

(11) 特許出願公開番号

特開2009-141555

(P2009-141555A)

(43) 公開日 平成21年6月25日(2009.6.25)

(51) Int. Cl.	F I	テーマコード (参考)
H04N 5/225 (2006.01)	H04N 5/225 F	5C053
H04N 5/91 (2006.01)	H04N 5/91 C	5C122

審査請求 未請求 請求項の数 15 O L (全 13 頁)

(21) 出願番号 特願2007-314454 (P2007-314454)
(22) 出願日 平成19年12月5日 (2007. 12. 5)

(71) 出願人 306037311
富士フイルム株式会社
東京都港区西麻布2丁目26番30号

(74) 代理人 100115107
弁理士 高松 猛

(74) 代理人 100132986
弁理士 矢澤 清純

(72) 発明者 穂山 俊文
埼玉県朝霞市泉水3丁目11番46号 富士フイルム株式会社内

Fターム(参考) 5C053 FA08 JA05 JA16 LA01
5C122 DA04 EA61 FJ01 FJ04 FJ15
FK37 FK42 GA20 GA24 GA31
HA04 HB01 HB05

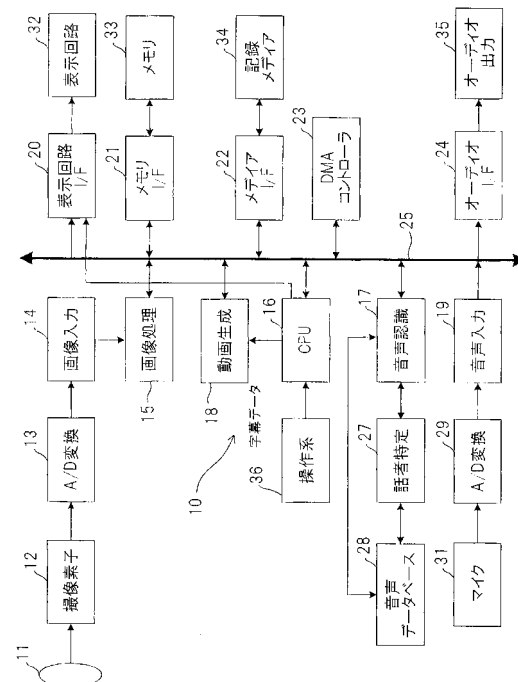
(54) 【発明の名称】 音声入力機能付き撮像装置及びその音声記録方法

(57) 【要約】

【課題】撮影後の編集によらずとも主要被写体の発する音声テキストデータで被写体画像に合成する。

【解決手段】被写体画像を撮像する撮像手段 1 2 と、撮像手段 1 2 が被写体画像を撮像するとき同時に音声を取り込む音声入力手段 3 1 と、該音声をテキストデータ化する音声認識手段 1 7 と、該テキストデータを被写体画像に合成する画像合成手段 1 8 とを備える。

【選択図】図 1



【特許請求の範囲】**【請求項 1】**

被写体画像を撮像する撮像手段と、該撮像手段が前記被写体画像を撮像するとき同時に音声を取り込む音声入力手段と、該音声をテキストデータ化する音声認識手段と、該テキストデータを前記被写体画像に合成する画像合成手段とを備えることを特徴とする音声入力機能付き撮像装置。

【請求項 2】

前記音声の話者を特定する話者特定手段を備え、前記画像合成手段は、特定された前記話者が予め登録された登録者である場合のみ前記合成を行うことを特徴とする請求項 1 に記載の音声入力機能付き撮像装置。

10

【請求項 3】

前記話者が特定され該話者の音声データを前記テキストデータで合成するとき該話者を特定する名前を該テキストデータと一緒に合成することを特徴とする請求項 2 に記載の音声入力機能付き撮像装置。

【請求項 4】

前記被写体画像中に前記登録者の顔が写っているか否かを識別する顔認識手段を備え、前記画像合成手段は、前記登録者の顔が写っており且つ該登録者が前記話者のとき前記テキストデータを前記被写体画像中の前記登録者の顔画像からの吹き出し形状で合成することを特徴とする請求項 1 に記載の音声入力機能付き撮像装置。

【請求項 5】

20

前記画像合成手段は、前記被写体画像中に複数の登録者が写っており前記音声は複数人の音声のときは各音声データに対応したテキストデータを各人の画像の近くに吹き出し形状で合成することを特徴とする請求項 4 に記載の音声入力機能付き撮像装置。

【請求項 6】

前記被写体画像中に前記登録者の顔が写っているか否かを識別する顔認識手段を備え、前記画像合成手段は、前記登録者の顔が写っておらず且つ該登録者が前記話者のとき前記テキストデータを字幕データとして合成することを特徴とする請求項 1 に記載の音声入力機能付き撮像装置。

【請求項 7】

前記登録者が前記被写体画像中に複数人写っている場合には該複数人の中の一人の登録者を指定する手段を備え、前記画像合成手段は、指定された登録者の音声を前記テキストデータで該被写体画像中に合成することを特徴とする請求項 1 に記載の音声入力機能付き撮像装置。

30

【請求項 8】

指向性を持つ前記音声入力手段の該指向性を前記被写体画像中の登録者の方向に制御する制御手段を備えることを特徴とする請求項 1 に記載の音声入力機能付き撮像装置。

【請求項 9】

被写体画像を撮像する撮像手段と、該撮像手段が前記被写体画像を撮像するとき同時に音声を取り込む音声入力手段と、該音声をテキストデータ化する音声認識手段とを備える音声入力機能付き撮像装置の音声記録方法において、前記テキストデータを前記被写体画像に合成することで前記音声を記録すること特徴とする音声入力機能付き撮像装置の音声記録方法。

40

【請求項 10】

前記音声の話者を特定し、特定された前記話者が予め登録された登録者である場合のみ前記テキストデータの合成を行うことを特徴とする請求項 9 に記載の音声入力機能付き撮像装置の音声記録方法。

【請求項 11】

前記話者が特定され該話者の音声を前記テキストデータで合成するとき該話者を特定する名前を該テキストデータと一緒に合成することを特徴とする請求項 10 に記載の音声入力機能付き撮像装置の音声記録方法。

50

【請求項 1 2】

前記被写体画像中に前記登録者の顔が写っているか否かを識別し、前記登録者の顔が写っており且つ該登録者が前記話者のとき前記テキストデータを前記被写体画像中の前記登録者の顔画像からの吹き出し形状で合成することを特徴とする請求項 9 に記載の音声入力機能付き撮像装置の音声記録方法。

【請求項 1 3】

前記被写体画像中に複数の登録者が写っており前記音声は複数人の音声のときは各音声データに対応したテキストデータを各人の画像の近くに吹き出し形状で合成することを特徴とする請求項 1 2 に記載の音声入力機能付き撮像装置の音声記録方法。

【請求項 1 4】

前記被写体画像中に前記登録者の顔が写っているか否かを識別し、前記登録者の顔が写っておらず且つ該登録者が前記話者のとき前記テキストデータを字幕データとして合成することを特徴とする請求項 9 に記載の音声入力機能付き撮像装置の音声記録方法。

【請求項 1 5】

前記登録者が前記被写体画像中に複数人写っている場合に、該複数人の中の一人の登録者の指定操作に従って、指定された登録者の音声を前記テキストデータで該被写体画像中に合成することを特徴とする請求項 9 に記載の音声入力機能付き撮像装置の音声記録方法。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明は、マイク等の音声入力機能を持った撮像装置及びその音声記録方法に関する。

【背景技術】

【0 0 0 2】

特許文献 1 には、会議などにおける音声をマイクで拾い、この音声をテキスト化して記録する技術が開示されている。また、特許文献 2 には、デジタルカメラ等で映像の他に音声をマイクで拾い、撮影された画像及び音声を解析して、音声をテキスト化したデータを画像ファイルに関連付けて記録する技術が開示されている。

【0 0 0 3】

また、特許文献 3 には、撮影された動画像データを解析し、顔検出機能により顔の特徴量と位置とを検出し、音声識別機能によって音声の特徴量を検出し、これらを基に動画像データ中の特定話者の位置を特定すると共に特定話者の音声をテキストデータ化し、特定話者が喋った内容を字幕として動画像データ中に合成する技術が開示されている。

【0 0 0 4】

また、特許文献 4 は、画像上の人物と音声との対応付けの精度を向上させ、音声をテキストデータ化して話者に対応させて表示させる技術を開示している。特許文献 5 は、外部から送信されてきた画像データおよび音声データを受信し、音声をテキストデータ化して画面上に表示する技術を開示する。

【0 0 0 5】

特許文献 6 は、入力された音声をテキストデータ化し、カメラ部で撮影された画像に対応させて記録させると共に、所望の画像に対応のテキストデータで検索できる様にする技術を開示する。

【特許文献 1】特開 2 0 0 6 - 1 8 9 6 2 6 号公報

【特許文献 2】特開 2 0 0 6 - 1 3 3 4 3 3 号公報

【特許文献 3】特開 2 0 0 7 - 2 7 9 9 0 号公報

【特許文献 4】特開 2 0 0 4 - 5 6 2 8 6 号公報

【特許文献 5】特開平 9 - 2 3 3 4 4 2 号公報

【特許文献 6】特開 2 0 0 5 - 2 0 4 4 0 号公報

【発明の開示】

【発明が解決しようとする課題】

【 0 0 0 6 】

上述した各特許文献に記載されている様に、マイクで拾った音声をテキストデータ化する技術が一般的となり、撮像された画像に対応付けて音声を字幕として、あるいは吹き出しとして、表示することが行われるようになってきている。

【 0 0 0 7 】

しかし、従来技術は、いずれも、撮影済みの画像を後でパソコン等を用いて解析して、話者を特定すると共に音声をテキストデータ化し、両者に対応付ける解析を行い、字幕表示や吹き出し表示を行うという、撮影後の編集処理であるため、面倒であるという問題がある。また、話者と音声との対応付けの精度が低く、撮影中の主要人物の音声を高精度に抽出して字幕表示や吹き出し表示を精度良く対応付けて行うことができないという問題がある。

10

【 0 0 0 8 】

本発明の目的は、撮影後の編集によらなくても主要被写体の発する音声をテキストデータで被写体画像に合成することができ、また、撮影中の主要人物の音声を高精度に対応付けて字幕表示や吹き出し表示することができる音声入力機能付き撮像装置及びその音声記録方法を提供することにある。

【課題を解決するための手段】

【 0 0 0 9 】

本発明の音声入力機能付き撮像装置は、被写体画像を撮像する撮像手段と、該撮像手段が前記被写体画像を撮像するとき同時に音声を取り込む音声入力手段と、該音声をテキストデータ化する音声認識手段と、該テキストデータを前記被写体画像に合成する画像合成手段とを備えることを特徴とする。

20

【 0 0 1 0 】

本発明の音声入力機能付き撮像装置は、前記音声の話者を特定する話者特定手段を備え、前記画像合成手段は、特定された前記話者が予め登録された登録者である場合のみ前記合成を行うことを特徴とする。

【 0 0 1 1 】

本発明の音声入力機能付き撮像装置は、前記話者が特定され該話者の音声データを前記テキストデータで合成するとき該話者を特定する名前を該テキストデータと一緒に合成することを特徴とする。

30

【 0 0 1 2 】

本発明の音声入力機能付き撮像装置は、前記被写体画像中に前記登録者の顔が写っているか否かを識別する顔認識手段を備え、前記画像合成手段は、前記登録者の顔が写っており且つ該登録者が前記話者のとき前記テキストデータを前記被写体画像中の前記登録者の顔画像からの吹き出し形状で合成することを特徴とする。

【 0 0 1 3 】

本発明の音声入力機能付き撮像装置の前記画像合成手段は、前記被写体画像中に複数の登録者が写っており前記音声は複数人の音声のときは各音声データに対応したテキストデータを各人の画像の近くに吹き出し形状で合成することを特徴とする。

【 0 0 1 4 】

40

本発明の音声入力機能付き撮像装置は、前記被写体画像中に前記登録者の顔が写っているか否かを識別する顔認識手段を備え、前記画像合成手段は、前記登録者の顔が写っておらず且つ該登録者が前記話者のとき前記テキストデータを字幕データとして合成することを特徴とする。

【 0 0 1 5 】

本発明の音声入力機能付き撮像装置は、前記登録者が前記被写体画像中に複数人写っている場合には該複数人の中の一人の登録者を指定する手段を備え、前記画像合成手段は、指定された登録者の音声を前記テキストデータで該被写体画像中に合成することを特徴とする。

【 0 0 1 6 】

50

本発明の音声入力機能付き撮像装置は、指向性を持つ前記音声入力手段の該指向性を前記被写体画像中の登録者の方向に制御する制御手段を備えることを特徴とする。

【0017】

本発明の音声入力機能付き撮像装置の音声記録方法は、被写体画像を撮像する撮像手段と、該撮像手段が前記被写体画像を撮像するとき同時に音声を取り込む音声入力手段と、該音声をテキストデータ化する音声認識手段とを備える音声入力機能付き撮像装置の音声記録方法において、前記テキストデータを前記被写体画像に合成することで前記音声を記録すること特徴とする。

【0018】

本発明の音声入力機能付き撮像装置の音声記録方法は、前記音声の話者を特定し、特定された前記話者が予め登録された登録者である場合のみ前記テキストデータの合成を行うことを特徴とする。

【0019】

本発明の音声入力機能付き撮像装置の音声記録方法は、前記話者が特定され該話者の音声を前記テキストデータで合成するとき該話者を特定する名前を該テキストデータと一緒に合成することを特徴とする。

【0020】

本発明の音声入力機能付き撮像装置の音声記録方法は、前記被写体画像中に前記登録者の顔が写っているか否かを識別し、前記登録者の顔が写っており且つ該登録者が前記話者のとき前記テキストデータを前記被写体画像中の前記登録者の顔画像からの吹き出し形状で合成することを特徴とする。

【0021】

本発明の音声入力機能付き撮像装置の音声記録方法は、前記被写体画像中に複数の登録者が写っており前記音声が複数人の音声のときは各音声データに対応したテキストデータを各人の画像の近くに吹き出し形状で合成することを特徴とする。

【0022】

本発明の音声入力機能付き撮像装置の音声記録方法は、前記被写体画像中に前記登録者の顔が写っているか否かを識別し、前記登録者の顔が写っておらず且つ該登録者が前記話者のとき前記テキストデータを字幕データとして合成することを特徴とする。

【0023】

本発明の音声入力機能付き撮像装置の音声記録方法は、前記登録者が前記被写体画像中に複数人写っている場合に、該複数人の中の一人の登録者の指定操作に従って、指定された登録者の音声を前記テキストデータで該被写体画像中に合成することを特徴とする。

【発明の効果】

【0024】

本発明によれば、撮像装置で被写体画像を撮影するときに撮影と同時に主要被写体の発する音声をテキストデータ化して被写体画像中に合成するため、使い勝手の優れた撮像装置を提供することが可能となる。また、登録者の場合にのみ音声をテキストデータで画像に合成するため、主要被写体と音声との対応付けの精度が向上する。

【発明を実施するための最良の形態】

【0025】

以下、本発明の一実施形態について、図面を参照して説明する。

【0026】

図1は、本発明の第1実施形態に係る撮像装置の機能ブロック図である。この撮像装置10は、撮影レンズ11と、この撮影レンズ11を通して結像された光像に応じた電気信号を出力する撮像素子12と、撮像素子12のアナログの撮像画像信号をデジタル信号に変換するアナログデジタル変換回路13と、デジタルの撮像画像信号を取り込む画像入力インタフェース14と、画像入力インタフェース14から取り込まれた撮像画像信号を処理する画像処理手段15とを備える。

【0027】

この撮像装置 10 は、更に、撮像装置 10 の全体を統括制御するシステム制御部（CPU）16 と、音声認識手段 17 と、動画像生成手段 18 と、音声入力インタフェース 19 と、表示回路インタフェース 20 と、メモリインタフェース 21 と、メディアインタフェース 22 と、DMA（ダイレクトメモリアクセス）コントローラ 23 と、オーディオインタフェース 24 と、これらを相互接続するバス 25 とを備える。

【0028】

音声認識手段 17 には、話者特定手段 27 と、この特定手段 27 に接続された音声データベース 28 とが接続され、音声入力インタフェース 19 には、アナログデジタル変換回路 29 を介して音声入力手段であるマイク 31 が接続される。

【0029】

表示回路インタフェース 20 には撮像装置 10 の背面等に取り付けられる液晶等の表示回路 32 が接続され、メモリインタフェース 21 には撮像装置 10 のメインメモリとなるフレームメモリ 33 が接続され、メディアインタフェース 22 には、着脱自在の外部メモリ（記録メディア）34 が接続され、オーディオインタフェース 24 にはスピーカ等のオーディオ出力手段 35 が接続される。

【0030】

CPU 16 は、音声認識手段 17 が音声認識しテキストデータ化した字幕データや吹き出しデータを、動画生成手段 18 や表示回路インタフェース 20 に出力する様になっている。また、CPU 16 には、シャッターボタンなどの操作系 36 の信号が入力される。CPU 16 が音声認識手段 17 の機能を実行する構成としても良い。

【0031】

図 1 に示す音声データベース 28 には、例えば、撮像装置 10 のユーザが良く撮る家族等の各人の音声、各人の名前と対応付けて登録されている。

【0032】

図 2 は、図 1 に示す撮像装置 10 が実行する処理手順を示すフローチャートである。動画像の記録が開始される（ステップ S1）と、撮像素子 12 から被写体画像が取り込まれると共に、マイク 31 から音声を取り込まれる（ステップ S2）。音声認識手段 17 はこの音声を解析すると共にデータベース 28 を参照し（ステップ S3）、この音声データベース 28 に登録されている音声であるか否かを判定する（ステップ S4）。

【0033】

マイク 31 から取り込んだ音声登録者の音声でないと判定した場合には、そのままステップ S8 に進む。登録者の音声であると判定した場合には、ステップ S4 からステップ S5 に進み、話者特定手段 27 は、データベース 28 から話者が誰であることを特定し、次のステップ S6 で、音声認識手段 17 は、その音声をテキストデータに変換する。

【0034】

そして、次のステップ S7 で、CPU 16 は、テキストデータ化された音声データを動画生成手段 18 と表示回路インタフェース 20 とに渡し、動画生成手段 18 は、図 3 に示す様に、動画像上に、特定された話者の名前“A”と、その音声データの字幕データとを合成して、記録メディア 34 に書き込み、ステップ S8 に進む。

【0035】

尚、音声の字幕データを合成した画像データを記録メディア 34 に書き込むのではなく、無線で外部記憶手段に伝送する構成としても良い。これは以下の実施形態でも同様である。

【0036】

ステップ S8 では、動画記録が終了したか否かを判定し、終了した場合にはこの図 2 の処理を終了し、動画記録が終了しない場合にはステップ S2 に戻り、ステップ S1～ステップ S7 の処理を繰り返す。

【0037】

表示回路 32 には、動画像データが表示されるが、このとき、字幕データと話者の名前とが重ねて表示される。

10

20

30

40

50

【 0 0 3 8 】

尚、音声を字幕データとして記録する部分についてのみ説明したが、マイク 3 1 から取り込んだ音声のままのデータでも記録することはいうまでもない。音声のままのデータを記録する場合には、音声データベース上に登録されている話者であるか否かに関係なく全て記録する。これは以下の実施形態でも同様である。

【 0 0 3 9 】

図 4 は、本発明の第 2 実施形態に係る撮像装置の機能ブロック図である。この撮像装置 4 0 は、図 1 に示す撮像装置 1 0 と殆どの機能が重複し、静止画像撮像用である点のみ異なる。このため、重複する部分には同一符号を付してその説明は省略し、異なる部分についてのみ説明する。

10

【 0 0 4 0 】

本実施形態の撮像装置 4 0 は、静止画像撮像専用であり、図 1 に示した動画生成手段 1 8 が設けられておらず、CPU 1 6 は、字幕データを画像処理部 1 5 に渡す様になっている。

【 0 0 4 1 】

図 5 は、本実施形態に係る撮像装置が実行する処理手順を示すフローチャートである。先ず、2 段シャッターボタンが半押し状態になっているか否かを判定する（ステップ S 1 1）。半押し状態になっていなければ繰り返しステップ S 1 1 を実行し、半押し状態になっている場合には、音声をマイク 3 1 から取り込む（ステップ S 1 2）。次に、シャッターボタンが全押しになかったか否かを判定し（ステップ S 1 3）、全押しになっていない場合にはステップ S 1 2 を繰り返し実行する。

20

【 0 0 4 2 】

シャッターボタンが全押しになった場合にはステップ S 1 3 からステップ S 1 4 に進み、マイク 3 1 から取り込んだ音声データを記録し、次に、音声データベース 2 8 を参照する（ステップ S 1 5）。また、シャッターボタンが全押しされた場合には、ステップ S 1 4 以下の処理と並行して、被写体画像の撮像処理が行われる。

【 0 0 4 3 】

そして、音声認識手段 1 7 は、この音声データベース 2 8 に登録されている音声であるか否かを判定し（ステップ S 1 6）、マイク 3 1 から取り込んだ音声に登録者の音声でないと判定した場合には、この図 5 の処理を終了する。登録者の音声であると判定した場合には、ステップ S 1 6 からステップ S 1 7 に進み、話者特定手段 2 7 は、データベース 2 8 から話者が誰であるかを特定し、次のステップ S 1 8 で、音声認識手段 1 7 は、その音声をテキストデータに変換する。

30

【 0 0 4 4 】

そして、次のステップ S 1 9 で、CPU 1 6 は、テキストデータ化された音声データを画像処理部 1 5 と表示回路インタフェース 2 0 とに渡し、画像処理部 1 5 は、図 3 に示す様に、静止画像上に、特定された話者の名前“ A ”と、その音声データの字幕データとを合成して、記録メディア 3 4 に書き込み、図 5 の処理を終了する。表示回路 3 2 には、静止画像データが表示されるが、このとき、字幕データと話者の名前とが重ねて表示される。

40

【 0 0 4 5 】

図 6 は、本発明の第 3 実施形態に係る撮像装置 5 0 の機能ブロック図である。本実施形態の撮像装置 5 0 は、図 1 に示す撮像装置 1 0 と殆どの機能が重複するため、同一機能ブロックには同一符号を付してその説明は省略し、異なる部分についてのみ説明する。

【 0 0 4 6 】

本実施形態に係る撮像装置 5 0 は、画像処理部 1 5 の処理結果を受け取り撮像画像中に人間の「顔」が存在するか否かを検出すると共に、検出した「顔」が登録されている特定者の「顔」であるか否かを検出する顔検出／顔認識処理手段 4 1 と、この顔検出／顔認識処理手段 4 1 が人間の「顔」として検出するとき使用する顔検出用データ及び登録者（例えば家族等）の顔画像を登録者の名前と対応付けて格納した顔データベース 4 2 とを備え

50

る。顔検出／顔認識処理手段４１はバス２５に接続されている。

【００４７】

図７は、本実施形態に係る撮像装置が実行する処理手順を示すフローチャートである。動画像の記録が開始される（ステップＳ２１）と、撮像素子１２から被写体画像が取り込まれると共に、マイク３１から音声を取り込まれる（ステップＳ２２）。音声認識手段１７はこの音声を解析すると共にデータベース２８を参照し（ステップＳ２３）、この音声

【００４８】

マイク３１から取り込んだ音声登録者の音声でないと判定した場合には、ステップＳ２２に戻って次の音声入力を待機し、登録者の音声であると判定した場合には、ステップＳ２４からステップＳ２５に進み、話者特定手段２７は、データベース２８から話者が誰であるかを特定すると共に、音声認識手段１７はその音声をテキストデータに変換する。

【００４９】

次に、顔検出／顔認識処理手段４１は、撮像された画像の中に人間の「顔」が存在するか否かを解析し（ステップＳ２６）、「顔」が検出されなかった場合には、ステップＳ２５でテキストデータ化された音声データを図３に示す様に字幕として合成し（ステップＳ３０）、ステップＳ３１に進む。

【００５０】

人間の「顔」が検出された場合にはステップＳ２６からステップＳ２７に進んで顔データベース４２を参照する。そして、次にステップＳ２８では、検出された「顔」が登録者の「顔」であるか否か、及び、登録者である場合にはその登録者が、ステップＳ２５でテキストデータ化した音声の特定話者であるか否かを判定する（ステップＳ２８）。

【００５１】

登録者の顔で無い場合、あるいは特定話者と一致しない登録者の顔である場合には、上記のステップＳ３０に進み、テキストデータ化された音声データを字幕として合成する。

【００５２】

ステップＳ２８の判定の結果、登録者の顔であり、且つ特定話者と一致すると判定した場合には、次にステップＳ２９に進み、テキストデータ化した音声データを、特定話者と一致する撮像画像中の登録者の顔画像の近くに、図８に示す様に、吹き出し形状で合成し、ステップＳ３１に進む。

【００５３】

ステップＳ３１では、動画記録が終了したか否かを判定し、終了した場合にはこの図７の処理を終了し、動画記録が終了しない場合にはステップＳ２２に戻り、ステップＳ２２～ステップＳ３０の処理を繰り返す。

【００５４】

吹き出し形状で音声のテキストデータを画像に合成し表示する場合、音声データが画像データと重なることになる。このため、下の画像が見づらくなならないように、吹き出し位置、字の大きさ、透明度等を設定可能にするのが好ましい。

【００５５】

この様にすることで、動画記録時に画面上に複数人が撮影されている場合でも自動的に登録者、話者を特定して各人に対応した吹き出し合成を行うため、画像を見れば誰が喋った内容か直ぐ分かるようになる。

【００５６】

尚、図７では、音声による話者特定とテキストデータ化を先に行い、その後に顔検出、顔認識を行ったが、これを逆に行っても、また同時並行的に行っても良いことは言うまでもない。

【００５７】

図９は、本発明の第４実施形態に係る撮像装置６０の機能ブロック図である。本実施形態の撮像装置６０は、図６に示す撮像装置５０と殆どの機能が重複するため、同一機能ブロックには同一符号を付してその説明は省略し、異なる部分についてのみ説明する。

【 0 0 5 8 】

本実施形態の撮像装置 6 0 は、マイク 3 1 が指向性を持ち、且つその指向性の制御（集音範囲，集音方向，集音距離（感度）等の制御）が可能なマイクであり、マイク制御手段 4 3 がマイク 3 1 を C P U 1 6 からの指示により制御する構成になっている。また、図 6 に示す話者特定手段 2 7 と音声データベース 2 8 とがこの撮像装置 6 0 には設けられていない。

【 0 0 5 9 】

図 1 0 は、本実施形態に係る撮像装置が実行する処理手順を示すフローチャートである。動画像の記録が開始される（ステップ S 4 1 ）と、撮像素子 1 2 から被写体画像が取り込まれる。次のステップ S 4 2 では、撮像画像中に登録者の「顔」が存在するか否かを判定し、登録者の顔が検出されるまでステップ S 4 2 を繰り返し実行する。

10

【 0 0 6 0 】

画面中に登録者の顔が検出された場合には、ステップ S 4 2 からステップ S 4 3 に進み、今度は 1 画面の中に複数の登録者の顔が存在するか否かを判定する。複数の登録者の顔が検出された場合には、ステップ S 4 3 からステップ S 4 4 に進み、ユーザの顔選択処理を行った後、ステップ S 4 5 に進む。1 画面の中に一人の登録者の顔しか検出されない場合にはステップ S 4 4 を飛び越してステップ S 4 5 に進む。

【 0 0 6 1 】

図 1 1 は、1 画面の中に 3 人の人間が撮像されている状態を示しており、そのうちの二人の顔が名前 “ A ” “ B ” の登録者であり、もう一人が登録者でない人の場合を示している。図 9 の顔検出 / 顔認識処理手段 4 1 は、検出した顔部分を矩形枠で示すため、図 1 0 のステップ S 4 4 では、ユーザは、登録者 A , B のいずれか一方を操作系 3 6 のボタン操作により指定することになる。図 1 1 に示す例では、登録者 A が指定されたため矩形枠を二重枠で表示したところを示している。

20

【 0 0 6 2 】

ステップ S 4 5 では、画面中の一人の登録者あるいはステップ S 4 4 で指定された登録者に対して、マイク制御を行う。この登録者の撮像画像から撮像装置と登録者との間の距離，登録者が居る方向がズーム倍率等で判別できるため、この登録者が喋る音声を精度良く集音できるように、C P U 1 6 はマイク 3 1 の指向性制御を行う。

【 0 0 6 3 】

次のステップ S 4 6 では、マイク 3 1 から集音した音声データを取り込み、ステップ S 4 7 でこれをテキストデータ化し、ステップ S 4 8 でテキストデータした字幕（勿論、吹き出しでも良い。）を図 1 1 に示す様に画像上に合成し、ステップ S 4 9 に進む。

30

【 0 0 6 4 】

ステップ S 4 9 では、動画記録が終了したか否かを判定し、終了した場合にはこの図 1 0 の処理を終了し、動画記録が終了していない場合にはステップ S 4 2 に戻り、ステップ S 4 3 ~ ステップ S 4 8 を繰り返し実行する。この繰り返し時には、ステップ S 4 4 のユーザ指定は、何らかのユーザによるボタン操作が無い限り実行しない構成とすることで、指定した登録者をずっと追って集音することが可能となる。

【 0 0 6 5 】

尚、図 1 0 の実施形態では、ユーザによる顔選択処理（ステップ S 4 4 ）を動画記録時に行ったが、動画記録前に予め選択する構成としても良い。このステップ S 4 4 における選択時に、図 1 1 に示したように、選択対象者の顔を二重枠で区別したり、人物名を画像中表示することで、選択が容易且つ確実にできる様にするのが良い。

40

【 0 0 6 6 】

この様に、本実施形態によれば、ユーザが記録したい被写体の音声のみを取り込んで字幕化することが可能となる。

【 0 0 6 7 】

以上述べた様に、本発明の各実施形態によれば、画像の撮像中に、音声をテキストデータ化し、該当する被写体画像に対応付けて合成するため、話者と音声データ（テキストデ

50

ータ)との対応付けの精度が向上すると共に、後で編集する手間が省け、デジタルカメラ等の撮像装置の使い勝手が向上する。

【産業上の利用可能性】

【0068】

本発明に係る撮像装置は、話者と音声データとの対応付けの精度が向上し字幕や吹き出しとして音声のテキストデータを撮像中の画像に合成できるため、デジタルスチルカメラやビデオカメラ等に適用すると有用である。

【図面の簡単な説明】

【0069】

【図1】本発明の第1実施形態に係る撮像装置の機能ブロック図である。

10

【図2】図1に示す撮像装置の処理手順を示すフローチャートである。

【図3】図2の処理手順により音声データが字幕表示された画像を示す図である。

【図4】本発明の第2実施形態に係る撮像装置の機能ブロック図である。

【図5】図4に示す撮像装置の処理手順を示すフローチャートである。

【図6】本発明の第3実施形態に係る撮像装置の機能ブロック図である。

【図7】図6に示す撮像装置の処理手順を示すフローチャートである。

【図8】図7に示す処理手順により音声データが吹き出し表示された画像を示す図である。

。

【図9】本発明の第4実施形態に係る撮像装置の機能ブロック図である。

20

【図10】図9に示す撮像装置の処理手順を示すフローチャートである。

【図11】図10に示す処理手順により音声データが字幕表示された画像を示す図である。

。

【符号の説明】

【0070】

10, 40, 50, 60 撮像装置

12 撮像素子

15 画像処理手段(画像合成手段)

16 CPU

17 音声認識手段

18 動画生成手段(画像合成手段)

30

27 話者特定手段

28 音声データベース

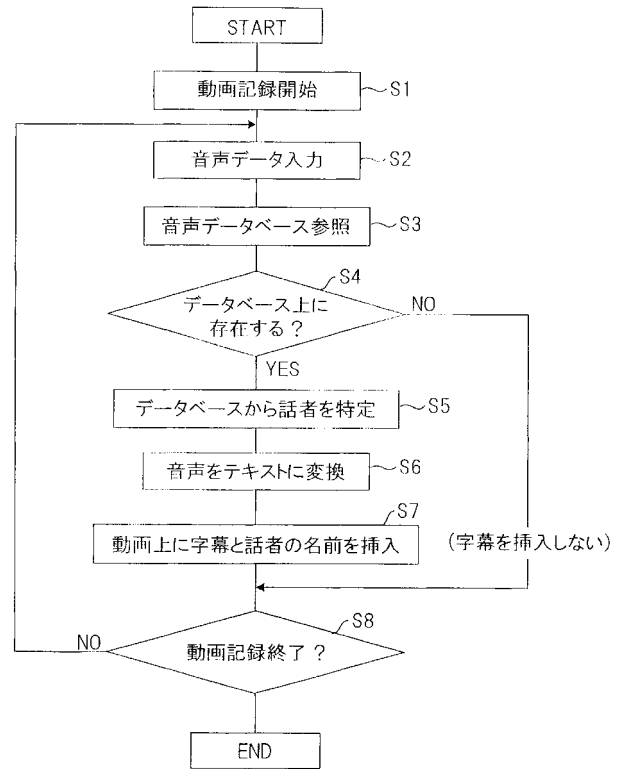
31 マイク(音声入力手段)

41 顔検出/顔認識処理手段

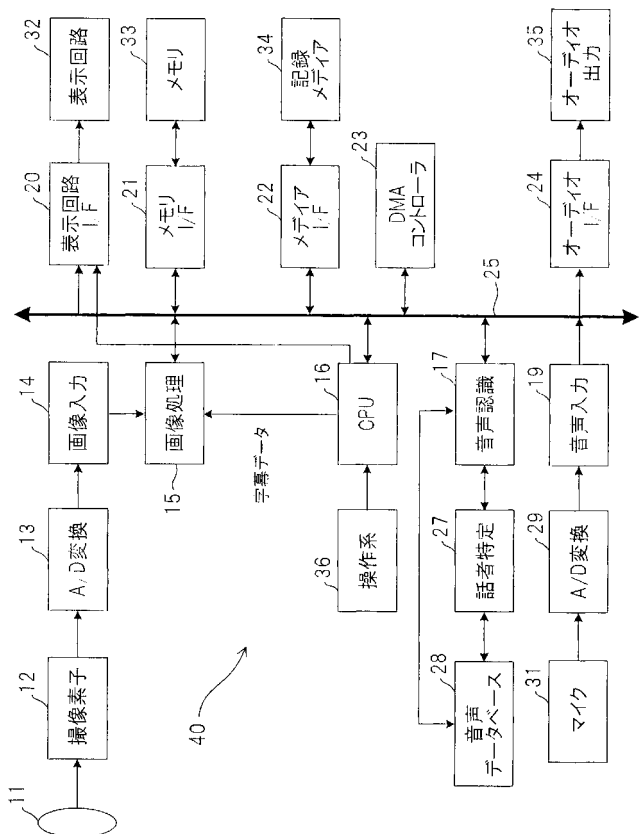
42 顔データベース

43 マイク制御手段

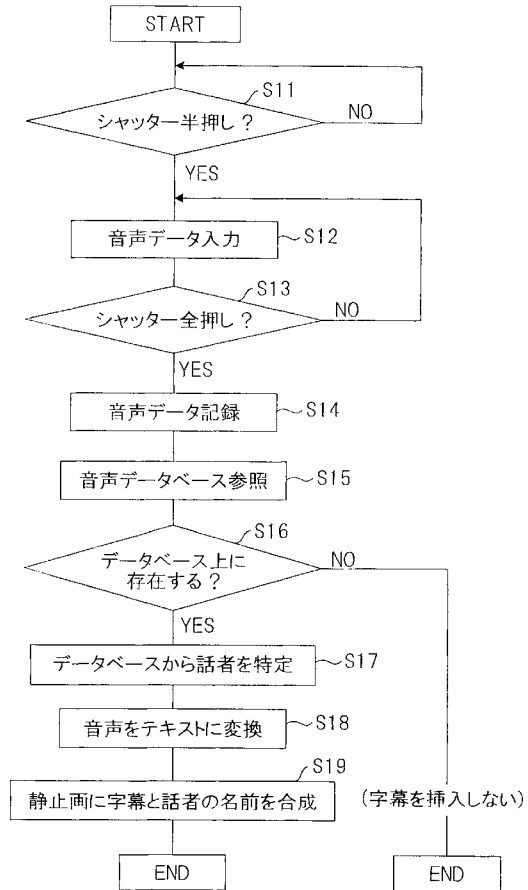
【 図 2 】



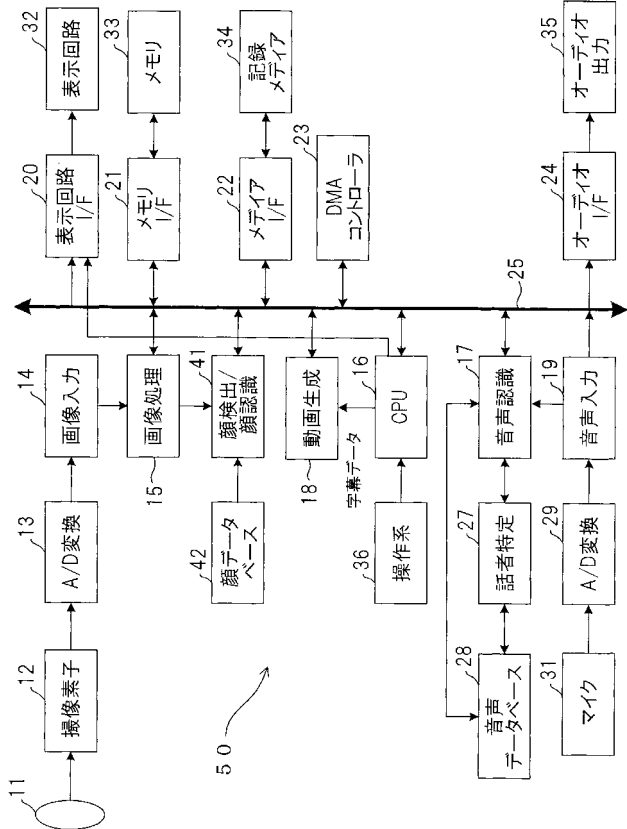
【 図 4 】



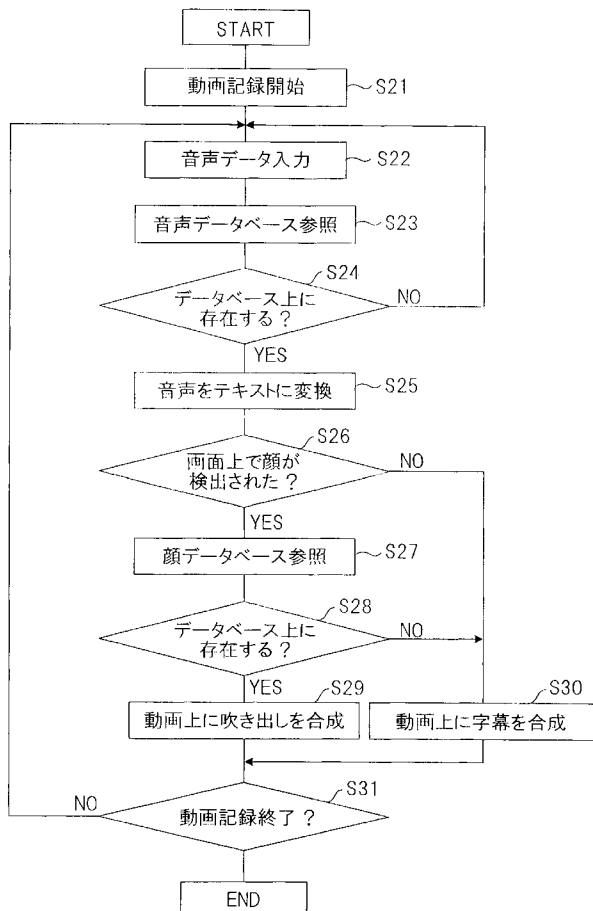
【図 5】



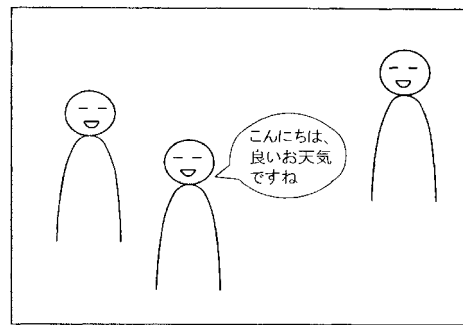
【図 6】



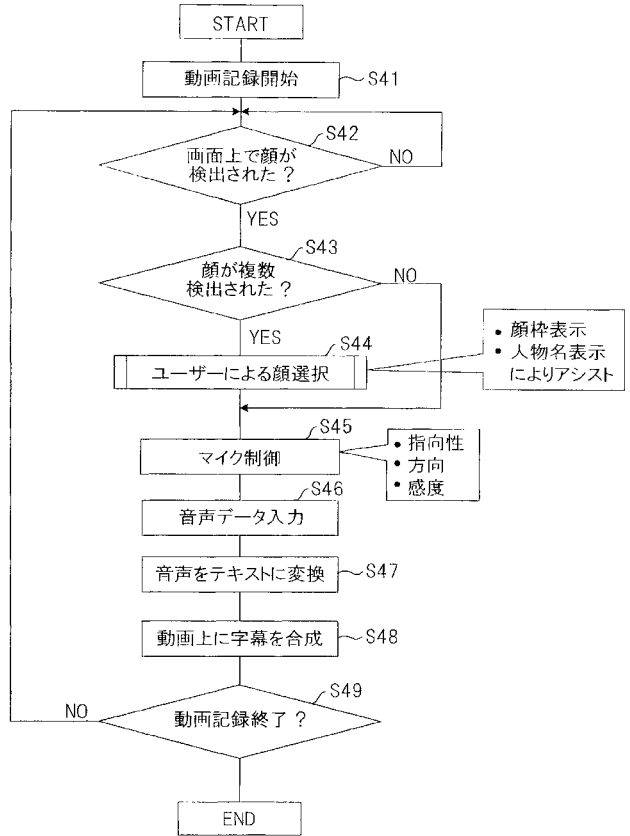
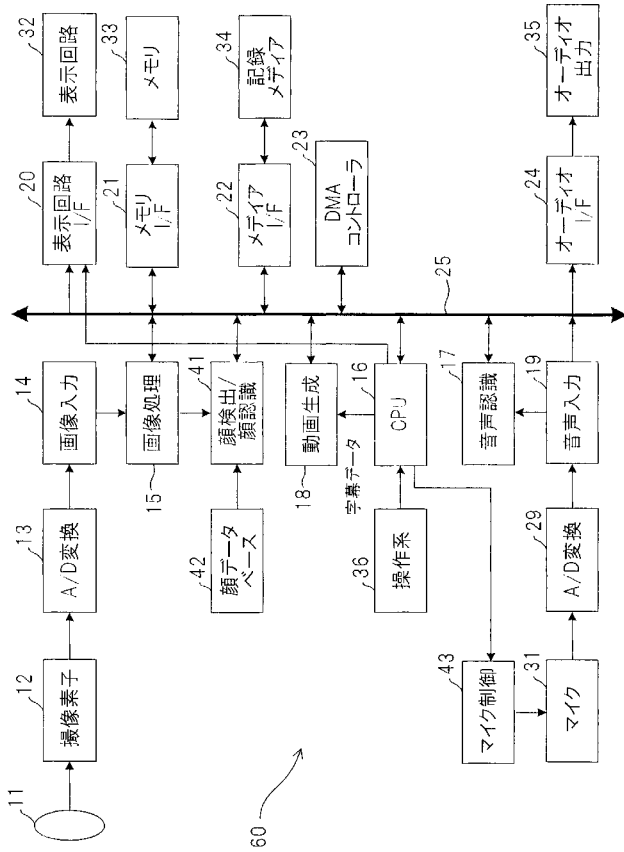
【図 7】



【図 8】



【 図 1 0 】



【 図 1 1 】

