

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4775066号
(P4775066)

(45) 発行日 平成23年9月21日(2011.9.21)

(24) 登録日 平成23年7月8日(2011.7.8)

(51) Int.Cl.	F I	
HO4N 5/232 (2006.01)	HO4N 5/232	Z
HO4N 5/91 (2006.01)	HO4N 5/91	R
HO4N 1/387 (2006.01)	HO4N 1/387	
GO6T 1/00 (2006.01)	GO6T 1/00	340A
G1OL 15/00 (2006.01)	G1OL 15/00	200G
請求項の数 11 (全 17 頁) 最終頁に続く		

(21) 出願番号 特願2006-86787(P2006-86787)
 (22) 出願日 平成18年3月28日(2006.3.28)
 (65) 公開番号 特開2007-266793(P2007-266793A)
 (43) 公開日 平成19年10月11日(2007.10.11)
 審査請求日 平成20年12月19日(2008.12.19)

(73) 特許権者 000001443
 カシオ計算機株式会社
 東京都渋谷区本町1丁目6番2号
 (74) 代理人 100088100
 弁理士 三好 千明
 (72) 発明者 栗山 祐司
 東京都羽村市栄町3丁目2番1号 カシオ
 計算機株式会社羽村技術センター内
 審査官 田村 誠治

最終頁に続く

(54) 【発明の名称】 画像加工装置

(57) 【特許請求の範囲】

【請求項1】

画像を取得する画像取得手段と、
 この画像取得手段により取得された画像から人の口を識別する画像識別手段と、
 この画像取得手段に対応して音を取得する音取得手段と、
 この音取得手段により取得された音を認識し、この認識した音を表示データに変換する音認識手段と、

前記画像識別手段により識別された前記人の口の数が複数ある時は、識別されたこれら複数の人の口から動きを検出し、前記音認識手段により認識された音に対応する動きのある人の口を判定する判定手段と、

前記判定手段による判定に基づいて、前記画像識別手段により識別された前記人の口に対応する位置に、前記音認識手段により変換された表示データを、合成する画像合成手段とを備え、

前記画像合成手段は、前記画像識別手段により識別された前記人の口が一つの場合、前記判定手段による判定に基づかずに、この一つの人の口に対応する位置に、前記音認識手段により変換された表示データを合成することを特徴とする画像加工装置。

【請求項2】

前記判定手段により検出された画像中における人の口の動きに基づき、前記表示データを訂正処理することを特徴とする請求項1記載の画像加工装置。

【請求項3】

前記画像合成手段は、前記表示データを吹き出しとともに合成することを特徴とする請求項 1 又は 2 に記載の画像加工装置。

【請求項 4】

前記画像識別手段は、識別した前記人の口の人物が誰であることを識別し、

前記音認識手段は、前記画像識別手段が識別した人物に応じて、変換する表示データの表示形態を変化させることを特徴とする請求項 1 から 3 にいずれか記載の画像加工装置。

【請求項 5】

前記画像識別手段は、更に前記人の口の人物の種別を識別し、

前記音認識手段は、前記画像識別手段が識別した人物の種別に応じて、変換する表示データの表示形態を変化させることを特徴とする請求項 1 から 4 にいずれか記載の
画像加工装置

10

【請求項 6】

前記画像識別手段は、更に前記画像の内容を識別し、

この画像識別手段が識別した画像の内容に応じて表示データを生成する内容表示データ生成手段を更に備え、

前記画像合成手段は、前記内容表示データ生成手段により生成された表示データを前記画像中に合成することを特徴とする請求項 1 から 5 にいずれか記載の画像加工装置。

【請求項 7】

前記画像取得手段は、前記画像とともに当該画像に付随する情報を取得し、

この画像取得手段が取得した前記情報に基づき、表示データを生成する情報表示データ生成手段を更に備え、

前記画像合成手段は、前記情報表示データ生成手段により生成された表示データを前記画像中に合成することを特徴とする請求項 1 から 6 にいずれか記載の画像加工装置。

20

【請求項 8】

前記画像合成手段は、前記画像識別手段により前記人の口の識別ができなかった場合、前記表示データを前記画像中における背景部分に合成することを特徴とする請求項 1 から 7 にいずれか記載の画像加工装置。

【請求項 9】

前記画像合成手段は、前記表示データを前記画像中における識別された人の口の人物と重ならない位置に合成することを特徴とする請求項 1 から 8 にいずれか記載の画像加工装置。

30

【請求項 10】

前記画像合成手段により前記表示データが合成された画像を記録する記録手段及び / 又は前記画像合成手段により前記表示データが合成された画像を表示する表示手段を更に備えることを特徴とする請求項 1 から 9 にいずれか記載の画像加工装置。

【請求項 11】

画像加工装置が備えるコンピュータを、

画像を取得する画像取得手段と、

この画像取得手段により取得された画像から人の口を識別する画像識別手段と、

この画像取得手段に対応して音を取得する音取得手段と、

この音取得手段により取得された音を認識し、この認識した音を表示データに変換する音認識手段と、

40

前記画像識別手段により識別された前記人の口の数が複数ある時は、識別されたこれら複数の人の口から動きを検出し、前記音認識手段により認識された音に対応する動きのある人の口を判定する判定手段と、

前記判定手段による判定に基づいて、前記画像識別手段により識別された前記人の口に対応する位置に、前記音認識手段により変換された表示データを、合成する画像合成手段として機能させ、

前記画像合成手段は、前記画像識別手段により判別された前記人の口が一つの場合、前記判定手段による判定に基づかずに、この一つの人の口に対応する位置に、前記音認識手

50

段により変換された表示データを合成することを特徴とする画像加工プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音認識結果に基づく表示データを生成して画像に合成する画像加工装置に関する。

【背景技術】

【0002】

従来の画像加工装置として特許文献1記載のものが知られている。この画像加工装置は、撮影の際に被写体から発せられた音声を検出して認識し、文字コードに変換する。また、音声の検出はステレオマイクで行い、このステレオマイクで検出した音声に基づき三角法を用いて音声の発生位置を推測する。そして、画像中の推測した音声発生位置に、前記変換された文字コードに基づく文字イメージと吹き出しとからなる表示データを合成する。

10

【特許文献1】特開平11-55614号公報

【発明の開示】

【発明が解決しようとする課題】

【0003】

しかしながら、前述した従来の画像加工装置にあっては、画像中における文字イメージ等の合成位置を、当該画像から直接的に得ることなく、画像との関係においては間接的な要素である音声に基づき推測する。このため、推測された文字イメージ等の合成位置が実際に音声を発生している被写体と一致しない場合が生じ、画像中の音声発生被写体に対応する位置に精度よく文字イメージ等を合成することができない。

20

【0004】

本発明は、かかる従来の課題に鑑みてなされたものであり、画像中の適正位置に音に対応する表示データを合成することのできる画像加工装置を提供することを目的とする。

【課題を解決するための手段】

【0005】

前記課題を解決するため請求項1記載の発明に係る画像加工装置にあっては、画像を取得する画像取得手段と、この画像取得手段により取得された画像から人の口を識別する画像識別手段と、この画像取得手段に対応して音を取得する音取得手段と、この音取得手段により取得された音を認識し、この認識した音を表示データに変換する音認識手段と、前記画像識別手段により識別された前記人の口の数が複数ある時は、識別されたこれら複数の人の口から動きを検出し、前記音認識手段により認識された音に対応する動きのある人の口を判定する判定手段と、前記判定手段による判定に基づいて、前記画像識別手段により識別された前記人の口に対応する位置に、前記音認識手段により変換された表示データを、合成する画像合成手段とを備え、前記画像合成手段は、前記画像識別手段により識別された前記人の口が一つの場合、前記判定手段による判定に基づかずに、この一つの人の口に対応する位置に、前記音認識手段により変換された表示データを合成する。

30

【0011】

また、請求項2記載の発明に係る画像加工装置にあっては、前記判定手段により検出された画像中における人の口の動きに基づき、前記表示データを訂正処理する。

40

【0012】

また、請求項3記載の発明に係る画像加工装置にあっては、前記画像合成手段は、前記表示データを吹き出しとともに合成する。

【0013】

また、請求項4記載の発明に係る画像加工装置にあっては、前記画像識別手段は、識別した前記人の口の人物が誰であることを識別し、前記音認識手段は、前記画像識別手段が識別した人物に応じて、変換する表示データの表示形態を変化させる。

【0014】

50

また、請求項 5 記載の発明に係る画像加工装置にあっては、前記画像識別手段は、更に前記人の口の人物の種別を識別し、前記音認識手段は、前記画像識別手段が識別した人物の種別に応じて、変換する表示データの表示形態を変化させる。

【 0 0 1 5 】

また、請求項 6 記載の発明に係る画像加工装置にあっては、前記画像識別手段は、更に前記画像の内容を識別し、この画像識別手段が識別した画像の内容に応じて表示データを生成する内容表示データ生成手段を更に備え、前記画像合成手段は、前記内容表示データ生成手段により生成された表示データを前記画像中に合成する。

【 0 0 1 6 】

また、請求項 7 記載の発明に係る画像加工装置にあっては、前記画像取得手段は、前記画像とともに当該画像に付随する情報を取得し、この画像取得手段が取得した前記情報に基づき、表示データを生成する情報表示データ生成手段を更に備え、前記画像合成手段は、前記情報表示データ生成手段により生成された表示データを前記画像中に合成する。

【 0 0 1 7 】

また、請求項 8 記載の発明に係る画像加工装置にあっては、前記画像合成手段は、前記画像識別手段により前記人の口の識別ができなかった場合、前記表示データを前記画像中における背景部分に合成する。

【 0 0 1 8 】

また、請求項 9 記載の発明に係る画像加工装置にあっては、前記画像合成手段は、前記表示データを前記画像中における識別された人の口の人物と重ならない位置に合成する。

【 0 0 1 9 】

また、請求項 10 記載の発明に係る画像加工装置にあっては、前記画像合成手段により前記表示データが合成された画像を記録する記録手段及び/又は前記画像合成手段により前記表示データが合成された画像を表示する表示手段を更に備える。

【 0 0 2 0 】

また、請求項 11 記載の発明に係る画像加工プログラムにあっては、画像加工装置が備えるコンピュータを、画像を取得する画像取得手段と、この画像取得手段により取得された画像から人の口を識別する画像識別手段と、この画像取得手段に対応して音を取得する音取得手段と、この音取得手段により取得された音を認識し、この認識した音を表示データに変換する音認識手段と、前記画像識別手段により識別された前記人の口の数が複数ある時は、識別されたこれら複数の人の口から動きを検出し、前記音認識手段により認識された音に対応する動きのある人の口を判定する判定手段と、前記判定手段による判定に基づいて、前記画像識別手段により識別された前記人の口に対応する位置に、前記音認識手段により変換された表示データを、合成する画像合成手段として機能させ、前記画像合成手段は、前記画像識別手段により判別された前記人の口が一つの場合、前記判定手段による判定に基づかずに、この一つの人の口に対応する位置に、前記音認識手段により変換された表示データを合成する。

【発明の効果】

【 0 0 2 1 】

本発明によれば、画像中の音発生被写体を該画像に基づいて識別することから、画像から直接的に音発生被写体を識別することにより、精度よく音発生被写体を識別することができる。よって、精度よく画像中の音発生被写体に対応する位置に、音を変換した表示データを表示することが可能となる。

【発明を実施するための最良の形態】

【 0 0 2 2 】

以下、本発明の一実施の形態を図に従って説明する。図 1 は、本発明の一実施の形態を適用したデジタルカメラ 10 の回路構成を示すブロック図であり、このデジタルカメラ 10 は、後述する AF 機能とともに AE、AWB 等の一般的な機能をも有するものである。すなわち、レンズブロック 11 には、ズームレンズ、フォーカスレンズ等の光学系、及び光学系を駆動するための駆動機構が含まれており、前記光学系は、駆動機構に設けられて

10

20

30

40

50

いるモーター12によって光軸方向に駆動される。なお、本実施の形態において、前記AFは、フォーカスレンズを光軸方向に移動させながら、各位置で撮像した画像のAF評価値(コントラスト値)を検出し、AF評価値のピーク位置を合焦位置とするコントラスト検出方式である。

【0023】

デジタルカメラ10全体を制御するCPU13には、バス14及びタイミング発生器(TG:Timing Generator)15を介してモータードライバ16が接続されており、モータードライバ16は、CPU13の命令に従いタイミング発生器15が発生するタイミング信号に基づき、モーター12を駆動する。なお、ストロボ17もタイミング発生器15が発生するタイミング信号により駆動される。

10

【0024】

また、このデジタルカメラ10は撮像素子としてCCD18を有している。CCD18は、レンズブロック11の光軸上に配置されており、被写体は、レンズブロック11によってCCD18の受光面に結像される。CCD18は、CPU13の命令に従いタイミング発生器15が生成するタイミング信号に基づき垂直及び水平ドライバ19によって駆動され、被写体の光学像に応じたアナログの撮像信号をユニット回路20に出力する。ユニット回路20は、CCD18の出力信号に含まれるノイズを相関二重サンプリングによって除去するCDS回路や、ノイズが除去された撮像信号をデジタル信号に変換するA/D変換器等から構成され、デジタルに変換した撮像信号を画像処理部21へ出力する。

【0025】

20

画像処理部21は、入力した撮像信号に対しペダスタルクランプ等の処理を施し、それを輝度(Y)信号及び色差(UV)信号に変換するとともに、オートホワイトバランス、輪郭強調、画素補間などの画品質向上のためのデジタル信号処理を行う。画像処理部21で変換されたYUVデータは順次SDRAM22に格納されるとともに、RECスルー・モードでは1フレーム分のデータ(画像データ)が蓄積される毎にビデオ信号に変換され、バックライト(BL)24を備える液晶モニタ(LCD)23へ送られてスルー画像として画面表示される。

【0026】

そして、スチル撮影モードにおいては、シャッターキー操作をトリガとして、CPU13は、CCD18、垂直及び水平ドライバ19、ユニット回路20、及び画像処理部21に対してスルー画撮影モード(RECスルー・モード)から静止画撮影モードへの切り替えを指示し、この静止画撮影モードによる撮影処理により得られ、SDRAM22に一時記憶された画像データは、CPU13により圧縮され、最終的には所定のフォーマットの静止画ファイルとして外部メモリ25に記録される。また、ムービー録画モードにおいては、1回目のシャッターキーと2回目のシャッターキー操作との間に、SDRAM22に順次記憶される複数の画像データがCPU13により順次圧縮されて、圧縮動画データが生成され動画ファイルとして外部メモリ25に記録される。この外部メモリ25に記録された静止画ファイル及び動画ファイルは、PLAY・モードにおいてユーザーの選択操作に応じてCPU13に読み出されるとともに伸張され、YUVデータとしてSDRAM22に展開された後、液晶モニタ(LCD)23に表示される。

30

40

【0027】

フラッシュメモリ26には、CPU13に前記各部を制御させるための各種のプログラム、例えばAE、AF、AWB制御用のプログラムや、さらには、後述するフローチャートに示す処理を実行するためのプログラム等の各種のプログラムが格納されている。

【0028】

また、デジタルカメラ10は、電源スイッチ、モード選択キー、シャッターキー、ズームキー、後述するピント枠を手動選択するためのピント枠選択キー等の複数の操作キー及びスイッチを含むキー入力部(KEY)27、ニッケル水素電池等の充電可能なバッテリー28、このバッテリー28の電力を各部に供給するための電源制御回路29、及びこれらを制御するマイコン30を有している。マイコン30は、キー入力部27における前記

50

操作キーの操作の有無を定常的にスキャンしており、ユーザーによっていずれかの操作キーが操作されると、その操作内容に応じた操作信号をCPU13へ送る。なお、シャッターキーは、半押しと全押しとが可能な所謂ハーフシャッター機能を有するものである。

【0029】

また、このデジタルカメラ10は、前記ムービー録画モードにおいて、周囲音を記録する録音機能を備えており、CPU13には、音声処理回路を有する音声チップ32を介して、スピーカ(SP)33と、マイクロホン(MIC)34とが接続されている。音声チップ32は、ムービー録画モード時には、マイクロホン34から入力された音声波形を処理して、音声波形データをCPU13に入力する。そして、CPU13は、ムービー録画モードにおいて1回目と2回目のシャッターキー操作間に、音声チップ32から入力された音声波形データを圧縮し、この圧縮周囲音データと前記圧縮動画データとを含む音声付き動画ファイルを生成して外部メモリ25に記録する。この外部メモリ25に記録された音声付き動画ファイルは、PLAY・モードにおいて動画データが再生される際に、周囲音データが音声チップ32で音声波形に変換されてスピーカ33により再生される。

10

【0030】

さらに、バス14にはGPS35が接続されており、前記フラッシュメモリ26には前記プログラム等とともに地図データが記憶されている。したがって、CPU13はGPS35により検出された現在位置の緯度・経度と、フラッシュメモリ26内の地図データとに現在位置の地名を取得することが可能である。フラッシュメモリ26には、音声をテキストデータに変換するための音声-テキストデータ変換テーブルや、音声以外の音を擬音表示データ(例えば、クラッカーの破裂音を擬音表示データ「パン」、自動車の音を擬音表示データ「ブー」)に変換する音-表示データ変換テーブルが記憶されている。また、画像の動きを擬音表示するための「ビュー」や表情を擬音表示するための「プン」「ニコニコ」等の表示データ、あるいは画像の弧の動きを強調表示するための「((」等、暑さや寒さや擬音表示するための「ジリジリ」「ヒュー」、汗マーク等の画像内容-表示データ変換テーブル、「もうかりまっか」「ぼちぼちでんな」等の複数の慣用語からいずれかをランダムに選択するためのランダムテキストデータ、口の動きをテキストデータに変換するための口の動き-テキストデータ変換テーブル、テキストデータを対応する方言テキストデータに変換するための方言変換テーブル等が記憶されている。

20

【0031】

加えて、前記フラッシュメモリ26には、図2に示す被写体種別判定テーブル261が格納されている。被写体種別判定テーブル261には、「人」、「人の口」、・、「自動車」・・・等の被写体となり得る被写体種別毎にその画像の特徴を示す特徴量データDが記憶されている。さらに、被写体種別判定テーブル261には、顔に関しては、「怒っている顔」、「泣いている顔」等の顔の表情種別毎に特徴量データDが記憶されているとともに、「個人名A」、「個人名B」等の個人名に対応して画像の特徴を示す特徴量データDも記憶されている。これら各画像の特徴量データDは、色相=HHH、彩度=SSS、明度=VVV、輪郭形状=FFF、大きさ=LLL・・・等の複数種の特徴量で構成されている。

30

【0032】

以上の構成に係る本実施の形態において、前述のようにムービー録画モードにおいては、1回目のシャッターキーと2回目のシャッターキー操作との間に、SDRAM22に順次記憶される複数の画像データがCPU13により順次圧縮される。また、1回目と2回目のシャッターキー操作間に、音声チップ32から入力された音声波形データが圧縮され、この圧縮周囲音データと圧縮動画データとを含む音声付き動画ファイルを生成されて外部メモリ25に記録される。さらに、この音声付き動画ファイルの記録に際してCPU13は、GPS35により検出された緯度・経度と前記地図データとに基づき、撮影地域を検出して、動画ファイルのヘッダーに記憶するとともに、撮影日時、撮影時の明るさ等の撮影条件データもヘッダーに記録する。したがって、音声付き動画ファイルには、圧縮周囲音データと圧縮動画データが記憶されているとともに、付加情報として撮影地域、撮影

40

50

日時、撮影条件等が記憶されている。

【 0 0 3 3 】

そして、PLAY・モードにおいて画像加工モードを設定し、外部メモリ25からいずれかの音声付き動画ファイルを選択すると、CPU13は図3～図10に示すフローチャートに従って処理を実行する。すなわち、図3に示すように、選択された音声付き動画ファイルからの画像データ及び周囲音データの読み出しを開始する(ステップS101)。この読み出した周囲音データに関しては、再生することなく後述する周囲音認識処理を実行する(ステップS102)。なお、周囲音に関しても、音声チップ32で再生しスピーカ33から放音するようにしてもよい。引き続き、後述する画像加工処理を実行し(ステップS103)、この画像加工処理された画像データを含む動画データをSDRAM22に順次記憶するとともに、この画像加工処理された画像データを含む動画データを再生して、液晶モニタ23に表示させる(ステップS104)。

10

【 0 0 3 4 】

しかる後に、前記動画ファイルから読み出している動画データの再生を終了したか、又は動作再生を停止させるキー操作がなされたか否かの終了判断を行い(ステップS105)、終了と判断したならば、ステップS104で順次記憶した複数の画像データからなる動画データを圧縮し、別動画ファイルとして外部メモリ25に記録する(ステップS106)。したがって、後日これら動画ファイルに基づく加工動画を再生することもできるし、加工動画中の任意のフレームを選択して静止画としてプリントアウトすることもできる。

20

【 0 0 3 5 】

図4は、前記周囲音認識処理(ステップS102)の処理手順を示すフローチャートである。先ず、前記音声付き動画ファイルから動画データと同期して順次読み出される音声データに周囲音が含まれているか否かを判断する(ステップS201)。周囲音が含まれている場合には、その波形、スペクトル等の音声データの特徴と読み出された周囲音の特徴とを比較することにより、該読み出された周囲音が音声であるか否かを判断する(ステップS202)。つまり、周囲音を音声認識し、音声認識不可能であれば、音声ではないと判断する。この判断の結果、読み出された周囲音が音声以外の音であった場合には、当該音を擬音表示データに変換する(ステップS203)。例えば、周囲音がクラッカーの破裂音であれば、「パン」の文字からなる擬音データに変換し、周囲音が音楽であれば音

30

【 0 0 3 6 】

また、読み出された周囲音が音声認識可能であれば、これを音声であると判断し、この音声を認識処理してテキストデータに変換する処理を開始する(ステップS204)。また、音声と同期して順次読み出される動画中における人間の口の動きを認識する(ステップS205)。このステップ205での処理に際しては、図10において後述するように先ずフレーム画像中における人間の口の存在を検出する。そして、この検出したフレーム画像中における口の変化を時系列的に検出することにより、口の動きを認識する。この認識した口の動きに対応するテキストデータを前記口の動き-テキストデータ変換テーブルから読み出すことにより、口の動きに対応するテキストデータを得る。なお、言うまでもなく、動画中に人間の口が存在しない場合や人間の口が存在しても口が動いていない場合にはステップS205～S208の処理をスキップすることになる。

40

【 0 0 3 7 】

次に、この口の動きに対応するテキストデータと、ステップS204で音声からの変換を開始しているテキストデータとを照合し(ステップS206)、両者に不一致があるか否かを判断する(ステップS207)。両者に不一致がある場合には、音声から変換しているテキストデータの不一致部分を、口の動きに対応するテキストデータに訂正する(ステップS208)。なお、これとは逆に、口の動きに対応するテキストデータの不一致部分を、音声から変換しているテキストデータに訂正するようにしてもよい。

【 0 0 3 8 】

50

また、音声を終了したか否かを判断し（ステップS209）、音声を終了するまでステップS205からの処理を繰り返す。音声を終了したならば、音声が強く終わったか否かを判断し（ステップS210）、強く終わった場合にはテキストデータの末尾に感嘆符“！”を追加する（ステップS211）。さらに、音声が上がって終わったか否かを判断し（ステップS212）、上がって終わった場合にはテキストデータの末尾に疑問符“？”を追加する（ステップS213）。

【0039】

図5～図7は、前記画像加工処理（ステップS103）の処理手順を示す一連のフローチャートである。まず、図5に示すように、動きの早い被写体があるか否かを判断する（ステップS301）。この判断に際しては、予め動画における画像変化速度の基準値 A_{mm}/s を定めておき、動画中にこの基準値 A_{mm}/s よりも速い速度で動いた被写体があるか否かを判断する。そして、この判断した被写体の動画を構成するフレーム画像中における位置（位置座標）を検出する（ステップS302）。また、フラッシュメモリ26から動きの早い被写体に対応する擬音を示す表示データ（本例では前記「ビュー」）を読み出し（ステップS303）、この読み出した擬音を示す表示データを前記ステップS302で検出した位置の近傍に合成する（ステップS304）。したがって、このステップS301～S304での処理により、例えば投げられたボールの近傍に擬音表示データ「ビュー」が合成される。

【0040】

また弧の動きの被写体があるか否かを判断する（ステップS305）。この判断に際しては、動画を構成するフレームの前後の関係から、弧の動きの被写体の有無を判断する。そして、弧の動きの被写体があった場合には、フレーム画像中における位置を検出する（ステップS306）。また、フラッシュメモリ26から弧の動きを線を示す表示データ（本例では前記「（（」）を読み出し（ステップS307）、この読み出した擬音を示す表示データを前記ステップS306で検出した位置の近傍に合成する（ステップS308）。したがって、このステップS305～S308での処理により、例えば尻尾を振る犬の尻尾の近傍に「（（」を合成することができる。

【0041】

引き続き、周囲音があるか否か（周囲音を読み出されたか否か）を判断し（図6ステップS309）、周囲音がない場合には、再生画像中に人の顔があるか否かを判断する（ステップS310）。

【0042】

この判断に際しては、図10のフローチャートに示すように、動画を構成するフレーム内の抽出領域を検出する（ステップS1）。この抽出領域の検出は、フレーム画像の画像データの輝度信号及び色差信号から、近い輝度又は色差信号別に、同系色の色相別等に領域を分割し、さらに、領域の境界線となる輪郭線を抽出し、この輪郭線で囲まれた部分を一つの抽出領域として検出する。引き続き、この検出した抽出領域を順次選択し（ステップS2）、この選択した抽出領域におけるフレーム画像の特徴抽出処理を実行する（ステップS3）。つまり、選択した抽出領域において、前記特徴量データDが有する特徴種別の特徴量を抽出する。したがって、本例においては、特徴量データDは、色相、彩度、明度、輪郭形状、大きさ・・・であったことから、抽出領域にこれら色相、彩度、明度、輪郭形状、大きさ・・・の特徴量を抽出する。

【0043】

そして、このステップS3で抽出した特徴量と、被写体種別判定テーブル261に記憶されている比較対照となっている被写体種別（ステップS310の場合「人の顔」）の特徴量データDの色相 = H H H、彩度 = S S S、明度 = V V V、輪郭形状 = f f f、大きさ = L L L・・・と各々比較し類似度を各々算出する（ステップS4）。つまり、被写体種別判定テーブル261に記憶されている判断対象の被写体種別の特徴量データDの各値と抽出した特徴量の各値との比率を算出する。次に、この算出した比率である類似度が所定値以上である否かを判断し（ステップS5）、類似度が所定値以上である場合には、当該

10

20

30

40

50

被写体があると判断する(ステップS6)。そして、あると判断した被写体の画像上における位置を検出し、この検出した位置をその被写体種別と共にSDRAM22に記憶する(ステップS7)。

【0044】

また、類似度が所定値未満である場合には、最後の抽出領域まで以上のステップS2～ステップS5の処理を実行したか否かを判断し(ステップS8)、最後の抽出領域となるまでステップS2からの処理を繰り返す。したがって、後述するように画像中に複数の口が存在する場合には、各口に対応してステップS6とステップS7の処理が実行されて、複数の各口に対応してその位置がSDRAM22に記憶されることとなる。よって、最後の抽出領域となるまで、ステップS5の判断がNOであって、類似度が所定値以上の抽出領域がない場合には、SDRAM22には被写体の画像上における位置、及び被写体種別が記憶されない。したがって、SDRAM22に被写体の画像上における位置、及び被写体種別が記憶されているか否かにより、当該被写体があるか否かを判断することができる。

10

【0045】

そして、ステップS310の判断がNOであって、人の顔の被写体がない場合には、前記ヘッダーに記憶されている撮影条件データ等に基づきフラッシュメモリ26から表示データを読み出し(ステップS311)、この読み出した表示データを画像の任意の位置に合成する(ステップS312)。したがって、このステップS311及びS312での処理により、周囲音がない場合であっても、ヘッダーに記憶されている明るさや撮影日時に応じて、「ジリジリ」や「ヒュー」の擬音表示データを、画像の適宜の位置に合成することができる。

20

【0046】

また、ステップS310での判断の結果、人の顔があった場合には、フレーム画像中におけるその位置を前記図10のステップS7においてSDRAM22に人の顔と共に記憶された検出位置を取得する(ステップS313)。次に、この検出された位置の画像である顔に表情があるか否かを判断する(ステップS314)。この判断も図10に示したフローチャートに従って行い、表情がある場合には、フラッシュメモリ26から表情に応じた表示データ(本例では前記「ニコニコ」「ブンブン」)を読み出し(ステップS315)、この読み出した擬音を示す表示データを前記ステップS313で取得した位置の近傍に合成する(ステップS316)。したがって、このステップS313～S316での処理により、周囲音がない場合であっても、被写体の顔の近傍に「ニコニコ」「ブンブン」等を合成して表示することができる。

30

【0047】

また、ステップS314での判断の結果、表情がないと判断された場合には、フラッシュメモリ26から前記ランダムテキストデータのいずれかをランダムに選択する(ステップS317)。引き続き、ステップS313で取得した検出位置に最も近い背景領域を検出する(ステップS308)。この背景領域の検出は、図10に示したフローチャートのステップを利用して行うことができる。

【0048】

すなわち、前述したように、図10のステップS1においては、動画を構成するフレーム内の抽出領域を検出する。この抽出領域の検出は、フレーム画像の画像データの輝度信号及び色差信号から、近い輝度又は色差信号別に、例えば同系色の色相別等に領域を分割し、さらに、領域の境界線となる輪郭線を抽出し、この輪郭線で囲まれた部分を一つの抽出領域として検出する。したがって、このように、抽出領域と検出された領域以外の領域を背景領域であるとして検出することができる。

40

【0049】

そして、検出位置に最も近い背景領域を検出したならば、この検出した検出領域内に吹き出しを合成し(ステップS319)、この吹き出し内に前記ステップS317で選択したテキストデータを合成する(ステップS320)。したがって、音声がない場合であっ

50

ても、画像内の人物が「もうかりまっか」等を発言しているかのような画像を合成して表示することができる。

【0050】

他方、ステップS309での判断の結果、周囲音がある場合には、前記図10に示したフローチャートに従って処理を実行することにより、フレーム画像中に人が存在するか否かを判断する(図7のステップS321)。人が存在する場合には、同様の処理により被写体種別判定テーブル261に個人名がある被写体であるか否かを判断する(ステップS322)。ある場合には、フレーム画像中における前記図10のステップS7においてSDRAM22に人と共に記憶された検出位置を取得し(ステップS323)、この取得した位置の被写体に個人名を合成する(ステップS324)。

10

【0051】

また、音声があるか否かを判断し(ステップS325)、音声がない場合にはステップS333に進む。音声がある場合には、前記同様の処理により人の口が存在するか否かを判断し(ステップS326)、人の口が存在しない場合、つまり音声があり(ステップS325; YES)、人も写っているが(ステップS321; YES)、口は写っていない場合には(ステップS326; NO)、後述する第1の吹き出し合成処理を実行する(ステップS327)。

【0052】

また、口が存在する場合には、複数の口が存在するか否かを判断する(ステップS328)。つまり、前述のように図10のフローチャートに従った処理より、複数の口が存在する場合には、フレーム画像中における各口は特定されていることから、これに基づき複数の口の有無を判断する。この判断がNOであって単一の口のみが写っている場合には、次のステップS329の判断を行うことなく、後述する第2の吹き出し合成処理を実行する(ステップS332)。また、複数の口が写っている場合には、動いている口があるか否かを判断する(ステップS329)。つまり、前述のように図10のフローチャートに従った処理より、複数の口が存在する場合には、フレーム画像中における各口は特定されていることから、このフレーム画像中における各口の変化の有無を時系列的に検出することにより、動いている口があるか否かを判断することができる。

20

【0053】

そして、動いている口がない場合には、後述する第1の吹き出し合成処理(ステップS327)を実行する。また、動いている口がある場合には、該動いている口は1つであるか否かを判断し(ステップS330)、1つである場合には後述する第2の吹き出し合成処理を実行する(ステップS332)。しかし、動いている口が1つではなく、複数ある場合には、前記ステップS325でYES(音声あり)と判断された音声に対応する口を検出する(ステップS331)。

30

【0054】

すなわち、前述の図4のフローチャートにおいては、ステップS204で音声を認識処理してテキストデータに変換する処理を開始し、また、ステップS205では音声とともに順次読み出される動画中における人間の口の動きを認識する。したがって、動いている複数の口において、音声認識により順次変換されるテキストデータと前記ステップS205で認識される動きとが同期する口を検出することにより、音声に対応する口、つまりテキストデータに変換されている音声に対応して動いている口を検出することができる。したがって、このステップS331の処理は、図4のフローチャートに示した周囲音認識処理で実行されるテキストデータ変換処理と口の動き認識処理とを利用して、判断を行う。

40

【0055】

なお、第2の吹き出し合成処理は、後述するようにテキストデータに基づき実行される処理、つまりは音声の存在を前提として実行される処理である。したがって、本実施の形態においては、単一の口が写っているか又は動いている口が写っている場合には、音声も録音されていることが前提となる。

【0056】

50

そして、前記ステップS 3 2 5で音声がないと判断された場合、第1の吹き出し合成処理(ステップS 3 3 2)又は第2の吹き出し合成処理(ステップS 3 3 2)を実行した後、同様に図10のフローチャートに従った処理を実行することにより、人以外の他の音発生被写体があるか否かを判断する(ステップS 3 3 3)。ある場合には、前記図10のステップS 7においてS D R A M 2 2に人以外の他の音発生被写体と共に記憶された検出位置を取得し(ステップS 3 3 4)、この取得した位置の近傍に、前記ステップS 2 0 3で変換された擬音表示データを合成する(ステップS 3 3 5)。したがって、図11(A)に示すように、加工前の画像においてクラッカーP 1が検出されると、同図(B)の加工後の画像に示すように、クラッカーP 1の近傍に擬音表示データP 2「パン」を合成することができる。

10

【0057】

さらに、前記人又は音発生被写体以外の背景に前記ステップS 2 0 3で変換された擬音表示データを合成する(ステップS 3 3 6)。したがって、周囲音が例えば拍手であれば、図11(B)に示すように、「パチパチ」なる表示データP 3が合成される。また、音楽が流れていれば、音符からなる表示データP 4を合成される。

【0058】

他方、前記ステップS 3 2 1で人が存在しないと判断された場合には、音声があるか否かを判断する(ステップS 3 3 7)。そして、音声がある場合には第1の吹き出し処理を実行し(ステップS 3 3 8)、音声がない場合にはステップS 3 3 3に進む。

【0059】

20

図8は、前記第1の吹き出し合成処理(ステップS 3 2 7、ステップS 3 3 8)の処理手順を示すフローチャートである。まず、前記ステップS 3 1 8での説明と同様の処理を行うことにより、フレーム画像中において背景領域を検出する(ステップS 4 0 1)。なお、ステップS 3 2 7でこの第1の吹き出し合成処理を実行する場合には、ステップS 3 2 1で人が存在すると判断されているので、このステップS 3 2 1で存在すると判断された人の近傍に背景領域を検出する。

【0060】

そして、背景領域を検出したならば、この検出した検出領域内に収まるような吹き出しを生成する(ステップS 4 0 2)。しかる後に、前記図4のフローチャートに従った処理により得られているテキストデータを方言に変換する(ステップS 4 0 3)。つまり前述のように、この音声付き動画ファイルの記録に際しては、G P S 3 5により検出された緯度・経度と地図データとに基づき検出された撮影地域が、当該動画ファイルのヘッダーに記憶されている。したがって、この撮影地域を読み出し、前記テキストデータを、フラッシュメモリ26内の方言変換テーブルを用いて、前記撮影地域に対応する方言のテキストデータに変換する。

30

【0061】

さらに、この変換したテキストデータをステップS 4 0 2で生成した吹き出し内に合成して、この吹き出しとテキストデータとからなる吹き出しテキストデータを生成する(ステップS 4 0 4)。引き続き、表示色変更処理を実行して、この吹き出しテキストデータの表示色を、ステップS 3 0 9で検出された周囲音(音声)の高さに応じて変更する(ステップS 4 0 5)。また、表示サイズ変更処理を実行して、この吹き出しテキストデータの表示サイズを、ステップS 3 0 9で検出された周囲音(音声)の音量に応じて変更する(ステップS 4 0 6)。また、前記ステップS 3 2 1で存在が検出された人に対応する個人名が被写体種別判定テーブル261にあるか否かを判断する(ステップS 4 0 7)。ある場合には、フォント変更処理を実行して、この吹き出しテキストデータにおけるテキストデータのフォントを、前記ステップS 4 0 7で個人名ありと判断された個人名(あるいは性別)に応じて変更する(ステップS 4 0 8)。そして、以上の処理により確定した吹き出しテキストデータを前記ステップS 4 0 1で検出した検出領域内に、合成する(ステップS 4 0 9)。

40

【0062】

50

したがって、この図 8 に示した第 1 の吹き出し合成処理により、人の口が写っていない場合であって、音声を検出された場合には、吹き出し内に音声に対応するテキストデータが合成された吹き出しテキストデータが、背景に合成されることとなる。

【 0 0 6 3 】

図 9 は、前記第 2 の吹き出し合成処理（ステップ S 3 3 2）の処理手順を示すフローチャートである。まず、前記ステップ S 3 2 8 で単一の口であると判断された口、又は前記ステップ S 3 3 0 で動いている口は 1 つであると判断された当該口、又はステップ S 3 3 1 で取得された口の位置を S D R A M 2 2 から取得する（ステップ S 5 0 1）。引き続き、前記ステップ S 3 1 8 での説明と同様の処理を行うことにより、検出位置に最も近い背景領域を検出する（ステップ S 5 0 2）。以下は、前記ステップ S 4 0 2 ~ S 4 0 9 と同様の処理であり、この検出した検出領域内に収まるような吹き出しを生成し（ステップ S 5 0 3）、テキストデータを方言に変換する（ステップ S 5 0 4）。この変換したテキストデータをステップ S 5 0 3 で生成した吹き出し内に合成して、この吹き出しとテキストデータとからなる吹き出しテキストデータを生成する（ステップ S 5 0 5）。引き続き、表示色変更処理を実行して、この吹き出しテキストデータの表示色を周囲音（音声）の高さに応じて変更する（ステップ S 5 0 6）。また、表示サイズ変更処理を実行して、この吹き出しテキストデータの表示サイズを、ステップ S 3 0 9 で検出された周囲音（音声）の音量に応じて変更する（ステップ S 5 0 7）。

10

【 0 0 6 4 】

また、前記ステップ S 3 2 1 で存在が検出された人に対応する個人名が被写体種別判定テーブル 2 6 1 にあるか否かを判断する（ステップ S 5 0 8）。ある場合には、フォント変更処理を実行して、この吹き出しテキストデータにおけるテキストデータのフォントを、個人名（あるいは性別）に応じて変更する（ステップ S 5 0 9）。そして、以上の処理により確定した吹き出しテキストデータを前記ステップ S 5 0 1 で検出した検出領域内に、合成する（ステップ S 5 1 0）。

20

【 0 0 6 5 】

したがって、この図 9 に示した第 2 の吹き出し合成処理により、人の口が写っている場合であってその動きも検出され、音声も検出された場合には、吹き出し内に音声に対応するテキストデータが合成された吹き出しテキストデータが、動きのある人の口の近傍であって背景に表示されることとなる。これにより、図 1 1（B）に示すように、加工後の画像には、口 P 5 を動かしている人 P 6 の、該口 5 の近傍であって、他の被写体とは重ならない背景に、吹き出し内に音声に対応するテキストデータ（「おめでとう」）を有する吹き出しテキスト表示データ P が合成される。

30

【 0 0 6 6 】

なお、本実施の形態においては、予め撮影して記録した音声付き動画ファイルを再生する際に本発明を適用する場合を示したが、音声付き静止画を再生する際、音声付き静止画又は音声付き動画を撮影する際のスルー画像表示時、音声付き静止画又は音声付き動画を撮影記録する際に本発明を適用するようにしてもよい。また、テキストデータを方言に変換するようにしたが、方言に変換することなく合成するようにしてもよい。また、実施の形態においては、本発明をデジタルカメラに適用した場合について示したが、これに限ることなく、撮影機能のみを有するビデオカメラ等の撮影装置、再生機能のみを有するビデオデッキ等の映像機器、画像加工機能のみを有する画像加工機器、撮影機能と再生機能とを併有する各種映像機器に本発明を適用するようにしてもよい。

40

【 図面の簡単な説明 】

【 0 0 6 7 】

【 図 1 】 本発明の一実施の形態に係るデジタルカメラの電氣的構成を示すブロック図である。

【 図 2 】 被写体種別判定テーブルを示す概念図である。

【 図 3 】 画像加工モードの処理手順を示すフローチャートである。

【 図 4 】 周囲音認識処理の処理手順を示すフローチャートである。

50

【図5】画像加工処理の処理手順を示すフローチャートである。

【図6】図5に続くフローチャートである。

【図7】図6に続くフローチャートである。

【図8】第1の吹き出し合成処理の処理手順を示すフローチャートである。

【図9】第2の吹き出し合成処理の処理手順を示すフローチャートである。

【図10】被写体存在判定処理の処理手順を示すフローチャートである。

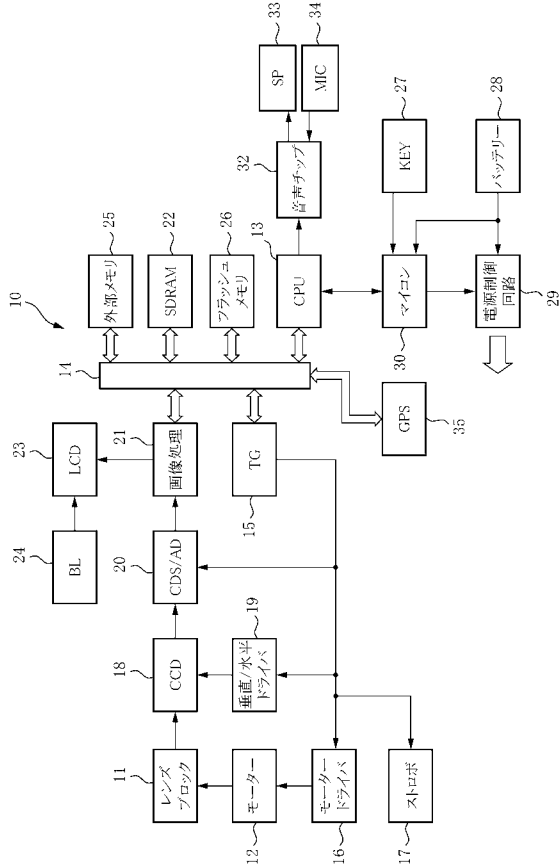
【図11】(A)は加工前、(B)は加工後の画像を示す図である。

【符号の説明】

【0068】

10	デジタルカメラ	10
11	レンズブロック	
14	バス	
15	タイミング発生器	
18	CCD	
19	水平ドライバ	
20	ユニット回路	
21	画像処理部	
22	SDRAM	
23	液晶モニタ	
25	外部メモリ	20
26	フラッシュメモリ	
27	キー入力部	
32	音声チップ	
33	スピーカ	
34	マイクロホン	
35	GPS	
261	被写体種別判定テーブル	

【図1】

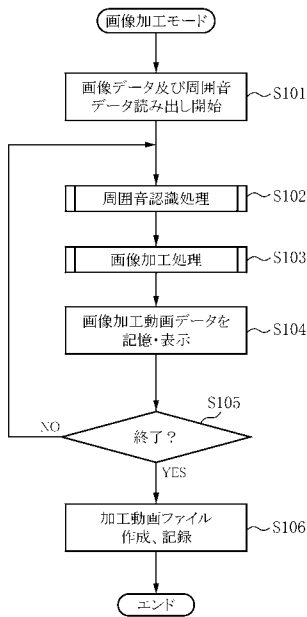


【図2】

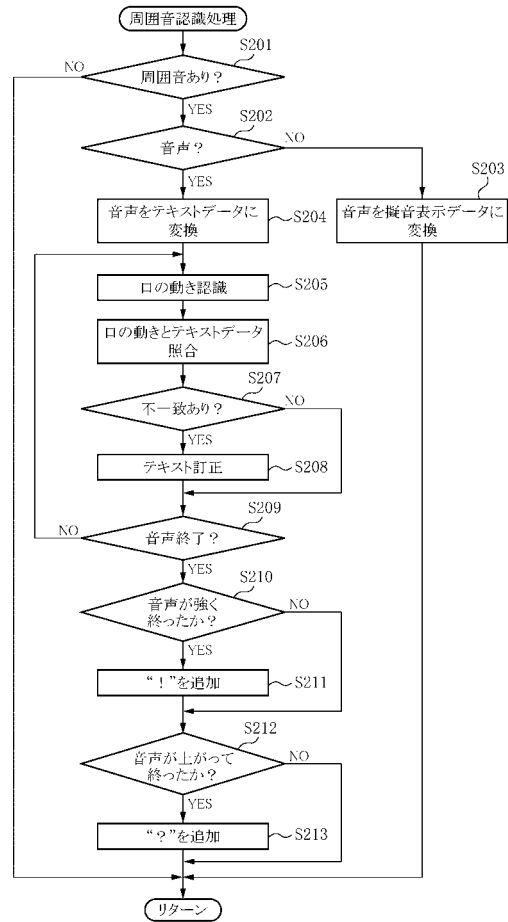
261

被写体種別	特徴量データD
人	--
人の口	--
人の顔	--
自動車	--
電話	--
クラッカー	--
起っている顔	--
泣いている顔	--
個人名A	--
個人名B	--

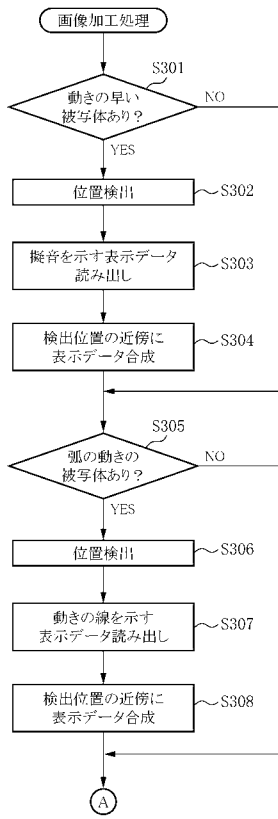
【図3】



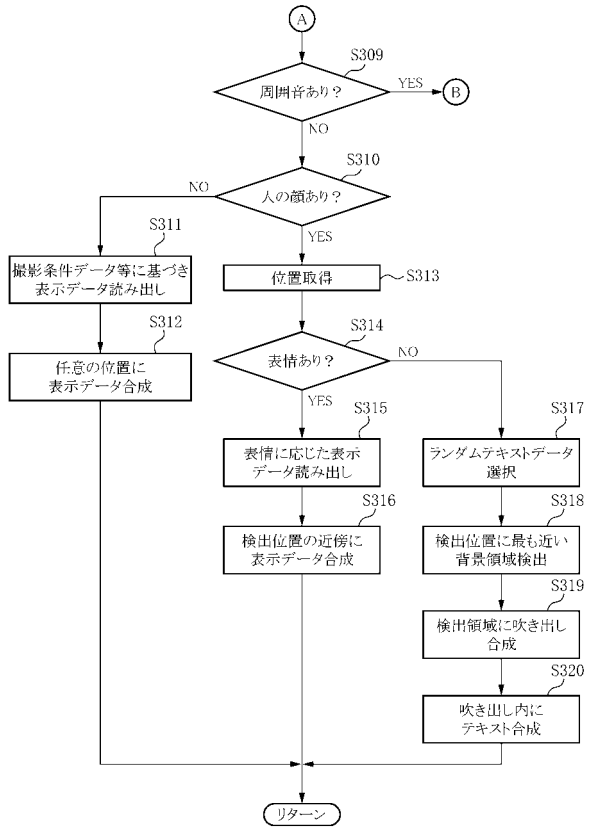
【図4】



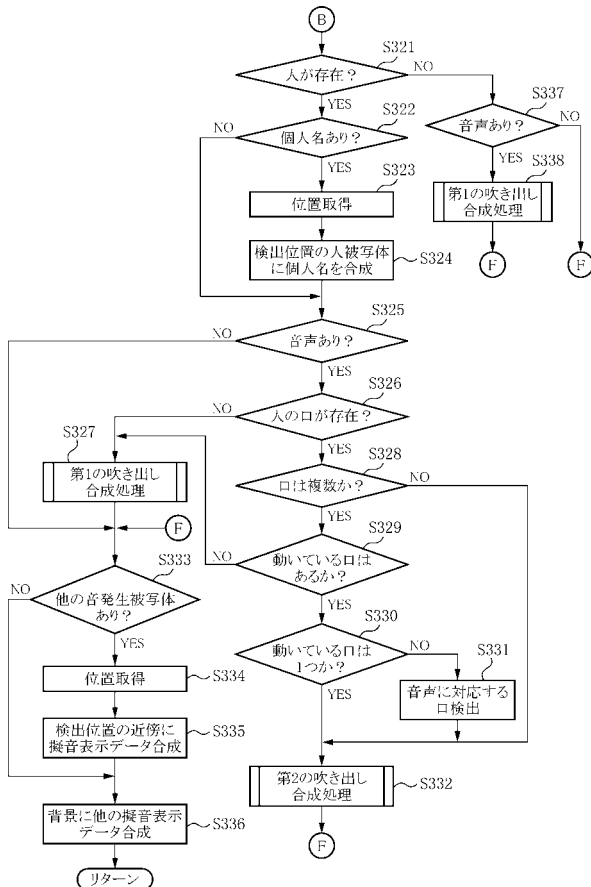
【図5】



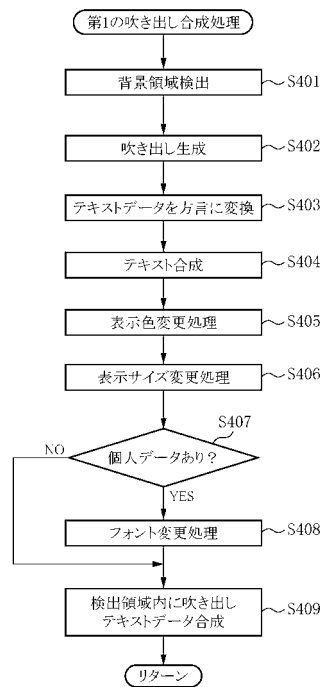
【図6】



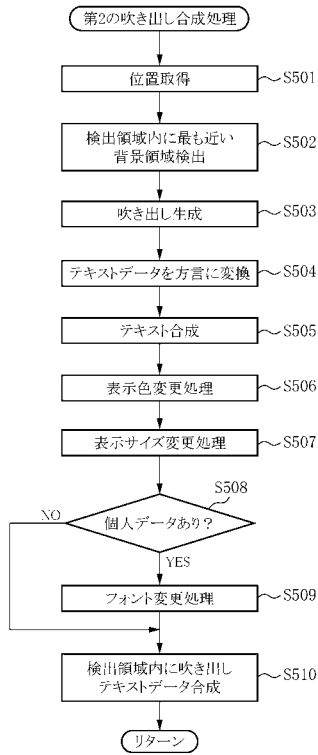
【図7】



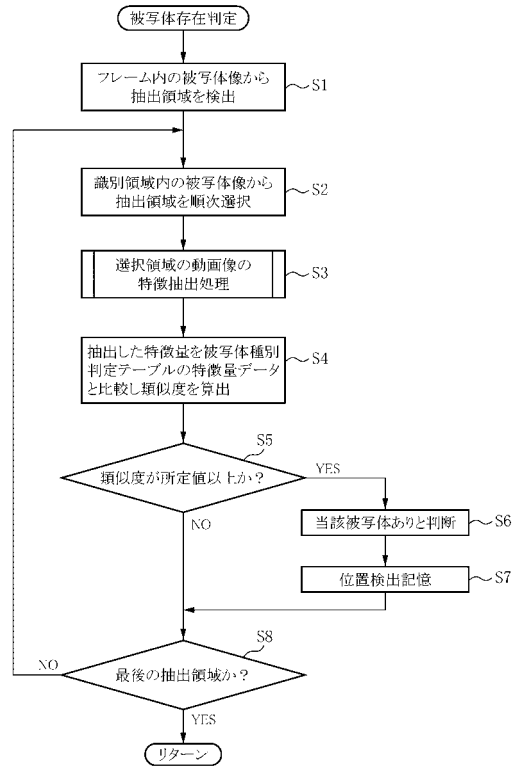
【図8】



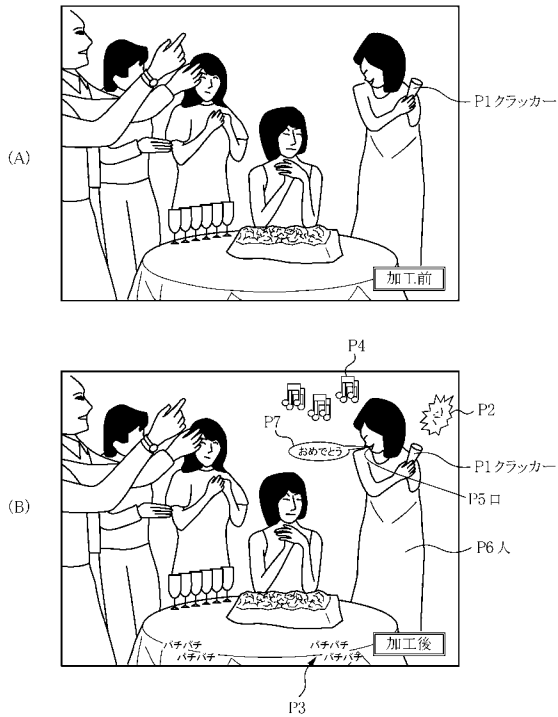
【図9】



【図10】



【図11】



フロントページの続き

(51) Int.Cl. F I
H 0 4 N 101/00 (2006.01) H 0 4 N 101:00

(56) 参考文献 特開 2 0 0 4 - 0 5 6 2 8 6 (J P , A)
特開平 0 8 - 3 1 7 3 6 3 (J P , A)
特開 2 0 0 0 - 1 9 6 9 3 5 (J P , A)
特開 2 0 0 5 - 1 2 4 1 6 9 (J P , A)
特開 2 0 0 0 - 3 3 3 0 8 0 (J P , A)
特開 2 0 0 1 - 0 5 1 3 3 8 (J P , A)
特開 2 0 0 5 - 0 5 1 2 7 8 (J P , A)

(58) 調査した分野(Int.Cl. , DB名)

H 0 4 N 5 / 2 3 2
G 0 6 T 1 / 0 0
G 1 0 L 1 5 / 0 0
H 0 4 N 1 / 3 8 7
H 0 4 N 5 / 9 1
H 0 4 N 1 0 1 / 0 0