

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第5053950号
(P5053950)

(45) 発行日 平成24年10月24日 (2012.10.24)

(24) 登録日 平成24年8月3日 (2012.8.3)

(51) Int.Cl.

F I

HO 4 N 5/225 (2006.01)

HO 4 N 5/225 F

HO 4 N 101/00 (2006.01)

HO 4 N 101:00

請求項の数 14 (全 37 頁)

(21) 出願番号	特願2008-194800 (P2008-194800)	(73) 特許権者	000001007
(22) 出願日	平成20年7月29日 (2008.7.29)		キヤノン株式会社
(65) 公開番号	特開2010-34841 (P2010-34841A)		東京都大田区下丸子3丁目30番2号
(43) 公開日	平成22年2月12日 (2010.2.12)	(74) 代理人	100126240
審査請求日	平成23年7月12日 (2011.7.12)		弁理士 阿部 琢磨
		(74) 代理人	100124442
			弁理士 黒岩 創吾
		(72) 発明者	山本 寛樹
			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
		審査官	深沢 正志
			最終頁に続く

(54) 【発明の名称】 情報処理方法、情報処理装置、プログラムおよび記憶媒体

(57) 【特許請求の範囲】

【請求項 1】

情報処理装置が行う情報処理方法であって、
第 1 の検出手段が、予め設定された基準を満たす音の開始を検出する第 1 の検出工程と
、
第 1 の取得手段が、前記開始の検出に応答して第 1 の画像データを取得する第 1 の取得
工程と、
第 1 の記憶手段が、前記第 1 の画像データをメモリに記憶する第 1 の記憶工程と、
第 2 の検出手段が、前記音の終了を検出する第 2 の検出工程と、
第 2 の取得手段が、前記終了の検出に応答して第 2 の画像データを取得する第 2 の取得
工程と、
第 2 の記憶手段が、前記第 2 の画像データを前記メモリに記憶する第 2 の記憶工程と、
決定手段が、前記音に含まれる意味に応じて、前記第 1 の画像データまたは前記第 2 の
画像データのいずれかを保存する対象のデータとして決定する決定工程とを有することを
特徴とする情報処理方法。

【請求項 2】

情報処理装置が行う情報処理方法であって、
第 1 の検出手段が、予め設定された基準を満たす音の開始を検出する第 1 の検出工程と
、
第 2 の検出手段が、前記音の終了を検出する第 2 の検出工程と、

取得手段が、前記開始または前記終了の検出に応答して画像データを取得する取得工程と、

記憶手段が、前記取得工程で取得した前記画像データをメモリに記憶する記憶工程と、
決定手段が、前記音に含まれる意味に応じて、前記メモリに記憶した前記画像データを
保存する対象のデータとして決定する決定工程とを有することを特徴とする情報処理方法。

【請求項 3】

更に、消去手段が、前記決定工程で、保存する対象のデータとして決定されなかった画
像データを前記メモリから消去する消去工程を有することを特徴とする請求項 1 または請
求項 2 に記載の情報処理方法。

10

【請求項 4】

更に、保存手段が、前記保存する対象のデータとして決定された前記画像データを第 2
のメモリに保存する保存工程を有することを特徴とする請求項 1 乃至請求項 3 のいずれか
1 項に記載の情報処理方法。

【請求項 5】

前記画像を取得する工程は、前記開始を検出した時点または前記終了を検出した時点に
実行されることを特徴とする請求項 1 乃至請求項 4 のいずれか 1 項に記載の情報処理方法。

【請求項 6】

前記開始を検出した時点で画像データを取得し、前記開始を検出した時点から前記音が
予め設定された時間継続しなかった場合、更に、前記メモリから取得した前記画像データ
を消去する第 2 の消去工程と、

20

前記第 2 の消去工程に続いて、前記第 1 の検出工程に相当する、第 1 の再検出工程とを
実行することを特徴とする請求項 1 乃至請求項 5 のいずれか 1 項に記載の情報処理方法。

【請求項 7】

前記終了を検出した時点で画像データを取得し、前記終了を検出した時点から予め設定
された時間に再び予め設定した基準を満たす音を検出した場合、更に、前記メモリが取得
したら前記画像データを消去する第 3 の消去工程と、

前記第 3 の消去工程に続いて、前記第 2 の検出工程に相当する、第 2 の再検出工程とを
実行することを特徴とする請求項 1 乃至請求項 6 のいずれか 1 項に記載の情報処理方法。

30

【請求項 8】

前記画像を取得する工程は、前記開始を検出した時点から予め設定した遅延時間が経過
した時点または前記終了を検出した時点から予め設定した遅延時間が経過した時点に実行
されることを特徴とする請求項 1 乃至請求項 7 のいずれか 1 項に記載の情報処理方法。

【請求項 9】

前記予め設定された基準とは、一定以上の音量を有するか否かであることを特徴とする
請求項 1 乃至請求項 8 のいずれか 1 項に記載の情報処理方法。

【請求項 10】

前記音を音声認識することによって前記意味を特定することを特徴とする請求項 1 乃至
請求項 9 のいずれか 1 項に記載の情報処理方法。

40

【請求項 11】

予め設定された基準を満たす音の開始を検出する第 1 の検出手段と、
前記開始の検出に₍₁₎応答して第 1 の画像データを取得する第 1 の取得手段と、
前記第 1 の画像データをメモリに記憶する第 1 の記憶手段と、
前記音の終了を検出する第 2 の検出手段と、
前記終了の検出に₍₂₎応答して第 2 の画像データを取得する第 2 の取得手段と、
前記第 2 の画像データを前記メモリに記憶する第 2 の記憶手段と、
前記音に含まれる意味に応じて、前記第 1 の画像データまたは前記第 2 の画像データの
いずれかを保存する対象のデータとして決定する決定手段とを有することを特徴とする情
報処理装置。

50

【請求項 1 2】

予め設定された基準を満たす音の開始を検出する第 1 の検出手段と、
前記音の終了を検出する第 2 の検出手段と、
前記開始または前記終了の検出に応答して画像データを取得する取得手段と、
前記取得手段で取得した前記画像データをメモリに記憶する記憶手段と、
前記音に含まれる意味に応じて、前記メモリに記憶した前記画像データを保存する対象のデータとして決定する決定手段とを有することを特徴とする情報処理装置。

【請求項 1 3】

コンピュータを、請求項 1 1 又は請求項 1 2 に記載の情報処理装置が有する各手段として機能させるためのプログラム。

10

【請求項 1 4】

請求項 1 3 に記載のプログラムを記憶したコンピュータ読み取り可能な記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音で撮像を指示する技術に関する。

【背景技術】

【0002】

従来、一定以上の音量を検知すると撮像動作を実行する機能（以下、音量検知シャッターとする）を備えたカメラが知られている（例えば、特許文献 1）。この機能を利用すると、発声のタイミングに合わせて撮像することが可能となる。

20

【0003】

また、撮像を指示する声を認識すると撮像動作を実行する機能（以下、音声認識シャッターとする）を備えたカメラが知られている（例えば、特許文献 2）。この機能を利用すると、ユーザが撮像を所望して発声した場合に撮像することが可能となる。なお、音声認識シャッターを利用して撮像する場合、ユーザが撮像を指示する声を発しても、音声コマンドの発声が完了するまではカメラの撮像動作は実行されない。よって、所望する撮像のタイミングを逃してしまうことがある。

【特許文献 1】特開平 1 1 - 1 9 4 3 9 2 号公報

【特許文献 2】特開 2 0 0 6 - 1 8 4 5 8 9 号公報

30

【発明の開示】

【発明が解決しようとする課題】

【0004】

従来の音量検知シャッターを利用して撮像する場合には、音の発声タイミングに連動して撮像動作を実行できる。しかし、この場合には、例えば大きな雑音等、目的の声以外を検知した場合にも撮像動作を実行してしまうため、不要な画像を保存してしまうという課題がある。

【0005】

例えば、“撮影”という発声に基づいてユーザが所望するタイミングで撮像する工程と、“消去”という音声コマンドに基づいて既に撮像した画像を消去する工程を備えることにより、上記課題を解決できる。しかしながら、2 種類の音声コマンドを入力する作業は効率がよくない。

40

【0006】

本発明は係る従来例を鑑みてなされたものであり、単一の音声コマンドに基づいて、音が入力されたタイミングを反映した撮像で得られた画像であって、かつユーザが所望する画像を効率良く保存することを主な目的とする。

【課題を解決するための手段】

【0007】

上記目的を達成するための情報処理方法の 1 つとして、情報処理装置が行う情報処理方法であって、第 1 の検出手段が、予め設定された基準を満たす音の開始を検出する第 1 の

50

検出工程と、第１の取得手段が、前記開始の検出に応答して第１の画像データを取得する第１の取得工程と、第１の記憶手段が、前記第１の画像データをメモリに記憶する第１の記憶工程と、第２の検出手段が、前記音の終了を検出する第２の検出工程と、第２の取得手段が、前記終了の検出に応答して第２の画像データを取得する第２の取得工程と、第２の記憶手段が、前記第２の画像データを前記メモリに記憶する第２の記憶工程と、決定手段が、前記音に含まれる意味に応じて、前記第１の画像データまたは前記第２の画像データのいずれかを保存する対象のデータとして決定する決定工程とを有することを特徴とする。

【発明の効果】

10

【０００８】

本発明によれば、単一の音声コマンドに基づいて、音が入力されたタイミングを反映した撮像で得られた画像であって、かつユーザが所望する画像を効率良く得ることができる。

【発明を実施するための最良の形態】

【０００９】

以下、本発明に好適な実施形態について、図面を参照しながら説明していく。

【００１０】

（第１の実施形態）

図１は第１の実施形態に係る情報処理装置の構成の一例であるデジタルカメラを示す機能ブロック図である。

20

【００１１】

図１においてデジタルカメラ２００は、制御部１０１、操作部１０２、撮像部１０３、メモリ（画像記憶用）１１０、記憶媒体（画像記憶用）１１１を備える。

【００１２】

また、デジタルカメラ２００は、マイク１１２、メモリ（音声認識データ用）１１３、メモリ（認識結果制御テーブル用）１１４、ディスプレイ１１５を備える。

【００１３】

（各部の説明）

制御部１０１は、操作部１０２、撮像部１０３、メモリ（画像記憶用）１１０、記憶媒体（画像記憶用）１１１、マイク１１２、メモリ（音声認識データ用）１１３、メモリ（認識結果制御テーブル用）１１４、ディスプレイ１１５の動作を制御する。

30

【００１４】

尚、制御部１０１における処理は後述する。

【００１５】

また、制御部１０１は、CPU（中央演算装置）、ROM（Read Only Memory）、RAM（Random Access Memory）等によって構成される。

【００１６】

また、制御部１０１は、ソフトウェアモジュールとして操作制御部１２２、撮像制御部１２３、画像記憶制御部１０４、音声入力部１０５、音声検出部１０６、音声認識部１０７、認識結果処理部１０８、表示制御部１０９を有する。

40

【００１７】

操作制御部１２２は、ユーザが操作部１０２に対して行った操作を検知するための部分である。

【００１８】

撮像制御部１２３は、撮像部１０３に撮像動作を実行させるための部分である。

【００１９】

画像記憶制御部１０４は、メモリ（画像記憶用）１１０および記憶媒体（画像記憶用）１１１へのデータの書込み、メモリ（画像記憶用）１１０および記憶媒体（画像記憶用）

50

1 1 1 に記憶されているデータの読み出し、消去等を制御する。

【 0 0 2 0 】

音声入力部 1 0 5 は、マイク 1 1 2 を介して入力される音をデジタルの音声信号に変換して出力する部分である。

【 0 0 2 1 】

音声検出部 1 0 6 は、音声入力部 1 0 5 が変換したデジタルの音声信号をフレーム単位で順次処理し、基準を満たす音の開始および終了を検出する。

【 0 0 2 2 】

尚、ユーザが発声した区間（時間帯）は、基準を満たす音の開始を検出してから基準を満たす音の終了を検出するまでの時間帯を発声区間とする。

10

【 0 0 2 3 】

尚、フレームとは、時間的に変化する音声信号をほぼ定常とみなせる固定時間長（例えば、25.6 ミリ秒とする）毎に区分するために設けた処理単位である。なお、このフレーム数によって時間を表現することも可能である。

【 0 0 2 4 】

音声認識部 1 0 7 は、ソフトウェアモジュールとして音響分析部、探索部を有し、ユーザが発声した区間に含まれるコマンド（いわゆる音声コマンド）を認識する。

【 0 0 2 5 】

尚、コマンドとは音声認識部 1 0 7 が認識可能な音のまとまりであり、例えば、“Shoot” 等である。

20

【 0 0 2 6 】

音響分析部は、音声信号をフレーム単位で分析して、例えば MFCC (Mel Frequency Cepstrum Coefficient) 等の特徴量のデータを出力する。

【 0 0 2 7 】

探索部は、例えば、Viterbi アルゴリズム等の周知のアルゴリズムを用いた探索処理を行い、所定個数のコマンドと、各々のコマンドに対応する認識スコアとを認識結果として出力する。

【 0 0 2 8 】

また、探索部は、探索処理を実行する際、メモリ（音声認識データ用）1 1 3 に含まれる音響モデルと言語モデルとを用いる。

30

【 0 0 2 9 】

尚、音響モデル、言語モデルの詳細は後述する。

【 0 0 3 0 】

尚、認識スコアとは、音響的な類似度を示す周知の音響スコア、言語モデルから求まる周知の言語スコア、またはこれら 2 つの重みつき和であってもよい。また、認識結果の確からしさを示す周知の信頼度スコアでもよい。

【 0 0 3 1 】

尚、異なるスコアまたは複数のスコアの用いることで、種々の音の応じた最適な探索処理を実行することが可能となる。

40

【 0 0 3 2 】

認識結果処理部 1 0 8 は、音声認識部 1 0 7 が出力した認識結果のデータを取得し、メモリ（認識結果制御テーブル用）1 1 4 に記憶された認識結果制御テーブルを参照して、認識結果に含まれるコマンドに対応する制御を決定する。

【 0 0 3 3 】

尚、本実施形態に利用する認識結果制御テーブルの一例は後述する。

【 0 0 3 4 】

表示制御部 1 0 9 は、ディスプレイ 1 1 5 に表示する表示内容を制御する。

操作部 1 0 2 は、ユーザがデジタルカメラ 2 0 0 を手動で操作するため部分である。

尚、操作部 1 0 2 は、ボタン、スイッチ等によって構成される。

50

【 0 0 3 5 】

撮像部 1 0 3 は、レンズによって結像した像の撮像信号を生成し、生成された撮像信号に A / D 変換等の画像処理を施す。

尚、撮像部 1 0 3 は、レンズ、撮像センサ等によって構成される。

【 0 0 3 6 】

メモリ（画像記憶用）1 1 0 は、撮像部 1 0 3 が撮像した画像の画像データを一時的に記憶する。尚、メモリ（画像記憶用）1 1 0 は、R A M 等である。

【 0 0 3 7 】

記憶媒体（画像記憶用）1 1 1 は、撮像部 1 0 3 が撮像した画像の画像データを最終的に蓄積する。尚、記憶媒体（画像記憶用）1 1 1 は、不揮発性メモリである。

10

【 0 0 3 8 】

メモリ（画像記憶用）1 1 0 は第 1 のメモリとして機能し、記憶媒体（画像記憶用）は第 2 のメモリとして機能する。

【 0 0 3 9 】

マイク 1 1 2 は、ユーザの音声入力を受け付け、入力された音声データを音声入力部 1 0 4 に出力する。

尚、マイク 1 1 2 は、周知のモノラルマイク、ステレオマイク等である。

【 0 0 4 0 】

メモリ（音声認識データ用）1 1 3 は、音声認識の実行に必要なデータと、例えば H M M (H i d d e n M a r k o v M o d e l) 等の周知の音響モデルと、N - g r a m 、形態素解析等の周知の言語モデルとを記憶する。

20

【 0 0 4 1 】

尚、N - g r a m とは、語の連鎖確率等を用いて言語の統計的な情報によって構成された言語モデルである。

【 0 0 4 2 】

また、音声認識で受理可能な特定の語や語の接続規則を記述した音声認識文法を言語モデルとして利用してもよい。尚、本実施形態に利用する音声認識文法の一例は、後述する。

【 0 0 4 3 】

また、メモリ（音声認識データ用）1 1 3 は、不揮発性メモリ等である。

30

メモリ（認識結果制御テーブル用）1 1 4 は、認識結果制御テーブルを格納する。また、メモリ（認識結果制御テーブル用）1 1 4 は、不揮発性メモリである。

尚、本実施形態に利用する認識結果制御テーブルの一例は後述する。

【 0 0 4 4 】

尚、不揮発性メモリとは、周知のハードディスク、コンパクトフラッシュ（登録商標）、S D (S e c u r e D i g i t a l) カード等でもよい。

【 0 0 4 5 】

また、不揮発性メモリとは、C D (C o m p a c t D i s k)、D V D (D i g i t a l V e r s a t i l e D i s k) 等でもよい。

【 0 0 4 6 】

また、不揮発性メモリとは、L A N (L o c a l A r e a N e t w o r k) アダプタ、U S B (U n i v e r s a l S e r i a l B u s) アダプタ等のインタフェースを介して情報処理装置 1 0 0 と接続可能な外部の記憶媒体であってもよい。

40

【 0 0 4 7 】

ディスプレイ 1 1 5 は、撮像部 1 0 3 で撮像された画像、メモリ（画像記憶用）1 1 0 、記憶媒体（画像記憶用）1 1 1 等に記憶された画像等を表示する。

【 0 0 4 8 】

また、ディスプレイ 1 1 5 は、例えば L C D (L i q u i d C r y s t a l D i s p l a y) や有機 E L (E l e c t r o - L u m i n e s c e n c e) 等である。

【 0 0 4 9 】

50

(カメラ本体の外観の説明)

図2は、本実施形態で想定されるデジタルカメラの外観を示す図である。尚、図2(A)はデジタルカメラ200の前面の外観、図2(B)はデジタルカメラ200の背面の外観である。

【0050】

尚、図1と共通の要素には同一の符号を付し、その説明を省略する。

【0051】

図2において、デジタルカメラ200は、シャッターボタン201、音声シャッター切替えスイッチ202、モードダイヤル203、四方向選択ボタン204、決定ボタン205、電源ボタン206、録音ボタン207を備える。これらは、図1の操作部102に相当する。

10

【0052】

(デジタルカメラ200の各部の説明)

201は、撮像を指示するための操作に用いるシャッターボタンである。

202は、音声指示によって撮像動作を実行する機能を使用するか否かを切り替える音声シャッター切替えスイッチである。

203は、回転することにより、デジタルカメラ200の動作モードを周知の撮影モード、再生モード等に切り替えるモードダイヤルである。

204は、上下左右の任意方向の指示を入力する四方向選択ボタンである。

205は、各種の操作の確定を指示する決定ボタンである。

20

206は、デジタルカメラ200の電源のON/OFFを切り替えるための電源ボタンである。

207は、音声入力の開始および終了を指示する手動操作に用いる録音ボタンである。

【0053】

(音声検出部106の説明)

次に、音声検出部106の機能の詳細を説明する。

【0054】

音声検出部106は、所定の基準(開始条件)を満たした音を検出し、所定の基準(終了条件)を満たした音を検出する。

【0055】

30

続いて所定の基準を満たした音を検出した時点から予め設定された時間が経過した時点で所定の基準を満たした音であることを確定する。

【0056】

また、入力される音声信号の変化によっては、所定の基準を満たした音ではないと判断する、すなわち所定の基準を満たした音の検出を取り消す。

【0057】

(音声検出部106によって判定される検出状態を示す図)

図3は、音声検出部106によって判定された検出状態の一例を示す図である。

【0058】

音声検出部106は、音声信号の検出状況によって仮想的に4つの状態のいずれかに遷移する。

40

第一状態301は、音の入力を開始した直後の状態、すなわち音声信号を検出していない状態(以下、SILENCEとする)とする。

第二状態302は、所定の基準を満たす音の開始を検出し、音の開始の検出を確定していない状態(以下、POSSIBLE SPEECHとする)とする。

第三状態303は、所定の基準を満たす音の開始が確定した状態(以下、SPEECHとする)とする。

第四状態304は、音の入力を終了した直後の状態、すなわち音の開始の検出を確定していない状態(以下、POSSIBLE SILENCEとする)とする。

【0059】

50

尚、本実施形態では音の検出状況を仮想的に4つの状態に分類する例を示すが、第二状態302と第四状態304をまとめて、3つの状態に分類して音の検出状況を判断しても本実施形態と同様の効果が得られる。

【0060】

(検出状態の遷移についての説明)

第一状態301において、音の開始(マイク112からの所定の基準を満たす音の入力の開始)を検出すると第二状態302に遷移する(305)。

第二状態302において、音の開始を取り消すと第一状態301に遷移する(306)。また、第二状態302において、音の開始を確定すると第三状態303に遷移する(307)。

10

第三状態303において、音の終了(マイク112からの所定の基準を満たす音の入力の終了)を検出すると第四状態304に遷移する(308)。

第四状態304において、音の終了を取り消すと第三状態303に遷移する(309)。また、第四状態304において、所定の基準を満たす音の終了を確定すると音の検出を終了する(310)。

第四状態304から所定の基準を満たす音の終了を確定すると音の検出を終了させることで、後述する音声認識の処理の際に、音声検出の処理による計算資源、電力等の消費を抑えることが可能となる。

【0061】

尚、第四状態304において所定の基準を満たす音の終了を確定した場合に、第一状態301に遷移するようにしてもよい。

20

第四状態304から第一状態301に遷移させることで、続けて次の発声を検出することが可能となる。

【0062】

(音声検出部106による処理の概念図)

図4は、音声検出部106による処理の一例を示す概念図である。

【0063】

図4は、ユーザが“Shoot”という言葉が発声した場合の様子を示している。

【0064】

尚、“Shoot”は撮像を指示するコマンドの一例であり、コマンドの種類については後述する。

30

【0065】

図4において420は音声信号である。

また、音声信号のうち421に示した区間の音声信号はユーザの発声ではなく雑音を検出したものである。

また、音声信号のうち422に示した区間の音声信号はユーザが“Shoot”と発声した音を検出したものである。

【0066】

本実施形態の音声検出部106は、所定の基準を満たす音か否かの判断として、音量を検出する。

40

【0067】

尚、音量が所定の閾値以上になると発声の開始を検出し、音量が所定の閾値未満になると発声の終了を検出する。

【0068】

(図4中のパラメータの説明)

図4において、401は音声信号420から周知の方法で求めた音量($E(t)$)、402が発声開始を検出するための閾値($TH1$)、403が発声終了を検出するための閾値($TH2$)である。

【0069】

尚、 $E(t)$ は時刻 t を始点とするフレームにおける音量を表す。

50

【 0 0 7 0 】

即ち、第一状態 3 0 1 で $E(t) \geq TH1$ となると発声開始を検出し、第三状態 3 0 3 で $E(t) < TH2$ となると発声終了を検出する。

また、発声開始の検出と発声終了の検出に同じ閾値を用いてもよい ($TH1 = TH2$)。また、発声開始の検出条件 ($E(t) \geq TH1$) となるフレームが所定数検出された場合に発声開始を確定する。

【 0 0 7 1 】

同様に、発声終了の検出条件 ($E(t) < TH2$) となるフレームが所定数検出された場合に発声終了を確定する。

【 0 0 7 2 】

本実施形態では、発声開始、発声終了を確定するまでのフレーム数をそれぞれ $D1$ (例えば、4 フレーム)、 $D2$ (例えば、6 フレーム) とする。

【 0 0 7 3 】

したがって、第二状態 3 0 2 に遷移してから $E(t) \geq TH1$ となるフレームが $D1$ 回検出された場合、発声開始を確定して第三状態 3 0 3 に遷移する。

また、第二状態 3 0 2 に遷移してからのフレームが $D1$ 回検出される前に、音量が $E(t) < TH1$ となった場合、第一状態 3 0 1 に遷移する。

【 0 0 7 4 】

尚、第二状態 3 0 2 から第一状態 3 0 1 に遷移する処理は発声開始を取り消す処理に相当する。

【 0 0 7 5 】

同様に、第四状態 3 0 4 に遷移してから $E(t) < TH2$ となるフレームが $D2$ 回検出された場合、発声終了を確定し音声検出を終了する。

【 0 0 7 6 】

また、第四状態 3 0 4 に遷移してからのフレームが $D2$ 回検出される前に、音量が $E(t) \geq TH2$ となった場合、第三状態 3 0 3 に遷移する。

尚、第四状態 3 0 4 から第三状態 3 0 3 に遷移する処理は発声終了を取り消す処理に相当する。

【 0 0 7 7 】

尚、発声の開始を確定するまでに必要なフレーム数 $D1$ は発声の終了を確定するまでに必要なフレーム数 $D2$ よりも小さき場合が一般的であるが、同じ数 ($D1 = D2$) であってもよい。

【 0 0 7 8 】

4 3 0 は音声信号 4 2 0 に対する音声検出部 1 0 6 が判定した認識状態の様子を示している。

音声入力開始後は第一状態 3 0 1 である。

【 0 0 7 9 】

音量 4 0 1 が閾値 $TH1$ 以上となる時点 $t1$ を始点とするフレームで発声開始を検出 (4 0 4) して第二状態 3 0 2 に遷移する。

続いて、時点 $t2$ を始点とするフレームでは、第二状態 3 0 2 に遷移してからのフレーム数が $D1$ 回となる前に音量 4 0 1 が閾値 $TH1$ 未満となるので発声開始を取り消し (4 0 5)、第一状態 3 0 1 に遷移する。

続いて、時点 $t3$ を始点とするフレームでは、再び音量 4 0 1 が閾値 $TH1$ 以上になるので発声開始を検出 (4 0 6) して、第二状態 3 0 2 に遷移する。

【 0 0 8 0 】

第二状態 3 0 2 に遷移してから音量 4 0 1 が $TH1$ 以上となるフレーム数が $D1$ 回となる時点 $t4$ で発声開始を確定 (4 0 7) して第三状態 3 0 3 に遷移する。

【 0 0 8 1 】

第三状態 3 0 3 において、音量 4 0 2 が発声終了を検出するための閾値 $TH2$ 未満となる時点 $t5$ を始点とするフレームで発声終了を検出 (4 0 8) して第四状態 3 0 4 に遷移す

10

20

30

40

50

る。

続く時点 t_6 を始点とするフレームにおいて音量 401 が閾値 TH_2 以上となるので発声終了を取り消し (409)、第三状態 303 に遷移する。

続く時点 t_7 を始点とするフレームで再び音量 401 が閾値 TH_2 未満となるので発声終了を検出 (410) して第四状態 304 に遷移する。

【0082】

以後、第四状態 304 に遷移してから音量 401 が閾値 TH_2 未満となるフレームの数が D_2 回となる時点 t_8 で発声終了を確定する (411)。

【0083】

また、発声開始の確定、および発声終了の確定は、フレームの数ではなく、音量が閾値以上あるいは閾値未満である状態が一定の時間継続するか否かで判断してもよい。

10

【0084】

すなわち、閾値 (TH_1) 以上の音量が、発声開始を確定するまでのフレーム数 D_1 (例えば、4 フレーム) に相当する時間 S_1 (102 . 4 ミリ秒間) 検出された場合、発声の開始を確定する。

【0085】

同様に、閾値 (TH_2) 以下の音量が、発声終了を確定するまでのフレーム数 D_2 (例えば、6 フレーム) に相当する時間 S_1 (153 . 6 ミリ秒間) 検出された場合、発声の終了を確定する。

【0086】

20

尚、所定の音量が断続して検知された場合でも、継続とみなして継続時間を判断してもよい。

【0087】

このような構成とすることで、音声検出部 106 は、検出すべき音が一瞬途切れ、対応するフレームの音量が取得できない場合があっても、音が一瞬途切れた後すぐに発声される場合には、適切な処理を実行することが可能となる。

【0088】

(音声検出部 106 による処理動作を示すフローチャート)

図5は、音声検出部 106 による処理動作を示すフローチャートである。

ステップ $S501$ で、発声開始を検出した時にフレーム番号を初期化する。

30

【0089】

以下、フレーム単位で音声の検出を行う。

即ち、音声検出部 106 はフレーム単位に処理を行う際に、当該フレーム毎に音量を計算する。

尚、音量は例えば、対数パワー等信号強度に係る値を周知の方法で音声信号から算出する。

【0090】

尚、短時間の対数パワーは例えば次式で算出する。

$$E(t) = \log \{ (x(t, i)^2) / N \} \quad (1 \leq i \leq N) \cdots (\text{数} 1)$$

ここで、 N はフレームあたりの音声信号のサンプル数、 i はフレーム内の音声サンプルのインデックスである。

40

【0091】

また、 $x(t, i)$ は時点 t を始点とするフレーム内の i 番目サンプルの音声信号を表している。

また、 $x(t, i)^2$ は $x(t, i)$ の2乗を意味する。

次に、ステップ $S502$ で、第一状態 301 における処理を開始する。

【0092】

次に、ステップ $S503$ で、時点 t を始点とするフレームにおける音量 $E(t)$ が発声の開始を検出するために用いる閾値 TH_1 以上であるか判断する。

音量 $E(t)$ が TH_1 以上の場合 (ステップ $S503$ において YES)、ステップ $S50$

50

5で第二状態302に遷移する。

音量E(t)がTH1未満の場合(ステップS503においてNO)、次のフレームの処理(ステップS504)を繰り返す。

【0093】

次に、ステップS506で、第二状態302に遷移したフレームを発声開始フレームTsと設定する。

【0094】

次に、ステップS507で、音量E(t)がTH1未満であるか判断する。

音量E(t)がTH1未満の場合(ステップS507においてYES)、第一状態301に遷移する。

10

音量E(t)がTH1以上の場合(ステップS507においてNO)、ステップS508で、第二状態302に遷移してからのフレーム数がD1回未満であるか判断する。

【0095】

第二状態302に遷移してからのフレーム数がD1回未満である場合(ステップS508においてYES)、次のフレームの処理(ステップS509)を繰り返す。

【0096】

第二状態302に遷移してからのフレーム数がD1回以上である場合(ステップS508においてNO)、ステップS510で、第三状態303に遷移する。

【0097】

次に、ステップS512で、音量E(t)が発声終了の検出に用いる閾値TH2未満であるか判断する。

20

音量E(t)がTH2未満の場合(ステップS512においてYES)、ステップS514で第四状態304に遷移する。

音量E(t)がTH2以上の場合(ステップS512においてNO)、ステップS513で次のフレームの処理を行う。

【0098】

次に、ステップS515で、第四状態304に遷移したフレームを発声終了フレームTeと設定する。

【0099】

次に、ステップS516で、音量E(t)がTH2以上であるか判断する。

30

音量E(t)が閾値TH2以上の場合(ステップS516においてYES)、第三状態303に遷移する。

音量E(t)がTH2未満の場合(ステップS516においてNO)、ステップS517で第四状態304に遷移してからのフレーム数がD2回未満であるか判断する。

【0100】

第四状態304に遷移してからのフレーム数がD2回未満である場合(ステップS517においてYES)、ステップS518で次のフレームの処理を行う。

【0101】

第四状態304に遷移してからのフレーム数がD2回以上である場合(ステップS517においてNO)、ステップS519で音声の検出を終了するか判断する。

40

音声の検出を終了する場合(ステップS519においてYES)、ステップS520で音声の検出を終了する。

音声検出を終了しない場合(ステップS519においてNO)、次の発声の検出に備える場合は第一状態301に遷移する。

【0102】

以上の処理によって音声検出部106はフレームTsからフレームTeまでを発声区間として検出する。

【0103】

音声認識部107は音声検出部106が検出した発声区間(フレームTsからフレームTeまで)の音声信号を処理して音声認識結果を求める。

50

【 0 1 0 4 】

尚、フローチャートと用いた上記の説明では音量の変化に基づいて音声区間を検出する場合を説明したが、これに限定しなくてもよい。

【 0 1 0 5 】

また、音声検出を行う場合、零交差回数、ピッチ、音声モデルと非音声モデルが出力する尤度比等の周知の特徴量やこれらを組み合わせた特徴量を用いてもよい。

【 0 1 0 6 】

このような特徴量を用いることで、例えば、周囲から入力される音が大きような環境下においても発声開始および発声終了を効率良く検出することが可能となる。

【 0 1 0 7 】

尚、発声開始および発声終了を確定する条件は、以下に示すようにフレーム数以外の条件を用いてもよい。

【 0 1 0 8 】

例えば、発声開始を検出するための閾値 T H 1 よりさらに大きな音量である所定の閾値 T H 3 を設け、発声開始を検出後、音量が所定の閾値 T H 3 に達したフレームで発声開始を確定してもよい。

【 0 1 0 9 】

また、発声終了の確定に対しては、発声終了を検出する閾値 T H 2 よりも小さな音量である所定の閾値 T H 4 を設けて、発声終了を検出後、音量が所定の閾値 T H 4 よりも小さくなったフレームで発声終了を確定してもよい。

【 0 1 1 0 】

このような条件用いて判定することで、発声開始および発声終了を確定するまでの時間を短縮させることが可能となる。

【 0 1 1 1 】

次に、以上の構成を備えたデジタルカメラ 2 0 0 において、音声の指示によって撮像動作を実行する場合について説明する。

【 0 1 1 2 】

(音声検出部 1 0 6 、撮像制御部 1 2 3 、画像記憶制御部 1 0 4 の処理の対応関係)

図 1 7 は、音声検出部 1 0 6 、撮像制御部 1 2 3 、画像記憶制御部 1 0 4 の処理の一例を示す図である。

尚、図 3 と共通の要素には同一符号を付し、その説明を省略する。

【 0 1 1 3 】

図 1 7 において、発声開始を検出すると (3 0 5) 、撮像制御部 1 2 3 は撮像部 1 0 3 に撮像動作を実行させる。

尚、発声開始を検出した場合 (3 0 5) とは、図 5 のステップ S 5 0 3 の Y E S と判断された場合に相当する。

【 0 1 1 4 】

また、発声終了を検出すると (3 0 8) 、撮像制御部 1 2 3 は撮像部 1 0 3 に撮像動作を実行させる。

尚、発声終了を検出した場合 (3 0 8) とは、図 5 のステップ S 5 1 2 の Y E S と判断された場合に相当する。

【 0 1 1 5 】

即ち、撮像部 1 0 3 は、音声検出処理の内部状態が第一状態 3 0 1 から第二状態 3 0 2 に遷移する時点および第三状態 3 0 3 から第四状態 3 0 4 に遷移する時点で撮像する。

【 0 1 1 6 】

また、画像記憶制御部 1 0 4 は、一旦撮像した画像を、発声開始を取り消した場合 (3 0 6) 、発声終了を取り消した場合 (3 0 9) で消去する。

【 0 1 1 7 】

尚、発声開始を取り消した場合 (3 0 6) とは、図 5 のステップ S 5 0 7 の Y E S と判断された場合である。

10

20

30

40

50

【 0 1 1 8 】

また、発声終了を取り消した場合（ 3 0 9 ）とは、図 5 のステップ S 5 1 6 の Y E S と判断された場合である。

【 0 1 1 9 】

即ち、図 1 7 において、発声開始を取り消すと、画像制御部 1 0 4 は発声開始を検出した場合（ 3 0 5 ）に撮像した画像を消去する。

【 0 1 2 0 】

同様に、発声終了を取り消すと、画像制御部 1 0 4 は発声終了を検出した場合（ 3 0 8 ）に撮像した画像を消去する。

【 0 1 2 1 】

即ち、第二状態 3 0 2 から第一状態 3 0 1 に遷移する時、第四状態 3 0 4 から第三状態 3 0 3 に遷移する時に直前の遷移で撮像した画像を消去する。

【 0 1 2 2 】

（音声認識文法）

図 9 は、本実施形態で利用する音声認識文法の一例を示す図である。

【 0 1 2 3 】

この例では音声認識文法 9 0 0 はルールを記述する部分 9 1 0 と認識するコマンドおよび発音を記述する部分 9 2 0 で構成される。

【 0 1 2 4 】

コマンドおよび発音を記述する部分 9 2 0 には一行毎に単語の I D 9 2 1、コマンド 9 2 2、発音 9 2 3 が記述されている。

【 0 1 2 5 】

尚、ルールを記述する部分 9 1 0 には、9 2 2 に記載された計 9 語を認識するための方法が音声認識部 1 0 7 に読み取り可能なプログラムコードの形態で記述されている。

【 0 1 2 6 】

“ S h o o t ”、“ G o ”、“ チーズ (C h e e s e) ”、“ はい、ちーず (S a y C h e e s e) ”、“ F i v e F o u r T h r e e ” は後述する撮像を音声で指示するためのコマンドである。

【 0 1 2 7 】

“ S p o t M e t e r i n g ”（スポット測光）、“ C e n t e r M e t e r i n g ”（中央部重点測光）、“ U s e a f l a s h ”（ストロボ発光）、“ N o F l a s h ”（ストロボ発光禁止）は音声で撮影条件を設定するためのコマンドである。

【 0 1 2 8 】

以下の説明において、本実施形態のデジタルカメラ 2 0 0 では、図 9 に示した音声認識文法 9 0 0 を言語モデルとして用いる。

【 0 1 2 9 】

なお本実施例においては、音声を好適な例として説明するが、本発明はこれに限らない。例えば、各音声コマンドの代わりに、何かしらの意味に置換できる音を適用することでもできるであろう。例えば、笑い声、列車が通過する際に発生する音等を適用可能である。なお、この場合には、音声認識の技術ではなく、公知の音の種類の検知技術を代用することになるであろう。

【 0 1 3 0 】

このような構成とすることで、音声に限らず、特徴のある音がマイク 1 1 2 を介して入力された場合にも、ユーザは特徴のある各種の音に応答したタイミングで撮像された画像を得ることが可能となる。

【 0 1 3 1 】

（認識結果制御テーブル）

認識結果制御テーブルは、認識結果に対応する撮像、測光、ストロボ発光等の処理を記述したテーブル形式のデータであり、認識結果処理部 1 0 8 が認識結果に対応するカメラ制御を決定する際に参照する。

10

20

30

40

50

【 0 1 3 2 】

尚、認識結果制御テーブルは、認識結果処理部 1 0 8 が読み取り可能なプログラムコードの形態でメモリ（認識結果制御テーブル用）1 1 4 に格納されている。

【 0 1 3 3 】

図 1 0 は、認識結果制御テーブルの一例を示す図である。

図 1 0 において、1 0 0 0 は認識結果処理データである。

9 1 0 に認識に利用するコマンドが記述されており、1 0 0 2 に 9 1 0 のコマンドに対応するデジタルカメラ 2 0 0 の制御内容が記述されている。

【 0 1 3 4 】

（音声によって撮像を指示する場合のデジタルカメラ 2 0 0 における処理の一例を示すフローチャート） 10

図 6 ~ 図 8 は、音声によって撮像を指示する場合のデジタルカメラ 2 0 0 における処理の一例を示すフローチャートである。

【 0 1 3 5 】

まず図 6 のフローチャートを参照して説明する。

【 0 1 3 6 】

ステップ S 6 0 1 で、音声シャッター機能がオンに設定されているか否か判断する。

【 0 1 3 7 】

音声シャッター機能がオンに設定されている場合（ステップ S 6 0 1 において Y E S ）、ステップ S 6 0 2 で録音ボタン 2 0 7 を押下されて音声入力を開始する操作が行われたか否か判断する。 20

【 0 1 3 8 】

音声シャッター機能がオフに設定されている場合（ステップ S 6 0 1 において N O ）、ステップ S 6 9 9 で音声シャッター機能以外の処理を行う。

【 0 1 3 9 】

尚、ユーザは操作部 1 0 2 が備える音声シャッター切替えスイッチ 2 0 2 を操作して、音声シャッター機能のオン・オフを切り替える。

【 0 1 4 0 】

また、音声シャッター機能のオン・オフの判断は制御部 1 0 1 が行う。

【 0 1 4 1 】

音声入力を開始する操作が行われた場合（ステップ S 6 0 2 の Y E S ）、ステップ S 6 0 3 で音声入力部 1 0 5 は音声入力の処理を開始し、音声検出部 1 0 6 は音声を検出する処理を開始する。 30

【 0 1 4 2 】

音声入力を開始する操作以外の操作が行われた場合（ステップ S 6 0 2 の N O ）、ステップ S 6 9 9 で音声シャッター機能以外の処理を行う。

【 0 1 4 3 】

尚、音声入力を開始する操作は、録音ボタン 2 0 7 を押下する操作以外の操作でもよい。 40

【 0 1 4 4 】

例えば、オートフォーカス機能を備えたデジタルカメラでは、シャッターボタン 2 0 1 を半押しすると、焦点を合わせる動作をするものがある。

【 0 1 4 5 】

この時、オートフォーカス機能の動作に連動して音声入力の処理を開始するようにしてもよい。即ち、ユーザがシャッターボタン 2 0 1 を半押しすると音声入力および音声検出の処理を開始する。

【 0 1 4 6 】

このような構成とすることで、手動作による操作が簡略化される。従って、ユーザは音声入力の処理を素早く開始することができる。

【 0 1 4 7 】

また、手動作によって音声検出の開始することなく、音声入力部 105 に音声信号が入力された時点で音声検出を開始するようにしてもよい。

【0148】

このような構成とすることで、音声検出の処理を素早く開始することができる。また、ユーザがカメラを手動作で操作できない場合にも音声検出を開始することができる。従って監視カメラ、防犯カメラ、高所に据え置きされたカメラ等に利用することができる。

【0149】

ステップ S604 で、音声検出部 106 が発声開始を検出したか否か判断する。

【0150】

尚、ステップ S604 において発声開始を検出したか否かの判断は、音声検出部 106 が第一状態 301 から第二状態 302 に遷移させる処理を実行したか否かという判断に基づく。

【0151】

発声開始を検出した場合（ステップ S604 において YES）、ステップ S605 で撮像部 103 が撮像動作を実行する。

【0152】

ステップ S606 で、画像記憶制御部 104 は直前のステップ S605 で撮像された画像の第 1 の画像データをメモリ 110 に記憶させる。

【0153】

尚、ステップ S605 で撮像した画像、即ち音声検出部 106 が発声開始を検出した時点で撮像した画像を画像 A とする。

【0154】

発声開始を検出しなかった場合（ステップ S604 において NO）、発声開始の検出を繰り返す。

【0155】

ステップ S607 で、音声検出部 106 が発声開始を取り消すか否か判断する。

【0156】

尚、ステップ S605 において発声開始を取り消すか否かの判断は、音声検出部 106 が第二状態 302 から第一状態 301 に遷移させる処理を実行したか否かという判断に基づく。

発声開始を取り消す場合（ステップ S607 において YES）、ステップ S608 で画像記憶制御部 104 はメモリ 110 に記憶した画像 A を消去する。

発声開始を取り消さない場合（ステップ S607 において NO）、ステップ S609 で音声検出部 106 は発声開始を確定したか否か判断する。

【0157】

尚、ステップ S609 において発声開始を確定したか否かの判断は、音声検出部 106 が第二状態 302 から第三状態 303 に遷移させる処理を実行したか否かという判断に基づく。

発声開始を確定した場合（ステップ S609 の YES）、ステップ S610 で音声認識部 107 が音声認識処理を開始する。

発声開始を確定しなかった場合（ステップ S609 の NO）、発声開始を取り消すか否かの判定を繰り返す。

【0158】

以降の処理は図 7 のフローチャートを参照して説明する。

ステップ S711 で、音声検出部 106 は発声終了を検出したか否か判断する。

【0159】

尚、ステップ S711 において発声終了を検出したか否かの判断は、音声検出部 106 が第三状態 303 から第四状態 304 に遷移させる処理を実行したか否かという判断に基づく。

【0160】

発声終了を検出した場合（ステップS 7 1 1においてYES）、ステップS 7 1 2で撮像部1 0 3が撮像を行う。

【0 1 6 1】

次に、ステップS 7 1 3で、画像記憶制御部1 0 4が直前のステップS 7 1 2で撮像した画像の第2の画像データをメモリ1 1 0に記憶する。

尚、ステップS 7 1 2で撮像した画像、即ち音声検出部1 0 6が発声終了を検出した時点で撮像した画像を画像Bとする。

尚、一般に「はい、チーズ（Say Cheese）」等の掛け声をかけてに、その発声が終了した後（「ず」を発声した後）に一拍（例えば、0.5秒間）遅延して撮像する場合がある。

10

【0 1 6 2】

これを考慮して、本実施例では、音声検出部1 0 6が「はい、チーズ（Say Cheese）」の発声終了を検出した時点から一定の遅延時間が経過してから、撮像部1 0 3が撮像を行う。

【0 1 6 3】

このような構成とすることで、ユーザが希望する撮像タイミングの種類を増やすことができる。

【0 1 6 4】

次に、ステップS 7 1 5で、音声検出部1 0 6は発声の終了を取り消すか否か判断する。

20

【0 1 6 5】

尚、ステップS 7 1 5において発声終了を取り消すか否かの判断は、音声検出部1 0 6の第四状態3 0 4から第三状態3 0 3に遷移したと認識したか否かという判断に基づく。

【0 1 6 6】

発声終了を取り消す場合（ステップS 7 1 5においてYES）、ステップS 7 1 4で画像記憶制御部1 0 4はメモリ1 1 0に記憶された画像Bを消去する。

【0 1 6 7】

次に、ステップS 7 1 6で、音声検出部1 0 6が発声終了を確定するか否か判断する。

【0 1 6 8】

尚、ステップS 7 1 6において発声終了を確定したか否かの判断は、音声検出部1 0 6の第四状態3 0 4から状態遷移を終了したか否かという判断に基づく。

30

【0 1 6 9】

発声終了を確定した場合（ステップS 7 1 6においてYES）、ステップS 7 1 7で音声入力部1 0 5および音声検出部1 0 6による処理を終了する。

【0 1 7 0】

次に、音声検出終了後のステップS 7 1 8で、音声認識部1 0 7は音声検出部1 0 6が検出した発声区間の音声信号を全て処理するまで音声認識処理を行う。

【0 1 7 1】

音声認識の処理が終了した場合（ステップS 7 1 8においてYES）、ステップS 7 1 9で認識結果処理部1 0 8は音声認識部1 0 7が求めた認識結果を取得する。

40

【0 1 7 2】

以降の処理は図8のフローチャートを参照して説明する。

【0 1 7 3】

ステップS 8 2 1で、認識結果処理部1 0 8は取得した認識結果の認識スコアに基づいて対応するコマンドを受理するか棄却するか判断する。

【0 1 7 4】

尚、コマンドを受理するとは、制御部1 0 1が認識されたコマンドに対応する制御を決定することをいう。また、コマンドを棄却するとは、制御部1 0 1が認識されたコマンドに対応する制御が決定されないことをいう。

【0 1 7 5】

50

取得した認識スコアが所定の閾値以上であり、対応するコマンドを受理した場合（ステップS 8 2 1においてYES）、ステップS 8 2 2で認識結果制御テーブルを参照して認識結果に含まれるコマンドに対応するカメラの制御を決定する。

【0176】

認識されたコマンドが発声開始時点で撮像を指示する語（“Shoot”または“Go”）の場合（ステップS 8 2 2においてYES）、ステップS 8 2 3で画像記憶制御部104がメモリ110に記憶されている画像Aの画像データを記憶媒体111に保存する。

【0177】

尚、ステップS 8 2 3の処理は、認識結果処理部108の決定にしたがった処理である。

10

【0178】

次に、ステップS 8 2 4で、撮像した画像をユーザが確認できるように表示制御部109は画像Aをディスプレイ115に表示する。

認識されたコマンドが発声開始時点で撮像を指示する語（“Shoot”または“Go”）でない場合（ステップS 8 2 2においてNO）、ステップS 8 2 6で、発声終了時点で撮像を指示する語（“チーズ（Cheese）”）であるか否か判断する。

認識されたコマンドが発声終了時点で撮像を指示する語（“Cheese”）の場合（ステップS 8 2 6においてYES）、ステップS 8 2 7で画像記憶制御部104が画像Bの画像データを記憶媒体（画像記憶用）111に保存する。

【0179】

尚、ステップS 8 2 7の処理は、認識結果処理部108の決定にしたがった処理である。

20

【0180】

ステップS 8 2 8で、撮像した画像をユーザが確認できるように表示制御部109は画像Bをディスプレイ115に表示する。

【0181】

認識されたコマンドが撮像を指示する語以外の語（“Spot Metering”等）の場合（ステップS 8 2 6においてNO）、ステップS 8 2 9で、認識結果処理部108が認識結果制御テーブル114を参照して、撮像以外のカメラの制御を行う。

【0182】

ステップS 8 2 5で画像記憶制御部104がメモリ110に記憶している全ての画像（画像Aおよび画像B）の画像データを消去する。

30

【0183】

即ち、所定のコマンドが認識されず、認識結果が棄却されると撮像部103が撮像した画像は消去される。

【0184】

この処理により、周囲の雑音や認識対象語以外の発声、ユーザ以外の話し声等カメラ操作を意図しない音声の認識結果を棄却し、これらの音を誤って検出して撮像した画像を自動的に消去する。

【0185】

尚、判定に用いる閾値は、あらかじめ決めた固定値でもよいし、ガーベッジモデルが出力する認識スコアを r 倍（ $0 < r$ ）した値を用いてもよい。

40

【0186】

ガーベッジモデルとは、音声以外の雑音区間や想定される複数の未知語を用いて作成した音響モデルであり、音声認識用データ113に含まれる。

【0187】

尚、ステップS 8 2 2～ステップS 8 2 9の処理は、発声開始時点で撮像した画像と発声終了時点で撮像した画像とから認識結果に従って保存する画像を決定している。

【0188】

したがって、ユーザは発声内容によって保存する画像の撮像タイミングを自由に変える

50

ことができる。

【 0 1 8 9 】

尚、上記の説明では、ステップ S 8 2 5 の後に処理を終了するように説明したが、引き続き次の音声入力を行うため、ステップ S 6 0 2 の処理へ進んでも良い。

【 0 1 9 0 】

このように構成し、シャッターボタン 2 0 1 を半押しすることで音声入力開始を操作する場合は、シャッターボタン 2 0 1 を半押ししている間は何度でも音声入力によるカメラ制御が可能になる。

【 0 1 9 1 】

例えば、シャッターボタン 2 0 1 を半押ししたまま、“ C e n t e r M e t e r i n g ”等の発声で撮影条件を設定し、つづく発声で撮像指示を出す、といったことが可能になる。

10

【 0 1 9 2 】

(“ S h o o t ” で撮像する場合の説明)

図 1 1 は本実施形態に係るデジタルカメラ 2 0 0 を利用して、“ S h o o t ”という音声指示で撮像する場合の動作示す図である。

【 0 1 9 3 】

図 1 1 の横軸 1 1 5 0 は時間軸であり、左から右に時刻が推移する。 t 1 ~ t 7 は時点を示している。

1 1 1 0 は音声入力部 1 0 5 が A / D 変換した音声信号である。

20

1 1 1 1 はユーザが“ S h o o t ”と発声した区間の音声信号（音声波形）である。

1 1 2 0 は音声信号 1 1 1 0 に対応する音量の変化を示す。

1 1 2 1 は音声検出部 1 0 6 で用いる発声開始検出用の閾値（ T H 1 ）、1 1 2 2 は発声終了検出用の閾値（ T H 2 ）である。

1 1 3 0 は、音声検出部 1 0 6 が認識した状態の変化を視覚的に示したものである。

1 1 4 0 はデジタルカメラ 2 0 0 の動作の内容を示している。

【 0 1 9 4 】

続いて、時点 t 1 から時点 t 7 までの時間経過に沿ってデジタルカメラ 2 0 0 の動作を説明する。

【 0 1 9 5 】

30

(時点 t 1)

音量 1 1 2 0 が閾値 T H 1 以上になる時点 t 1 を始点とするフレームで音声検出部 1 0 6 が発声開始を検出する。これは、上述した所定の基準（開始条件）を満たした音を検出する工程に相当する。

この時、音声検出部 1 0 6 は第一状態 3 0 1 から第二状態 3 0 2 に遷移させる処理を実行する（ 1 1 3 0 の時点 t 1 の部分）。

発声開始を検出した時点で撮像部 1 0 3 が同時点の被写体（ I M G 0 0 3 ）を撮像し、続いて、画像記憶制御部 1 0 4 が撮像した画像の画像データをメモリ（画像記憶用） 1 1 0 に記憶する（以上、 1 1 4 1 ）。

【 0 1 9 6 】

40

(時点 t 2)

音声検出部 1 0 6 は、発声開始を検出した時点 t 1 を始点とするフレームから D 1 番目のフレームである時点 t 2 を始点とするフレームで発声開始を確定する。

同時に、音声認識部 1 0 7 による音声認識の処理を開始する（以上 1 1 4 2 ）。

この時、音声検出部 1 0 6 は、第二状態 3 0 2 から第三状態 3 0 3 に遷移させる処理を実行する（ 1 1 3 0 の時点 t 2 の部分）。

【 0 1 9 7 】

(時点 t 3)

続いて、音量 1 1 2 0 が閾値 T H 2 未満になる時点 t 3 を始点とするフレームで音声検出部 1 0 6 が発声終了を検出する。これは、上述した所定の基準（終了条件）を満たした

50

音を検出する。

この時、音声検出部 106 は、第三状態 303 から第四状態 304 へ遷移させる処理を実行する(1130の時点 t3の部分)。

音声検出部 106 が発声終了を検出した時点 t3 で、撮像部 103 がこの時点の被写体(IMG005)を撮像し、続いて画像記憶制御部 104 が撮像した画像を撮像した画像の画像データをメモリ(画像記憶用)110に記憶する(以上、1143)。

【0198】

(時点 t4)

音声検出部 106 が発声終了を検出した時点 t3 を始点とするフレームから D2 番目のフレームとなる前の、時点 t4 を始点とするフレームで音量 1120 が閾値 TH2 以上になると、音声検出部 106 は発声終了を取り消す。

この時、音声検出部 106 は第四状態 304 から第三状態 303 に遷移させる処理を実行する(1130の時点 t4の部分)。

発声終了が取り消された時点 t4 で、画像記憶制御部 104 は発声終了を検出した時点 t3 で撮像した画像 IMG005 の画像データをメモリ(画像記憶用)110から消去する(以上、1144)。

【0199】

(時点 t5)

続く時点 t5 を始点とするフレームで音量 1120 が閾値 TH2 未満になるので、音声検出部 106 が発声終了を検出する。

この時、音声検出部 106 は第三状態 303 から第四状態 304 に遷移させる処理を実行する(1130の時点 t5の部分)。

また、撮像部 103 はこの時点 t5 の被写体(IMG006)を撮像し、画像記憶制御部 104 が撮像した画像の画像データをメモリ(画像記憶用)110に記憶する(以上1145)。

【0200】

(時点 t6)

発声終了を検出した時点 t5 を始点とするフレームから、音量 1120 が閾値 TH2 以上になることなく D2 番目のフレームである時点 t6 を始点とするフレームで、音声検出部 106 は発声の終了を確定する(1146)。

この時、前述したように、音声検出部 106 は第四状態 304 から第1状態に遷移させる処理を実行してもよく、音声検出部 106 は状態を遷移させる処理を終了させてもよい。

【0201】

(時点 t7)

その後、音声認識部 107 の処理が終了した時点 t7 で認識結果処理部 108 がデジタルカメラ 200 の制御方法を決定する。

ここで“Shoot”が認識結果として得られた場合は、認識結果制御テーブルを参照して“Shoot”に対応する処理を決定する。

図10に示したように“Shoot”は発声の開始を検出した時点での撮像動作と対応付けられたコマンドである。

認識結果処理部 108 の決定にしたがって、画像記憶制御部 104 が発声開始を検出した時点 t1 で撮像した画像(IMG003)の画像データを画像記憶媒体 111 に保存する。

同時に、画像記憶制御部 104 は発声終了時点で撮像した画像(IMG006)を保存せずにメモリ(画像記憶用)110から消去する。

【0202】

(“チーズ(Cheese)”で撮像する場合の説明)

図12は、本実施形態に係るデジタルカメラ 200 を利用して、“チーズ(Cheese)”という音声指示で撮像する場合の動作を示す図である。

【0203】

図 1 1 と同様に 1 2 5 0 は時間軸であり、1 2 1 0 は音声信号、1 2 2 0 は音量、1 2 3 0 は音声検出部 1 5 0 が認識した状態、1 2 4 0 はデジタルカメラ 2 0 0 の動作を示す。

1 2 1 1 はユーザの発声前に混入した雑音を検知した区間であり、1 2 1 2 はユーザが発した “ C h e e s e ” という音声を検知した区間である。

1 2 2 1 は音声検出部 1 0 6 で用いる発声区間を検出するための閾値 (T H 1) である。

【 0 2 0 4 】

尚、図 1 2 では、発声開始、発声終了の検出に同じ閾値 T H 1 を用いる。

【 0 2 0 5 】

以下、時間経過に沿ってデジタルカメラ 2 0 0 の動作を説明する。

10

【 0 2 0 6 】

(時点 t 1)

時点 t 1 を始点とするフレームで音声検出部 1 0 6 が発声開始を検出すると、撮像部 1 0 3 が時点 t 1 を始点とするフレームに対応する被写体 1 2 0 2 (I M G 0 0 2) を撮像する。また、画像記憶制御部 1 0 4 が撮像した画像の画像データをメモリ (画像記憶用) 1 1 0 に一時的に記憶する (1 2 4 1) 。

【 0 2 0 7 】

(時点 t 2)

時点 t 2 を始点とするフレームで、発声開始を検出してからのフレーム数が D 1 回となる前に音量が閾値 T H 1 未満になるため、音声検出部 1 0 6 が発声開始を取り消す。この時、画像記憶制御部 1 0 4 が 1 2 4 1 で撮像した I M G 0 0 2 を消去する。

20

【 0 2 0 8 】

(時点 t 3)

時点 t 3 を始点とするフレームで音声検出部 1 0 6 が再び発声開始を検出すると、撮像部 1 0 3 が時点 t 3 を始点とするフレームに対応する被写体 1 2 0 3 (I M G 0 0 3) を撮像する。また、画像記憶制御部 1 0 4 が撮像した画像の画像データをメモリ (画像記憶用) 1 1 0 に一時的に記憶する (1 2 4 3) 。

【 0 2 0 9 】

(時点 t 4)

時点 t 4 を始点とするフレームで音声検出部 1 0 6 が発声開始を確定すると、音声認識部 1 0 7 が音声認識の処理を開始する (1 2 4 4) 。

30

【 0 2 1 0 】

(時点 t 5)

時点 t 5 を始点とするフレームで音声検出部 1 5 0 が発声終了を検出すると、撮像部 1 0 3 が時点 t 5 始点とするフレームに対応する被写体の画像 1 2 0 5 (I M G 0 0 5) を撮像する。また、続いて画像記憶制御部 1 0 4 が撮像した画像の画像データをメモリ (画像記憶用) 1 1 0 に一時的に記憶する (1 2 4 5) 。

【 0 2 1 1 】

(時点 t 6)

時点 t 6 を始点とするフレームで音声検出部 1 0 6 が発声終了を確定する (1 2 4 6)

40

【 0 2 1 2 】

(時点 t 7)

発声終了確定後、音声認識部 1 0 7 による音声認識処理が終了する時点 t 7 で、認識結果処理部 1 0 8 が得られた認識結果に基づいてカメラの制御を決定する。

尚、図 1 0 に示したように “ C h e e s e ” は発声の終了を検知した時点での撮像動作と対応付けられたコマンドである。

したがって、画像記憶制御部 1 0 4 は、発声終了を検出した時点 t 5 で撮像した画像 (I M G 0 0 5) の画像データを画像記憶媒体 1 1 1 に保存し、発声開始を検出した時点 t 3 で撮像した画像 (I M G 0 0 3) の画像データは保存せずに消去する。

50

【0213】

以上、図11、図12を用いて説明したように、本実施形態で説明したデジタルカメラ200では、発声を開始したタイミングの画像を得たい場合は“Shoot”（または“Go”）と発声すればよい。

【0214】

また、本実施形態で説明したデジタルカメラ200では、発声を完了したタイミングの画像を得たい場合は“Cheese”と発声すればよい。

【0215】

また、発声を開始したタイミングから一定時間後（すなわち、“Two One Zero”が発声されるべき時間後）の時間関係にあるタイミングで撮像された画像の画像データを得心たい場合は“Five Four Three”と発声すればよい。

10

【0216】

また、発声を終了したタイミングから一定時間後（例えば、0.5秒後）の画像の画像データを得心たい場合は“はい、チーズ（Say Cheese）”と発声すればよい。

【0217】

“Shoot”（または“Go”）と発声した場合、音声認識の終了を待たずに撮像するため、乗り物等の動いている被写体を撮影する場合に好適である。

また、“Cheese”（または、“Say Cheese”）と発声した場合、発声後に撮影するため、集合写真や記念写真等、被写体に撮影タイミングを伝えて撮影する場合に好適である。

20

また“Five Four Three”と発声した場合、発声を開始してから一定時間後（すなわち、“Two One Zero”が発声されるべき時間後）の所望のタイミングで撮像された画像を得ることができる。

【0218】

従って、撮影シーンによって自由に撮影タイミングを変えた撮影が可能になり、ユーザの利便性が向上する。

【0219】

また、撮影後にユーザが意図しないタイミングで撮影した画像を手動作によって削除する手間も必要ないという利点がある。

【0220】

30

即ち、図12で説明したように、音声入力時に混入した周囲の雑音を発声と誤って撮像した場合でも、発声開始が確定されなければ自動的に消去する。

【0221】

また、雑音や撮像を意図しない発声により撮像した場合でも、図8のS821における処理で、撮像を意図しない語が認識された場合は認識結果を棄却して誤って撮像した画像を消去する。

【0222】

したがって、音声指示による撮影を行う場合に、周囲雑音による誤動作の影響を少なくするという効果がある。

【0223】

40

（第1の実施形態の変形例1）

本実施形態においては、発声の開始を検出したタイミングで撮像するか、発声の終了を検出したタイミングで撮像するようにしてもよい。

【0224】

（フローチャート）

図13に発声開始を検出した時点でのみ撮像する場合のフローチャートを示す。

図13に示したフローチャートは図6～図8で説明した処理と異なるフローチャートになるステップS811以降の処理のみ示している。

【0225】

また、図7～図8と同じ処理については同じ符号で示している。以下、図7～図8と図

50

3の相違点のみ説明する。

図13のフローチャートでは、図7のフローチャートで行っていた、発声終了を検出した時に撮像する処理（ステップS712、ステップS713）および撮像した画像を消去する処理（ステップS714）は行わない。

また図13のフローチャートでは、図8のフローチャートで行っていた、発声終了時に撮像を指示する語を認識した場合の認識結果処理部108が行う処理（ステップS826、ステップS827、ステップS828）を行わない。

【0226】

その他の処理については、図6～図8で説明した処理と同じである。

【0227】

なお、発声開始を検出した時点でのみ撮像する場合は、図9に示した音声認識文法から発声開始時点で撮像を指示する語（“Cheese”、“Say Cheese”等）を削除する。

【0228】

音声認識文法を変更しない場合は、図10に示した認識結果制御データを変更し、“Cheese”、“Say Cheese”等を認識した場合の処理を、発声の開始を検出した時点で撮像する処理に変更する。

【0229】

これにより、ユーザが“Cheese”、“Say Cheese”と発声すると、発声開始時点で撮像した画像の画像データが画像記憶媒体111に記憶される。

【0230】

同様にして、発声終了を検出した時点でのみ撮像するように変形することもできる。この場合は、発声開始を検出したときに撮像する処理（ステップS605、ステップS606）および発声開始が取り消されたときの処理（ステップS608）が省かれる。

【0231】

また、認識結果処理部108が行う処理のうち、ステップS822～ステップS824が省かれる。

【0232】

このとき、ステップS821で認識結果を受理した場合（ステップS821においてYES）に、ステップS826以降の処理を行う。

【0233】

また、発声開始時に撮像を指示する語を音声認識文法900から削除するか、認識結果制御データに記述された処理内容を変更する。

【0234】

（第1の実施形態の変形例2）

本実施形態では、認識結果によって発声の開始を検出した時点および発声の終了を検出した時点の画像の画像データを記憶媒体（画像記憶用）111に記憶する構成してもよい。

【0235】

たとえば、“Say Cheese”に対して発声開始を検出した時点および発声終了を検出した時点の両方で撮像するように認識結果制御データに記述すれば、両方時点における画像の画像データが記憶媒体（画像記憶用）111に記憶される。

【0236】

このような構成とすることで、ユーザが指示することが可能な撮像タイミングの種類が増えて、ユーザの利便性が向上する。

【0237】

（第1の実施形態の変形例3）

本実施形態では、認識結果処理部108が行う処理において、認識結果を棄却した場合（ステップS821においてNO）、メモリ（画像記憶用）110に記憶した画像A、画像Bを消去（ステップS825）するか否かユーザに確認させるようにしてもよい。

10

20

30

40

50

【 0 2 3 8 】

また、ユーザが記憶媒体（画像記憶用）1 1 1 に記憶する画像を選択するようにしてもよい。

【 0 2 3 9 】

また、認識結果が棄却された場合は、画像 A、画像 B の両方の画像データを記憶媒体（画像記憶用）1 1 1 に記憶するようにしてもよい。

【 0 2 4 0 】

例えば、ディスプレイ 1 1 5 に画像 A、画像 B を表示し、画像データの消去の可否を四方向ボタン 2 0 4 で選べるようする。

【 0 2 4 1 】

また、四方向ボタン 2 0 4 で記憶する画像をユーザが選択し、決定ボタン 2 0 5 が押された時点で選択されている画像の画像データを記憶媒体（画像記憶用）1 1 1 に記憶するようにする。

【 0 2 4 2 】

撮影を指示する語以外が認識された場合（ステップ S 8 2 6 において N O ）についても同様に、画像消去の確認、記憶媒体（画像記憶用）1 1 1 に記憶する画像の選択をユーザが行えるように構成できる。

【 0 2 4 3 】

また、画像 A、画像 B の画像データをともに記憶媒体（画像記憶用）1 1 1 に記憶するようにしてもよい。

【 0 2 4 4 】

このように構成することで、音声認識の性能が劣化するような環境で音声指示による撮像機能を使用する場合に、音声認識の誤りによって所望の撮像画像を誤って消去することが防止でき、ユーザの利便性が向上する。

【 0 2 4 5 】

尚、メモリ（画像記憶用）1 1 0 の記憶容量に応じて、1 回の音声認識において保持する画像の数を決定してもよい。

【 0 2 4 6 】

このように構成することで、メモリ 1 1 0 が限られた記憶容量を考慮して、ユーザが希望する画像の候補をできるだけ一時的に保持しておくことができる。

【 0 2 4 7 】

（第 1 の実施形態の変形例 4 ）

本実施形態の認識結果処理部 1 7 0 の処理において、撮像タイミングの異なる語の認識スコアの差が所定の閾値より小さい場合は、発声開始時点および発声終了時点で撮像された画像の両方を記憶媒体（画像記憶用）1 1 1 に記憶するようにしてもよい。

【 0 2 4 8 】

例えば、発声開始時点での撮像を指示する“ S h o o t ”と発声終了時点での撮像を指示する“ C h e e s e ”の認識スコアの差が所定値未満の場合に、発声開始時点および発声終了時点で撮像された画像を両方とも記憶媒体（画像記憶用）1 1 1 に記憶する。

【 0 2 4 9 】

あるいは、二つの画像をディスプレイ 1 1 5 に表示して、ユーザが選択するようにしてもよい。

【 0 2 5 0 】

このように構成することで、音声認識の性能が劣化するような環境で音声指示によって撮像を実行する機能を使用する場合に、音声認識の認識誤りによって、所望の撮像画像を誤って消去することが防止でき、ユーザの利便性が向上する。

【 0 2 5 1 】

（第 1 の実施形態の変形例 5 ）

本実施形態では、撮像した画像の画像データをメモリ（画像記憶用）1 1 0 に一時的に記憶し、認識結果確定後に記憶媒体（画像記憶用）1 1 1 に記憶するよう説明したが、最

10

20

30

40

50

初から記憶媒体（画像記憶用）１１１に記憶するようにしてもよい。

【０２５２】

この場合、ステップＳ６０８、ステップＳ７１４における画像データの消去の処理は、記憶媒体（画像記憶用）１１１に記憶された画像データを消去することになる。

【０２５３】

また、ステップＳ８２３、ステップＳ８２７の処理は行わない。

【０２５４】

さらに、認識結果を棄却した場合（ステップＳ８２１においてＮＯ）および認識結果が撮影を指示する語でない場合（ステップＳ８２６においてＮＯ）は、記憶媒体（画像記憶用）１１１に記憶されている画像Ａおよび画像Ｂの画像データを消去する。

10

【０２５５】

さらに、認識結果が発声開始時点で撮像を指示する語である場合には画像Ｂの画像データを消去し、発声終了時点で撮像を指示する語である場合には画像Ａの画像データを消去する。

【０２５６】

（第１の実施形態の変形例６）

例えば、道路脇等、周囲の雑音の影響を受けやすい場所で本実施形態に係るデジタルカメラ２００を使用する場合、音声検出部１０６の内部状態が短時間に頻繁に変化する場合がある。

【０２５７】

20

短時間のうちに撮像と画像データの消去が繰り返されると、デジタルカメラ２００の連写機能が画像データを消去した直後の撮像に対応しきれず、メモリ（画像記憶用）１１０上に画像が記憶されないということが起こりうる。

【０２５８】

これに対処するため、例えば、ステップＳ１０８で発声開始を検出した時点では撮像した画像Ａの画像データを消去せずに、次の発声開始を検出するまで画像Ａの画像データをメモリ（画像記憶用）１１０に記憶しておいてもよい。

【０２５９】

この場合、次に発声開始を検出した時点で画像Ａの画像データを消去するか、画像Ａの画像データに新たに撮像された画像の画像データを上書きするようにする。

30

【０２６０】

同様に、ステップＳ７１５において発声終了を取り消した場合も、画像Ｂの画像データを消去せずに、次に発声終了を検出するまでメモリ（画像記憶用）１１０に記憶しておいてもよい。

【０２６１】

このように構成することで、カメラの連写が音声検出の状態変化の速度に間に合わない場合でも、少なくとも最初に撮像された画像は残しておくことができる。

【０２６２】

尚、上記各実施形態では、カメラについて説明したが、本発明はビデオカメラ等の他の撮像装置にも適用することができる。

40

【０２６３】

（第１の実施形態の変形例７）

本実施形態では、マイク１１２として、周知ステレオマイクを用いる。

【０２６４】

また、音声認識部１０７は、左右のマイク１１２を介して入力されるそれぞれの音声信号の音量、ピッチ等の関係を前述した特徴量として用いてもよい。

【０２６５】

このような特徴量を用いることで、例えば、デジタルカメラ２００に対して右側から迫る音源と左側から迫る音源を判別することができる。すなわち、撮像する際の状況を認識して撮像することが可能となる。

50

【0266】

(第1の実施形態の変形例8)

本実施形態では、認識結果制御テーブルに含まれるコマンドの一例として示した“チーズ(Cheese)”に換えて“ハイ、チーズ(Say Cheese)”というコマンドに発声終了時点で撮像する処理を対応付けてもよい。

【0267】

また、認識結果制御テーブルに含まれるコマンドの一例として示した“Go”に換えて“今”というコマンドに発声開始時点で撮像する処理を対応付けてもよい。

【0268】

(第2の実施形態)

図16は、本発明の第2の実施形態に係る情報処理装置1600の構成の一例を示す機能ブロック図である。

【0269】

尚、図1と共通の要素には同一符号を付し、その説明を省略する。

【0270】

情報処理装置1600は、入力装置1602、撮像装置1603、格納装置(画像記憶用)1610、記憶装置(画像記憶用)1611、集音装置1612と接続可能であることを特徴とする。

【0271】

また、情報処理装置1600は、格納装置(音声認識データ用)1613、格納装置(認識結果制御データ用)1614、表示制御装置1609と接続可能であることを特徴とする。

【0272】

尚、入力装置1602は操作部102に、撮像装置1603は撮像部103に、格納装置(画像記憶用)1610はメモリ(画像記憶用)110に、記憶装置(画像記憶用)1611は記憶媒体(画像記憶用)111に対応する機能を備える。

【0273】

また、集音装置1612はマイク112に、格納装置(音声認識データ用)1613はメモリ(音声認識データ用)113に対応する機能を備える。

【0274】

また、格納装置(認識結果制御データ用)1614は認識結果データ用メモリ114に、表示制御装置1609は表示制御部109に対応する機能を備える。

【0275】

情報処理装置1600としては、例えばマイクロプロセッサ等が想定できる。

【0276】

図14、図15は、情報処理装置1600における処理動作の一例を示したフローチャートである。

【0277】

まず、図14のフローチャートを参照して説明する。

ステップS1400で、音声入力部105は音声信号が入力されたか否か判断する。

【0278】

音声信号が入力された場合(ステップS1400においてYES)、ステップS1401で音声検出部106はフレームfを初期化する($f=0$)。

【0279】

次に、ステップS1402で音声検出部106は音声信号の検出状態を第一状態301に設定する。

【0280】

次に、ステップS1403で音声検出部106は検出の対象となるフレームを設定する。

【0281】

10

20

30

40

50

次に、ステップ S 1 4 0 4 で音声検出部 1 0 6 は音声入力部 1 0 5 に入力された音声信号の特徴量データを記憶する。

【 0 2 8 2 】

尚、特徴量データとは、音声認識部 1 0 7 が音声認識を行う場合に使用するデータである。

【 0 2 8 3 】

次に、ステップ S 1 4 0 5 で音声検出部 1 0 6 は音声の検出状態が第 1 の状態から第 4 の状態のいずれであるか判断する。

【 0 2 8 4 】

ステップ S 1 4 0 5 で音声検出部 1 0 6 が検出状態を第一状態 3 0 1 であると判断した場合、ステップ S 1 4 0 6 で、第 1 の検出として、音声検出部 1 0 6 は閾値 T H 1 以上の音量を検出したか否か判断する。

【 0 2 8 5 】

閾値 T H 1 以上の音量を検出した場合（ステップ S 1 4 0 6 において Y E S ）、ステップ S 1 4 0 7 で、音声検出部 1 0 6 は検出状態を第二状態 3 0 2 に遷移させる（このタイミングを第 1 の時刻とする）。

【 0 2 8 6 】

次に、ステップ S 1 4 0 8 で、撮像制御部 1 2 3 は撮像装置 1 6 0 3 に撮像動作を実行させる信号を出力する。

【 0 2 8 7 】

尚、ステップ S 1 4 0 8 で出力された信号によって撮像された画像を画像 A とする。

【 0 2 8 8 】

次に、ステップ S 1 4 0 9 で、画像記憶制御部 1 0 4 は、第 1 の取得として、直前のステップ S 1 4 0 8 によって撮像された画像 A を表す画像データを格納装置（画像記憶用）1 6 1 0 に記憶させる信号を出力する。

【 0 2 8 9 】

次に、ステップ S 1 4 1 0 で、第 1 の記憶として、音声検出部 1 0 6 は処理中のフレーム f を発声開始フレーム F s として記憶する。

【 0 2 9 0 】

次に、ステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

【 0 2 9 1 】

また、ステップ S 1 4 0 6 で閾値 T H 1 以上の音量を検出しなかった場合（ステップ S 1 4 0 6 において N O ）、同様にステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

【 0 2 9 2 】

また、ステップ S 1 4 0 5 で音声検出部 1 0 6 が検出状態を第二状態 3 0 2 であると判断した場合、ステップ S 1 4 1 1 で、処理中のフレーム f が発声開始フレーム F s から M 1 回目のフレーム以上であるか否か判断する。

【 0 2 9 3 】

また、処理中のフレーム f が発声開始フレーム F s から M 1 回目のフレーム未満である場合（ステップ S 1 4 1 1 において Y E S ）、ステップ S 1 4 1 3 で音声検出部 1 0 6 が閾値 T H 1 より大きい音量を検出したか否か判断する。

【 0 2 9 4 】

閾値 T H 1 より大きい音量を検出しなかった場合（ステップ S 1 4 1 3 において N O ）、ステップ S 1 4 1 4 で音声検出部 1 0 6 はカウンタ F a の値を初期化する。

【 0 2 9 5 】

次に、ステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

【 0 2 9 6 】

10

20

30

40

50

尚、カウンタ F a とは、発声開始フレーム F s を設定し直すか否か判定するために使用する。

【 0 2 9 7 】

また、閾値 T H 1 より未満の音量を検出した場合（ステップ S 1 4 1 3 において Y E S ）、ステップ S 1 4 1 5 で音声検出部 1 0 6 はカウンタ F a の値を 1 増やす。

【 0 2 9 8 】

次に、ステップ S 1 4 1 6 で音声検出部 1 0 6 はカウンタ F a の値が N 1 以上であるか判断する。

【 0 2 9 9 】

カウンタ F a の値が N 1 以上である場合（ステップ S 1 4 1 6 において Y E S ）、ステップ S 1 4 1 7 で画像記憶制御部 1 0 4 は格納装置（画像記憶用）1 6 1 0 に記憶された画像 A を表す画像データを消去するための信号を出力する。

10

【 0 3 0 0 】

尚、ステップ S 1 4 1 7 における処理は、音声認識後に画像データを消去する処理に対して、第 2 の消去に相当する。

【 0 3 0 1 】

次に、ステップ S 1 4 1 8 で、発声の開始を再検出する第 1 の再検出をおこなうために、音声検出部 1 0 6 は検出状態を第一状態 3 0 1 に遷移させる。

【 0 3 0 2 】

次に、ステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

20

【 0 3 0 3 】

また、カウンタ F a の値が N 1 未満である場合（ステップ S 1 4 1 6 において N O ）、同様にステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

【 0 3 0 4 】

また、ステップ S 1 4 1 1 で処理中のフレーム f が発声開始フレーム F s から M 1 回目のフレーム以上である場合（ステップ S 1 4 1 1 において N O ）、ステップ S 1 4 1 2 で音声検出部 1 0 6 は検出状態を第三状態 3 0 3 に遷移させる。

【 0 3 0 5 】

30

また、ステップ S 1 4 0 5 で音声検出部 1 0 6 が検出状態を第三状態 3 0 3 であると判断した場合、ステップ S 1 4 1 9 で、第 2 の検出として、音声検出部 1 0 6 は閾値 T H 2 以下の音量を検出したか否か判断する。

【 0 3 0 6 】

閾値 T H 2 以下の音量を検出した場合（ステップ S 1 4 1 9 において Y E S ）、ステップ S 1 4 2 0 で、音声検出部 1 0 6 は検出状態を第四状態 3 0 4 に遷移させる（このタイミングを第 2 の時刻とする）。

【 0 3 0 7 】

次に、ステップ S 1 4 2 1 で、撮像制御部 1 2 3 は撮像装置 1 6 0 3 に撮像動作を実行させるための信号を出力する。

40

【 0 3 0 8 】

尚、ステップ S 1 4 2 1 で出力された信号によって撮像された画像を画像 B とする。

【 0 3 0 9 】

次に、ステップ S 1 4 2 2 で、画像記憶制御部 1 0 4 は、第 2 の取得として、直前のステップ S 1 4 2 1 で撮像された画像 B を表す画像データを格納装置（画像記憶用）1 6 1 0 に記憶させる信号を出力する。

【 0 3 1 0 】

次に、ステップ S 1 4 2 3 で、第 2 の記憶として、音声検出部 1 0 6 は処理中のフレーム f を発声終了フレーム F e として設定する。

【 0 3 1 1 】

50

次に、ステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

【 0 3 1 2 】

また、ステップ S 1 4 1 9 で閾値 T H 1 以上の音量を検出なかった場合（ステップ S 1 4 1 9 において N O ）、同様にステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

【 0 3 1 3 】

また、ステップ S 1 4 0 5 で音声検出部 1 0 6 が検出状態を第四状態 3 0 4 であると判断した場合、ステップ S 1 4 2 4 で、処理中のフレーム f が発声終了フレーム F e から M 2 回目のフレーム以上であるか否か判断する。

【 0 3 1 4 】

また、処理中のフレーム f が発声終了フレーム F e から M 2 回目のフレーム未満である場合（ステップ S 1 4 2 4 において Y E S ）、ステップ S 1 4 2 6 で音声検出部 1 0 6 が閾値 T H 2 より大きい音量を検出したか否か判断する。

【 0 3 1 5 】

閾値 T H 2 より大きい音量を検出しなかった場合（ステップ S 1 4 2 6 において N O ）、ステップ S 1 4 2 7 で音声検出部 1 0 6 はカウンタ F b の値を初期化する。

【 0 3 1 6 】

次に、ステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は音声検出の対象となるフレームを設定する。

【 0 3 1 7 】

尚、カウンタ F b とは、発声終了フレーム F e を設定し直すか否か判定するために使用する。

【 0 3 1 8 】

また、閾値 T H 2 より大きい音量を検出した場合（ステップ S 1 4 2 6 において Y E S ）、ステップ S 1 4 2 8 で音声検出部 1 0 6 はカウンタ F b の値を 1 増やす。

【 0 3 1 9 】

次に、ステップ S 1 4 2 9 で音声検出部 1 0 6 はカウンタ F b の値が N 2 以上であるか判断する。

【 0 3 2 0 】

カウンタ F b の値が N 2 以上である場合（ステップ S 1 4 2 9 において Y E S ）、ステップ S 1 4 3 0 で画像記憶制御部 1 0 4 は格納装置（画像記憶用）1 6 1 0 に記憶された画像 B を表す画像データを消去するための信号を出力する。

【 0 3 2 1 】

尚、ステップ S 1 4 3 0 における処理は、音声認識後に画像データを消去する処理に対して、第 3 の消去に相当する。

【 0 3 2 2 】

次に、ステップ S 1 4 3 1 で音声検出部 1 0 6 は、発声の終了を再検出する第 2 の再検出をおこなうために、検出状態を第三状態 3 0 3 に遷移させる。

【 0 3 2 3 】

次に、ステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

【 0 3 2 4 】

また、カウンタ F b の値が N 2 未満である場合（ステップ S 1 4 2 9 において N O ）、同様にステップ S 1 4 0 3 に戻り、音声検出部 1 0 6 は次の音声検出の対象となるフレームを設定する。

【 0 3 2 5 】

また、ステップ S 1 4 2 4 で処理中のフレーム f が発声開始フレーム F e から M 2 回目のフレーム以上である場合（ステップ S 1 4 2 4 において N O ）、ステップ S 1 4 2 5 で音声検出部 1 0 6 は音声検出を終了する。

10

20

30

40

50

【 0 3 2 6 】

次に、図 1 5 のフローチャートを参照して説明する。

【 0 3 2 7 】

ステップ S 1 5 3 2 で音声認識部 1 0 7 はステップ S 1 5 0 4 で取得した各フレームの特徴量データと音声認識用データとに基づいて音声認識を行う。

【 0 3 2 8 】

次に、ステップ S 1 5 3 3 で音声認識部 1 0 7 による音声認識を終了する。

【 0 3 2 9 】

尚、ステップ S 1 5 3 3 の処理は、音声認識部 1 0 7 によって音声認識の結果が得られた後に実行する。

10

【 0 3 3 0 】

次に、ステップ S 1 5 3 4 で、認識結果処理部 1 0 8 は音声認識の結果が発声開始のタイミングで撮像を指示する内容であるか否か判断する。

【 0 3 3 1 】

発声開始のタイミングで撮像を指示する内容である場合（ステップ S 1 5 3 4 の Y E S）、ステップ S 1 5 3 5 で画像 B を消去するための信号を出力する。

【 0 3 3 2 】

発声開始のタイミングで撮像を指示する内容でない場合（ステップ S 1 5 3 4 の N O）、ステップ S 1 5 3 6 で、認識結果処理部 1 0 8 は音声認識の結果が発声終了のタイミングで撮像を指示する内容であるか否か判断する。

20

【 0 3 3 3 】

発声終了のタイミングで撮像を指示する内容である場合（ステップ S 1 5 3 6 の Y E S）、ステップ S 1 5 3 7 で画像 A を消去するための信号を出力する。

【 0 3 3 4 】

発声終了のタイミングで撮像を指示する内容でない場合（ステップ S 1 5 3 6 の N O）、ステップ S 1 5 3 8 で、画像 A、画像 B を消去するための信号を出力する。

【 0 3 3 5 】

次に、ステップ S 1 5 3 9 で、認識結果処理部 1 0 8 は音声認識の結果が発声開始のタイミングから一定時間経過後に撮像を指示する内容であるか否か判断する。

【 0 3 3 6 】

発声開始から一定時間経過後に撮像を指示する内容である場合（ステップ S 1 5 3 9 の Y E S）、ステップ S 1 5 4 0 で一定時間経過後（このタイミングを第 3 の時刻とする）に、撮像制御部 1 2 3 は撮像装置 1 6 0 3 に撮像動作を実行させるための信号を出力する。

30

【 0 3 3 7 】

尚、ステップ S 1 5 4 0 で出力された信号によって撮像された画像を画像 C とする。

【 0 3 3 8 】

次に、ステップ S 1 5 4 1 で、画像記憶制御部 1 0 4 は、第 3 の保持として、直前のステップ S 1 5 4 0 で撮像された画像 C を表す画像データを格納装置（画像記憶用）1 6 1 0 に記憶させる信号を出力して、処理を終了する。

40

【 0 3 3 9 】

また、発声開始のタイミングから一定時間経過後に撮像を指示する意味内容でない場合（ステップ S 1 4 3 9 の N O）、処理を終了する。

【 0 3 4 0 】

このような構成とすることで、発声区間に対して、第 1 の関係である発声開始のタイミングで撮像された第 1 の画像（画像 A）と、第 2 の関係である発声終了のタイミングで撮像された第 2 の画像（画像 B）とを得ることができる。

【 0 3 4 1 】

また、発声区間に対して、第 3 の関係である発声終了から一定時間のタイミングで撮像された第 3 の画像（画像 C）を得ることができる。

50

【0342】

さらに、音声区間の音声の意味内容に応じて、複数の画像からユーザが所望するタイミングで撮像された画像を選択することができる。

【0343】

また、このような構成とすることで、本実施形態の情報処理装置1600と、外部機器とを連動させて、ユーザが欲するタイミングで撮像された画像を効率良く取得することができる。

【0344】

また、本実施形態の情報処理装置1600によると、断続的に音声が入力された場合にも、1つのコマンドとして認識することが可能であるため、発声区間が長くなるような言葉コマンドとして利用した場合にも認識の誤りが軽減される。

10

【0345】

(プログラムCLのサポート)

尚、本発明の目的は、前述した実施形態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が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。

【0346】

尚、コンピュータは、CPU、MPU等であってもよい。

【0347】

この場合、記憶媒体から読み出されたコンピュータ読み取り可能な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。

20

【0348】

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等を用いることができる。

【0349】

また、コンピュータが読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、OS(オペレーティングシステム)等が実際の処理の一部または全部を行ってもよい。

30

【0350】

尚、こ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。

尚、OSはコンピュータ上で稼働している。

【0351】

また、まず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る。

【0352】

その後、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等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。

40

【図面の簡単な説明】

【0353】

【図1】本発明の第1の実施形態に係る情報処理装置の構成の一例を示す機能ブロック図である。

【図2】本発明の第1の実施形態で想定されるデジタルカメラの外観を示す図である。

【図3】音声検出部106が判定した認識状態の一例を示す図である。

【図4】音声検出部106の動作の一例を示す概念図である。

【図5】音声検出部106における処理動作を示すフローチャートである。

【図6】音声によって撮像を指示する場合のデジタルカメラ200における処理の一例を

50

示す第一のフローチャートである。

【図 7】音声によって撮像を指示する場合のデジタルカメラ 200 における処理の一例を示す第二のフローチャートである。

【図 8】音声によって撮像を指示する場合のデジタルカメラ 200 における処理の一例を示す第三のフローチャートである。

【図 9】本発明の第 1 の実施形態で利用する音声認識文法の一例を示す図である。

【図 10】認識結果制御テーブルの一例を示す図である。

【図 11】本発明の第 1 の実施形態に係るデジタルカメラ 200 を利用して、“ S h o o t ” という音声指示で撮像する場合の動作を示す図である。

【図 12】本発明の第 1 の実施形態に係るデジタルカメラ 200 を利用して、“ C h e e s e ” という音声指示で撮像する場合の動作を示す図である。 10

【図 13】発声開始を検出した時点でのみ撮像する場合のフローチャートである。

【図 14】情報処理装置 1600 における処理動作の一例を示した第 1 のフローチャートである。

【図 15】情報処理装置 1600 における処理動作の一例を示した第 2 のフローチャートである。

【図 16】本発明の第 2 の実施形態に係る情報処理装置 1600 の構成の一例を示す機能ブロック図である。

【図 17】音声検出部 106 が判定した認識状態と撮像部 103、画像記憶制御部 104 の動作の一例を示す図である。 20

【符号の説明】

【 0 3 5 4 】

101 制御部

104 画像記憶制御部

105 音声入力部

106 音声検出部

107 音声認識部

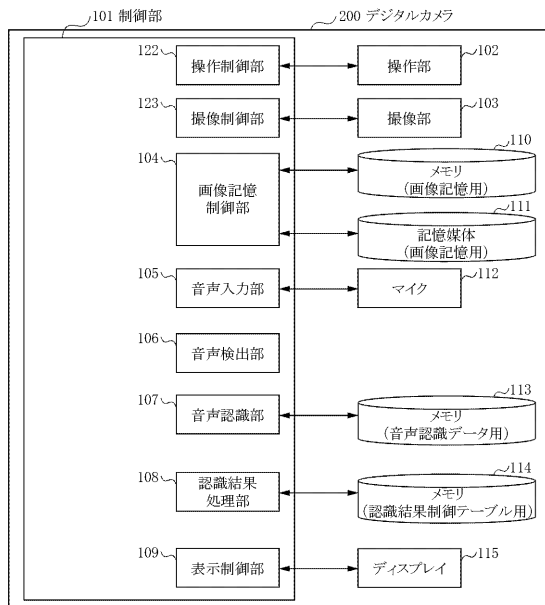
108 認識結果処理部

109 表示制御部

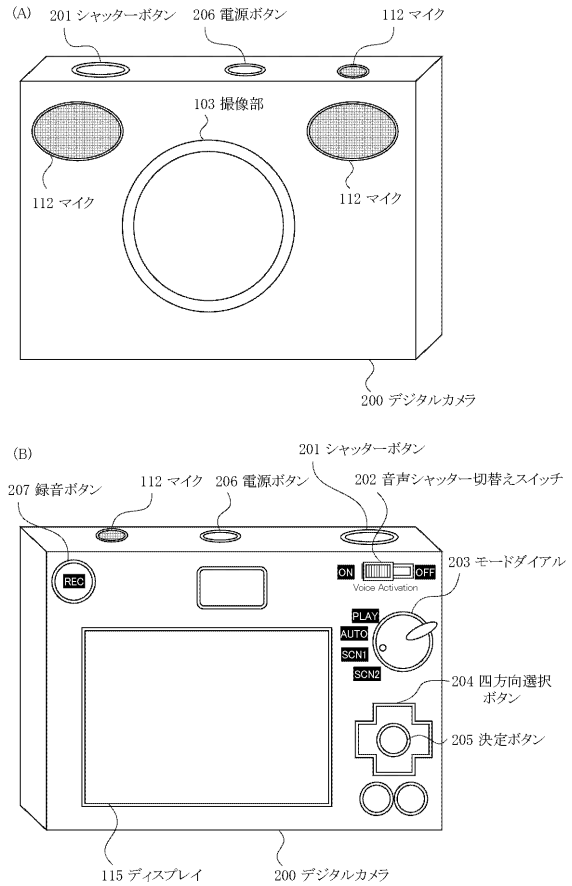
122 操作制御部

123 撮像制御部 30

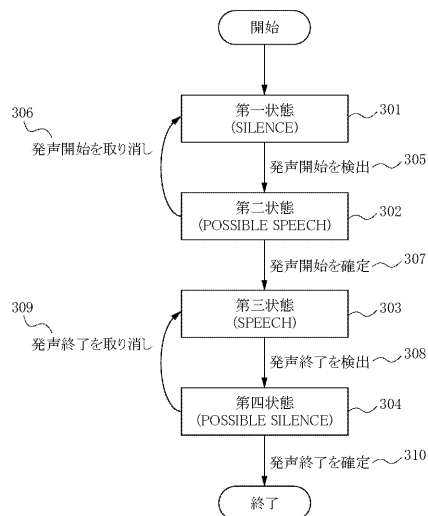
【図 1】



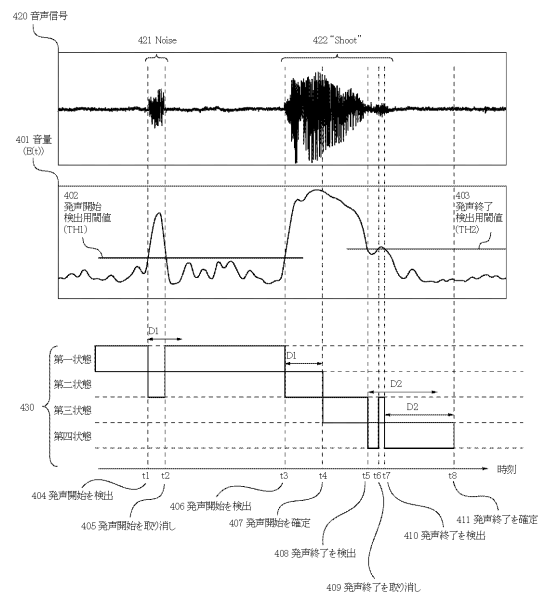
【図 2】



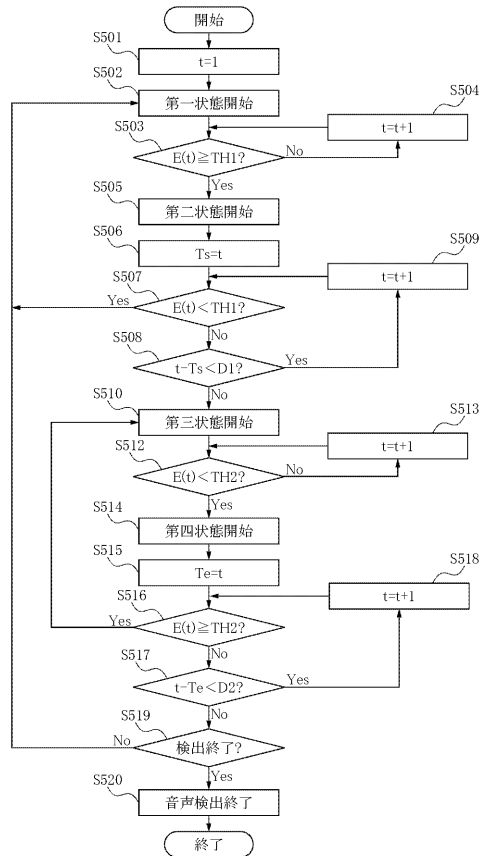
【図 3】



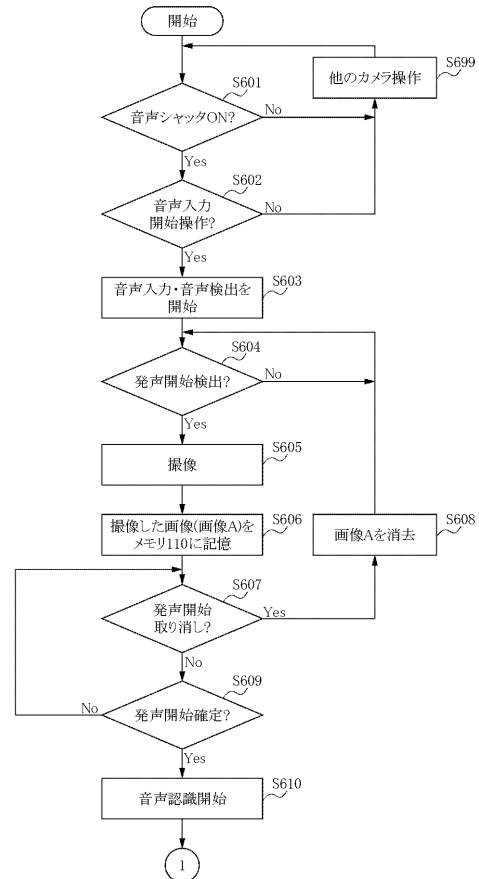
【図 4】



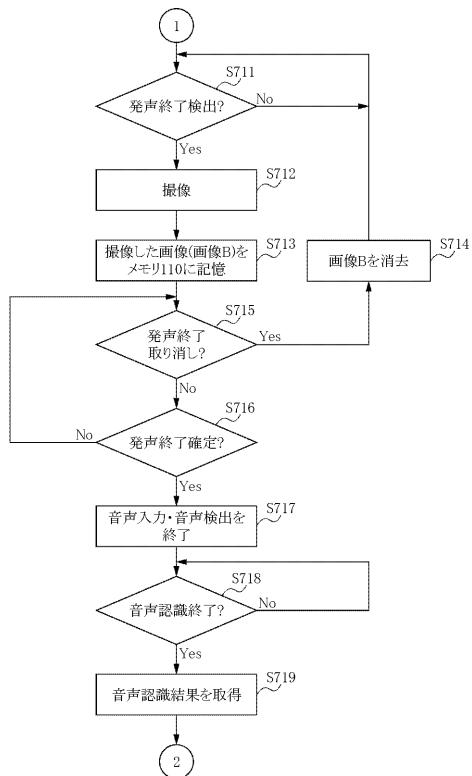
【図5】



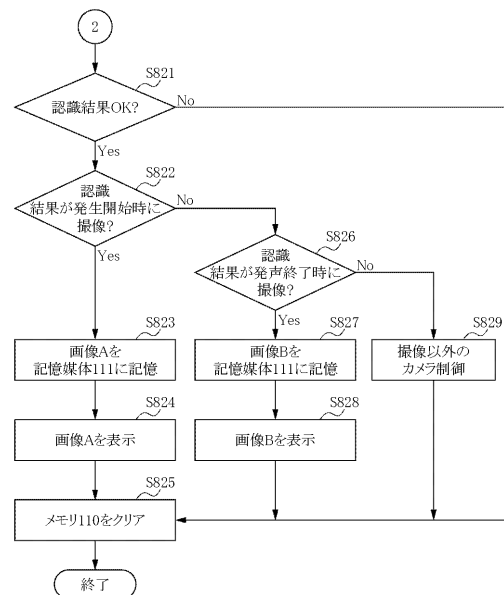
【図6】



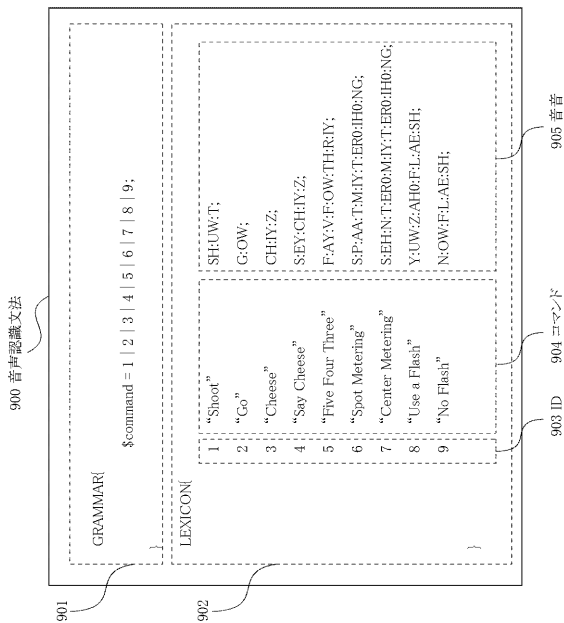
【図7】



【図8】



【図 9】



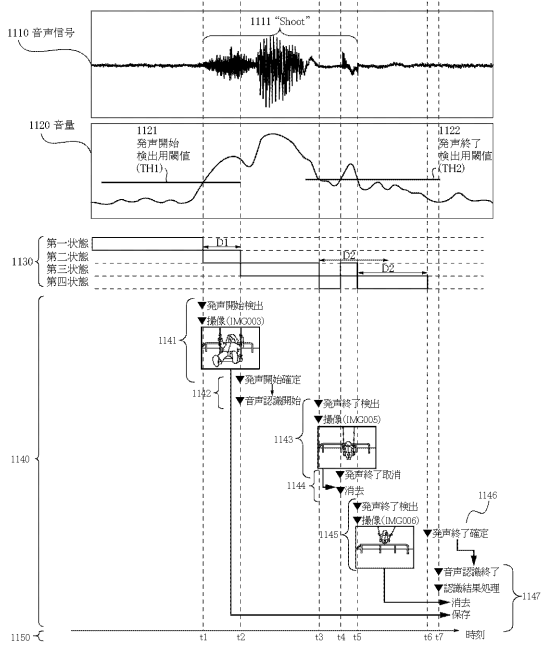
【図 10】

コマンド	処理
Shoot	発声開始時点で撮像された画像を保持
Go	発声開始時点で撮像された画像を保持
Cheese	発声終了時点で撮像された画像を保持
Say Cheese	発声終了から所定時間経過後に撮像された画像を保持
Five Four Three	発声開始から所定時間経過後に撮像された画像を保持
Spot Metering	スポット測光に設定
Center Metering	中央部重点測光に設定
Use a Flash	ストロボ発光に設定
No Flash	ストロボ発光禁止に設定

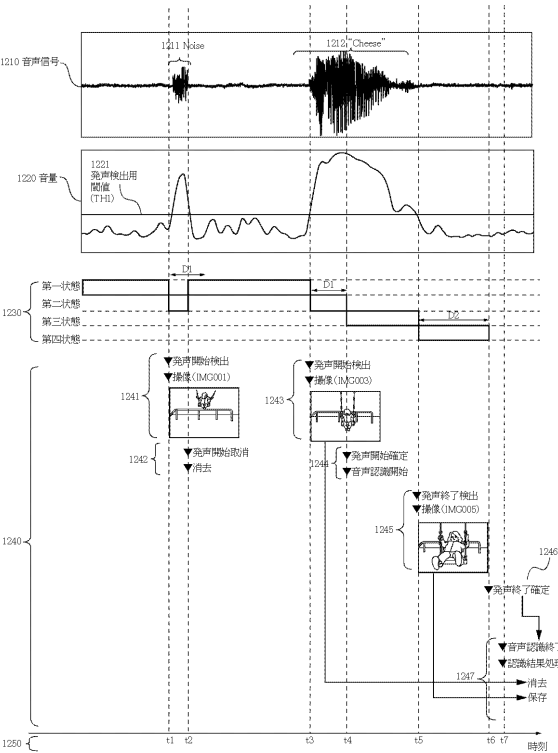
認識結果制御データ

904 コマンド 1002 制御内容

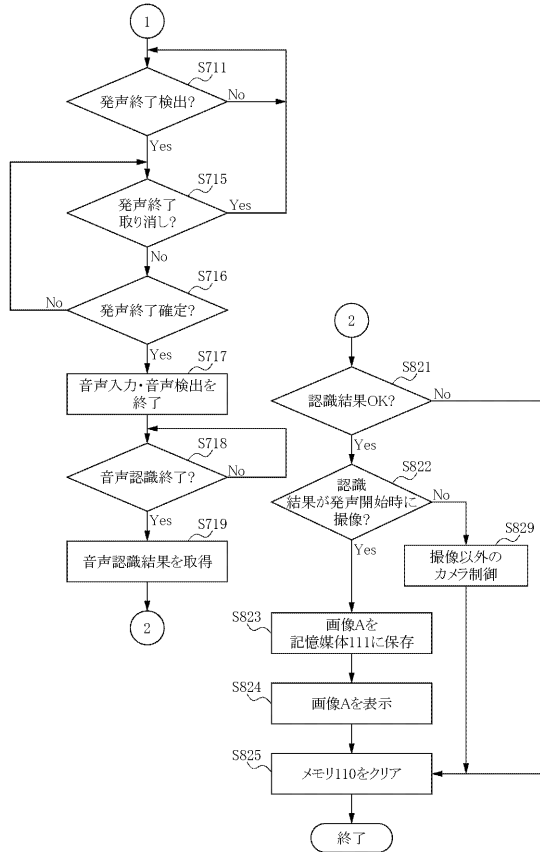
【図 11】



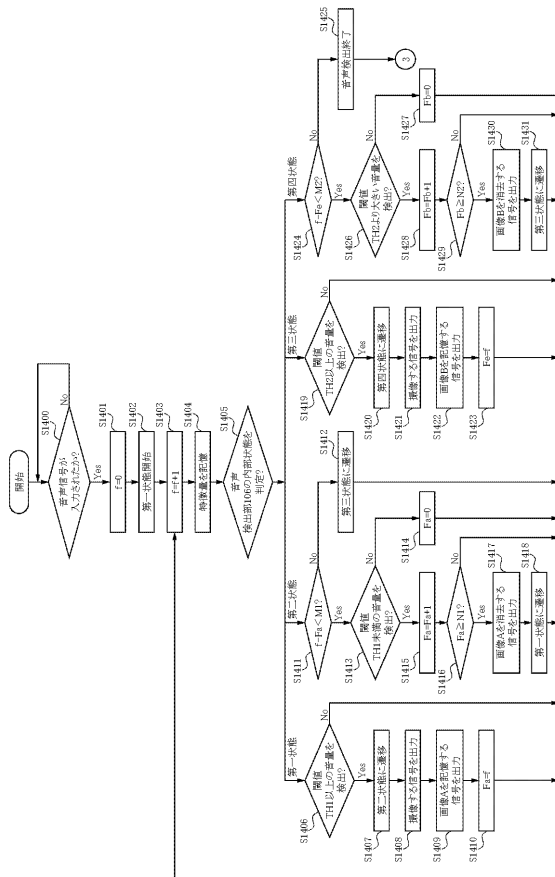
【図 12】



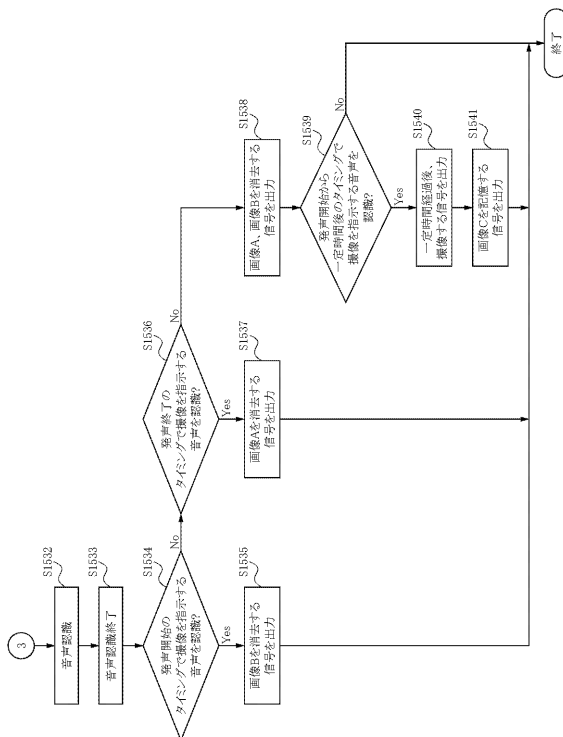
【 図 1 3 】



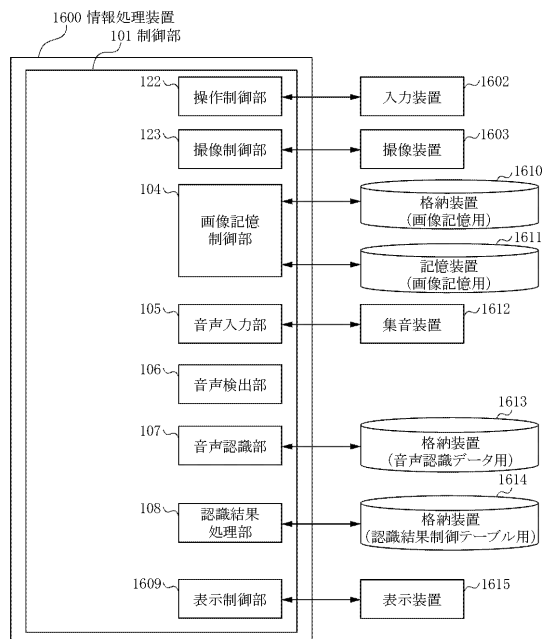
【 図 1 4 】



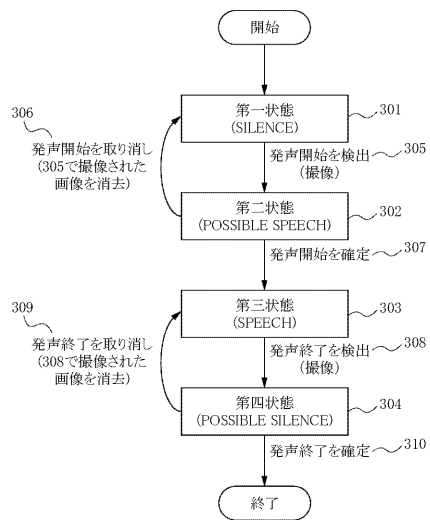
【 図 1 5 】



【 図 1 6 】



【図 17】



フロントページの続き

- (56)参考文献 特開2005-181365(JP,A)
特開2004-301895(JP,A)
特開2004-287063(JP,A)
特開2005-159558(JP,A)
特開2005-122139(JP,A)
特開平11-194392(JP,A)
特開2006-184589(JP,A)

- (58)調査した分野(Int.Cl., DB名)
H04N 5/225