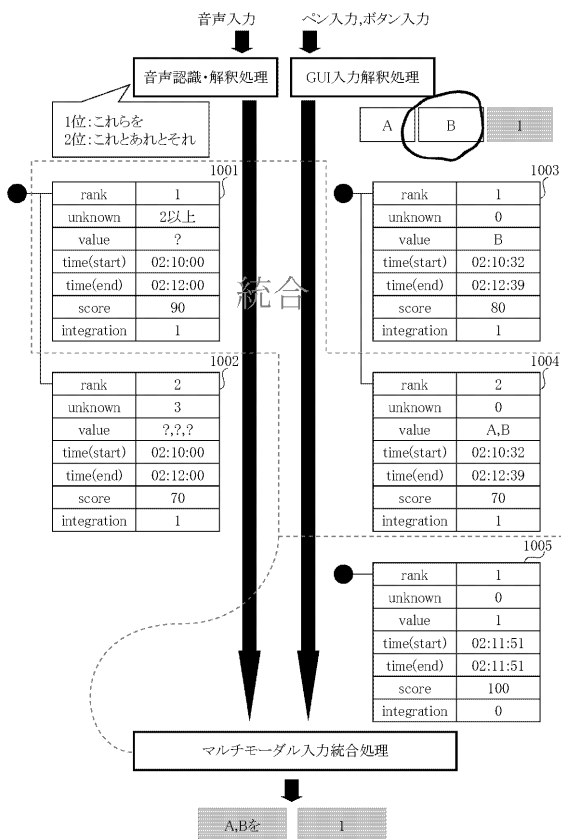
【図 8】



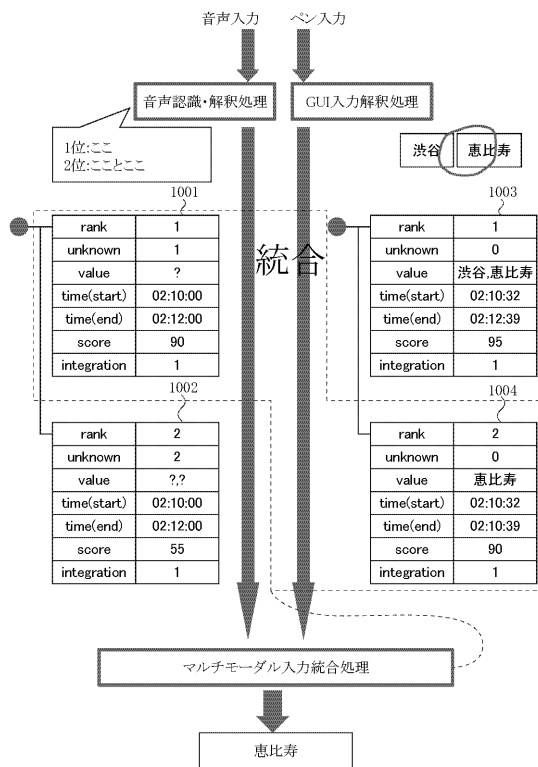
【図 9】



【図 10】



【図 11】



---

フロントページの続き

(56)参考文献 特開平08-063319(JP,A)  
特開2003-084900(JP,A)  
特開2000-209378(JP,A)  
特開2001-282413(JP,A)  
特開2003-140687(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/01  
G06F 3/16  
G10L 11/00 - 21/06