

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5998861号
(P5998861)

(45) 発行日 平成28年9月28日(2016.9.28)

(24) 登録日 平成28年9月9日(2016.9.9)

(51) Int.Cl.

F I

G 1 O L 15/22 (2006.01)

G 1 O L 15/22 2 O O H

G 1 O L 15/28 (2013.01)

G 1 O L 15/28 2 3 O K

請求項の数 18 (全 30 頁)

(21) 出願番号 特願2012-246118 (P2012-246118)
 (22) 出願日 平成24年11月8日(2012.11.8)
 (65) 公開番号 特開2014-95766 (P2014-95766A)
 (43) 公開日 平成26年5月22日(2014.5.22)
 審査請求日 平成27年1月23日(2015.1.23)

(73) 特許権者 000002185
 ソニー株式会社
 東京都港区港南1丁目7番1号
 (74) 代理人 100095957
 弁理士 亀谷 美明
 (74) 代理人 100096389
 弁理士 金本 哲男
 (74) 代理人 100101557
 弁理士 萩原 康司
 (74) 代理人 100128587
 弁理士 松本 一騎
 (72) 発明者 大村 淳己
 東京都港区港南1丁目7番1号 ソニー株
 式会社内

最終頁に続く

(54) 【発明の名称】 情報処理装置、情報処理方法及びプログラム

(57) 【特許請求の範囲】

【請求項 1】

入力画像を取得する画像取得部と、

発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させる制御部と、

前記入力画像に映るユーザの身体を認識する画像認識部と、

を備え、

前記制御部は、前記ユーザの音声について実行される音声認識を、前記画像認識部により認識される前記ユーザの身体の所定の部分と前記オブジェクトとの間の前記画面内の位置関係に基づいて、制御する、

情報処理装置。

10

【請求項 2】

前記所定の部分は、前記ユーザの口を含み、

前記制御部は、前記ユーザの口と前記オブジェクトとの間の距離に基づいて、前記音声認識のための音声入力をアクティブ化する、

請求項 1 に記載の情報処理装置。

【請求項 3】

前記所定の部分は、前記ユーザの手を含み、

前記制御部は、前記ユーザの手の動きに従って前記オブジェクトを前記画面内で移動させる、

請求項 2 に記載の情報処理装置。

20

【請求項 4】

前記制御部は、前記入力画像に映る前記ユーザのジェスチャに応じて、前記音声認識のための音声入力を非アクティブ化する、請求項 2 又は 3 に記載の情報処理装置。

【請求項 5】

前記制御部は、前記音声認識のための音声入力がアクティブ化されているか否かを、前記オブジェクトの表示属性を変化させることにより前記ユーザに通知する、請求項 1 ~ 4 のいずれか 1 項に記載の情報処理装置。

【請求項 6】

前記制御部は、前記音声認識において音声が発出されているか否かを、前記オブジェクトの表示属性を変化させ又は前記オブジェクトが重畳された出力画像の状態を変化させることにより、前記ユーザに通知する、請求項 1 ~ 5 のいずれか 1 項に記載の情報処理装置。

10

【請求項 7】

前記制御部は、前記音声認識において検出されている音声のレベルに応じて、前記オブジェクトの前記表示属性又は前記出力画像の前記状態の変化のレベルを変化させる、請求項 6 に記載の情報処理装置。

【請求項 8】

前記音声認識は、可変的な指向性を有するマイクロフォンにより取得される音声信号を用いて実行される、請求項 1 ~ 7 のいずれか 1 項に記載の情報処理装置。

【請求項 9】

20

前記制御部は、前記オブジェクトの位置を前記ユーザの動きに応じて変化させ、前記マイクロフォンの指向性は、前記オブジェクトの位置に応じて設定される、請求項 8 に記載の情報処理装置。

【請求項 10】

前記制御部は、前記オブジェクトの向きを前記ユーザの動きに応じて変化させ、前記マイクロフォンの指向性は、前記オブジェクトの向きに応じて設定される、請求項 8 又は 9 に記載の情報処理装置。

【請求項 11】

前記制御部は、前記音声認識において認識された音声の内容を表すテキストを含む第 1 の追加的なオブジェクトを、前記入力画像に映る前記ユーザの近傍にさらに重畳する、請求項 1 ~ 10 のいずれか 1 項に記載の情報処理装置。

30

【請求項 12】

前記制御部は、前記音声認識が失敗した場合に、前記第 1 の追加的なオブジェクトの表示属性を変化させ又は特別な文字列を前記テキストに挿入することにより、前記音声認識の失敗を前記ユーザに通知する、請求項 11 に記載の情報処理装置。

【請求項 13】

前記制御部は、前記音声認識において検出されている音声のレベルと、前記音声認識を有効に行うために求められる音声のレベルとを示す第 2 の追加的なオブジェクトを、前記入力画像にさらに重畳する、請求項 1 ~ 12 のいずれか 1 項に記載の情報処理装置。

【請求項 14】

40

前記制御部は、1 つ以上の音声コマンドの候補の各々を表すテキストオブジェクトを、前記入力画像にさらに重畳する、請求項 1 ~ 13 のいずれか 1 項に記載の情報処理装置。

【請求項 15】

前記情報処理装置は、テレビジョン装置であり、
前記音声コマンドは、前記情報処理装置を前記ユーザが遠隔的に制御するために発せられるコマンドである、

請求項 14 に記載の情報処理装置。

【請求項 16】

前記オブジェクトは、マイクロフォンを模したアイコンである、請求項 1 ~ 15 のいずれか 1 項に記載の情報処理装置。

50

【請求項 17】

情報処理装置により実行される情報処理方法であって、
入力画像を取得することと、
発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させることと、
前記入力画像に映るユーザの身体を認識することと、
前記ユーザの音声について実行される音声認識を、認識された前記ユーザの身体の所定の部分と前記オブジェクトとの間の前記画面内の位置関係に基づいて、制御することと、
を含む情報処理方法。

【請求項 18】

情報処理装置を制御するコンピュータを、
入力画像を取得する画像取得部と、
発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させる制御部と、
前記入力画像に映るユーザの身体を認識する画像認識部と、
として機能させ、
前記制御部は、前記ユーザの音声について実行される音声認識を、前記画像認識部により認識される前記ユーザの身体の所定の部分と前記オブジェクトとの間の前記画面内の位置関係に基づいて、制御する、

プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、情報処理装置、情報処理方法及びプログラムに関する。

【背景技術】

【0002】

従来、音声認識は、ユーザによる情報機器への入力を支援する技術として利用されている。例えば、下記特許文献1は、ユーザが発した音声の内容を音声認識によってテキストに変換し、得られたテキストを複数のユーザ間のコミュニケーションのために画面に表示する技術を開示している。

【先行技術文献】

【特許文献】

【0003】

【特許文献1】特開2012-58838号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、音声認識が機能し音声入力アクティブになっているタイミングと、ユーザが音声認識のために音声を発するタイミングとが整合しないケースが少なくない。これらタイミングが整合しなければ、認識すべき音声認識されず、又は予期しない音声認識されてしまうという不都合が生じ得る。

【0005】

従って、ユーザが適切なタイミングで音声認識のために音声を発することを支援する、改善された仕組みが提供されることが望ましい。

【課題を解決するための手段】

【0006】

本開示によれば、入力画像を取得する画像取得部と、発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させる制御部と、を備え、前記制御部は、ユーザの音声について実行される音声認識を、前記オブジェクトを用いて制御する、情報処理装置が提供される。

【0007】

また、本開示によれば、情報処理装置により実行される情報処理方法であって、入力画

10

20

30

40

50

像を取得することと、発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させることと、ユーザの音声について実行される音声認識を、前記オブジェクトを用いて制御することと、を含む情報処理方法が提供される。

【0008】

また、本開示によれば、情報処理装置を制御するコンピュータを、入力画像を取得する画像取得部と、発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させる制御部と、として機能させ、前記制御部は、ユーザの音声について実行される音声認識を、前記オブジェクトを用いて制御する、プログラムが提供される。

【発明の効果】

【0009】

10

本開示に係る技術によれば、ユーザが適切なタイミングで音声認識のために音声を発することを支援することができる。

【図面の簡単な説明】

【0010】

【図1】第1の実施形態に係る情報処理装置の概要について説明するための説明図である。

【図2】第2の実施形態に係る情報処理装置の概要について説明するための説明図である。

【図3】第1の実施形態に係る情報処理装置のハードウェア構成の一例を示すブロック図である。

20

【図4】第1の実施形態に係る情報処理装置の論理的機能の構成の一例を示すブロック図である。

【図5】画像認識の結果の一例について説明するための説明図である。

【図6】画像認識の結果の他の例について説明するための説明図である。

【図7】音声認識を制御するために使用される制御オブジェクトの第1の例について説明するための説明図である。

【図8】音声認識を制御するために使用される制御オブジェクトの第2の例について説明するための説明図である。

【図9】音声入力をアクティブ化するためのアクティブ化条件の第1の例について説明するための説明図である。

30

【図10】音声入力をアクティブ化するためのアクティブ化条件の第2の例について説明するための説明図である。

【図11】音声認識結果の視覚的なフィードバックの一例について説明するための説明図である。

【図12】認識された音声の内容を表す追加的な表示オブジェクトの一例について説明するための第1の説明図である。

【図13】認識された音声の内容を表す追加的な表示オブジェクトの一例について説明するための第2の説明図である。

【図14】音声認識を支援する追加的な表示オブジェクトの一例について説明するための説明図である。

40

【図15】マイクロフォンの指向性の制御の一例について説明するための第1の説明図である。

【図16】マイクロフォンの指向性の制御の一例について説明するための第2の説明図である。

【図17】マイクロフォンの指向性の制御の一例について説明するための第3の説明図である。

【図18】出力画像のウィンドウ構成の第1の例について説明するための説明図である。

【図19】出力画像のウィンドウ構成の第2の例について説明するための説明図である。

【図20】第1の制御シナリオについて説明するための説明図である。

【図21】第2の制御シナリオについて説明するための説明図である。

50

【図 2 2】第 3 の制御シナリオについて説明するための説明図である。

【図 2 3】第 4 の制御シナリオについて説明するための説明図である。

【図 2 4】第 1 の実施形態に係る処理の流れの一例を示すフローチャートの前半部である。

【図 2 5】第 1 の実施形態に係る処理の流れの一例を示すフローチャートの後半部である。

【図 2 6】第 2 の実施形態に係る情報処理装置のハードウェア構成の一例を示すブロック図である。

【図 2 7】第 2 の実施形態における制御シナリオの一例について説明するための説明図である。

10

【発明を実施するための形態】

【0011】

以下に添付図面を参照しながら、本開示の好適な実施の形態について詳細に説明する。なお、本明細書及び図面において、実質的に同一の機能構成を有する構成要素については、同一の符号を付することにより重複説明を省略する。

【0012】

また、以下の順序で説明を行う。

1. 概要
2. 第 1 の実施形態
 - 2-1. ハードウェア構成例
 - 2-2. 機能構成例
 - 2-3. 制御シナリオの例
 - 2-4. 処理の流れの例
3. 第 2 の実施形態
4. まとめ

20

【0013】

< 1. 概要 >

本節では、図 1 及び図 2 を用いて、本開示に係る技術が適用され得る情報処理装置の概要について説明する。本開示に係る技術は、ユーザインタフェースのための手段として音声認識を活用する様々な装置及びシステムに適用可能である。一例として、本開示に係る技術は、テレビジョン装置、デジタルスチルカメラ又はデジタルビデオカメラなどのデジタル家電機器に適用されてもよい。また、本開示に係る技術は、P C (Personal Computer)、スマートフォン、P D A (Personal Digital Assistant) 又はゲーム端末などの端末装置に適用されてもよい。また、本開示に係る技術は、カラオケシステム又はアミューズメント装置のような特殊な用途を有するシステム又は装置に適用されてもよい。

30

【0014】

図 1 は、第 1 の実施形態に係る情報処理装置 100 の概要について説明するための説明図である。図 1 を参照すると、情報処理装置 100 は、テレビジョン装置である。情報処理装置 100 は、カメラ 101、マイクロフォン 102 及びディスプレイ 108 を備える。カメラ 101 は、情報処理装置 100 のディスプレイ 108 を見るユーザを撮像する。マイクロフォン 102 は、ユーザが発する音声を集音する。ディスプレイ 108 は、情報処理装置 100 により生成される画像を表示する。ディスプレイ 108 により表示される画像は、コンテンツ画像に加えて、ユーザインタフェース (U I) 画像を含み得る。図 1 の例では、ユーザ U a 及び U b がディスプレイ 108 を見ている。ディスプレイ 108 には、U I 画像 W 0 1 が表示されている。U I 画像 W 0 1 は、カメラ 101 により撮像される撮像画像を用いて生成され、それによりいわゆるミラー表示が実現される。情報処理装置 100 は、音声認識機能を有する。ユーザ U a 及び U b は、マイクロフォン 102 を介して情報処理装置 100 へ音声を入力することにより、情報処理装置 100 を操作し又は情報処理装置 100 へ情報を入力することができる。

40

【0015】

50

図2は、第2の実施形態に係る情報処理装置200の概要について説明するための説明図である。図2を参照すると、情報処理装置200は、タブレットPCである。情報処理装置200は、カメラ201、マイクロフォン202及びディスプレイ208を備える。カメラ201は、情報処理装置200のディスプレイ208を見るユーザを撮像する。マイクロフォン202は、ユーザが発する音声を集音する。ディスプレイ208は、情報処理装置200により生成される画像を表示する。ディスプレイ208により表示される画像は、コンテンツ画像に加えて、UI画像を含み得る。図2の例では、ユーザUcがディスプレイ208を見ている。ディスプレイ208には、UI画像W02が表示されている。UI画像W02は、カメラ201により撮像される撮像画像を用いて生成され、それによりいわゆるミラー表示が実現される。情報処理装置200は、音声認識機能を有する。ユーザUcは、マイクロフォン202を介して情報処理装置200へ音声を入力することにより、情報処理装置200を操作し又は情報処理装置200へ情報を入力することができる。

10

【0016】

これら装置において、音声認識機能が動作し音声入力がアクティブになっている間、ユーザが音声認識のための音声のみを発するとは限らない。また、音声入力がアクティブになっていない時にユーザが音声認識のための音声を発する可能性もある。このようなタイミングの不整合は、認識しなくてもよい音声の認識又は音声認識の不成功などといった、ユーザにとって不都合な結果を招来し得る。そこで、情報処理装置100及び200は、次節より詳細に説明する仕組みに従って、ユーザが適切なタイミングで音声認識のために音声を発することを支援する。

20

【0017】

< 2. 第1の実施形態 >

[2-1. ハードウェア構成例]

図3は、情報処理装置100のハードウェア構成の一例を示すブロック図である。図3を参照すると、情報処理装置100は、カメラ101、マイクロフォン102、入力デバイス103、通信インタフェース(I/F)104、メモリ105、チューナ106、デコーダ107、ディスプレイ108、スピーカ109、遠隔制御I/F110、バス111及びプロセッサ112を備える。

【0018】

(1) カメラ

カメラ101は、CCD(Charge Coupled Device)又はCMOS(Complementary Metal Oxide Semiconductor)などの撮像素子を有し、画像を撮像する。カメラ101により撮像される画像(動画を構成する各フレーム)は、情報処理装置100による処理のための入力画像として扱われる。

【0019】

(2) マイクロフォン

マイクロフォン102は、ユーザにより発せられる音声を集音し、音声信号を生成する。マイクロフォン102により生成される音声信号は、情報処理装置100による音声認識のための入力音声として扱われる。マイクロフォン102は、無指向性マイクロフォンであってもよく、又は固定的な若しくは可変的な指向性を有していてもよい。あるシナリオにおいて、マイクロフォン102は可変的な指向性を有し、その指向性は動的に制御される。

40

【0020】

(3) 入力デバイス

入力デバイス103は、ユーザが情報処理装置100を直接的に操作するために使用されるデバイスである。入力デバイス103は、例えば、情報処理装置100の筐体に配設されるボタン、スイッチ及びダイヤルなどを含み得る。入力デバイス103は、ユーザ入力を検出すると、検出されたユーザ入力に対応する入力信号を生成する。

【0021】

50

(4) 通信インタフェース

通信 I / F 1 0 4 は、情報処理装置 1 0 0 による他の装置との間の通信を仲介する。通信 I / F 1 0 4 は、任意の無線通信プロトコル又は有線通信プロトコルをサポートし、他の装置との間の通信接続を確立する。

【 0 0 2 2 】

(5) メモリ

メモリ 1 0 5 は、半導体メモリ又はハードディスクなどの記憶媒体により構成され、情報処理装置 1 0 0 による処理のためのプログラム及びデータ、並びにコンテンツデータを記憶する。メモリ 1 0 5 により記憶されるデータは、例えば、後に説明する画像認識及び音声認識のための特徴データを含み得る。なお、本明細書で説明するプログラム及びデータの一部又は全部は、メモリ 1 0 5 により記憶されることなく、外部のデータソース（例えば、データサーバ、ネットワークストレージ又は外付けメモリなど）から取得されてもよい。

【 0 0 2 3 】

(6) チューナ

チューナ 1 0 6 は、アンテナ（図示せず）を介して受信される放送信号から、所望のチャンネルのコンテンツ信号を抽出し及び復調する。そして、チューナ 1 0 6 は、復調したコンテンツ信号をデコーダ 1 0 7 へ出力する。

【 0 0 2 4 】

(7) デコーダ

デコーダ 1 0 7 は、チューナ 1 0 6 から入力されるコンテンツ信号からコンテンツデータを復号する。デコーダ 1 0 7 は、通信 I / F 1 0 4 を介して受信されるコンテンツ信号からコンテンツデータを復号してもよい。デコーダ 1 0 7 により復号されるコンテンツデータに基づいて、コンテンツ画像が生成され得る。

【 0 0 2 5 】

(8) ディスプレイ

ディスプレイ 1 0 8 は、L C D (Liquid Crystal Display)、O L E D (Organic Light-Emitting Diode) 又は C R T (Cathode Ray Tube) などにより構成される画面を有し、情報処理装置 1 0 0 により生成される画像を表示する。例えば、図 1 及び図 2 を用いて説明したコンテンツ画像及び U I 画像が、ディスプレイ 1 0 8 の画面に表示され得る。

【 0 0 2 6 】

(9) スピーカ

スピーカ 1 0 9 は、振動板及びアンプなどの回路素子を有し、情報処理装置 1 0 0 により生成される出力音声信号に基づいて、音声を出力する。スピーカ 1 0 9 の音量は、変更可能である。

【 0 0 2 7 】

(10) 遠隔制御インタフェース

遠隔制御 I / F 1 1 0 は、ユーザにより使用されるリモートコントローラから送信される遠隔制御信号（赤外線信号又はその他の無線信号）を受信するインタフェースである。遠隔制御 I / F 1 1 0 は、遠隔制御信号を検出すると、検出された遠隔制御信号に対応する入力信号を生成する。

【 0 0 2 8 】

(11) バス

バス 1 1 1 は、カメラ 1 0 1、マイクロフォン 1 0 2、入力デバイス 1 0 3、通信 I / F 1 0 4、メモリ 1 0 5、チューナ 1 0 6、デコーダ 1 0 7、ディスプレイ 1 0 8、スピーカ 1 0 9、遠隔制御 I / F 1 1 0 及びプロセッサ 1 1 2 を相互に接続する。

【 0 0 2 9 】

(12) プロセッサ

プロセッサ 1 1 2 は、例えば、C P U (Central Processing Unit) 又は D S P (Dig

10

20

30

40

50

ital Signal Processor) などであってよい。プロセッサ 112 は、メモリ 105 又は他の記憶媒体に記憶されるプログラムを実行することにより、後に説明する情報処理装置 100 の様々な機能を動作させる。

【0030】

[2 - 2 . 機能構成例]

図 4 は、図 3 に示した情報処理装置 100 のメモリ 105 及びプロセッサ 112 により実現される論理的機能の構成の一例を示すブロック図である。図 4 を参照すると、情報処理装置 100 は、画像取得部 120、音声取得部 130、アプリケーション部 140、認識部 150、特徴データベース (DB) 160 及び制御部 170 を備える。認識部 150 は、画像認識部 152 及び音声認識部 154 を含む。制御部 170 は、認識制御部 172 及び表示制御部 174 を含む。なお、図 4 に示した機能ブロックの一部は、情報処理装置 100 の外部の (例えば、クラウドコンピューティング環境内の) 装置において実現されてもよい。例えば、画像認識部 152 は、以下に説明する画像認識処理を自ら実行する代わりに、当該処理を外部の画像認識機能に実行させてもよい。同様に、音声認識部 154 は、以下に説明する音声認識処理を自ら実行する代わりに、当該処理を外部の音声認識機能に実行させてもよい。

10

【0031】

(1) 画像取得部

画像取得部 120 は、カメラ 101 により撮像される画像を入力画像として取得する。入力画像は、典型的には、ユーザが映る動画を構成する一連のフレームの各々である。そして、画像取得部 120 は、取得した入力画像を認識部 150 及び制御部 170 へ出力する。

20

【0032】

(2) 音声取得部

音声取得部 130 は、マイクロフォン 102 により生成される音声信号を入力音声として取得する。そして、音声取得部 130 は、取得した入力音声を認識部 150 へ出力する。

【0033】

(3) アプリケーション部

アプリケーション部 140 は、情報処理装置 100 が有する様々なアプリケーション機能を実行する。例えば、テレビジョン番組再生機能、電子番組表表示機能、録画設定機能、写真再生機能、動画再生機能、音楽再生機能及びインターネットブラウジング機能などが、アプリケーション部 140 により実行されてよい。アプリケーション部 140 は、アプリケーション機能を通じて生成される (コンテンツ画像を含み得る) アプリケーション画像及び音声を、制御部 170 へ出力する。

30

【0034】

本実施形態において、アプリケーション部 140 により実行されるアプリケーション機能の少なくとも一部は、後述する音声認識部 154 と連携し、ユーザからの音声入力を受け付ける。例えば、テレビジョン番組再生機能は、音声認識部 154 により認識される音声コマンドに従って、再生されるチャンネル及び音量などの設定を変更し得る。電子番組表表示機能は、音声認識部 154 により認識される音声コマンドに従って、表示すべき電子番組表のチャンネル及び時間帯を変更し得る。写真再生機能は、音声認識部 154 により認識される指定日に撮像された写真を再生し得る。インターネットブラウジング機能は、音声認識部 154 により認識されるキーワードを用いたインターネット検索を実行し得る。

40

【0035】

(4) 画像認識部

画像認識部 152 は、画像取得部 120 から入力される入力画像に映るユーザの身体を認識する。例えば、画像認識部 152 は、入力画像から抽出される画像特徴量をユーザの

50

身体の所定の部分について特徴DB160により予め記憶される画像特徴量と照合することにより、当該所定の部分を認識する。所定の部分とは、例えば、ユーザの手、口及び顔のうちの少なくとも1つを含み得る。

【0036】

図5は、画像認識部152による画像認識の結果の一例について説明するための説明図である。図5を参照すると、入力画像W03にユーザUaが映っている。ユーザUaは、カメラ101の方向を向き、左手を挙げている。画像認識部152は、画像特徴量の照合又はその他の公知の手法を用いて、入力画像W03内の手領域A01、口領域A02及び顔領域A03を認識し得る。そして、画像認識部152は、認識したこれら領域の画像内の位置を示す位置データを、制御部170へ出力する。

10

【0037】

一例として、画像認識部152は、入力画像内で認識した顔領域の部分画像（顔画像）を特徴DB160により予め記憶される既知のユーザの顔画像データと照合することにより、ユーザを識別してもよい。画像認識部152によるユーザ識別結果は、例えば、音声認識の調整、UI画像に表示されるメニューの個人化又はアプリケーション部140によるコンテンツの推薦などの用途に使用され得る。なお、ユーザの識別（即ち、個人認識）は、入力画像ではなく、入力音声に基づいて行われてもよい。

【0038】

本実施形態において、画像認識部152は、入力画像に映るユーザのジェスチャをも認識し得る。なお、本明細書において、ジェスチャとの用語は、ユーザの身体の動的な動きを伴わないいわゆるポーズ（形状）をも含むものとする。

20

【0039】

図6は、画像認識部152による画像認識の結果の他の例について説明するための説明図である。図6を参照すると、入力画像W04にユーザUa及びUbが映っている。ユーザUaは、右手の人差し指を口に当てるジェスチャを行っている。画像認識部152は、入力画像W04内の手領域A04を認識し、ユーザUaの上記ジェスチャをさらに認識し得る。ユーザUbは、両手で口を塞ぐジェスチャを行っている。画像認識部152は、入力画像W04内の手領域A05を認識し、ユーザUbの上記ジェスチャをさらに認識し得る。画像認識部152は、ユーザのジェスチャを認識すると、認識したジェスチャの種類を示すジェスチャデータを、制御部170へ出力する。

30

【0040】

（5）音声認識部

音声認識部154は、音声取得部130から入力される入力音声に基づいて、ユーザの音声を認識する。本実施形態において、音声取得部130から音声認識部154への音声入力は、後述する認識制御部172によりアクティブ化され、又は非アクティブ化される。音声入力がアクティブである間、音声認識部154は、入力音声をその内容を示すテキストに変換する。実行中のアプリケーションがフリーテキストの入力を受け付ける場合には、音声認識部154は、認識した音声の内容を示すテキストを、アプリケーション部140へ出力し得る。その代わりに、実行中のアプリケーションが所定の音声コマンドセット内の音声コマンドの入力を受け付ける場合には、音声認識部154は、ユーザの音声から認識した音声コマンドを識別する識別子を、アプリケーション部140へ出力してもよい。音声入力为非アクティブである間、音声認識部154は、音声認識を実行しない。

40

【0041】

音声認識部154は、音声取得部130から入力される入力音声のレベルを判定し、判定したレベルを制御部170へ通知してもよい。後述する認識制御部172は、音声認識部154から通知される入力音声のレベルに応じて、画面上でのユーザへの様々なフィードバックを行い得る。

【0042】

上述したように、あるシナリオにおいて、マイクロフォン102は可変的な指向性を有する。この場合、後述する認識制御部172により、マイクロフォン102の指向性が設

50

定される。そして、音声認識部 154 は、設定された指向性に対応する方向に位置するユーザの音声を、マイクロフォン 102 により取得される音声信号を用いて認識する。

【0043】

(6) 特徴データベース

特徴 DB 160 は、画像認識部 152 により画像認識のために使用される画像特徴データ、及び音声認識部 154 により音声認識のために使用される音声特徴データを予め記憶する。画像特徴データは、例えば、ユーザの手、口又は顔などの所定の部分の既知の画像特徴量を含み得る。また、画像特徴データは、ユーザごとの顔画像データを含んでもよい。また、画像特徴データは、画像認識部 152 が認識すべきジェスチャを定義するジェスチャ定義データを含んでもよい。音声特徴データは、例えば、ユーザごとの発話の特徴を示す音声特徴量を含み得る。

10

【0044】

(7) 認識制御部 172

認識制御部 172 は、入力画像に重畳されるオブジェクトであって、発話に関連する当該オブジェクトを生成する。そして、認識制御部 172 は、生成した当該オブジェクトを用いて、音声認識部 154 により実行される音声認識を制御する。以下、音声認識を制御するために使用されるこのオブジェクトを、制御オブジェクトという。制御オブジェクトは、ユーザによる操作に従って画面上で移動してもよく、又は固定的な位置に表示されてもよい。

【0045】

20

図 7 は、制御オブジェクトの第 1 の例について説明するための説明図である。図 7 を参照すると、入力画像 W05 に制御オブジェクト IC1 が重畳されている。制御オブジェクト IC1 は、手持ち型のマイクロフォンを模したアイコンである。認識制御部 172 は、例えば、ユーザからの音声入力を受け付けるアプリケーション（以下、音声対応アプリケーションという）が起動されると、画面上の規定の表示位置又は画像認識部 152 により認識されるユーザの身体の近傍に、制御オブジェクト IC1 を表示させる。そして、認識制御部 172 は、ユーザの動き（例えば、手領域の動き）に応じて、制御オブジェクト IC1 の表示位置を変化させる。認識制御部 172 は、ユーザの動き（例えば、手領域の回転）に応じて、制御オブジェクト IC1 の向きを変化させてもよい。音声対応アプリケーションが終了すると、制御オブジェクト IC1 は画面から消去され、又は非アクティブ化されて既定の表示位置若しくは画面の端部へ移動し得る。

30

【0046】

図 8 は、制御オブジェクトの第 2 の例について説明するための説明図である。図 8 を参照すると、入力画像 W06 に制御オブジェクト IC2 が重畳されている。制御オブジェクト IC2 は、スタンド型のマイクロフォンを模したアイコンである。認識制御部 172 は、例えば、音声対応アプリケーションが起動されると、画面上の既定の表示位置に制御オブジェクト IC2 を表示させる。制御オブジェクト IC2 の表示位置は移動しない。音声対応アプリケーションが終了すると、制御オブジェクト IC2 は画面から消去され得る。

【0047】

なお、図 7 及び図 8 に示した制御オブジェクト IC1 及び IC2 は一例に過ぎない。例えば、口若しくは拡声器を模した他の種類のアイコン又はテキストラベルなどが制御オブジェクトとして使用されてもよい。また、制御オブジェクトの外観ではなく、制御オブジェクトの機能が発話に関連していてもよい。

40

【0048】

本実施形態において、認識制御部 172 は、画像認識部 152 により認識されるユーザの身体の所定の部分と制御オブジェクトとの間の画面内の位置関係に基づいて、音声認識部 154 により実行される音声認識を制御する。例えば、認識制御部 172 は、当該位置関係に基づくアクティブ化条件が満たされている場合に、音声認識部 154 への音声入力をアクティブ化する。認識制御部 172 は、アクティブ化条件が満たされていない場合に、音声認識部 154 への音声入力をアクティブ化しない。

50

【 0 0 4 9 】

図 9 は、音声入力をアクティブ化するためのアクティブ化条件の第 1 の例について説明するための説明図である。図 9 を参照すると、入力画像 W 0 7 a 及び W 0 7 b にユーザ U a が映っている。画像認識部 1 5 2 は、入力画像に映るユーザの口領域及び手領域を認識する。第 1 の例において、アクティブ化条件は、ユーザの口と制御オブジェクトとの間の距離が距離閾値 D 1 を下回る、という条件である。図中には、口領域の中心点 G 1 を中心とし半径が距離閾値 D 1 に等しい円が点線で示されている。認識制御部 1 7 2 は、認識される手領域 A 0 1 の動きに従って、制御オブジェクト I C 1 を画面内で移動させる。図 9 の上段では、ユーザの口と制御オブジェクト I C 1 との間の距離が距離閾値 D 1 を上回るため、音声入力は非アクティブである。即ち、ユーザが音声を発しても（又は近傍で雑音が発生しても）、音声認識部 1 5 4 は音声を認識しない。従って、その間、ユーザが意図しない音声認識に起因して、アプリケーションが予期しない動作をすることが防止される。図 9 の下段において、ユーザが手を動かした結果、ユーザの口と制御オブジェクト I C 1 との間の距離が距離閾値 D 1 を下回っている。そこで、認識制御部 1 7 2 は、アクティブ化条件が満たされていると判定し、音声入力をアクティブ化する。すると、ユーザにより発せられる音声が発せられる音声認識部 1 5 4 により認識されるようになる。なお、ユーザの身体の口以外の部分と制御オブジェクトとの間の距離が上記距離閾値と比較されてもよい。

10

【 0 0 5 0 】

図 1 0 は、音声入力をアクティブ化するためのアクティブ化条件の第 2 の例について説明するための説明図である。図 1 0 を参照すると、入力画像 W 0 8 a 及び W 0 8 b にユーザ U b が映っている。また、入力画像 W 0 8 a 及び W 0 8 b に制御オブジェクト I C 2 が重畳されている。画像認識部 1 5 2 は、入力画像に映るユーザの口領域 A 0 6 を認識する。第 2 の例において、アクティブ化条件は、ユーザの口と制御オブジェクトとの間の距離が距離閾値 D 2 を下回る、という条件である。図中には、制御オブジェクト上の基準点 G 2 を中心とし半径が距離閾値 D 2 に等しい円が点線で示されている。図 1 0 の上段では、ユーザの口と制御オブジェクト I C 2 との間の距離が距離閾値 D 2 を上回るため、音声入力は非アクティブである。即ち、ユーザが音声を発しても（又は近傍で雑音が発生しても）、音声認識部 1 5 4 は音声を認識しない。従って、その間、ユーザが意図しない音声認識に起因して、アプリケーションが予期しない動作をすることが防止される。図 1 0 の下段において、ユーザが移動した結果、ユーザの口と制御オブジェクト I C 2 との間の距離が距離閾値 D 2 を下回っている。そこで、認識制御部 1 7 2 は、アクティブ化条件が満たされていると判定し、音声入力をアクティブ化する。すると、ユーザにより発せられる音声が発せられる音声認識部 1 5 4 により認識されるようになる。

20

30

【 0 0 5 1 】

なお、図 9 及び図 1 0 を用いて説明したアクティブ化条件は一例に過ぎない。例えば、制御オブジェクトへのタッチ又は制御オブジェクトを高く掲げるなどといった、制御オブジェクトに関連する所定のジェスチャの検出が、アクティブ化条件として定義されてもよい。

【 0 0 5 2 】

音声入力が一度アクティブ化された後、認識制御部 1 7 2 は、所定の非アクティブ化条件が満たされるまで、音声入力のアクティブ状態を継続させる。非アクティブ化条件は、例えば、上記アクティブ化条件の単純な反対（例えば、ユーザの口と制御オブジェクトとの間の距離が距離閾値を上回る、など）であってもよい。その代わりに、非アクティブ化条件は、画像認識部 1 5 2 によるユーザの所定のジェスチャの認識などであってもよい。音声入力を非アクティブ化するためのジェスチャとは、例えば、人差し指を口に当てるジェスチャなどであってもよい。また、非アクティブ化条件は、一単位の音声コマンドの認識の成功、又はアクティブ化からの所定の期間の経過などを含んでもよい。

40

【 0 0 5 3 】

音声入力がアクティブである間、認識制御部 1 7 2 は、音声認識部 1 5 4 による音声認識に関連するユーザへの視覚的なフィードバックをも制御する。

50

【 0 0 5 4 】

例えば、認識制御部 1 7 2 は、制御オブジェクトの表示属性を変化させることにより、音声認識部 1 5 4 への音声入力がアクティブ化されていることをユーザに通知する。認識制御部 1 7 2 により変更される制御オブジェクトの表示属性は、例えば、色、輝度、透明度、サイズ、形状及びテクスチャのうち少なくとも 1 つを含み得る。図 9 及び図 1 0 の例では、音声入力がアクティブであるか否かが、制御オブジェクトのテクスチャの変化によって示されている。

【 0 0 5 5 】

また、例えば、認識制御部 1 7 2 は、音声認識部 1 5 4 から通知される入力音声のレベルをユーザへフィードバックする。入力音声のレベルのフィードバックは、制御オブジェクトの表示属性を変化させ、又は制御オブジェクトが重畳された U I 画像の状態を変化させることにより行われてよい。図 1 1 は、音声認識結果の視覚的なフィードバックの一例について説明するための説明図である。図 1 1 を参照すると、制御オブジェクト I C 1 が重畳された U I 画像 W 0 9 に、エフェクト F b 1 が適用されている。エフェクト F b 1 は、制御オブジェクト I C 1 (ユーザの口であってもよい) から波動が放出されているかのような U I 画像の状態を表現する。入力音声のレベルが所定の閾値を下回る場合には、エフェクト F b 1 は解除され得る。こうしたフィードバックによれば、ユーザは、自身が発した音声を情報処理装置 1 0 0 が適切に検出しているか否かを、直感的に把握することができる。認識制御部 1 7 2 は、上記所定の閾値を上回る入力音声のレベルに応じて、制御オブジェクトの表示属性の変化のレベル又は出力画像の状態の変化のレベルを変化させてもよい。例えば、入力音声のレベルがより大きいほどより広い画像領域に、エフェクト F b 1 が適用されてもよい。それにより、ユーザは、自身が発した音声について情報処理装置 1 0 0 が検出したレベルを、直感的に把握することができる。なお、認識制御部 1 7 2 は、エフェクト F b 1 の表示属性(例えば、色など)を、音声認識のステータス又はエラーの有無を示すように変化させてもよい。入力音声のレベルの所定の基準値との比較の結果が、U I 画像 W 0 9 においてテキストで示されてもよい。

【 0 0 5 6 】

また、例えば、認識制御部 1 7 2 は、音声認識部 1 5 4 により認識された音声の内容を表すテキストを含む追加的な表示オブジェクトを、入力画像に映るユーザの近傍にさらに重畳してもよい。図 1 2 及び図 1 3 は、認識された音声の内容を表す追加的な表示オブジェクトの一例について説明するための説明図である。図 1 2 を参照すると、制御オブジェクト I C 1 及び追加オブジェクト F b 2 が U I 画像 W 1 0 に重畳されている。追加オブジェクト F b 2 は、U I 画像 W 1 0 に映るユーザ U a が発した音声の内容を表すテキストを含む吹き出しである。こうしたフィードバックによれば、ユーザは、自身が発した音声を情報処理装置 1 0 0 が正しく認識したか否かを、即座に把握することができる。図 1 3 を参照すると、追加オブジェクト F b 2 は、ランダム文字列 S t r 1 を含む。ランダム文字列 S t r 1 は、所定の閾値を上回るレベルの入力音声を検出されたものの、当該入力音声に基づく音声認識が失敗した場合に、追加オブジェクト F b 2 に挿入され得る。こうしたフィードバックによれば、ユーザは、自身が発した音声のレベルが十分であったものの音声認識が失敗したことを、即座に把握することができる。音声認識の失敗は、追加オブジェクト F b 2 の表示属性を変化させることによりユーザに通知されてもよい。なお、追加オブジェクト F b 2 は、ランダム文字列の代わりに、空白を含んでもよい。ランダム文字列又は空白の長さは、音声認識が失敗した間の発話時間の長さに応じて決定されてもよい。

【 0 0 5 7 】

また、例えば、認識制御部 1 7 2 は、音声認識部 1 5 4 により検出されている音声のレベルと、音声認識を有効に行うために求められる音声のレベルとを示す追加的なオブジェクトを入力画像に重畳してもよい。音声認識を有効に行うために求められる音声のレベルは、メモリ 1 0 5 により予め記憶されてもよく、又は環境の雑音レベルに依存して動的に計算されてもよい。図 1 4 は、音声認識を支援する追加的な表示オブジェクトの一例につ

いて説明するための説明図である。図 1 4 を参照すると、U I 画像 W 1 2 に、制御オブジェクト I C 1、追加オブジェクト F b 2 及び追加オブジェクト F b 3 が重畳されている。追加オブジェクト F b 2 は、音声の内容を表すテキストを含む吹き出しである。ここでは、ユーザが発した音声のレベルが十分ではないことに起因して音声認識が失敗した結果、追加オブジェクト F b 2 の背景色が暗い色に変更されている。追加オブジェクト F b 3 は、音声のレベルを通知するインジケータである。追加オブジェクト F b 3 の外側の点線の円周の半径は、音声認識を有効に行うために求められる音声のレベルに対応する。塗りつぶされた円の半径は、音声認識部 1 5 4 から通知される入力音声のレベルに対応する。入力音声のレベルが高くなれば、塗りつぶされた円は大きくなる。なお、追加オブジェクト F b 3 は、図 1 4 の例に限定されず、例えば帯状のインジケータなどであってもよい。こうしたフィードバックによれば、ユーザは、自身が発した音声のレベルが不十分であった場合に、どの程度声を大きくすれば音声認識が成功し得るかを、直感的に把握することができる。なお、認識制御部 1 7 2 は、追加オブジェクト F b 3 の表示属性（例えば、色など）を、音声認識のステータス又はエラーの有無を示すように変化させてもよい。入力音声のレベルの所定の基準値との比較の結果が、U I 画像 W 1 2 においてテキストで示されてもよい。

【 0 0 5 8 】

マイクログフォン 1 0 2 が可変的な指向性を有する場合には、認識制御部 1 7 2 は、制御オブジェクトを用いてマイクログフォン 1 0 2 の指向性を設定することにより、音声認識の精度を向上させてもよい。例えば、認識制御部 1 7 2 は、制御オブジェクトの画面上の位置に応じて、マイクログフォン 1 0 2 の指向性を設定してもよい。また、認識制御部 1 7 2 は、制御オブジェクトの画面上の向きに応じて、マイクログフォン 1 0 2 の指向性を設定してもよい。

【 0 0 5 9 】

図 1 5 ~ 図 1 7 は、マイクログフォンの指向性の制御の一例について説明するための説明図である。図 1 5 の上段において、U I 画像 W 1 3 に、制御オブジェクト I C 1 が重畳されている。制御オブジェクト I C 1 の表示位置は、ユーザ U a の手領域の動きに応じて変化し得る。図示された時点において、制御オブジェクト I C 1 の表示位置は、画面の中央のやや左である。図 1 5 の下段には、ユーザ U a の頭上の視点から見た、情報処理装置 1 0 0 とユーザ U a との間の実空間における位置関係が示されている。認識制御部 1 7 2 は、例えば、カメラ 1 0 1 の画角と制御オブジェクト I C 1 の表示位置とに基づいて、マイクログフォン 1 0 2 の指向性を角度 R 1 に設定する。ユーザ U a は角度 R 1 の方向に存在するため、結果として、ユーザ U a が発する音声をマイクログフォン 1 0 2 がより高い品質で集音することが可能となる。

【 0 0 6 0 】

図 1 6 の上段において、U I 画像 W 1 4 に、制御オブジェクト I C 1 が重畳されている。また、U I 画像 W 1 4 には、ユーザ U a 及び U b が映っている。図示された時点において、制御オブジェクト I C 1 の表示位置は、ユーザ U a よりもむしろユーザ U b の顔の近傍である。図 1 6 の下段には、ユーザ U a 及び U b の頭上の視点から見た、情報処理装置 1 0 0 とユーザ U a 及び U b との間の実空間における位置関係が示されている。認識制御部 1 7 2 は、例えば、カメラ 1 0 1 の画角と制御オブジェクト I C 1 の表示位置とに基づいて、マイクログフォン 1 0 2 の指向性を角度 R 2 に設定する。角度 R 2 の方向にはユーザ U b が存在するため、結果として、ユーザ U b が発する音声をマイクログフォン 1 0 2 がより高い品質で集音することが可能となる。

【 0 0 6 1 】

図 1 7 の上段において、U I 画像 W 1 5 に、制御オブジェクト I C 1 が重畳されている。制御オブジェクト I C 1 の画面上での向きは、ユーザ U a の手領域の向きに応じて変化し得る。U I 画像 W 1 5 には、ユーザ U a 及び U b が映っている。図示された時点において、制御オブジェクト I C 1 は、ユーザ U a より操作され、ユーザ U b の顔領域 A 0 7 の方向に向けられている。図 1 7 の下段には、ユーザ U a 及び U b の頭上の視点から見た、

情報処理装置 100 とユーザ U a 及び U b との間の実空間における位置関係が示されている。認識制御部 172 は、例えば、制御オブジェクト IC 1 の表示位置及び向き、並びにユーザ U b の顔領域 A 07 の位置に基づいて、マイクロフォン 102 の指向性を角度 R 3 に設定する。角度 R 3 の方向にはユーザ U b が存在するため、結果として、ユーザ U b が発する音声をマイクロフォン 102 がより高い品質で集音することが可能となる。

【0062】

図 16 又は図 17 を用いて説明したような手法によれば、複数のユーザが存在する場合に、制御オブジェクト IC 1 をあたかも現実のマイクロフォンであるかのように使用して、音声認識についての発話権をユーザ間で受け渡すことが可能となる。

【0063】

ここまで説明した例以外にも、ユーザのジェスチャに基づく様々なユーザインタフェースが実現されてよい。例えば、認識制御部 172 は、ユーザが手で口を塞ぐジェスチャの認識に応じて、音声認識部 154 によるそれまでの音声認識結果をキャンセルしてもよい。それにより、ユーザが誤った内容の音声を発し又は音声認識部 154 が音声の内容を誤って認識した場合に、ユーザが簡易に音声入力をやり直すことができる。また、認識制御部 172 は、予め定義されるジェスチャの認識に応じて、スピーカ 109 からの音声出力のボリュームを増加させ又は減少させてもよい。

【0064】

また、認識制御部 172 は、1 つ以上の音声コマンド候補の各々を表すテキストオブジェクトを、入力画像にさらに重畳してもよい。それにより、ユーザは、アプリケーション機能が受け付ける音声コマンドを事前に知っていなくても、必要とされる音声コマンドを適切に発することができる。

【0065】

(8) 表示制御部 174

表示制御部 174 は、ディスプレイ 108 を介する画像の表示を制御する。例えば、表示制御部 174 は、アプリケーション部 140 から入力されるアプリケーション画像をディスプレイ 108 に表示させる。また、表示制御部 174 は、音声対応アプリケーションが起動された場合に、認識制御部 172 により生成される UI 画像を、ディスプレイ 108 に表示させる。表示制御部 174 は、UI 画像のみをディスプレイ 108 に表示させてもよく、又はアプリケーション画像及び UI 画像を合成することにより生成される 1 つの出力画像をディスプレイ 108 に表示させてもよい。

【0066】

図 18 及び図 19 は、本実施形態において採用され得る出力画像のウィンドウ構成の例をそれぞれ示している。これら図において、UI 用ウィンドウ W_{UI} 及びアプリケーション用ウィンドウ W_{APP} がディスプレイ 108 により表示される。UI 用ウィンドウ W_{UI} は、認識制御部 172 により生成される UI 画像を表示する。アプリケーション用ウィンドウ W_{APP} は、アプリケーション部 140 から入力されるアプリケーション画像（例えば、コンテンツ画像）を表示する。図 18 の第 1 の例では、アプリケーション用ウィンドウ W_{APP} は、UI 用ウィンドウ W_{UI} の右下のコーナーに合成されている。図 19 の第 2 の例では、UI 用ウィンドウ W_{UI} はアプリケーション用ウィンドウ W_{APP} の一部分にブレンディングされている。こうしたウィンドウ構成によれば、ユーザは、例えばコンテンツ画像を閲覧しながら、リモートコントローラが手元になくても、制御オブジェクトを用いて情報処理装置 100 を自らの音声で操作することができる。

【0067】

[2-3. 制御シナリオの例]

上述した情報処理装置 100 において行われ得るいくつかの制御シナリオの例について、図 20 ~ 図 23 を用いて説明する

【0068】

(1) 第 1 のシナリオ

図 20 は、第 1 の制御シナリオについて説明するための説明図である。図 20 を参照す

10

20

30

40

50

ると、5つのUI画像ST11～ST15が時間軸に沿って示されている。

【0069】

UI画像ST11にはユーザUdが映っており、ミラー表示が実現されている。

【0070】

次のUI画像ST12は、例えば音声対応アプリケーションが起動し、又はユーザが手を挙げるなどのジェスチャをした後に表示され得る。UI画像ST12には、制御オブジェクトIC1が重畳されている。但し、この時点では、音声認識部154への音声入力はアクティブ化されていない。

【0071】

次のUI画像ST13は、例えばユーザUdが制御オブジェクトIC1を口の近傍に移動させた後に表示され得る。認識制御部172は、アクティブ化条件が満たされた結果として、音声認識部154への音声入力をアクティブ化する。UI画像ST13において、制御オブジェクトIC1の表示属性は、アクティブ状態を示すように変化している。

【0072】

次のUI画像ST14は、ユーザUdが音声を発している間に表示され得る。UI画像ST14において、制御オブジェクトIC1の表示属性は、引き続きアクティブ状態を示している。また、UI画像ST14にはエフェクトFb1が適用されると共に、認識された音声の内容を示す追加オブジェクトFb2がUI画像ST14に重畳されている。

【0073】

次のUI画像ST15は、非アクティブ化条件が満たされた場合に表示され得る。ここでは、音声入力を非アクティブ化させるジェスチャとして人差し指を口に当てるジェスチャが定義されているものとする。認識制御部172は、当該ジェスチャの認識に応じて、音声認識部154への音声入力を非アクティブ化する。制御オブジェクトIC1の表示位置は例えば既定の表示位置に戻され、制御オブジェクトIC1の表示属性は非アクティブ状態を示すように変更される。

【0074】

(2) 第2のシナリオ

図21は、第2の制御シナリオについて説明するための説明図である。図21を参照すると、5つのUI画像ST21～ST25が時間軸に沿って示されている。

【0075】

UI画像ST21には、ユーザUdが映っている。また、UI画像ST21に制御オブジェクトIC1が重畳されている。但し、この時点では、音声認識部154への音声入力はアクティブ化されていない。

【0076】

次のUI画像ST22は、例えばユーザUdが制御オブジェクトIC1を口の近傍に移動させた後に表示され得る。認識制御部172は、アクティブ化条件が満たされた結果として、音声認識部154への音声入力をアクティブ化する。UI画像ST22において、制御オブジェクトIC1の表示属性は、アクティブ状態を示すように変化している。

【0077】

次のUI画像ST23は、ユーザUdが音声を発している間に表示され得る。UI画像ST23において、制御オブジェクトIC1の表示属性は、引き続きアクティブ状態を示している。第2の制御シナリオでは、ユーザUdが音声を発している間、手の動きに関わらず、制御オブジェクトIC1の表示位置は、ユーザUdの口の近傍に維持される。従って、ユーザは、例えば電子メールのメッセージのように長い文章を音声で入力するような場合に、手を挙げ続けることで疲れることなく、音声入力を継続することができる。

【0078】

次のUI画像ST24において、ユーザUdは、手で口を塞ぐジェスチャをしている。認識制御部172は、かかるジェスチャの認識に応じて、それまでの音声認識結果をキャンセルする。第2の制御シナリオにおいて、音声認識部154への音声入力のアクティブ状態は、その後も維持される。

10

20

30

40

50

【 0 0 7 9 】

次のUI画像ST 2 5において、ユーザU dは再び音声を発している。その結果、当初ユーザU dが発した音声の内容とは異なる内容の音声、音声認識部1 5 4により適切に認識されている。

【 0 0 8 0 】

(3) 第3のシナリオ

図2 2は、第3の制御シナリオについて説明するための説明図である。図2 2を参照すると、3つのUI画像ST 3 1～ST 3 3が時間軸に沿って示されている。

【 0 0 8 1 】

UI画像ST 3 1にはユーザU dが映っており、ミラー表示が実現されている。

10

【 0 0 8 2 】

次のUI画像ST 3 2は、例えばユーザが手を挙げるなどのジェスチャをした後に表示され得る。UI画像ST 3 2には、制御オブジェクトIC 2が重畳されている。また、UI画像ST 3 2には、音声対応アプリケーションが受け付ける音声コマンド候補(コマンドA～コマンドD)の各々を表す4つのテキストオブジェクトが重畳されている。

【 0 0 8 3 】

次のUI画像ST 3 3において、例えばユーザU dが制御オブジェクトIC 2の近傍に近付いた結果として音声入力がアクティブ化されている。そして、ユーザU dがコマンドBを読み上げる音声を発し、発せられたコマンドBを音声認識部1 5 4が適切に認識している。音声コマンド候補は、例えば、情報処理装置1 0 0をユーザが遠隔的に制御するために予め用意される1つ以上のコマンドであってよい。

20

【 0 0 8 4 】

このように、本実施形態では、ユーザの手元にリモートコントローラがなくても、ユーザが情報処理装置1 0 0を遠隔的に制御することが可能である。例えば、リモートコントローラが紛失した状況、又は他のユーザによりリモートコントローラが保持されている状況でも、ユーザは、ストレスを感じることなく、所望のタイミングで情報処理装置1 0 0を制御することができる。なお、UI画像ST 3 2が表示された後、所定の音声コマンド又はジェスチャの認識に応じて、音声コマンドA～Dを表すテキストオブジェクトが他の音声コマンド候補を表すテキストオブジェクトに置き換えられてもよい。

【 0 0 8 5 】

30

(4) 第4のシナリオ

第4のシナリオは、制御オブジェクトが介在しない補足的なシナリオである。図2 3は、第4の制御シナリオについて説明するための説明図である。図2 3を参照すると、3つのUI画像ST 4 1～ST 4 3が時間軸に沿って示されている。

【 0 0 8 6 】

UI画像ST 4 1にはユーザU dが映っており、ミラー表示が実現されている。

【 0 0 8 7 】

次のUI画像ST 4 2において、ユーザU dは、耳元で手を丸めるジェスチャをしている。認識制御部1 7 2は、かかるジェスチャの認識に応じて、スピーカ1 0 9からの音声出力のボリュームを増加させる。ボリュームの増加量は、ジェスチャが認識されている時間の長さに依存して変化してもよい。

40

【 0 0 8 8 】

次のUI画像ST 4 3において、ユーザU dは、人差し指を口に当てるジェスチャをしている。認識制御部1 7 2は、かかるジェスチャの認識に応じて、スピーカ1 0 9からの音声出力のボリュームを減少させる。ボリュームの減少量は、ジェスチャが認識されている時間の長さに依存して変化してもよい。

【 0 0 8 9 】

このように、本実施形態では、ユーザのジェスチャに基づく様々なユーザインタフェースが実現され得る。音声入力がアクティブか否か、又は音声対応アプリケーションが実行中であるか否かに依存して、同じ種類のジェスチャが互いに異なる意味に解釈されてもよ

50

い。なお、ユーザ独自のジェスチャをユーザに登録させるためのユーザインタフェースが提供されてもよい。例えば、“手で(制御オブジェクトを)払いのける”というジェスチャが登録され、当該ジェスチャが音声入力のアクティブ化/非アクティブ化のためのジェスチャとして定義されてもよい。個々のジェスチャのための動き、及びジェスチャと対応する処理との間のマッピングをユーザにカスタマイズさせるためのユーザインタフェースがさらに提供されてもよい。

【0090】

[2-4. 処理の流れの例]

図24及び図25のフローチャートは、本実施形態に係る情報処理装置100により実行され得る処理の流れの一例を示している。ここで説明する処理は、カメラ101により撮像される動画を構成する一連のフレームの各々について繰り返される。

10

【0091】

図24を参照すると、まず、画像取得部120は、カメラ101により撮像される画像を入力画像として取得する(ステップS100)。そして、画像取得部120は、取得した入力画像を認識部150及び制御部170へ出力する。

【0092】

次に、画像認識部152は、画像取得部120から入力される入力画像に映るユーザの身体を認識する(ステップS105)。例えば、画像認識部152は、入力画像内のユーザの手領域及び口領域を認識し、認識したこれら領域の位置を示す位置データを制御部170へ出力する。また、画像認識部152は、予め定義されるいくつかのユーザのジェスチャを追加的に認識してもよい。

20

【0093】

次に、認識制御部172は、音声対応アプリケーションが起動しているかを判定する(ステップS110)。音声対応アプリケーションが起動していない場合には、その後のステップS115～ステップS160の処理はスキップされる。音声対応アプリケーションが起動している場合(又はステップS105で認識されるジェスチャによって、音声対応アプリケーションが起動された場合)には、処理はステップS115へ進む。

【0094】

ステップS115において、認識制御部172は、発話に関連する制御オブジェクトの表示位置及び向きを決定する(ステップS115)。制御オブジェクトの表示位置は、既定の位置であってもよく、又は画像認識部152により認識されるユーザの手の動きに追従して移動してもよい。同様に、制御オブジェクトの向きは、既定の向きであってもよく、又はユーザの手の動きに追従して回転してもよい。

30

【0095】

次に、マイクロフォン102が可変的な指向性を有する場合には、認識制御部172は、ステップS115において決定した制御オブジェクトの表示位置及び向きに応じて、マイクロフォン102の指向性を設定する(ステップS120)。

【0096】

次に、認識制御部172は、入力画像をミラー表示するUI画像に、ステップS115において決定した表示位置及び向きを有する制御オブジェクトを重畳する(ステップS125)。ここでの制御オブジェクトの表示属性は、音声入力thatアクティブ化されていないことを示す値に設定され得る。

40

【0097】

図25に移り、次に、認識制御部172は、上述したアクティブ化条件及び非アクティブ化条件に従って、音声入力thatアクティブであるかを判定する(ステップS130)。例えば、ユーザの口領域と制御オブジェクトとの間の距離が距離閾値を下回る場合には、アクティブ化条件は満たされていると判定され得る。音声入力thatアクティブであると判定されない場合には、その後のステップS135～ステップS160の処理はスキップされる。音声入力thatアクティブであると判定された場合には、処理はステップS135へ進む。

【0098】

50

ステップS 1 3 5において、認識制御部 1 7 2は、音声認識部 1 5 4への音声入力が必要に応じてアクティブ化し、制御オブジェクトの表示属性を、音声入力が入力されたアクティブ化されていることを示す値に設定する(ステップS 1 3 5)。

【0099】

次に、音声取得部 1 3 0は、マイクロフォン 1 0 2から取得される入力音声を、音声認識部 1 5 4へ出力する(ステップS 1 4 0)。

【0100】

次に、音声認識部 1 5 4は、音声取得部 1 3 0から入力される入力音声に基づいて、ユーザの音声を認識する(ステップS 1 4 5)。そして、音声認識部 1 5 4は、音声認識の結果を、アプリケーション部 1 4 0及び認識制御部 1 7 2へ出力する。

10

【0101】

次に、認識制御部 1 7 2は、音声認識部 1 5 4から入力される音声認識結果についてのフィードバックを、UI画像に適用する(ステップS 1 5 0)。例えば、認識制御部 1 7 2は、図 1 1に例示したエフェクトF b 1をUI画像に適用してもよい。また、認識制御部 1 7 2は、図 1 2～図 1 4に例示した追加オブジェクトF b 2又はF b 3をUI画像に重畳してもよい。

【0102】

次に、認識制御部 1 7 2は、音声認識が成功したか否かを判定する(ステップS 1 5 5)。音声認識が成功していなければ、その後のステップS 1 6 0の処理はスキップされる。音声認識が成功していれば、処理はステップS 1 6 0へ進む。

20

【0103】

ステップS 1 6 0において、アプリケーション部 1 4 0は、音声認識結果に基づくアプリケーション処理を実行する(ステップS 1 6 0)。例えば、アプリケーション部 1 4 0は、認識された音声コマンドに対応する処理を実行してもよい。また、アプリケーション部 1 4 0は、認識された音声の内容を示すテキストを入力情報として受け付けてもよい。

【0104】

次に、表示制御部 1 7 4は、UI画像を含む出力画像をディスプレイ 1 0 8に表示させる(ステップS 1 6 5)。ここで表示される出力画像は、UI画像のみを含んでもよく、又はUI画像及びアプリケーション画像の双方を含んでもよい。その後、処理は図 2 4のステップS 1 0 0へ戻る。

30

【0105】

なお、ここまで、主にUI画像に1つの制御オブジェクトのみが重畳される例を説明した。しかしながら、かかる例に限定されず、UI画像に複数の制御オブジェクトが重畳されてもよい。例えば、入力画像に複数のユーザが映っている場合において、それぞれのユーザについて別個の制御オブジェクトを重畳すれば、制御オブジェクトをユーザ間で受け渡す作業を要することなく、各ユーザが所望のタイミングで音声コマンドを入力することが可能となる。

【0106】

< 3. 第2の実施形態 >

上述したように、本開示に係る技術は、テレビジョン装置に限定されず、様々な種類の装置に適用可能である。そこで、第2の実施形態として、本開示に係る技術がメッセージ交換用アプリケーションを有する情報処理装置 2 0 0に適用される例について説明する。図 2を用いて説明したように、情報処理装置 2 0 0は、タブレットPCである。

40

【0107】

(1) ハードウェア構成例

図 2 6は、情報処理装置 2 0 0のハードウェア構成の一例を示すブロック図である。図 2 6を参照すると、情報処理装置 2 0 0は、カメラ 2 0 1、マイクロフォン 2 0 2、入力デバイス 2 0 3、通信 I / F 2 0 4、メモリ 2 0 5、ディスプレイ 2 0 8、スピーカ 2 0 9、バス 2 1 1及びプロセッサ 2 1 2を備える。

【0108】

50

カメラ２０１は、ＣＣＤ又はＣＭＯＳなどの撮像素子を有し、画像を撮像する。カメラ２０１により撮像される画像（動画を構成する各フレーム）は、情報処理装置２００による処理のための入力画像として扱われる。

【０１０９】

マイクロフォン２０２は、ユーザにより発せられる音声を集音し、音声信号を生成する。マイクロフォン２０２により生成される音声信号は、情報処理装置２００による音声認識のための入力音声として扱われる。

【０１１０】

入力デバイス２０３は、ユーザが情報処理装置２００を操作し又は情報処理装置２００へ情報を入力するために使用されるデバイスである。入力デバイス２０３は、例えば、タッチパネル、ボタン及びスイッチなどを含み得る。入力デバイス２０３は、ユーザ入力を検出すると、検出されたユーザ入力に対応する入力信号を生成する。

【０１１１】

通信Ｉ／Ｆ２０４は、情報処理装置２００による他の装置との間の通信を仲介する。通信Ｉ／Ｆ２０４は、任意の無線通信プロトコル又は有線通信プロトコルをサポートし、他の装置との間の通信接続を確立する。

【０１１２】

メモリ２０５は、半導体メモリ又はハードディスクなどの記憶媒体により構成され、情報処理装置２００による処理のためのプログラム及びデータ、並びにコンテンツデータを記憶する。なお、プログラム及びデータの一部又は全部は、メモリ２０５により記憶されることなく、外部のデータソース（例えば、データサーバ、ネットワークストレージ又は外付けメモリなど）から取得されてもよい。

【０１１３】

ディスプレイ２０８は、ＬＣＤ又はＯＬＥＤなどにより構成される画面を有し、情報処理装置２００により生成される画像を表示する。例えば、第１の実施形態において説明したものと同様のＵＩ画像が、ディスプレイ２０８の画面に表示され得る。

【０１１４】

スピーカ２０９は、振動板及びアンプなどの回路素子を有し、情報処理装置２００により生成される出力音声信号に基づいて、音声を出力する。スピーカ２０９の音量は、変更可能である。

【０１１５】

バス２１１は、カメラ２０１、マイクロフォン２０２、入力デバイス２０３、通信Ｉ／Ｆ２０４、メモリ２０５、ディスプレイ２０８、スピーカ２０９及びプロセッサ２１２を相互に接続する。

【０１１６】

プロセッサ２１２は、例えば、ＣＰＵ又はＤＳＰなどであってよい。プロセッサ２１２は、メモリ２０５又は他の記憶媒体に記憶されるプログラムを実行することにより、第１の実施形態に係る情報処理装置１００のプロセッサ１１２と同様に、情報処理装置２００の様々な機能を動作させる。情報処理装置２００のメモリ２０５及びプロセッサ２１２により実現される論理的機能の構成は、アプリケーション機能が異なることを除き、図４に例示した情報処理装置１００の構成と同様であってよい。

【０１１７】

（２）制御シナリオの例

図２７は、第２の実施形態における制御シナリオの一例について説明するための説明図である。図２７を参照すると、４つの出力画像ＳＴ５１～ＳＴ５４が時間軸に沿って示されている。本シナリオにおいて、各出力画像は、上部のメッセージ交換用アプリケーションのアプリケーション画像と、下部のＵＩ画像とにより構成される。

【０１１８】

出力画像ＳＴ５１において、アプリケーション画像は、メッセージ入力ボックスを含む。メッセージ入力ボックスには、メッセージは入力されていない。ＵＩ画像にはユーザＵ

10

20

30

40

50

dが映っており、ミラー表示が実現されている。

【0119】

次の出力画像ST52は、例えばユーザが手を挙げるなどのジェスチャをした後に表示され得る。出力画像ST52において、UI画像に制御オブジェクトIC1が重畳されている。但し、この時点では、音声入力にはアクティブ化されていない。

【0120】

次の出力画像ST53は、例えばユーザUdが制御オブジェクトIC1を口の近傍に移動させた後に表示され得る。音声入力にはアクティブ化され、制御オブジェクトIC1の表示属性は、アクティブ状態を示すように変化している。メッセージ入力ボックスには、ユーザにより発せられた音声の内容が入力されている。

10

【0121】

次の出力画像ST54は、例えばユーザUdが制御オブジェクトIC1を口の近傍から離れた後に表示され得る。音声入力には非アクティブ化され、制御オブジェクトIC1の表示属性は、非アクティブ状態を示すように変化している。この状態でユーザが音声を発しても、メッセージ入力ボックスには音声の内容は入力されない。従って、ユーザは、手を動かす簡単な動作だけで、音声入力の状態を切り替えて、入力することを望む音声の内容だけをメッセージに含めることができる。

【0122】

<4.まとめ>

ここまで、図1～図27を用いて、本開示に係る技術の実施形態について詳細に説明した。上述した実施形態によれば、入力画像に重畳して表示される制御オブジェクトを用いて、情報機器により実行される音声認識が制御される。従って、ユーザは、画面上の制御オブジェクトの状態を手掛かりとして、音声認識のための適切なタイミングを判断することができる。

20

【0123】

また、上述した実施形態によれば、入力画像内で認識されるユーザの身体の所定の部分と制御オブジェクトとの間の位置関係に基づいて、音声認識が制御される。従って、ユーザは、画面に表示される自らの身体を動かすことにより、音声認識に関連する様々な機能性を扱うことができる。

【0124】

また、上述した実施形態によれば、ユーザの口と制御オブジェクトとの間の距離に基づいて、音声認識のための音声入力にアクティブ化され得る。また、制御オブジェクトは、ユーザの手の動きに従って画面内で移動し得る。従って、ユーザは、制御オブジェクトを移動させ又は自ら制御オブジェクトの方へ移動することにより、意図したタイミングで所望の音声のみを容易に認識させることができる。その際にユーザに求められる動きは、現実のマイクロフォンを扱う動きに類似しているため、こうした仕組みによって、ユーザにとって直感的なユーザインタフェースを実現することができる。

30

【0125】

また、上述した実施形態によれば、音声入力にアクティブ化されているか否かが、制御オブジェクトの表示属性の変化を通じてユーザに通知される。従って、ユーザは、画面上の制御オブジェクトのみに注意を払うだけで、適切なタイミングで発話することができる。

40

【0126】

なお、本明細書において説明した各装置による一連の処理は、典型的には、ソフトウェアを用いて実現される。一連の処理を実現するソフトウェアを構成するプログラムは、例えば、各装置の内部又は外部に設けられる記憶媒体（非一時的な媒体：non-transitory media）に予め格納される。そして、各プログラムは、例えば、実行時にRAM(Random Access Memory)に読み込まれ、CPUなどのプロセッサにより実行される。

【0127】

以上、添付図面を参照しながら本開示の好適な実施形態について詳細に説明したが、本

50

開示の技術的範囲はかかる例に限定されない。本開示の技術分野における通常の知識を有する者であれば、特許請求の範囲に記載された技術的思想の範疇内において、各種の変更例または修正例に想到し得ることは明らかであり、これらについても、当然に本開示の技術的範囲に属するものと了解される。

【0128】

なお、以下のような構成も本開示の技術的範囲に属する。

(1)

入力画像を取得する画像取得部と、

発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させる制御部と、
を備え、

前記制御部は、ユーザの音声について実行される音声認識を、前記オブジェクトを用いて制御する、

情報処理装置。

(2)

前記情報処理装置は、前記入力画像に映るユーザの身体を認識する画像認識部、をさらに備え、

前記制御部は、前記画像認識部により認識されるユーザの身体の所定の部分と前記オブジェクトとの間の前記画面内の位置関係に基づいて、前記音声認識を制御する、

前記(1)に記載の情報処理装置。

(3)

前記所定の部分は、ユーザの口を含み、

前記制御部は、ユーザの口と前記オブジェクトとの間の距離に基づいて、前記音声認識のための音声入力をアクティブ化する、

前記(2)に記載の情報処理装置。

(4)

前記所定の部分は、ユーザの手を含み、

前記制御部は、ユーザの手の動きに従って前記オブジェクトを前記画面内で移動させる、

前記(3)に記載の情報処理装置。

(5)

前記制御部は、前記入力画像に映るユーザのジェスチャに応じて、前記音声認識のための音声入力を非アクティブ化する、前記(3)又は前記(4)に記載の情報処理装置。

(6)

前記制御部は、前記音声認識のための音声入力がアクティブ化されているか否かを、前記オブジェクトの表示属性を変化させることによりユーザに通知する、前記(1)～(5)のいずれか1項に記載の情報処理装置。

(7)

前記制御部は、前記音声認識において音声が発出されているか否かを、前記オブジェクトの表示属性を変化させ又は前記オブジェクトが重畳された出力画像の状態を変化させることにより、ユーザに通知する、前記(1)～(6)のいずれか1項に記載の情報処理装置。

(8)

前記制御部は、前記音声認識において検出されている音声のレベルに応じて、前記オブジェクトの前記表示属性又は前記出力画像の前記状態の変化のレベルを変化させる、前記(7)に記載の情報処理装置。

(9)

前記音声認識は、可変的な指向性を有するマイクロフォンにより取得される音声信号を用いて実行される、前記(1)～(8)のいずれか1項に記載の情報処理装置。

(10)

前記制御部は、前記オブジェクトの位置をユーザの動きに応じて変化させ、

前記マイクロフォンの指向性は、前記オブジェクトの位置に応じて設定される、
前記（９）に記載の情報処理装置。

（１１）

前記制御部は、前記オブジェクトの向きをユーザの動きに応じて変化させ、
前記マイクロフォンの指向性は、前記オブジェクトの向きに応じて設定される、
前記（９）又は前記（１０）に記載の情報処理装置。

（１２）

前記制御部は、前記音声認識において認識された音声の内容を表すテキストを含む第１
の追加的なオブジェクトを、前記入力画像に映るユーザの近傍にさらに重畳する、前記（
１）～（１１）のいずれか１項に記載の情報処理装置。

10

（１３）

前記制御部は、前記音声認識が失敗した場合に、前記第１の追加的なオブジェクトの表
示属性を変化させ又は特別な文字列を前記テキストに挿入することにより、前記音声認識
の失敗をユーザに通知する、前記（１２）に記載の情報処理装置。

（１４）

前記制御部は、前記音声認識において検出されている音声のレベルと、前記音声認識を
有効に行うために求められる音声のレベルとを示す第２の追加的なオブジェクトを、前記
入力画像にさらに重畳する、前記（１）～（１３）のいずれか１項に記載の情報処理装置

。

（１５）

前記制御部は、１つ以上の音声コマンドの候補の各々を表すテキストオブジェクトを、
前記入力画像にさらに重畳する、前記（１）～（１４）のいずれか１項に記載の情報処理
装置。

20

（１６）

前記情報処理装置は、テレビジョン装置であり、

前記音声コマンドは、前記情報処理装置をユーザが遠隔的に制御するために発せられる
コマンドである、

前記（１５）に記載の情報処理装置。

（１７）

前記オブジェクトは、マイクロフォンを模したアイコンである、前記（１）～（１６）
のいずれか１項に記載の情報処理装置。

30

（１８）

情報処理装置により実行される情報処理方法であって、

入力画像を取得することと、

発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させることと、

ユーザの音声について実行される音声認識を、前記オブジェクトを用いて制御すること
と、

を含む情報処理方法。

（１９）

情報処理装置を制御するコンピュータを、

入力画像を取得する画像取得部と、

発話に関連するオブジェクトを前記入力画像に重畳して画面に表示させる制御部と、
として機能させ、

前記制御部は、ユーザの音声について実行される音声認識を、前記オブジェクトを用い
て制御する、

プログラム。

【符号の説明】

【０１２９】

１００，２００ 情報処理装置

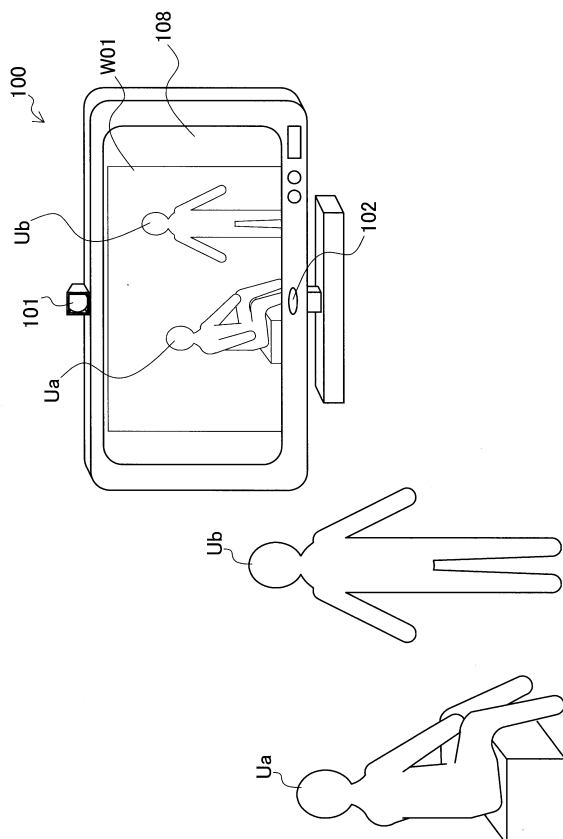
１２０ 画像取得部

40

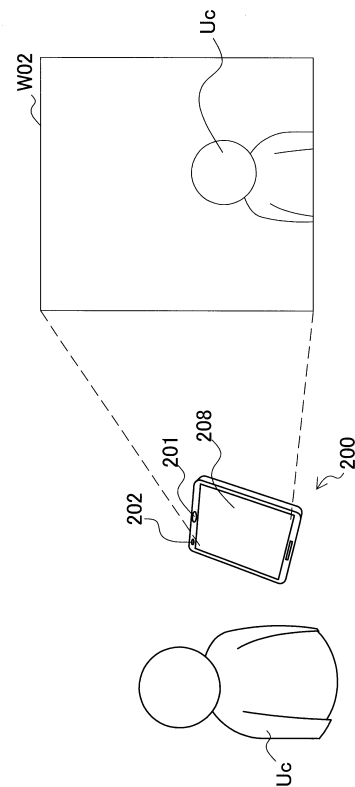
50

1 5 2	画 像 認 識 部
1 5 4	音 声 認 識 部
1 7 2	認 識 制 御 部
1 7 4	表 示 制 御 部
I C 1 , I C 2	制 御 オ ブ ジ ェ ク ト

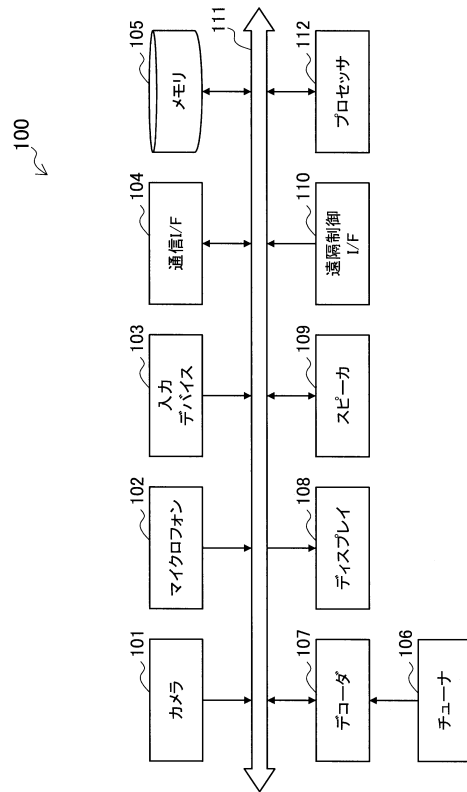
【 図 1 】



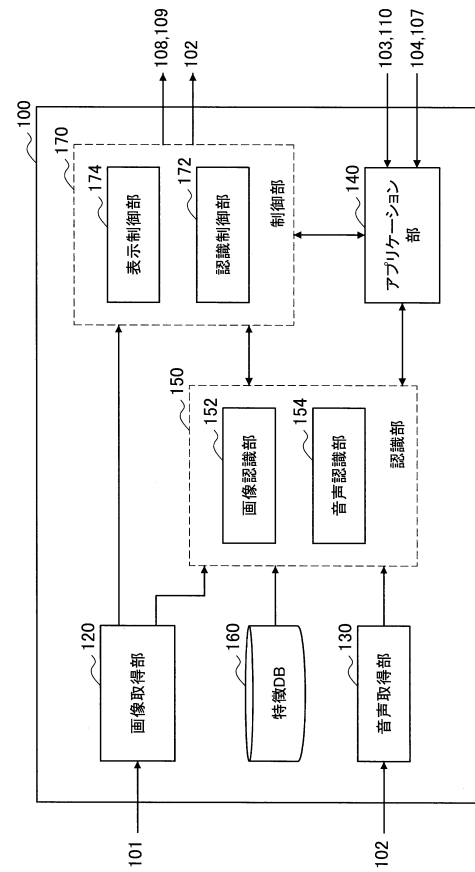
【 図 2 】



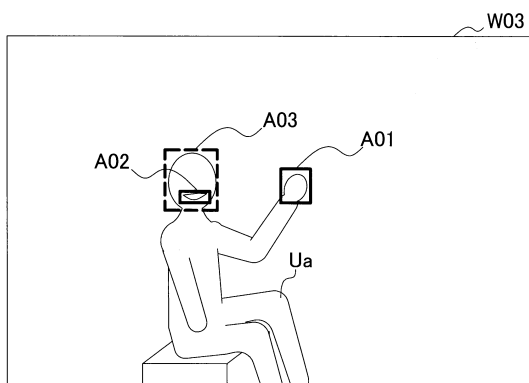
【図 3】



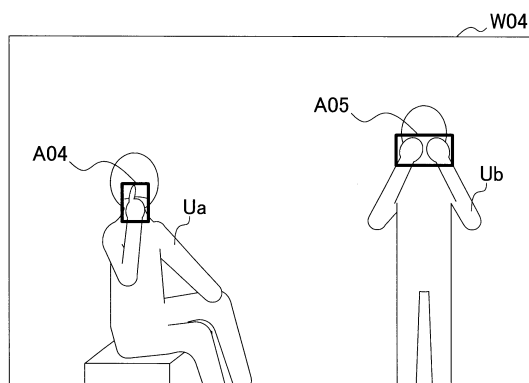
【図 4】



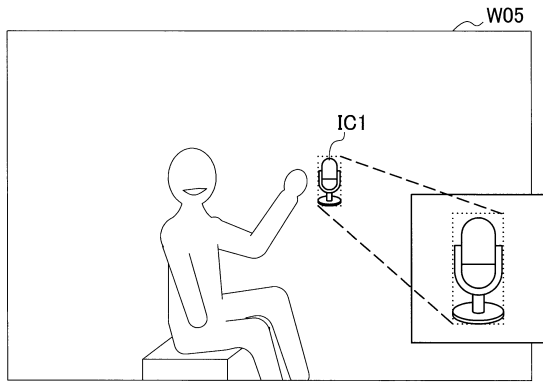
【図 5】



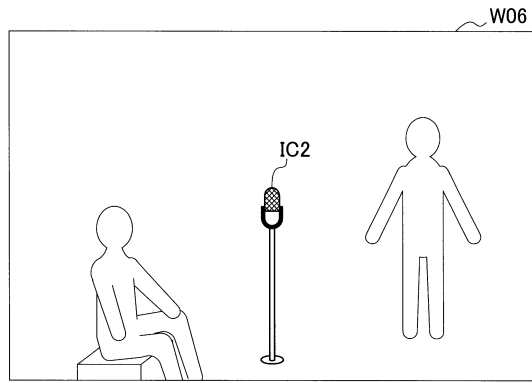
【図 6】



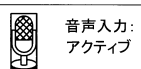
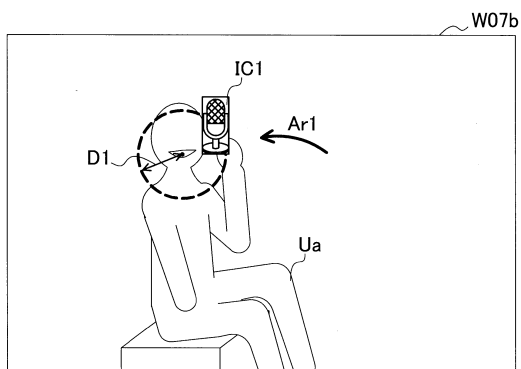
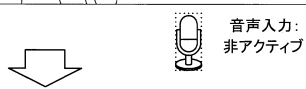
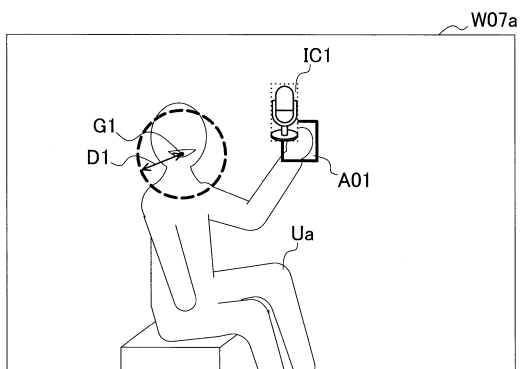
【図 7】



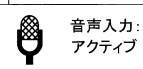
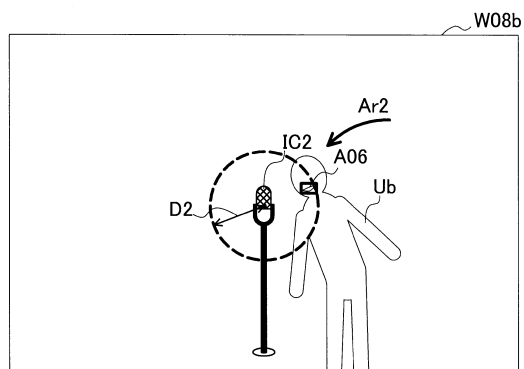
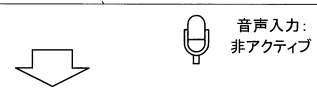
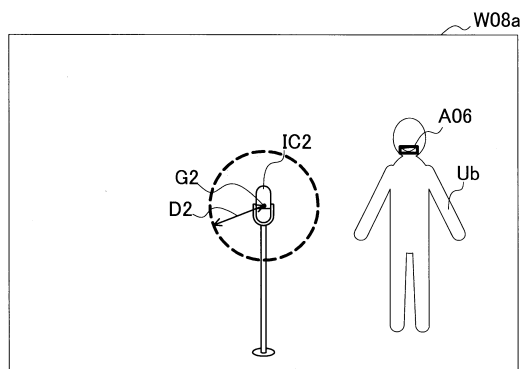
【図 8】



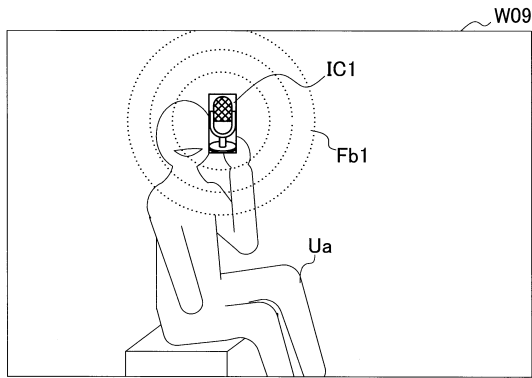
【図 9】



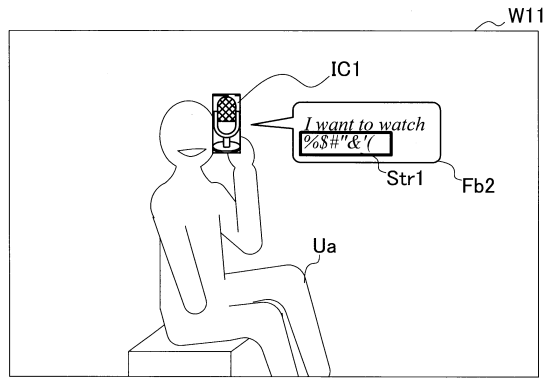
【図 10】



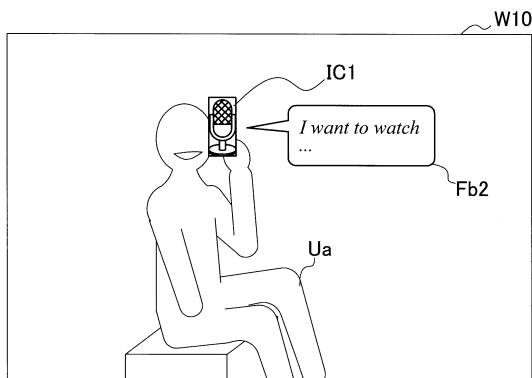
【図 1 1】



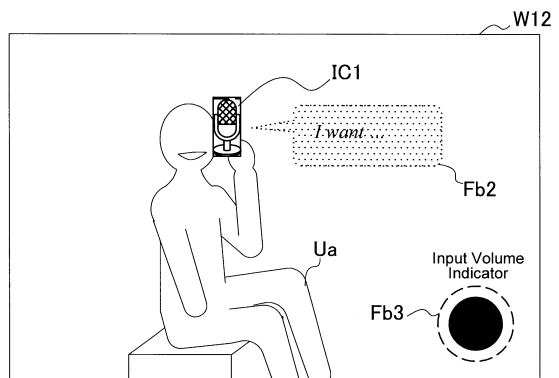
【図 1 3】



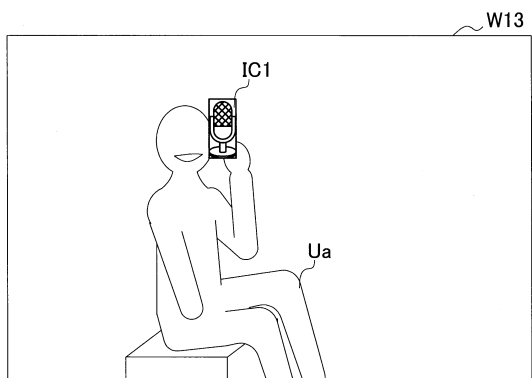
【図 1 2】



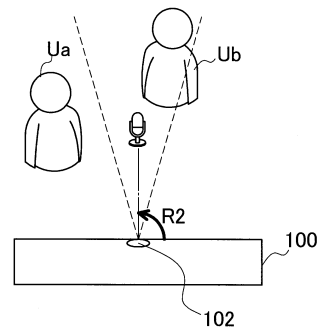
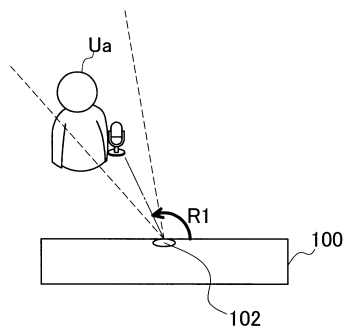
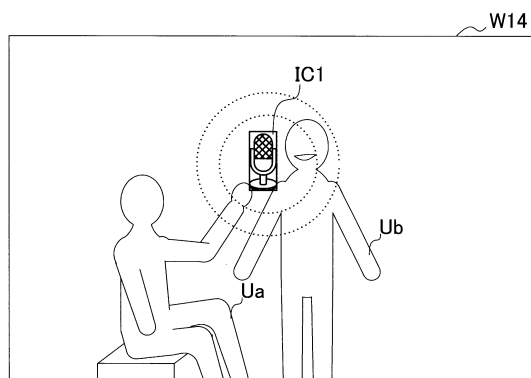
【図 1 4】



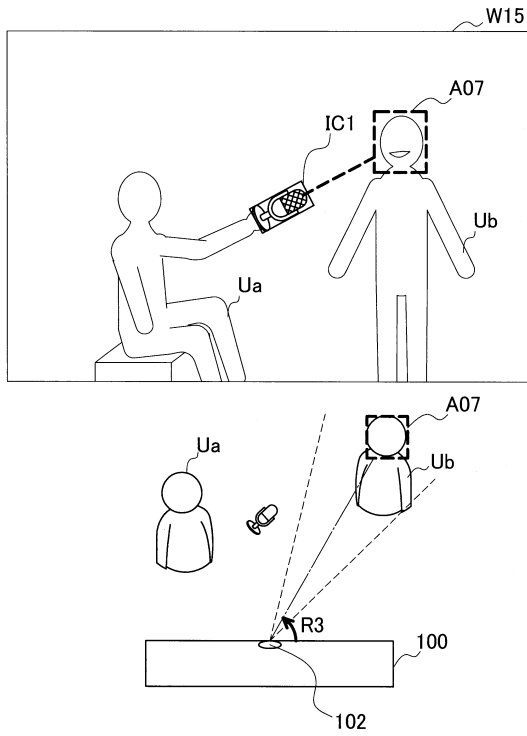
【図 1 5】



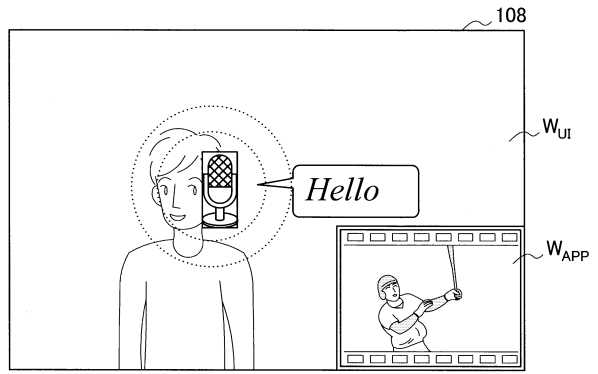
【図 1 6】



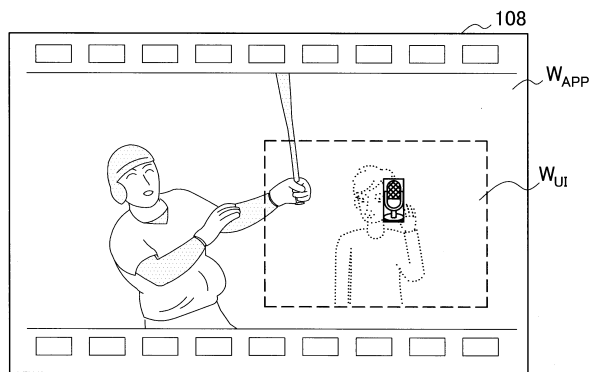
【図 17】



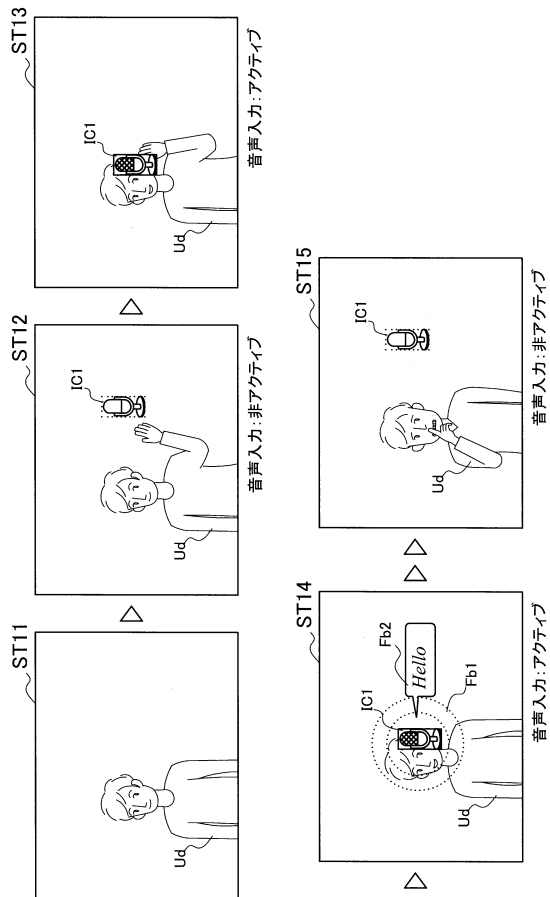
【図 18】



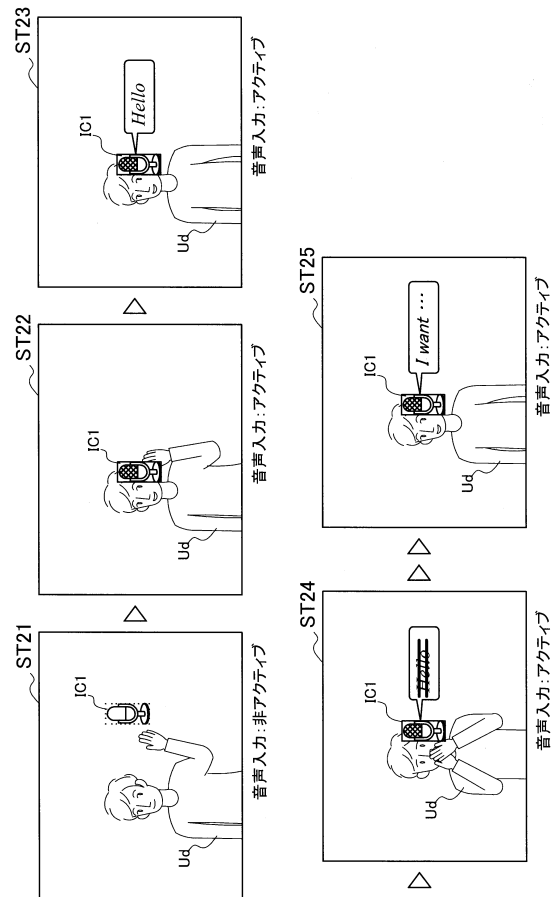
【図 19】



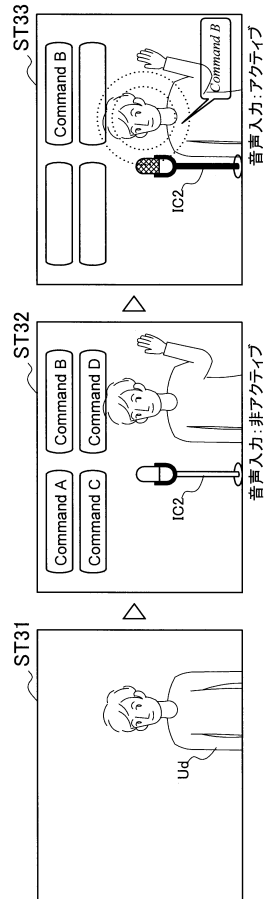
【図 20】



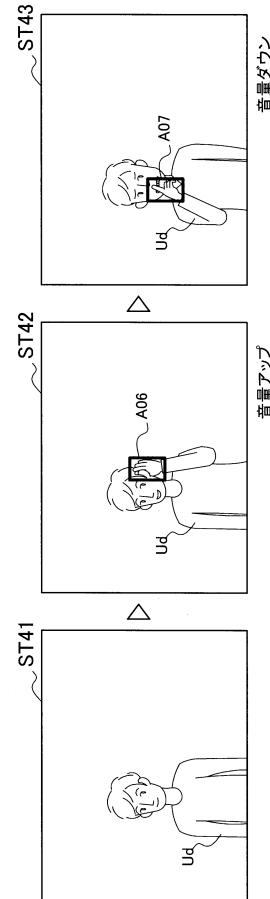
【図 21】



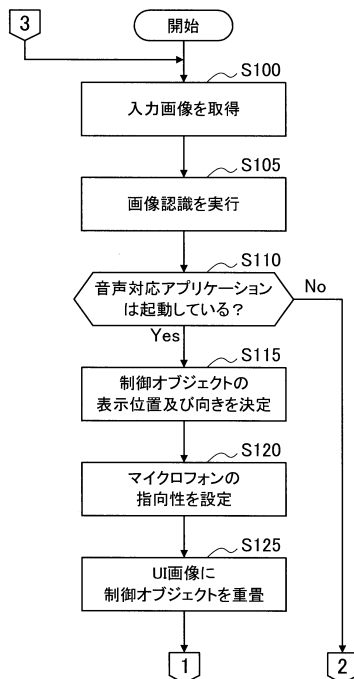
【図 2 2】



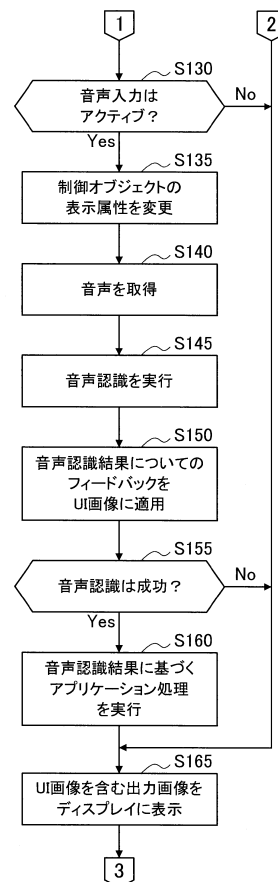
【図 2 3】



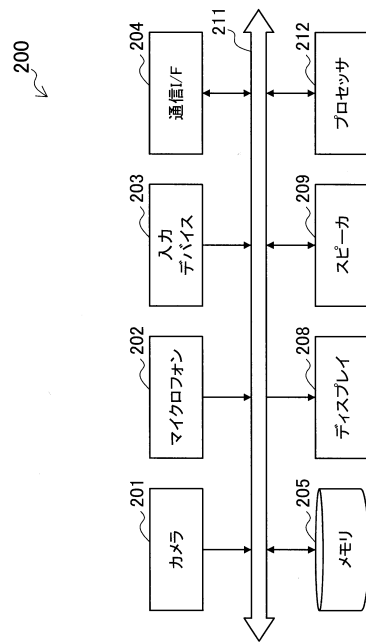
【図 2 4】



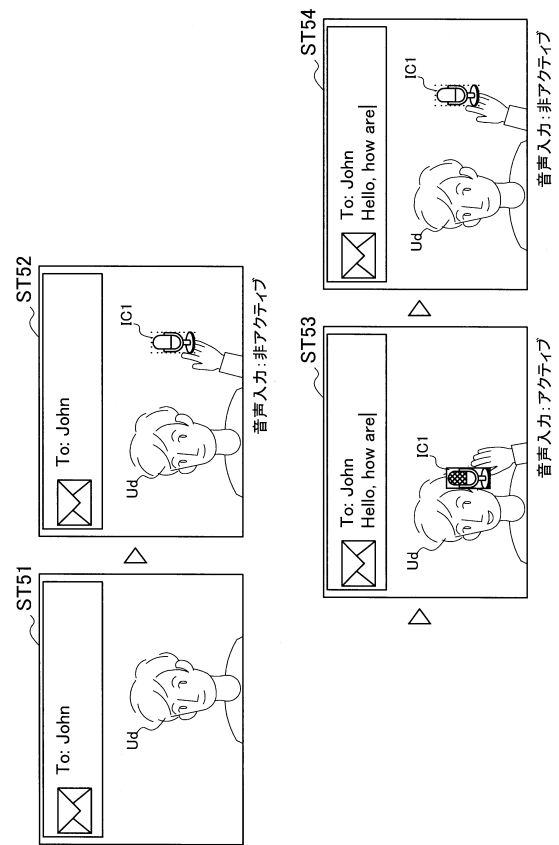
【図 2 5】



【図 26】



【図 27】



フロントページの続き

- (72)発明者 河野 道成
東京都港区港南1丁目7番1号 ソニー株式会社内
- (72)発明者 池田 卓郎
東京都港区港南1丁目7番1号 ソニー株式会社内
- (72)発明者 岡田 憲一
東京都港区港南1丁目7番1号 ソニー株式会社内

審査官 大野 弘

- (56)参考文献 特表2003-526120(JP,A)
特開2004-239963(JP,A)
特開2002-108390(JP,A)
特開2006-294066(JP,A)

- (58)調査した分野(Int.Cl., DB名)
- | | |
|------|-------|
| G10L | 15/22 |
| G10L | 15/28 |