

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4894741号
(P4894741)

(45) 発行日 平成24年3月14日(2012.3.14)

(24) 登録日 平成24年1月6日(2012.1.6)

(51) Int.Cl.

F I

G 0 6 T 7/20 (2006.01)

G 0 6 T 7/20 3 0 0 A

請求項の数 20 (全 33 頁)

(21) 出願番号 特願2007-312568 (P2007-312568)
 (22) 出願日 平成19年12月3日(2007.12.3)
 (65) 公開番号 特開2009-140009 (P2009-140009A)
 (43) 公開日 平成21年6月25日(2009.6.25)
 審査請求日 平成22年11月2日(2010.11.2)

(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 100082131
 弁理士 稲本 義雄
 (74) 代理人 100121131
 弁理士 西川 孝
 (72) 発明者 鈴木 洋貴
 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 佐藤 実

最終頁に続く

(54) 【発明の名称】 情報処理装置および情報処理方法、プログラム、並びに記録媒体

(57) 【特許請求の範囲】

【請求項1】

入力動画に、登録されているアクションが含まれているか否かを認識する情報処理装置において、

前記アクションを認識するためのモデルを含むモデル動画を画像平面および時間の3次元とした場合における、特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を記憶する記憶手段と、

前記入力動画を取得する第1の取得手段と、

前記第1の取得手段により取得された前記入力動画を画像平面および時間の3次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出する第1の特徴点抽出手段と、

前記第1の特徴点抽出手段により抽出された前記入力特徴点における特徴量である入力特徴量を抽出する第1の特徴量抽出手段と、

前記第1の特徴量抽出手段により抽出された前記入力特徴量と、前記記憶手段により記憶された前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成する特徴量比較手段と、

前記特徴量比較手段による比較の結果得られた前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求める姿勢推定手段と、

前記姿勢推定手段により得られる前記モデルの姿勢の推定結果、および、前記認識対応

10

20

特徴点ペア群に基づいて、認識結果を生成する認識結果生成手段と
を備える情報処理装置。

【請求項 2】

前記姿勢推定手段は、ランダムに選択した N 組の前記候補対応特徴点ペアにより決定される前記モデル動画の画像平面および時間の 3 次元における位置姿勢を決める画像変換パラメータをパラメータ空間に投射し、前記パラメータ空間上をクラスタリングすることにより形成されるクラスタのうち、最多メンバ数を有するクラスタを求め、前記最多メンバ数を有するクラスタのメンバである前記候補対応特徴点ペア群を前記認識対応特徴点ペア群とする

請求項 1 に記載の情報処理装置。

10

【請求項 3】

前記姿勢推定手段は、前記最多メンバ数を有するクラスタのセントロイドを検出し、前記セントロイドを、姿勢に対応するパラメータとして、前記モデルの姿勢を推定する

請求項 2 に記載の情報処理装置。

【請求項 4】

前記姿勢推定手段は、NN法により前記パラメータ空間上をクラスタリングする

請求項 2 に記載の情報処理装置。

【請求項 5】

前記画像変換パラメータは、アフィンパラメータである

請求項 2 に記載の情報処理装置。

20

【請求項 6】

前記姿勢推定手段は、前記アフィンパラメータのレンジを正規化し、正規化された前記アフィンパラメータをパラメータ空間に投射する

請求項 5 に記載の情報処理装置。

【請求項 7】

前記姿勢推定手段は、回転、拡大縮小、および、せん断変形のそれぞれを決定する 9 次元のパラメータの正規化係数を 1 . 0 とし、平行移動を決定するための 3 次元のパラメータのうち、横方向の平行移動に関するパラメータの正規化係数を想定される動画の横ピクセル数の逆数とし、縦方向の平行移動に関するパラメータの正規化係数を想定される動画の縦ピクセル数の逆数とし、時間方向の平行移動に関するパラメータの正規化係数を想定される動画の時間長の逆数とし、これらの正規化係数を前記アフィンパラメータに乘じることにより、前記アフィンパラメータのレンジを正規化する

30

請求項 6 に記載の情報処理装置。

【請求項 8】

前記姿勢推定手段は、回転、拡大縮小、および、せん断変形のそれぞれを決定する 9 次元のパラメータに対するクラスタリング規範となる距離の第 1 の閾値と、平行移動を決定するための 3 次元のパラメータに対するクラスタリング規範となる距離の第 2 の閾値を用いてクラスタリングを実行し、前記第 2 の閾値は前記第 1 の閾値よりも大きい

請求項 6 に記載の情報処理装置。

【請求項 9】

前記第 1 の特徴点抽出手段は、画像平面および時間の 3 次元に拡張された Harris 関数 H の極大および極小を与える画像平面および時間の 3 次元座標を、前記入力動画における前記入力特徴点として抽出する

請求項 1 に記載の情報処理装置。

40

【請求項 10】

前記第 1 の特徴量抽出手段は、画像平面および時間の 3 次元のそれぞれの次元について、4 次までの偏微分ガウスオペレーションをかけた画像情報から構成される特徴ベクトルを前記入力特徴量として抽出する

請求項 1 に記載の情報処理装置。

【請求項 11】

50

前記特徴量比較手段は、前記入力特徴量と、前記モデル特徴量とのノルムを、前記入力特徴量と前記モデル特徴量との非類似度の尺度に用いて、前記候補対応特徴点ペアを生成する

請求項 1 に記載の情報処理装置。

【請求項 1 2】

前記認識結果生成手段は、前記姿勢推定手段により得られた前記認識対応特徴点ペア群の要素数が所定の閾値以上である前記モデルを、登録されている前記アクションが含まれているモデルの認識結果とする

請求項 1 に記載の情報処理装置。

【請求項 1 3】

前記認識結果生成手段は、前記姿勢推定手段により得られた前記認識対応特徴点ペア群の要素数が所定の閾値以上である前記モデルを、要素数の多い順にソートし、検出されたモデル全てとそれらの順位とを、登録されている前記アクションが含まれているモデルの認識結果とする

請求項 1 に記載の情報処理装置。

【請求項 1 4】

前記認識結果生成手段は、前記姿勢推定手段により得られた前記認識対応特徴点ペア群の要素数が所定の閾値以上である前記モデルの要素数の総和に対する、それぞれのモデルの前記認識対応特徴点ペア群の要素数の割合を、前記認識対応特徴点ペア群の要素数が所定の閾値以上であるそれぞれの前記モデルの信頼度とする

請求項 1 に記載の情報処理装置。

【請求項 1 5】

前記認識結果生成手段は、前記姿勢推定手段により得られる前記モデルの姿勢の推定結果を認識結果とする

請求項 1 に記載の情報処理装置。

【請求項 1 6】

前記認識結果生成手段は、前記姿勢推定手段により得られた前記認識対応特徴点ペア群の要素数が所定の閾値以上である前記モデルの前記画像変換パラメータの最小二乗推定結果を認識結果とする

請求項 2 に記載の情報処理装置。

【請求項 1 7】

前記第 1 の取得手段により取得された前記入力動画を、前記モデルに対応する領域と背景に対応する領域とに分割する分割手段を更に備え、

前記第 1 の特徴点抽出手段は、前記分割手段によって分割された前記入力動画中の前記モデルに対応する領域から、前記入力特徴点を抽出する

請求項 1 に記載の情報処理装置。

【請求項 1 8】

アクションを認識するためのモデルを含むモデル動画を画像平面および時間の 3 次元とした場合における、特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を記憶する記憶部を有し、入力動画に、登録されている前記アクションが含まれているか否かを認識する情報処理装置の情報処理方法において、

前記入力動画を取得し、

前記入力動画を画像平面および時間の 3 次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出し、

前記入力特徴点における特徴量である入力特徴量を抽出し、

前記入力特徴量と、前記記憶部に記憶されている前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成し、

前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求め、

前記モデルの姿勢の推定結果、および、前記認識対応特徴点ペア群に基づいて、認識結

10

20

30

40

50

果を生成する

ステップを含む情報処理方法。

【請求項 19】

所定の記憶部に記憶されているアクションを認識するためのモデルを含むモデル動画を画像平面および時間の3次元とした場合における、特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を用いて、入力動画に、登録されている前記アクションが含まれているか否かを認識する処理をコンピュータに実行させるためのプログラムであって、

前記入力動画を取得し、

前記入力動画を画像平面および時間の3次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出し、

前記入力特徴点における特徴量である入力特徴量を抽出し、

前記入力特徴量と、前記記憶部に記憶されている前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成し、

前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求め、

前記モデルの姿勢の推定結果、および、前記認識対応特徴点ペア群に基づいて、認識結果を生成する

ステップを含む処理をコンピュータに実行させるプログラム。

【請求項 20】

請求項 19 に記載のプログラムが記録されている記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理装置および情報処理方法、プログラム、並びに記録媒体に関し、特に、物体の動きを検出する場合に用いて好適な、情報処理装置および情報処理方法、プログラム、並びに記録媒体に関する。

【背景技術】

【0002】

動画像データを取得し、物体の動きを検出するために、従来、さまざまな手法が用いられてきた。

【0003】

取得された動画像データから、認識すべき物体の、例えば、手や足などの特定のボディパーツを位置同定し、その動き情報から、アクションの認識を行う方法が多数提案されている。具体的には、例えば、ボディパーツ同定に特殊器具を用いる方法（例えば、特許文献 1 または特許文献 2 参照）、テンプレートマッチングを行う方法（例えば、特許文献 3 または特許文献 4 参照）、色情報や輪郭情報でボディパーツの同定を行う方法（例えば、特許文献 5 乃至特許文献 12 参照）などがある。

【0004】

【特許文献 1】特許 2 5 5 8 9 4 3 号公報

【特許文献 2】特許 3 1 4 4 4 0 0 号公報

【特許文献 3】特許 2 7 8 1 7 4 3 号公報

【特許文献 4】特開平 8 - 2 7 9 0 4 4 号公報

【特許文献 5】米国特許 6 2 5 6 0 0 B 1 号公報

【特許文献 6】特許 2 8 6 8 4 4 9 号公報

【特許文献 7】特許 2 9 3 4 1 9 0 号公報

【特許文献 8】特許 3 4 4 0 6 4 4 号公報

【特許文献 9】特開平 1 0 - 2 1 4 3 4 6 号公報

【特許文献 10】特開 2 0 0 3 0 3 9 3 6 5 号公報

【特許文献 11】特開 2 0 0 3 2 1 6 9 5 5 号公報

10

20

30

40

50

【特許文献 1 2】特許 2 8 6 8 4 4 9 号公報

【0 0 0 5】

また、時間差分法やオプティカルフローにより、動き領域を抽出し、その領域の重心時間変化パターンからアクションを認識する方法が提案されている（例えば、特許文献 1 3 参照）。

【0 0 0 6】

【特許文献 1 3】米国特許 U S 6 6 8 1 0 3 1 B 2 号公報

【0 0 0 7】

また、認識させたいアクション（モデルアクション）に対して、それが撮像されている大量の学習用動画を用意し、各動画から時空間イベントを記述する特徴量群を抽出し、例えば、サポートベクターマシンなどの統計学習手法を用いて学習することにより、その特徴量群のなかから、モデルアクションをそれ以外の時空間パターンと良く分離する特徴量を求め、認識処理時に入力動画からモデルアクションを認識する際には、学習により求めた特徴量のみを用いて検出有無の判定を行う方法が提案されている（例えば、非特許文献 1 参照）。

【0 0 0 8】

【非特許文献 1】C. Schuldt, I. Laptev, and B. Caputo. Recognizing human actions: a local SVM approach. In ICPR, pages III: 3236, 20

【発明の開示】

【発明が解決しようとする課題】

【0 0 0 9】

しかしながら、取得された動画像データから、認識すべき物体の、例えば、手や足などの特定のボディパーツを位置同定し、その動き情報からアクションの認識を行う手法は、認識可能なボディパーツに限定して行われるものであり、各アクションに特化した認識アルゴリズムを必要とする。すなわち、ボディパーツおよびアクションごとに検出アルゴリズムが全く異なるものとなるため、システム設計時に想定していないようなアクションを後から認識することができるようにすることはできない。具体的には、例えば、物体を用いたアクションや他のボディパーツを用いたアクションや複数人による協調的アクションを、ユーザが後から任意に登録して、それらのアクションを認識することはできない。

【0 0 1 0】

また、時間差分法やオプティカルフローにより、動き領域を抽出する方法を用いたとしても、重心時間変化パターンのみでは、様々なジェスチャを切り分けるのに十分な情報であるとはいえない。また背景などを含む画像が取得されてしまう実環境では、動き領域の抽出の精度を高くするのは困難であった。さらに、認識すべき動き領域が部分遮蔽されてしまった場合、重心位置が本来の位置からずれてしまい、認識精度が出ないことが予想される。

【0 0 1 1】

そして、統計学習を用いる場合、システム設計時に想定していないようなアクションを後から認識することは可能であるが、認識に適した特徴量を学習するために大量の学習用動画が必要になる。例えば、ユーザが新しいジェスチャをシステムに登録させたいと思った場合、ユーザは、登録させるジェスチャをシステムに学習させるため、大量の学習データを用意する必要がある。このようなシステムで多くのジェスチャを認識させるには、ユーザに、大変な労力を課してしまう。

【0 0 1 2】

本発明はこのような状況に鑑みてなされたものであり、入力画像の部分隠れなどに頑強で、かつ、多くの学習用データを必要とせず、物体の動きを検出することができるようにするものである。

【課題を解決するための手段】

【0 0 1 3】

本発明の第 1 の側面の情報処理装置は、入力動画に、登録されているアクションが含ま

10

20

30

40

50

れているか否かを認識する情報処理装置であって、前記アクションを認識するためのモデルを含むモデル動画を画像平面および時間の3次元とした場合における、特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を記憶する記憶手段と、前記入力動画を取得する第1の取得手段と、前記第1の取得手段により取得された前記入力動画を画像平面および時間の3次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出する第1の特徴点抽出手段と、前記第1の特徴点抽出手段により抽出された前記入力特徴点における特徴量である入力特徴量を抽出する第1の特徴量抽出手段と、前記第1の特徴量抽出手段により抽出された前記入力特徴量と、前記記憶手段により記憶された前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成する特徴量比較手段と、前記特徴量比較手段による比較の結果得られた前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求める姿勢推定手段と、前記姿勢推定手段により得られる前記モデルの姿勢の推定結果、および、前記認識対応特徴点ペア群に基づいて、認識結果を生成する認識結果生成手段とを備える。

10

【0014】

前記姿勢推定手段には、ランダムに選択したN組の前記候補対応特徴点ペアにより決定される前記モデル動画の画像平面および時間の3次元における位置姿勢を決める画像変換パラメータをパラメータ空間に投射させ、前記パラメータ空間上をクラスタリングすることにより形成されるクラスタのうち、最多メンバ数を有するクラスタを求めさせ、前記最

20

【0015】

前記姿勢推定手段には、前記最多メンバ数を有するクラスタのセントロイドを検出させ、前記セントロイドを、姿勢に対応するパラメータとして、前記モデルの姿勢を推定させるようにすることができる。

【0016】

前記姿勢推定手段には、NN法により前記パラメータ空間上をクラスタリングさせるようにすることができる。

【0017】

前記画像変換パラメータは、アフィンパラメータであるものとすることができる。

30

【0018】

前記姿勢推定手段には、前記アフィンパラメータのレンジを正規化させ、正規化された前記アフィンパラメータをパラメータ空間に投射させるようにすることができる。

【0019】

前記姿勢推定手段には、回転、拡大縮小、および、せん断変形のそれぞれを決定する9次元のパラメータの正規化係数を1.0とし、平行移動を決定するための3次元のパラメータのうち、横方向の平行移動に関するパラメータの正規化係数を想定される動画の横ピクセル数の逆数とし、縦方向の平行移動に関するパラメータの正規化係数を想定される動画の縦ピクセル数の逆数とし、時間方向の平行移動に関するパラメータの正規化係数を想定される動画の時間長の逆数とし、これらの正規化係数を前記アフィンパラメータに乘じることにより、前記アフィンパラメータのレンジを正規化させるようにすることができる。

40

【0020】

前記姿勢推定手段には、回転、拡大縮小、および、せん断変形のそれぞれを決定する9次元のパラメータに対するクラスタリング規範となる距離の第1の閾値と、平行移動を決定するための3次元のパラメータに対するクラスタリング規範となる距離の第2の閾値を用いてクラスタリングを実行させるようにすることができ、前記第2の閾値は前記第1の閾値よりも大きいものとすることができる。

【0021】

50

前記第 1 の特徴点抽出手段には、画像平面および時間の 3 次元に拡張された H a r r i s 関数 H の極大および極小を与える画像平面および時間の 3 次元座標を、前記入力動画における前記入力特徴点として抽出させるようにすることができる。

【 0 0 2 2 】

前記第 1 の特徴量抽出手段には、画像平面および時間の 3 次元のそれぞれの次元について、4 次までの偏微分ガウスオペレーションをかけた画像情報から構成される特徴ベクトルを前記入力特徴量として抽出させるようにすることができる。

【 0 0 2 3 】

前記特徴量比較手段には、前記入力特徴量と、前記モデル特徴量とのノルムを、前記入力特徴量と前記モデル特徴量との非類似度の尺度に用いて、前記候補対応特徴点ペアを生成させるようにすることができる。

10

【 0 0 2 4 】

前記認識結果生成手段には、前記姿勢推定手段により得られた前記認識対応特徴点ペア群の要素数が所定の閾値以上である前記モデルを、登録されている前記アクションが含まれているモデルの認識結果とさせるようにすることができる。

【 0 0 2 5 】

前記認識結果生成手段には、前記姿勢推定手段により得られた前記認識対応特徴点ペア群の要素数が所定の閾値以上である前記モデルを、要素数の多い順にソートさせ、検出されたモデル全てとそれらの順位とを、登録されている前記アクションが含まれているモデルの認識結果とさせるようにすることができる。

20

【 0 0 2 6 】

前記認識結果生成手段には、前記姿勢推定手段により得られた前記認識対応特徴点ペア群の要素数が所定の閾値以上である前記モデルの要素数の総和に対する、それぞれのモデルの前記認識対応特徴点ペア群の要素数の割合を、前記認識対応特徴点ペア群の要素数が所定の閾値以上であるそれぞれの前記モデルの信頼度とさせるようにすることができる。

【 0 0 2 7 】

前記認識結果生成手段には、前記姿勢推定手段により得られる前記モデルの姿勢の推定結果を認識結果とさせるようにすることができる。

【 0 0 2 8 】

前記認識結果生成手段には、前記姿勢推定手段により得られた前記認識対応特徴点ペア群の要素数が所定の閾値以上である前記モデルの前記画像変換パラメータの最小二乗推定結果を認識結果とさせるようにすることができる。

30

【 0 0 2 9 】

前記第 1 の取得手段により取得された前記入力動画を、前記モデルに対応する領域と背景に対応する領域とに分割する分割手段を更に備えさせるようにことができ、前記第 1 の特徴点抽出手段は、前記分割手段によって分割された前記入力動画中の前記モデルに対応する領域から、前記入力特徴点を抽出させるようにすることができる。

【 0 0 3 3 】

本発明の第 1 の側面の情報処理方法は、アクションを認識するためのモデルを含むモデル動画を画像平面および時間の 3 次元とした場合における、特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を記憶する記憶部を有し、入力動画に、登録されている前記アクションが含まれているか否かを認識する情報処理装置の情報処理方法において、前記入力動画を取得し、前記入力動画を画像平面および時間の 3 次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出し、前記入力特徴点における特徴量である入力特徴量を抽出し、前記入力特徴量と、前記記憶部に記憶されている前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成し、前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求め、前記モデルの姿勢の推定結果、および、前記認識対応特徴点ペア群に基づいて、認識結果を生成するステップを含む。

40

50

【0034】

本発明の第1の側面のプログラムは、所定の記憶部に記憶されているアクションを認識するためのモデルを含むモデル動画を画像平面および時間の3次元とした場合における、特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を用いて、入力動画に、登録されている前記アクションが含まれているか否かを認識する処理をコンピュータに実行させるためのプログラムであって、前記入力動画を取得し、前記入力動画を画像平面および時間の3次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出し、前記入力特徴点における特徴量である入力特徴量を抽出し、前記入力特徴量と、前記記憶部に記憶されている前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成し、前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求め、前記モデルの姿勢の推定結果、および、前記認識対応特徴点ペア群に基づいて、認識結果を生成するステップを含む処理をコンピュータに実行させる。

10

【0035】

本発明の第1の側面においては、入力動画が取得され、入力動画を画像平面および時間の3次元として、入力動画からアクションを認識するための特徴点である入力特徴点が抽出され、入力特徴点における特徴量である入力特徴量が抽出され、入力特徴量と、予め記憶されているモデル特徴量とが比較され、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアが生成され、候補対応特徴点ペアから、アウトライヤが除去され、入力動画上でのモデルの姿勢が推定されるとともに、モデルの姿勢に対応する認識対応特徴点ペア群が求められ、モデルの姿勢の推定結果、および、認識対応特徴点ペア群に基づいて、認識結果が生成される。

20

【0044】

ネットワークとは、少なくとも2つの装置が接続され、ある装置から、他の装置に対して、情報の伝達をできるようにした仕組みをいう。ネットワークを介して通信する装置は、独立した装置どうしであっても良いし、1つの装置を構成している内部ブロックどうしであっても良い。

【0045】

また、通信とは、無線通信および有線通信は勿論、無線通信と有線通信とが混在した通信、即ち、ある区間では無線通信が行われ、他の区間では有線通信が行われるようなものであっても良い。さらに、ある装置から他の装置への通信が有線通信で行われ、他の装置からある装置への通信が無線通信で行われるようなものであっても良い。

30

【0046】

認識処理装置は、独立した装置であっても良いし、情報処理装置の認識処理を行うブロックであっても良い。

【発明の効果】

【0047】

以上のように、本発明の第1の側面によれば、認識処理を行うことができ、特に、画像平面および時間の3次元を用いて処理を行うことにより、入力画像の部分隠れなどに頑強で、かつ、多くの学習用データを必要とせずに、物体の動きを検出することができる。

40

【発明を実施するための最良の形態】

【0049】

以下に本発明の実施の形態を説明するが、本発明の構成要件と、明細書または図面に記載の実施の形態との対応関係を例示すると、次のようになる。この記載は、本発明をサポートする実施の形態が、明細書または図面に記載されていることを確認するためのものである。従って、明細書または図面中には記載されているが、本発明の構成要件に対応する実施の形態として、ここには記載されていない実施の形態があったとしても、そのことは、その実施の形態が、その構成要件に対応するものではないことを意味するものではない。逆に、実施の形態が構成要件に対応するものとしてここに記載されていたとしても、そ

50

のことは、その実施の形態が、その構成要件以外の構成要件には対応しないものであることを意味するものでもない。

【 0 0 5 0 】

本発明の第 1 の側面の情報処理装置は、入力動画に、登録されているアクションが含まれているか否かを認識する情報処理装置（たとえば、図 1 の認識処理装置 1 1、または、図 1 の認識処理部 2 2 が有する機能を有する装置、もしくは、図 1 0 のパーソナルコンピュータ 5 0 0 ）であって、前記アクションを認識するためのモデルを含むモデル動画を画像平面および時間の 3 次元とした場合における、特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を記憶する記憶手段（例えば、図 1 の辞書登録部 6 1 ）と、前記入力動画を取得する第 1 の取得手段（例えば、図 1 の入力動画バッファ部 6 2 ）と、前記第 1 の取得手段により取得された前記入力動画を画像平面および時間の 3 次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出する第 1 の特徴点抽出手段（例えば、図 1 の特徴点抽出部 6 4 ）と、前記第 1 の特徴点抽出手段により抽出された前記入力特徴点における特徴量である入力特徴量を抽出する第 1 の特徴量抽出手段（例えば、図 1 の特徴量抽出部 6 5 ）と、前記第 1 の特徴量抽出手段により抽出された前記入力特徴量と、前記記憶手段により記憶された前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成する特徴量比較手段（例えば、図 1 の特徴量比較部 6 6 ）と、前記特徴量比較手段による比較の結果得られた前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢（例えば、モデル姿勢）を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求める姿勢推定手段（例えば、図 1 の姿勢パラメータ推定部 6 7 ）と、前記姿勢推定手段により得られる前記モデルの姿勢の推定結果、および、前記認識対応特徴点ペア群に基づいて、認識結果を生成する認識結果生成手段（例えば、図 1 の認識結果生成部 6 8 ）とを備える。

【 0 0 5 1 】

前記第 1 の取得手段により取得された前記入力動画を、前記モデルに対応する領域と背景に対応する領域とに分割する分割手段（例えば、図 1 の前処理実行部 6 3 ）を更に備えることができ、前記第 1 の特徴点抽出手段は、前記分割手段によって分割された前記入力動画中の前記モデルに対応する領域から、前記入力特徴点を抽出することができる。

【 0 0 5 4 】

本発明の第 1 の側面の情報処理方法は、アクションを認識するためのモデルを含むモデル動画を画像平面および時間の 3 次元とした場合における、特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を記憶する記憶部を有し、入力動画に、登録されている前記アクションが含まれているか否かを認識する情報処理装置（たとえば、図 1 の認識処理装置 1 1、または、図 1 の認識処理部 2 2 が有する機能を有する装置、もしくは、図 1 0 のパーソナルコンピュータ 5 0 0 ）の情報処理方法であって、前記入力動画を取得し（例えば、図 7 のステップ S 4 1 の処理）、前記入力動画を画像平面および時間の 3 次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出し（例えば、図 7 のステップ S 4 3 の処理）、前記入力特徴点における特徴量である入力特徴量を抽出し（例えば、図 7 のステップ S 4 4 の処理）、前記入力特徴量と、前記記憶部に記憶されている前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成し（例えば、図 7 のステップ S 4 5 の処理）、前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢（例えば、モデル姿勢）を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求め（例えば、図 7 のステップ S 4 6 の処理）、前記モデルの姿勢の推定結果、および、前記認識対応特徴点ペア群に基づいて、認識結果を生成する（例えば、図 7 のステップ S 4 7 の処理）ステップを含む。

【 0 0 5 5 】

本発明の第 1 の側面のプログラムは、所定の記憶部に記憶されているアクションを認識するためのモデルを含むモデル動画を画像平面および時間の 3 次元とした場合における、

特徴点であるモデル特徴点と、前記モデル特徴点における特徴量であるモデル特徴量の情報を用いて、入力動画に、登録されている前記アクションが含まれているか否かを認識する処理をコンピュータに実行させるためのプログラムであって、前記入力動画を取得し（例えば、図7のステップS41の処理）、前記入力動画を画像平面および時間の3次元として、前記入力動画から前記アクションを認識するための特徴点である入力特徴点を抽出し（例えば、図7のステップS43の処理）、前記入力特徴点における特徴量である入力特徴量を抽出し（例えば、図7のステップS44の処理）、前記入力特徴量と、前記記憶部に記憶されている前記モデル特徴量とを比較し、類似する特徴量を有する特徴点の組としての候補対応特徴点ペアを生成し（例えば、図7のステップS45の処理）、前記候補対応特徴点ペアから、アウトライヤを除去し、前記入力動画上での前記モデルの姿勢（例えば、モデル姿勢）を推定するとともに、前記モデルの姿勢に対応する認識対応特徴点ペア群を求め（例えば、図7のステップS46の処理）、前記モデルの姿勢の推定結果、および、前記認識対応特徴点ペア群に基づいて、認識結果を生成する（例えば、図7のステップS47の処理）ステップを含む処理をコンピュータに実行させる。

10

【0060】

以下、図を参照して、本発明の実施の形態について説明する。

【0061】

図1に、本発明を適用した認識処理装置11の構成を示す。

【0062】

認識処理装置11は、モデルアクションの登録を行う特徴抽出処理部21と、認識対象となる入力動画を取得して、認識処理を実行する認識処理部22とで構成されている。認識処理装置11は、ユーザが登録した画像シーケンス中のアクション、ジェスチャ、イベントなどの時空間パターンに対して、入力画像シーケンスから、類似した時空間パターンを検出し、検出された場合には、対応点の情報や対応する時刻および対応箇所、対応姿勢やそのパラメータ、または、それらの類似度合いなどを出力することができる。

20

【0063】

以下、画像シーケンス中のアクション、ジェスチャ、イベントなどの時空間パターンを総称してアクションと称するものとする。また、画像シーケンスは、動画または動画像とも称するものとする。

【0064】

ここでは、認識処理装置11として1つの装置であるものとして説明するが、特徴抽出処理部21および認識処理部22が、それぞれ1つの装置として構成されていても良いことは言うまでもない。

30

【0065】

まず、特徴抽出処理部21の各部について説明する。

【0066】

特徴抽出処理部21は、モデル動画記録部41、前処理実行部42、特徴点抽出部43、および、特徴量抽出部44を含んで構成されている。

【0067】

モデル動画記録部41は、認識処理のモデルとなる特徴量を取得するための動画像データを取得し、時間情報とともに記録する。記録される動画像データは、システムに認識させたいアクション（以下、モデルアクションと称する）を含むモデル画像シーケンス（以下、モデル動画とも称する）である。

40

【0068】

モデル動画記録部41は、動画像を撮像可能なカメラを内蔵し、例えば、録画開始終了ボタンのようなユーザインタフェースを用いて、ユーザの指示により、動画像データを取得するものとしても良いし、外部の装置から、有線または無線を介して、モデル動画として用いられる動画像データを取得するものとしても良い。そして、モデル動画記録部41は、例えば、図示しない操作入力部により入力されるユーザの操作入力に基づいて、取得された動画像データのうち、モデルとして用いる部分の開始および終了時刻を設定し、その

50

部分の動画像データをモデル動画として記録するとともに、前処理実行部 4 2 に供給する。

【 0 0 6 9 】

認識処理装置 1 1 においては、画像平面 $x - y$ に対して、時間 t を空間的奥行き方向の次元と見立てることにより、例えば、図 2 に示されるように、画像シーケンスを 3 次元画像として扱うものとする。すなわち、画像シーケンスは、時間（タイムスタンプ t ）と平面（ t 時刻における画像の $x - y$ 平面）とによる 3 次元座標系として捉えることができるため、画像シーケンスの数学的表現として、 $I(x, y, t)$ という表現を用いるものとする。したがって、以下、複数のモデル動画のうちの i 番目のモデル動画を、 $I_{\text{MODEL}}^{[i]}(x, y, t)$ のように表すものとする。

10

【 0 0 7 0 】

前処理実行部 4 2 は、モデル動画中のアクション部分と背景部とを分離する。具体的には、前処理実行部 4 2 は、例えば、モデル動画 $I_{\text{MODEL}}^{[i]}(x, y, t)$ からアクション部と背景部とを分離し、アクション部のピクセルが 1、背景部のピクセルが 0 となったマスク動画を生成することができる。アクション部分と背景部との分離の方法は任意の方法でよい。

【 0 0 7 1 】

前処理実行部 4 2 は、例えば、図示しない操作入力部のマウスポインタデバイスやタッチパッドといった入力インターフェースを用いて、ユーザから、モデル動画の各フレームにおけるアクションの領域を直接選択することができるようにし、その選択領域、すなわち、アクションピクセルが 1、非選択領域、すなわち、背景領域のピクセルが 0 という 2 値画像シーケンスを得ることができるようにしても良い。

20

【 0 0 7 2 】

また、前処理実行部 4 2 は、例えば、図 3 のモデル動画の時刻 t の画像 1 0 1 に対する背景画像 1 0 2、すなわち、モデル動画を撮像するカメラ等の設置場所において、アクションが撮像されていない、環境のみ、つまり背景のみが撮像された画像を取得し、動画の各時刻 t の画像から背景画像を差し引いて得られた背景差分画像 1 0 3 を算出し、この背景差分画像シーケンスを、所定の閾値により 2 値化して、2 値画像シーケンスを得ることができるようにしても良い。

【 0 0 7 3 】

30

また、図 4 に示されるように、上述した閾値を用いた 2 値化処理のみを行った後の画像 1 1 1 には、例えば、図中、白い部分の内側の黒い部分のように、ノイズ部分が残り、アクション部分となるべき範囲内に、背景部分と判定される領域が発生してしまう可能性が高い。そのような場合、例えば、形態学的膨張処理（例えば、8 近傍膨張処理）を行うようにすると好適である。

【 0 0 7 4 】

8 近傍膨張処理は、2 値画像上のある画素 P_0 を注目画素として、注目画素の近傍の 8 画素 P_1 乃至 P_8 中に、少なくとも 1 つアクション部分と判断された画素、すなわち、画素値が 1 である画素がある場合、注目画素 P_0 の画素値を 1、すなわち、アクション部分に変換する処理を行うことにより、画素値 1 の部分、すなわち、アクション部分を膨張させる処理である。この処理を数式で示すと、次の式 (1) となる。

40

【 0 0 7 5 】

$$f(P_0) = f(P_1) \vee f(P_2) \vee \cdots \vee f(P_8) \vee \cdots (1)$$

【 0 0 7 6 】

ここで、 \vee は、論理和（オア）を示し、 $f(P_x)$ は、画素 P_x における画素値を示すものである。

【 0 0 7 7 】

この処理における膨張の度合い、換言すれば、アクション部分としての認識領域の外周を膨らませる度合いを増すためには、以上の処理を所定の複数回繰り返すようにしても良いし、注目画素の近傍 8 画素よりも広い範囲の近傍画素のうち、少なくとも 1 つアクショ

50

ン部分と判断された画素、すなわち、画素値が1である画素がある場合、注目画素の画素値を1、すなわち、アクション部分に変更するものとしても良い。

【0078】

上述したように、ユーザの選択により得られた、または、背景差分画像103により得られた2値画像、もしくは、その2値画像に対して必要に応じて8近傍膨張処理を施した2値画像を、マスク画像1と称し、画像 $I_{MASK1}^{[i]}(x, y, t)$ と表すものとする。

【0079】

前処理実行部42は、図5に示されるように、マスク画像1に対応する画像 $I_{MASK1}^{[i]}(x, y, t)$ 、すなわち、中央値フィルタリング後の画像121を用いて、モデル動画とマスク動画1を乗算することで、背景がマスクアウトされ、背景のピクセルが0、それ以外が元の画素値となった、すなわち、背景分離処理が行われたモデル動画を得ることができる。このとき得られるマスキング済みの画像は、例えば、同じく図5に示される、背景分離処理後のモデル動画の時刻tの画像122となる。

【0080】

ここでは、特に区別する必要がない場合には、背景分離処理が行われたモデル動画についても、 $I_{MODEL}^{[i]}(x, y, t)$ のように表すものとする。

【0081】

特に、上述した8近傍膨張処理のように、アクション部分に対応する領域の膨張処理が施されたとき、背景分離処理後のモデル動画の時刻tの画像122は、実際のアクション部分（ここでは、人物の顔部分）の周辺の、本来ならば背景として分離されなくてはならない部分を含んでしまう可能性がある。この背景分離処理後のモデル動画の時刻tの画像122を用いて後述する特徴点抽出処理および特徴量抽出処理が行われた場合、背景に対応する部分の特徴点および特徴量が抽出されてしまう恐れがある。

【0082】

そこで、前処理実行部42は、上述したようにして得られた2値画像シーケンスのそれぞれのフレームに対して、形態学的収縮処理（例えば、8近傍収縮処理）を施すようにしても良い。形態学的収縮処理が施されたマスク画像を用いてマスキングした画像を利用して特徴点の抽出を行うことにより、得られる特徴点は少なくなってしまう可能性があるが、背景に対応する部分の特徴点および特徴量が抽出されてしまう恐れを限りなく除去することができ、認識の精度が向上する。

【0083】

形態学的収縮処理の具体的な例として、8近傍収縮処理について説明する。

【0084】

8近傍収縮処理とは、2値画像上のある画素 P_0 を注目画素とし、注目画素 P_0 に対する近傍8画素 P_1 乃至 P_8 中に少なくとも1つの背景画素、すなわち、画素値が0である画素がある場合、 P_0 の画素値を背景画素値、すなわち0に変更する処理である。この処理を数式で示すと、次の式(2)となる。

【0085】

$$f(P_0) = f(P_1) \cdot f(P_2) \cdot \dots \cdot f(P_8) \cdot \dots \quad (2)$$

【0086】

ここで、 \cdot は、論理積（アンド）を示し、 $f(P_x)$ は、画素 P_x における画素値を示すものである。

【0087】

この処理における収縮の度合い、換言すれば、アクション部分としての認識領域の外周を狭める度合いを増すためには、以上の処理を所定回数繰り返すようにしても良いし、注目画素の近傍8画素よりも広い範囲の近傍画素のうち少なくとも1つ背景部分と判断された画素（画素値が0である画素）がある場合、注目画素の画素値を0（すなわち、背景部分）に変えるものとしても良い。

【0088】

前処理実行部42は、マスク画像1（ $I_{MASK1}^{[i]}(x, y, t)$ ）に対して、形態学的

収縮処理（例えば 8 近傍収縮処理）を施し、マスク画像 2（ $I_{\text{MASK}2}^{[i]}(x, y, t)$ と表される）を生成する。

【0089】

前処理実行部 42 は、上述した 8 近傍膨張処理のような、アクション部分に対応する領域の膨張処理と、形態学的収縮処理（例えば、8 近傍収縮処理）との、一件相反する処理を画像に対して施す。これは、例えば、図 4 を用いて説明した閾値処理後の画像 111 において、図中、白い部分の内側の黒い部分のように発生するノイズ部分を除去しつつ、背景に対応する部分の特徴点および特徴量が抽出されてしまう恐れを限りなく除去するために、非常に有用な処理である。例えば、図 4 を用いて説明した閾値処理後の画像 111 に直接形態学的収縮処理を施した場合、図 4 中、白い部分の内側の黒い部分のように発生するノイズを除去することができず、誤検出の原因となってしまう。これに対して、図 4 の形態学的膨張後の画像 112 のように、白い部分の内側の黒い部分のように発生するノイズを除去した、換言すれば、白い部分の内側の黒い部分の値を全て 1（白い部分）としたあとに形態学的収縮処理を施した場合、収縮によって再びノイズ部分が発生ことはない。

【0090】

なお、前処理実行部 42 が省略され、前処理が実行されない場合、後述する処理が実行不可になり、モデル動画の特徴点および特徴量の抽出処理ができなくなるものではない。すなわち、前処理実行部 42 が省略され、前処理が実行されない場合には、背景部分に対しても、モデル動画の特徴点および特徴量の抽出処理が行われてしまうため、処理時間が長くなり、これらを用いた認識処理の認識精度も落ちてしまうことが考えられるが、背景部分を含む動画像に対して、モデル動画の特徴点および特徴量の抽出処理を行うことは可能である。

【0091】

図 1 に戻って、認識処理装置 11 の各部の説明を続ける。

【0092】

特徴点抽出部 43 は、モデル動画から特徴点の抽出を行う。特徴点抽出部 43 は、既に公知であるさまざまな手法のうちのいずれの手法を用いて特徴点を抽出するものとしても良いが、時空間、すなわち、図 2 を用いて説明した画像平面 $x - y$ に対して、時間 t を空間的奥行き方向の次元とした 3 次元の変形に対して頑強な特徴点抽出法を利用すると、認識の精度があがり、好適である。

【0093】

特徴点抽出部 43 が行う特徴点抽出のために用いる手法の具体的な例として、"I. Laptev, "On Space-Time Interest Points", in International Journal of Computer Vision, vol 64, number 2/3, 2005" に記載された技術を用いた特徴点抽出方法について説明する。以下、この特徴点抽出法で抽出される特徴点を ST（Spatio-Temporal）特徴点、または単に特徴点と称する。

【0094】

ST 特徴点は、時空間、すなわち、図 2 を用いて説明した画像平面 $x - y$ に対して、時間 t を空間的奥行き方向の次元とした 3 次元に拡張を行った一般化 Harris 尺度に基づいて検出される特徴点であり、式 (3) に示す 3 次元拡張 Harris 関数 H の極大および極小を与える 3 次元座標 (x, y, t) として定義される。

【0095】

$$H = \det(\mu) - k \cdot \text{trace}^3(\mu) \quad \dots (4)$$

【0096】

式 (3) における $\det(\mu)$ は、正方行列の行列式を示し、 $\text{trace}^3(\mu)$ は、行列の対角成分の和の 3 乗を示し、 k は、定数である。

【0097】

そして、式 (3) の μ は、次の式 (4) で与えられる。

【0098】

$$\mu(x, y, t; \sigma_x^2, \sigma_y^2) = G(x, y, t; \sigma_x^2, \sigma_y^2) * (L(L)^T)$$

10

20

30

40

50

・・・(4)

【0099】

式(4)中の $G(x, y, t; \sigma^2, \tau^2)$ は、3次元ガウスフィルタであり、次の式(5)で与えられる。

【0100】

$$G(x, y, t; \sigma^2, \tau^2) = \frac{1}{\sqrt{\pi}(\sigma^2)^3} \exp\left(-\frac{(x^2 + y^2)}{2\sigma^2} - \frac{t^2}{2\tau^2}\right) \quad \dots (5)$$

【0101】

そして、式(4)および式(5)において、 σ は、空間領域におけるガウス形状(裾野の広がり)を決めるパラメータであり、 τ は、時間領域におけるガウス形状を決めるパラメータである。換言すれば、 σ は、空間領域におけるローパスフィルタによる値のぼかし度合いに対応するパラメータであり、 τ は、時間領域におけるローパスフィルタによる値のぼかし度合いに対応するパラメータである。ここで、 σ および τ は、画像を取得するカメラデバイスなどの解像度によって最適な値を用いると好適であり、例えば、 $\sigma = 8$ 、 $\tau = 4$ 程度の値とすることができる。

【0102】

また、式(4)の $*$ は内積を示す。すなわち、式(4)の右辺は、 $(L_x(x, y, t; \sigma_L^2, \tau_L^2))^T$ を3次元ガウスフィルタ G でぼかすオペレーションであり、式(4)中の L は、次の式(6)で表される。

【0103】

【数1】

$$\nabla L = (L_x(x, y, t; \sigma_L^2, \tau_L^2), L_y(x, y, t; \sigma_L^2, \tau_L^2), L_t(x, y, t; \sigma_L^2, \tau_L^2))^T \quad \dots (6)$$

【0104】

そして、式(6)において、以下の式(7)乃至式(9)が成立する。

【0105】

$$L_x(x, y, t; \sigma_L^2, \tau_L^2) = x(G(x, y, t; \sigma_L^2, \tau_L^2)) * I(x, y, t) \quad \dots (7)$$

$$L_y(x, y, t; \sigma_L^2, \tau_L^2) = y(G(x, y, t; \sigma_L^2, \tau_L^2)) * I(x, y, t) \quad \dots (8)$$

$$L_t(x, y, t; \sigma_L^2, \tau_L^2) = t(G(x, y, t; \sigma_L^2, \tau_L^2)) * I(x, y, t) \quad \dots (9)$$

【0106】

すなわち、式(6)の L は時空間画像グラディエントを表しており、式(4)の μ は時空間画像グラディエントの2次モーメントマトリクスを示している。

【0107】

そして、パラメータ σ_L は、空間スケールパラメータであり、 τ_L は、時間スケールパラメータであり、それぞれ、特徴点を抽出する際に考慮に入れるべき時間または空間(画像の x - y 平面)方向のサイズを決めるパラメータとなっている。なお、 σ_L と τ_L とは、それぞれ独立に決められる。ここでは、例えば、 $\sigma_L = \{2, 4, 8\}$ 、 $\tau_L = \{2, 4, 8\}$ の、全9通りの組み合わせで S 、 T 特徴点の検出を行うこともできる。また、例えば、 $\sigma_L = \{2, 4, 8\}$ 、 $\tau_L = \{2, 4, 8\}$ などとするにより、時間的にも、 x - y 平面内の空間的にもスケールに幅を持つことができ、広い範囲から特徴点を拾うことが可能となる。

【0108】

特徴点抽出部43は、供給されたモデル動画 $I_{\text{MODEL}}^{[i]}(x, y, t)$ を用いて式(3)を計算して、その極大値および極小値を与える座標 (x, y, t) を検出し、 S 、 T 特徴

10

20

30

40

50

点とする。

【 0 1 0 9 】

このとき、特徴点抽出部 4 3 は、前処理として背景がマスクされた動画に対して S T 特徴点の抽出処理を行う。このとき、上述したマスク動画 1 を利用するものとしても良いが、特に、マスク動画 2 を用いてマスクされた動画を用いるようにした場合、または、抽出された S T 特徴点のうち、特徴点位置 (x , y , t) におけるマスク動画 2 の値が 1 の特徴点のみを、モデルアクションを記述する有効な特徴点とした場合、上述したように、背景部分から特徴点を抽出する可能性を大幅に排除することができ、好適である。

【 0 1 1 0 】

取得されたモデル動画において、アクションを起こした人体や物体によって背景は不連続的に隠される。この部分は画像情報が極端に変化するため、上述した S T 特徴点が検出されてしまう可能性が高い。また、ある特徴点に対して求められる特徴量は、その特徴点近傍領域の画像情報から計算されるため、特徴点が背景領域に近い点である場合、背景部分の情報を含むことになり、入力動画とアクション部において、画像平面 x - y と時間 t との 3 次元座標位置が同じ対応特徴点であっても、背景が微妙に異なることで特徴量のマッチングが取れなくなってしまう。すなわち、このような特徴点は、視点変化や背景変化に頑強でない特徴点である。

【 0 1 1 1 】

そこで、マスク動画 2 を用いたマスク処理により、背景部とアクション部の境界付近に検出された、視点変化や背景変化に頑強でない特徴点を排除する処理を入れることで、認識性能が向上する。

【 0 1 1 2 】

以上の処理により特徴点抽出部 4 3 において求められたモデル動画 $I_{MODEL}^{[i]}(x, y, t)$ の N 個の S T 特徴点を、 $P_{MODEL}^{[i]} = \{P_1^{[i]}, P_2^{[i]}, \dots, P_N^{[i]}\}$ と表するものとする。特徴点抽出部 4 3 は、N 個の S T 特徴点 $P_{MODEL}^{[i]} = \{P_1^{[i]}, P_2^{[i]}, \dots, P_N^{[i]}\}$ を、特徴量抽出部 4 4 に供給する。

【 0 1 1 3 】

特徴量抽出部 4 4 は、特徴点抽出部 4 3 により抽出されて供給されたモデル動画 $I_{MODEL}^{[i]}(x, y, t)$ の特徴点における特徴量を抽出する。一般的にアクションイベント認識において利用される特徴量は、例えば、“I. Laptev and T. Lindeberg, “Local Descriptors for Spatio-Temporal Recognition”, in ECCV Workshop “Spatial Coherence for Visual Motion Analysis”, Springer LNCS Vol.3667, pp. 91-103, 2004 ” など、いくつか提案されており、さまざまな方法の特徴量抽出部 4 4 の特徴量抽出処理に適用することが可能である。ここでは、その一例として、“I. Laptev, “On Space-Time Interest Points”, in International Journal of Computer Vision, vol 64, number 2/3, 2005 ” に記載されている技術を用いた特徴量の抽出について説明する。

【 0 1 1 4 】

特徴点抽出部 4 3 により抽出された S T 特徴点のうちのある S T 特徴点 $P = (x_p, y_p, t_p)$ における時空間特徴量 V_p は、次の式 (1 0) で定義される。

【 0 1 1 5 】

$$V_p = \{ \begin{matrix} L_{x, p} & L_{y, p} & L_{t, p} & L_{xx, p}^2 & \dots & L_{y, p}^3 & L_{y, p}^4 & L_{ttt, p} \\ L_{x, p} & L_{y, p} & L_{t, p} & L_{xx, p}^2 & \dots & L_{y, p}^3 & L_{y, p}^4 & L_{ttt, p} \end{matrix} \dots (10)$$

【 0 1 1 6 】

ここで、L に下付の x y z がついているものは、次の式 (1 1) に対応するものであり、パラメータ p は P が検出されたときのスケールパラメータである。

【 0 1 1 7 】

$$L_{x, y, t}^{m, n, k}(x_p, y_p, t_p; p^2, p^2) = \dots (11)$$

10

20

30

40

50

【0118】

式(11)において、 m は、式(10)における x の次数であり、 n は、式(10)における y の次数であり、 k は、式(10)における t の次数である。

【0119】

すなわち、時空間特徴量 V_p は、 x 、 y 、 t のそれぞれの次元について、4次までの偏微分ガウスオペレーションをかけた画像情報から構成される特徴ベクトルであり、その次元数は、 ${}_3C_2 + {}_4C_2 + {}_5C_2 + {}_6C_2 = 34$ の34次元となる。

【0120】

以上の処理で求めた i 番目のモデル動画 $I_{MODEL}^{[i]}(x, y, t)$ の各ST特徴点(ST特徴点の総数を N とする)で求めた特徴量を、ST特徴点 $P_j^{[i]}$ の特徴量が $V_j^{[i]}$ であるものとして、 $V_{MODEL}^{[i]} = \{V_1^{[i]}, V_2^{[i]}, \dots, V_N^{[i]}\}$ と表記する。

10

【0121】

特徴量抽出部44は、モデル動画 $I_{MODEL}^{[i]}$ から抽出されたST特徴点 $P_{MODEL}^{[i]}$ とその特徴量 $V_{MODEL}^{[i]}$ を、モデルアクション $MODEL^{[i]} = (P_j^{[i]}, V_j^{[i]})$ として、認識処理部22の辞書登録部61に供給する。

【0122】

次に、認識処理部22の各部について説明する。

【0123】

認識処理部22は、辞書登録部61、入力動画バッファ部62、前処理実行部63、特徴点抽出部64、特徴量抽出部65、特徴量比較部66、姿勢パラメータ推定部67、および、認識結果生成部68を含んで構成されている。

20

【0124】

認識処理部22の辞書登録部61は、特徴抽出処理部21の特徴量抽出部44から供給されたモデルアクション $MODEL^{[i]} = (P_j^{[i]}, V_j^{[i]})$ を、認識処理時に参照可能な形で保存する。

【0125】

入力動画バッファ部62は、入力動画データ(以下、入力動画とも称する)を取得し、バッファリングする。入力動画データは、辞書登録部61に保存されているモデルアクション $MODEL^{[i]} = (P_j^{[i]}, V_j^{[i]})$ を用いて所定のアクションの有無を認識するための認識対象である。入力動画バッファ部62にバッファされる入力動画を、 $I_{INPUT}(x, y, t)$ のように表すものとする。

30

【0126】

入力動画バッファ部62は、動画を撮像可能なカメラを内蔵し、例えば録画開始終了ボタンのようなユーザインタフェース用いて、ユーザの指示により、動画データを取得するものとしても良いし、外部の装置から、有線または無線を介して、動画データを取得するものとしても良い。

【0127】

認識処理部22における認識処理が逐次認識の形態をとる場合には、入力動画バッファ部62は、少なくとも、最新フレームから認識対象時間長(所定フレーム数)さかのぼったフレームまでの画像シーケンスを入力動画としてバッファする。また、入力動画バッファ部62は、例えば、録画開始終了ボタンのようなユーザインタフェース用いたユーザの指示により、所定時間の認識対象画像シーケンスを入力動画としてバッファする構成であっても良い。

40

【0128】

前処理実行部63は、入力動画中のアクション部分と背景部とを分離する。分離の方法は任意の方法でよく、例えば、逐次認識ではない場合においては、特徴抽出処理部21の前処理実行部42と同様の方法を用いてアクション部分と背景部を分離するものであっても良い。また、前処理実行部63を省略しても、認識処理部22における認識処理は実行可能である。

【0129】

50

すなわち、前処理を行わない背景部を含んだ入力動画から特徴点および特徴量の抽出を行っても、特徴点および特徴量の抽出時間、並びに、特徴量の比較処理にかかる時間が増えてしまうが、辞書登録部 6 1 に登録されているモデルアクション $MODEL^{[i]} = (P_j^{[i]}, V_j^{[i]})$ と最終的には一致しないので、認識処理は実行可能である。具体的には、後述する処理により特徴量の比較が行われ、対応する特徴量のペアが生成されるが、前処理が行われない場合、多くの誤った特徴量のペアが生成されてしまう恐れがある。しかしながら、後述する処理によって誤った特徴量のペアのほとんどはアウトライヤとして除去されることが期待されるため、前処理が行われなくても、認識処理は正しく実行される。もちろん、前処理を行ったほうが、処理時間が短縮され、認識精度も更になるので好適である。

10

【0130】

また、ここでは、前処理実行部 6 3 により前処理が行われているか否かにかかわらず、特徴点抽出部 6 4 に供給される入力動画も、 $I_{INPUT}(x, y, t)$ のように表すものとする。

【0131】

特徴点抽出部 6 4 は、特徴点抽出部 4 3 と同様の方法を用いて、入力動画 $I_{INPUT}(x, y, t)$ から特徴点の抽出を行い、抽出された特徴点の情報を、特徴量抽出部 6 5 に供給する。特徴点抽出部 6 4 により抽出された入力動画 $I_{INPUT}(x, y, t)$ の M 個の特徴点 (ST 特徴点) を、 $Q_{INPUT} = \{Q_1, Q_2, \dots, Q_M\}$ と表す。

20

【0132】

特徴量抽出部 6 5 は、特徴点抽出部 6 4 により抽出された入力動画 $I_{INPUT}(x, y, t)$ の特徴点 $Q_{INPUT} = \{Q_1, Q_2, \dots, Q_M\}$ の各点において、上述した特徴量抽出部 4 4 と同様の方法を用いて、特徴量の抽出を行う。特徴量抽出部 6 5 により抽出された入力動画 $I_{INPUT}(x, y, t)$ の M 個の特徴点 $Q_{INPUT} = \{Q_1, Q_2, \dots, Q_M\}$ の各点における特徴量を、特徴点 Q_k の特徴量が W_k であるものとして、 $W_{INPUT} = \{W_1, W_2, \dots, W_M\}$ と表記する。

【0133】

特徴量抽出部 6 5 は、この入力動画の特徴点特徴量情報 $INPUT = (Q_k, W_k)$ (ここで、k は、1 以上 M 以下の整数) を参照可能な形でバッファするとともに、特徴量比較部 6 6 に供給する。

30

【0134】

特徴量比較部 6 6 は、特徴量抽出部 6 5 から供給された特徴点特徴量情報 $INPUT = (Q_k, W_k)$ と、辞書登録部 6 1 に登録されているモデルアクション $MODEL^{[i]} = (P_j^{[i]}, V_j^{[i]})$ とのマッチング処理を行う。

【0135】

例えば、辞書登録部 6 1 に L 個のモデルアクションが登録されているものとする。すなわち、モデルアクション $MODEL^{[i]} = (P_j^{[i]}, V_j^{[i]})$ において、 $1 \leq i \leq L$ であるものとする。

【0136】

特徴量比較部 6 6 は、特徴量抽出部 6 5 から供給された特徴点特徴量情報 $INPUT = (Q_k, W_k)$ と、辞書登録部 6 1 に登録されているモデルアクション $MODEL^{[i]} = (P_j^{[i]}, V_j^{[i]})$ とで、類似度の高い特徴量のペア群の抽出を行う。抽出された特徴量のペア群に対応する特徴点のペア群を、候補対応特徴点ペア群と称するものとする。

40

【0137】

類似度の高い特徴点のペア群の抽出を行うために用いる類似尺度、または、非類似尺度には、様々なものを用いることができる。ここでは、その一例として、非類似度として任意のノルムを用いる場合について説明する。

【0138】

i 番目のモデル動画の j 番目の特徴点の特徴量 $V_j^{[i]}$ と入力動画の k 番目の特徴点の特徴量 W_k との非類似度 $D(V_j^{[i]}, W_k)$ を、次の式 (12) で定義する。

50

【0139】

$$D(V_j^{[i]}, W_k) = \text{norm}(V_j^{[i]}, W_k)$$

(12)

【0140】

特徴量比較部66は、全ての $V_j^{[i]}$ と W_k との組み合わせにおいて、式(12)で定義される非類似度 $D(V_j^{[i]}, W_k)$ を演算する。そして、特徴量比較部66は、式(12)で定義される非類似度 $D(V_j^{[i]}, W_k)$ の値に基づいて、特徴点 $P_j[i]$ に対する候補対応特徴点群を、例えば、 $D(V_j^{[i]}, W_k)$ が最も小さい K 個(例えば、 $K=3$ 程度の、複数であって、かつ、あまり大きくない値であると好適である)の特徴量 W_k に対応する特徴点 Q_k を、 $P_j^{[i]}$ の候補対応特徴点群とすることができる。また、特徴量比較部66は、式(12)で定義される非類似度 $D(V_j^{[i]}, W_k)$ の値が、所定の閾値を下回る全ての特徴量 W_k に対応する特徴点 Q_k を、 $P_j^{[i]}$ の候補対応特徴点群とすることも可能である。

10

【0141】

ここで、特徴量比較部66により得られる候補対応特徴点群における $P_j^{[i]}$ と各対応特徴点 Q_a をペアにして、候補対応特徴点ペア $[P_j^{[i]}, Q_a]$ のように表現し、入力動画と i 番目のモデルアクションに関する候補対応特徴点ペア群を $CM P^{[i]}$ と表記するものとする。すなわち、ペア群 $CM P^{[i]} = \{ (P_j^{[i]}, Q_a) \mid Q_a : P_j^{[i]} \text{の対応特徴点} \}$ となり、このとき、 i は1から L までの整数、 j は1から N までの整数となる。

20

【0142】

特徴量比較部66は、上述したような処理により得られた候補対応特徴点ペア群 $CM P^{[i]}$ の情報を、姿勢パラメータ推定部67に供給する。

【0143】

姿勢パラメータ推定部67は、特徴量比較部66により得られる候補対応特徴点ペア群 $CM P^{[i]}$ のアウトライヤ除去を施した後、各モデルアクション検出の有無の判定、および、検出有りのモデルに対するモデルアクションの姿勢パラメータ推定を行う。

【0144】

特徴量比較部66により得られる候補対応特徴点ペア群 $CM P^{[i]}$ の抽出処理においては、特徴量を抽出した特徴点の位置情報は使っていないため、巨視的に見ると、候補対応特徴点ペア群 $CM P^{[i]}$ には、対応特徴点間の位置関係が、モデルアクションの入力動画上での姿勢(モデル姿勢)と矛盾しない真の対応特徴点ペア(インライヤ)だけでなく、局所的な画像情報から得られたいずれかの特徴量に関して類似しているが、時空間的な幾何学的配置という視点から見ると対応しないような偽の対応特徴点ペア(アウトライヤ)も多数混在している。また、上述したように、前処理実行部63による前処理が省略される場合、アウトライヤの混在可能性が高くなる。すなわち、特徴量比較部66により得られる候補対応特徴点ペア群 $CM P^{[i]}$ を全て利用して、モデルアクションの入力動画中の存在有無の判定とモデルアクションの入力動画中の姿勢推定を行うようにした場合、アウトライヤの混在により、認識結果が著しく悪くなる。

30

【0145】

そこで、「モデルアクションは、入力動画中に時空間画像変換(つまり3次元画像変換)されて出現する」という時空間的な変換仮定を立てることにより、姿勢パラメータ推定部67は、候補対応特徴点ペア群 $CM P^{[i]}$ の中からもっとも正しそうな時空間画像変換パラメータ、および、それを決める候補対応特徴点ペア群 $CM P^{[i]}$ の部分集合を求め、その部分集合を、最終的な認識結果を計算する認識対応特徴点ペア群 $RM P^{[i]}$ とする。

40

【0146】

2次元静止画像における変換仮説の概念は、例えば、特開2006-065399に記載されているが、ここでは、それを、単純に2次元(平面)から3次元(空間)に拡張するのではなく、2次元と時間とによる時空間、すなわち、画像平面 $x-y$ に対して、時間 t を空間的奥行き方向の次元とした3次元に拡張することにより、動画に適用する。

【0147】

50

姿勢パラメータ推定部 67 が実行する画像変換としては、いずれも 3 次元画像変換に拡張された、ユークリッド変換、相似変換、アフィン変換、または、射影変換などを用いることができる、ここでは、その一例として、姿勢パラメータ推定部 67 が、3 次元アフィン変換の拘束の下、姿勢推定を行う場合を例として、詳細に説明する。

【0148】

姿勢パラメータ推定部 67 により実行される 3 次元アフィン変換は、 $x - y$ の 2 次元と時間 t とによる 3 次元において、平行移動および回転変換（ユークリッド変換）に、拡大縮小変換を加えた相似変換に、せん断変形を許すような変換であり、元の図形で直線上に並ぶ点は変換後も直線上に並び、平行線は変換後も平行線であるなど、幾何学的性質が保たれる変換である。変換前の点の座標を (x, y, t) 、変換後の点の座標を (x', y', t') とすると、3 次元アフィン変換は、次の式 (13) で示される。

10

【0149】

【数 2】

$$\begin{bmatrix} x' & y' & t' & 1 \end{bmatrix} = \begin{bmatrix} x & y & t & 1 \end{bmatrix} \begin{bmatrix} a_1 & a_4 & a_7 & 0 \\ a_2 & a_5 & a_8 & 0 \\ a_3 & a_6 & a_9 & 0 \\ b_1 & b_2 & b_3 & 1 \end{bmatrix} \quad \cdots (13)$$

【0150】

20

候補対応特徴点ペア群 $CM P^{[i]}$ からその一部を抽出することにより、候補対応特徴点ペア群 $CM P^{[i]}$ の部分集合である、ある対応特徴点ペア群 MP を抽出したとき、抽出された対応特徴点ペア群 MP について、対応特徴点ペア群 MP に含まれるモデルアクション特徴量の特徴点 P の座標を (x_s, y_s, t_s) 、それに対応した入力動画特徴量の特徴点 Q の座標を (x'_s, y'_s, t'_s) （ここで、 s は、 $1 \leq s \leq$ 対応特徴点ペア群 MP 中の対応特徴点ペア数）とすると、対応特徴点ペア群 MP から推定される、 $x - y - t$ の 3 次元のアフィンパラメータは、次の式 (14) から求められる。

【0151】

【数 3】

$$\begin{bmatrix} x_1 & y_1 & t_1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & x_1 & y_1 & t_1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & x_1 & y_1 & t_1 & 0 & 0 & 1 \\ x_2 & y_2 & t_2 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & x_2 & y_2 & t_2 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & x_2 & y_2 & t_2 & 0 & 0 & 1 \\ & & & & & & & & & & & \vdots \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_6 \\ a_7 \\ a_8 \\ a_9 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} x'_1 \\ y'_1 \\ t'_1 \\ x'_2 \\ y'_2 \\ t'_2 \\ \vdots \end{bmatrix} \quad \cdots (14)$$

30

40

【0152】

式 (14) において、 $a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9$ は、回転、拡大縮小、および、せん断変形のそれぞれを決定するパラメータを表し、 b_1, b_2, b_3 は、平行移動を決定するためのパラメータを表す。

【0153】

50

そして、式(14)の変数を、次の式(15)に示されるように置き換えたとき、その最小二乗解は、次の式(16)で表される。

【0154】

【数4】

$$C\Theta = d \quad \dots(15)$$

【数5】

$$\Theta = [C^T C]^{-1} C^T d \quad \dots(16)$$

【0155】

10

ここで、式(14)に示される決定パラメータ数が12であることから、2次元と時間とによる3次元アフィン変換パラメータを決定するためには、対応特徴点ペア群MP中に対応特徴点ペアが4組以上必要となる。よって、候補対応特徴点ペア群CMP^[i]に対応特徴点ペアが3組以下しか存在しない場合、姿勢パラメータ推定部67は、対応特徴点ペアが3組以下しか存在しないことを認識結果生成部68に通知するので、認識結果生成部68は、i番目のモデルアクションに対して非検出という認識をする。

【0156】

これに対して、候補対応特徴点ペア群CMP^[i]に対応特徴点ペアが4組以上存在する場合、i番目のモデルアクションは、入力動画に対して少なくともある程度は対応するという結果を得ることができる。

20

【0157】

候補対応特徴点ペア群CMP^[i]から、ランダムに対応特徴点ペア群Rを選択し、その対応特徴点ペア群Rにアウトライヤが1つ以上混入していた場合、その3次元画像変換パラメータは、パラメータ空間上に散らばって投射される。一方、ランダムに対応特徴点ペア群Rを選択し、その対応特徴点ペア群Rがインライヤのみから構成されていた場合、その3次元画像変換パラメータは、パラメータ空間上で距離の近い範囲にまとまって投影される。すなわち、インライヤである対応特徴点ペアでは、何れもモデルアクションの入力動画中の姿勢の真のアフィン変換パラメータに極めて類似したものであるため、その3次元画像変換パラメータの投影先は、パラメータ空間上で距離の近いものとなる。

【0158】

30

したがって、候補対応特徴点ペア群CMP^[i]から、ランダムに対応特徴点ペア群Rを選択し、その3次元画像変換パラメータをパラメータ空間上に投射していく処理を繰り返すと、インライヤはパラメータ空間上で密度の高い(メンバ数の多い)クラスタを形成し、アウトライヤは散らばって出現することになる。すなわち、パラメータ空間上でクラスタリングを行えば、最多メンバ数を持つクラスタの要素を認識することができる。そして、姿勢パラメータ推定部67は、このクラスタ内の要素をインライヤと認識することができる。

【0159】

姿勢パラメータ推定部67は、パラメータ空間上におけるクラスタリング手法として、NN(Nearest Neighbor)法を用いることができる。

40

【0160】

姿勢パラメータ推定部67は、候補対応特徴点ペア群CMP^[i]からランダムに4組以上のペアに対応特徴点ペア群R₁として選択し、上述した式(14)乃至式(16)を用いて、3次元アフィン変換パラメータ_{R1}を求め、パラメータ空間に投射する。姿勢パラメータ推定部67は、クラスタ数を表す変数NZをNZ=1とし、3次元アフィン変換パラメータ空間上で3次元アフィン変換パラメータ_{R1}をセントロイドとするクラスタZ₁を作る。具体的には、姿勢パラメータ推定部67は、このクラスタZ₁のセントロイドC₁をC₁=_{R1}とし、クラスタのメンバ数を表す変数nz₁をnz₁=1とする。

【0161】

そして、姿勢パラメータ推定部67は、候補対応特徴点ペア群CMP^[i]からランダム

50

に4組以上のペアを次の対応特徴点ペア群 R_2 として選択し、上述した式(14)乃至式(16)を用いて、3次元アフィン変換パラメータ R_2 を求め、パラメータ空間に投射する。そして、姿勢パラメータ推定部67は、NN法によりアフィン変換パラメータ空間をクラスタリングする。姿勢パラメータ推定部67は、クラスタリングの結果、新たなクラスタが発生した場合、そのクラスタを新たなクラスタ Z_2 とし、新たなクラスタが発生しなかった場合、クラスタ C_1 のメンバ数を $n_{z_1} = 2$ とする。

【0162】

そして、姿勢パラメータ推定部67は、所定の条件が満たされるまで、候補対応特徴点ペア群 $CM P^{[i]}$ からランダムに4組以上のペアを選択し、上述した式(14)乃至式(16)を用いて、3次元アフィン変換パラメータを求め、パラメータ空間に投射し、NN法によりアフィン変換パラメータ空間をクラスタリングする。

10

【0163】

クラスタリングについて具体的に説明すると、姿勢パラメータ推定部67は、次の式(17)に従って、3次元アフィン変換パラメータ R_{cnt} (cntは何回目の処理であることを示す変数)と、各クラスタ Z_g (g は、 $1 \leq g \leq N_Z$ となる値)のセントロイド C_g (g は、 $1 \leq g \leq N_Z$ となる値)との距離 $d(R_{cnt}, C_g)$ のうち、最小の距離 d_{min} を求める。

【0164】

$$d_{min} = \min \{ d(R_{cnt}, C_g) \} \quad (17)$$

【0165】

20

ここで、クラスタリング規範となる距離 $d(R_{cnt}, C_g)$ として、例えば、ユークリッド距離を用いることができ、セントロイド C_g として、クラスタメンバの平均ベクトルを用いることができる。

【0166】

そして、所定の閾値 E に対して $d_{min} < E$ であれば、姿勢パラメータ推定部67は、 d_{min} を与えるクラスタ Z_g に R_{cnt} を属させ、 R_{cnt} を含めた全メンバでクラスタ Z_g のセントロイド C_g を更新する(クラスタ Z_g のメンバ数 n_{z_g} は、1インクリメントされる)。一方、 $d_{min} \geq E$ であれば、姿勢パラメータ推定部67は、3次元アフィン変換パラメータ空間上で3次元アフィン変換パラメータ R_{cnt} をセントロイド C_{g+1} とする新しいクラスタ Z_{g+1} を作り、そのクラスタのメンバ数 n_{z_g} を $n_{z_{g+1}} = 1$ とし、クラスタ数 N_Z を $N_Z = N_Z + 1$ とする。

30

【0167】

そして、姿勢パラメータ推定部67は、所定の条件が満たされたか否かを判断する。所定の条件とは、例えば、最多メンバ数が所定の閾値(例えば15)を超え、かつ、最多メンバ数と2番目に多いメンバ数との差が所定の閾値(例えば3)を超える場合、または、処理の繰り返し回数が、所定の閾値(例えば5000回)を超える場合などである。所定の条件が満たされなかった場合、姿勢パラメータ推定部67は、繰り返し回数を計数するカウンタを1インクリメントし、候補対応特徴点ペア群 $CM P^{[i]}$ からランダムに4組以上のペアを選択し、上述した式(14)乃至式(16)を用いて、3次元アフィン変換パラメータを求め、パラメータ空間に投射し、NN法によりアフィン変換パラメータ空間をクラスタリングする処理を繰り返す。一方、所定の条件が満たされた場合、姿勢パラメータ推定部67は、最多メンバ数を持つクラスタ Z_{max} のメンバである対応特徴点ペア群を i 番目のモデルアクションに対する認識対応特徴点ペア群 $RM P^{[i]}$ として保持するとともに、クラスタ Z_{max} のセントロイド C_{max} を認識姿勢パラメータ $\tilde{p}^{[i]}$ として保持する。

40

【0168】

姿勢パラメータ推定部67は、 L 個のモデルアクションに対して、順次、処理を実行する。そして、 L 個のモデルアクションそれぞれの認識対応特徴点ペア群 $RM P^{[i]}$ と認識姿勢パラメータ $\tilde{p}^{[i]}$ とを、認識結果生成部68に供給する。なお、いずれかのモデルアクションにおいて、候補対応特徴点ペア群 $CM P^{[i]}$ に対応特徴点ペアが3組以下しか存在しない場合、姿勢パラメータ推定部67は、認識対応特徴点ペア群 $RM P^{[i]} = NU$

50

LL、認識姿勢パラメータ $\sim^{[i]}$ = NULLとして、認識結果生成部68に供給する。

【0169】

なお、上述したクラスタリング手法では、クラスタリング規範となる距離尺度 $d(\text{Rcnt}, C_g)$ として、たとえばユークリッド距離を用い、セントロイドとして、クラスタメンバーの平均ベクトルを用いる構成として説明したが、距離尺度として、クラスタの分散を考慮したマハラノビス距離を用い、セントロイドとしてクラスタメンバーの平均ベクトル及び分散の情報を用いる構成を用いることも可能であることはいうまでもない。

【0170】

さらに、上記クラスタリング手法では、3次元アフィン変換パラメータ Rcnt の12次元全ての次元について同じ重みでクラスタとパラメータベクトルとの距離 $d(\text{Rcnt}, C_g)$ を計算し、クラスタ更新かクラスタ新規作成かの判別を、全ての次元において同一の閾値 E を用いて行っている。しかしながら、3次元アフィン変換パラメータ Rcnt の12次元のうち、式(14)において a_1 乃至 a_9 で構成される最初の9次元と、 b_1 乃至 b_3 で構成される残りの3次元とでは、 b_1 乃至 b_3 が時空間内の平行移動を表すパラメータであり、 a_1 乃至 a_9 が回転、拡大縮小、せん断変形など平行移動以外の空間変形写像を表すパラメータであるため、レンジが非常に異なる。

【0171】

そこで、姿勢パラメータ推定部67は、例えば、3次元アフィン変換パラメータ Rcnt の各次元に対して独立に所定の正規化係数を乗ずることでレンジの正規化を行い、正規化後の3次元アフィン変換パラメータ Rcnt をパラメータ空間に投射し、クラスタリングを行うようにしてもよい。正規化係数 nf としては、例えば、 a_1 乃至 a_9 で構成される最初の9次元に対して $nf = 1/10$ 、10次元目の b_1 に対する正規化係数 nf を、想定される動画の横ピクセル数の逆数(例えば、動画のサイズがVGA(Video Graphics Array)サイズなら $nf = 1/640$)とし、11次元目の b_2 に対する正規化係数 nf を、想定される動画の縦ピクセル数の逆数(例えば、動画のサイズがVGAサイズなら $1/480$)とし、12次元目の b_3 に対する正規化係数 nf を、想定される動画の時間長の逆数とすることができる。

【0172】

また、姿勢パラメータ推定部67は、例えば、クラスタとパラメータベクトルとの距離 $d(\text{Rcnt}, C_g)$ を、 a_1 乃至 a_9 で構成される最初の9次元と、 b_1 乃至 b_3 で構成される残りの3次元とで独立に計算し、距離 $d(\text{Rcnt}^{<1-9>}, C_g^{<1-9>})$ と距離 $d(\text{Rcnt}^{<10-12>}, C_g^{<10-12>})$ のそれぞれに対して、閾値 $E^{<1-9>}$ と閾値 $E^{<10-12>}$ を別に設け、いずれの閾値判定も満たされるクラスタがあれば、そのクラスタの更新を、無ければ新規クラスタを生成するものとしてもよい。閾値の設定例としては、例えば、閾値 $E^{<1-9>} = 1$ とした場合、閾値 $E^{<10-12>} = 5$ とすることができる。

【0173】

このようにして、姿勢パラメータ推定部67は、L個のモデルアクションそれぞれにおいて、インライヤと認識された候補対応特徴点ペア群 $\text{CMP}^{[i]}$ の部分集合、すなわち、最終的な認識結果を計算する認識対応特徴点ペア群 $\text{RMP}^{[i]}$ と認識姿勢パラメータ $\sim^{[i]}$ とを、認識結果生成部68に供給する。

【0174】

認識結果生成部68は、姿勢パラメータ推定部67から供給された認識対応特徴点ペア群 $\text{RMP}^{[i]}$ と認識姿勢パラメータ $\sim^{[i]}$ とに基づいて、最終的なモデルアクションの認識結果を生成する。

【0175】

認識結果を利用するユーザまたはアプリケーションは、その目的によって、最も認識結果が高いと考えられる唯一のモデルアクションだけを出力して欲しい場合や、信頼度付きで複数の認識モデルアクションの候補を出力して欲しい場合、また、検出有無のみが知りたい場合や、対応するアクションが検出された場合にはその結果のみならず入力動画中の検出モデルアクションの姿勢パラメータ(3次元画像変換パラメータ)を出力して欲しい

10

20

30

40

50

場合などが考えられる。

【0176】

認識結果生成部68は、上述したように、全モデルアクションの認識対応特徴点ペア群 $RMP^{[i]}$ について、その要素数が4を上回るものが無かった場合、すなわち、認識対応特徴点ペア群 $RMP^{[i]} = NULL$ 、認識姿勢パラメータ $\sim^{[i]} = NULL$ であるとき、認識結果を「非検出」として出力する。そして、それ以外の場合は、いずれかのモデルアクションが認識されたこととなるので、認識結果生成部68は、認識結果を利用するユーザまたはアプリケーションの要求に基づいた形式で検出結果を生成し、出力する。

【0177】

認識結果生成部68は、例えば、認識対応特徴点ペア群 $RMP^{[i]}$ の要素数（対応特徴点ペア数）が所定の閾値以上となっているモデルアクション i 全てを、検出されたモデルアクションとして出力することができる。

10

【0178】

また、認識結果生成部68は、例えば、認識対応特徴点ペア群 $RMP^{[i]}$ の要素数が最大であるモデルアクション i を、検出されたモデルアクションとして出力することができる。

【0179】

また、認識結果生成部68は、例えば、認識対応特徴点ペア群 $RMP^{[i]}$ の要素数（対応特徴点ペア数）が所定の閾値以上となっているモデルアクション i 全てを要素数の多い順にソートし、検出されたモデルアクション i 全てとそれらの順位とを、検出結果として出力することができる。

20

【0180】

また、認識結果生成部68は、例えば、認識対応特徴点ペア群 $RMP^{[i]}$ の要素数（対応特徴点ペア数）が所定の閾値以上となっているモデルアクション i 全てを、検出されたモデルアクションとし、検出されたモデルアクション i の要素数の総和に対する、それぞれのモデルアクションの認識対応特徴点ペア群 $RMP^{[i]}$ の要素数の割合を信頼度として、検出されたモデルアクション i 全てとそれらの信頼度とを、検出結果として出力することができる。

【0181】

また、検出されたモデルアクションの姿勢パラメータ（3次元画像変換パラメータ）の出力が求められている場合、認識結果生成部68は、認識結果として、検出されたモデルアクション i 全てとそれらの姿勢パラメータとを出力する。

30

【0182】

認識結果生成部68は、例えば、検出されたモデルアクション i の $RMP^{[i]}$ の要素すべてを用いて、式(14)乃至式(16)を用いてパラメータの最小二乗推定を行い、その結果を検出モデルアクションの認識姿勢パラメータ $\sim^{[i]}$ として出力することができる。

【0183】

また、認識結果生成部68は、例えば、検出されたモデルアクション i の認識姿勢パラメータ $\sim^{[i]}$ を検出モデルアクションの認識姿勢パラメータ $\sim^{[i]}$ として出力するようにしても良い。

40

【0184】

次に、図6のフローチャートを参照して、認識処理装置11の特徴抽出処理部21が実行する特徴抽出処理について説明する。

【0185】

ステップS11において、モデル動画記録部41は、認識処理のモデルとなる特徴量を取得するための動画データを記録する。

【0186】

ステップS12において、モデル動画記録部41は、図示しない操作入力部により入力されるユーザの操作入力に基づいて、登録に利用するための動画の開始時刻と終了時刻の

50

指定を受け、その部分の動画像データをモデル動画として記録するとともに、前処理実行部 4 2 に供給する。

【 0 1 8 7 】

ステップ S 1 3 において、前処理実行部 4 2 は、モデル動画中のアクション部分と背景部とを分離する前処理を実行する。

【 0 1 8 8 】

上述したように、アクション部分と背景部との分離の方法は任意の方法でよいが、例えば、上述したように、アクション部分に対応する領域の形態学的膨張処理（例えば、8 近傍膨張処理）が行われると、アクション部分と背景部との分離のための 2 値画像から、アクション部分として検出するべき領域内に発生するノイズを除去することができるので好適であり、さらに、膨張処理後の 2 値画像に対して、形態学的収縮処理（例えば、8 近傍収縮処理）が施されると、アクション部分周辺の背景に対応する部分の特徴点および特徴量が抽出されてしまう恐れを限りなく除去することができ、認識の精度が向上するので、好適である。

【 0 1 8 9 】

ステップ S 1 4 において、特徴点抽出部 4 3 は、前処理が実行されたモデル動画から特徴点の抽出を行い、抽出された特徴点の情報を、特徴量抽出部 4 4 に供給する。

【 0 1 9 0 】

特徴点抽出部 4 3 は、既に公知であるさまざまな手法のうちのいずれの手法を用いて特徴点を抽出するものとしても良いが、時空間、すなわち、図 2 を用いて説明した画像平面 $x - y$ に対して、時間 t を空間的奥行き方向の次元とした 3 次元の変形に対して頑強な特徴点抽出法を利用すると、認識の精度があがり、好適である。特徴点抽出部 4 3 は、例えば、上述した ST 特徴点を抽出することができる。

【 0 1 9 1 】

ステップ S 1 5 において、特徴量抽出部 4 4 は、特徴点抽出部 4 3 により抽出されて供給されたモデル動画の特徴点における特徴量を抽出する。

【 0 1 9 2 】

特徴量抽出部 4 4 は、既に公知であるさまざまな手法のうちのいずれの手法を用いて特徴量を抽出するものとしても良い。特徴量抽出部 4 4 は、例えば、上述した式 (1 0) および式 (1 1) で示される時空間特徴量 V_p を用いて特徴量を抽出することができる。

【 0 1 9 3 】

ステップ S 1 6 において、特徴量抽出部 4 4 は、抽出されたモデル動画の特徴点における特徴量を認識処理部 2 2 の辞書登録部 6 1 に供給し、モデル動画ごとに特徴点と特徴量を記憶させ、処理が終了される。

【 0 1 9 4 】

このような処理により、大量のモデルデータを用意したり、複雑な学習処理を行うことなく、モデル動画の特徴点と特徴量を抽出し、認識用に記憶させることができる。

【 0 1 9 5 】

次に、図 7 のフローチャートを参照して、認識処理装置 1 1 の認識処理部 2 2 において実行される認識処理について説明する。この処理が実行されるとき、認識処理部 2 2 の辞書登録部 6 1 には、特徴抽出処理部 2 1 の特徴点抽出部 4 3 から供給されたモデル動画の特徴点および特徴量が、認識処理時に参照可能な形で保存されている。

【 0 1 9 6 】

ステップ S 4 1 において、入力動画バッファ部 6 2 は、入力動画像データを取得し、バッファリングする。

【 0 1 9 7 】

ステップ S 4 2 において、前処理実行部 6 3 は、入力動画中のアクション部分と背景部とを分離する前処理を実行する。

【 0 1 9 8 】

前処理実行部 6 3 は、例えば、特徴抽出処理部 2 1 の前処理実行部 4 2 と同様の方法を

10

20

30

40

50

用いてアクション部分と背景部を分離することができる。また、上述したように、前処理実行部 63 を省略しても、認識処理部 22 における認識処理は実行可能であるので、ステップ S 42 の処理は、省略することができる。

【0199】

ステップ S 43 において、特徴点抽出部 64 は、特徴点抽出部 43 と同様の方法を用いて、入力動画から特徴点を抽出し、抽出された特徴点の情報を、特徴量抽出部 65 に供給する。

【0200】

ステップ S 44 において、特徴量抽出部 65 は、特徴点抽出部 64 により抽出された入力動画の特徴点における特徴量を、上述した特徴量抽出部 44 と同様の方法を用いて抽出し、特徴量比較部 66 に供給する。

10

【0201】

ステップ S 45 において、図 8 のフローチャートを用いて後述する特徴量比較処理が実行される。

【0202】

そして、ステップ S 46 において、図 9 のフローチャートを用いて後述する姿勢パラメータ推定処理が実行される。

【0203】

ステップ S 47 において、認識結果生成部 68 は、姿勢パラメータ推定部 67 から供給された認識対応特徴点ペア群と認識姿勢パラメータ $\theta^{[i]}$ とに基づいて、認識結果を利用するユーザまたはアプリケーションの要求に基づいた形式で、最終的なモデルアクションの認識結果を生成し、認識結果を出力して処理が終了される。

20

【0204】

このような処理により、入力動画に含まれるアクションが、登録されているモデルアクションと一致しているか否かを認識する処理が行われ、認識結果が、認識結果を利用するユーザまたはアプリケーションの要求に基づいた形式で出力される。

【0205】

なお、認識処理において、上述したように、前処理実行部 63 を省略しても、認識処理部 22 における認識処理は実行可能であり、ステップ S 42 の処理は、省略することができる。ステップ S 42 の処理が省略された場合においても、アウトライヤが除去されることにより、背景部分から抽出された特徴点を含むペアについては除去されることとなるため、認識結果を著しく悪化させることはない。ステップ S 42 の処理が省略された場合、候補対応特徴点ペア群を生成する際の閾値を、ステップ S 42 の処理がある場合と比べて弱く設定すると好適である。

30

【0206】

次に、図 8 のフローチャートを参照して、図 7 のステップ S 45 において実行される特徴量比較処理について説明する。

【0207】

ステップ S 71 において、特徴量比較部 66 は、処理対象のモデルアクションを示す変数 i を、 $i = 1$ とする。

40

【0208】

ステップ S 72 において、特徴量比較部 66 は、処理中のモデルアクションの特徴点特徴量情報の ST 特徴点とその特徴量を示す変数 j を $j = 1$ とする。

【0209】

ステップ S 73 において、特徴量比較部 66 は、処理中の入力動画の特徴点特徴量情報の ST 特徴点とその特徴量を示す変数 k を $k = 1$ とする。

【0210】

ステップ S 74 において、特徴量比較部 66 は、 i 番目のモデル動画の j 番目の特徴点の特徴量 $V_j^{[i]}$ と入力動画の k 番目の特徴点の特徴量 W_k との非類似度 D を、例えば、上述した式 (12) を用いて算出する。

50

【0211】

ステップS75において、特徴量比較部66は、 k = 処理中の入力動画の特徴点特徴量情報のST特徴点の数であるか否かを判断する。

【0212】

ステップS75において、 k = 処理中の入力動画の特徴点特徴量情報のST特徴点の数ではないと判断された場合、ステップS76において、特徴量比較部66は、 $k = k + 1$ とし、処理は、ステップS74に戻り、それ以降の処理が繰り返される。

【0213】

ステップS75において、 k = 処理中の入力動画の特徴点特徴量情報のST特徴点の数であると判断された場合、ステップS77において、特徴量比較部66は、 j = 処理中のモデルアクションの特徴点特徴量情報のST特徴点の数であるか否かを判断する。

10

【0214】

ステップS77において、 j = 処理中のモデルアクションの特徴点特徴量情報のST特徴点の数でないと判断された場合、ステップS78において、特徴量比較部66は、 $j = j + 1$ とし、処理は、ステップS73に戻り、それ以降の処理が繰り返される。

【0215】

ステップS77において、 j = 処理中のモデルアクションの特徴点特徴量情報のST特徴点の数であると判断された場合、ステップS79において、特徴量比較部66は、 i = 登録モデルアクション数であるか否かを判断する。

【0216】

20

ステップS79において、 i = 登録モデルアクション数でないと判断された場合、ステップS80において、特徴量比較部66は、 $i = i + 1$ とし、処理は、ステップS72に戻り、それ以降の処理が繰り返される。

【0217】

ステップS79において、 i = 登録モデルアクション数であると判断された場合、ステップS81において、特徴量比較部66は、得られた全ての非類似度Dの値に基づいて、モデル特徴点のそれぞれに対応する入力動画の特徴点である対応特徴点の候補となる候補対応特徴点群を求め、モデル特徴点と対応特徴点の候補のペアを、候補対応特徴点ペアとして姿勢パラメータ推定部67に供給し、処理は、図7のステップS45に戻り、ステップS46に進む。

30

【0218】

このような処理により、モデル動画と入力動画の全ての特徴点の組み合わせにおいて、例えば、上述した式(12)を用いて説明した非類似度Dが演算されて、対応するモデル特徴点と対応特徴点の候補のペアが求められる。

【0219】

次に、図7のステップS46において実行される姿勢パラメータ推定処理について説明する。

【0220】

ステップS111において、姿勢パラメータ推定部67は、処理対象のモデルアクションを示す変数 i を、 $i = 1$ とする。

40

【0221】

ステップS112において、姿勢パラメータ推定部67は、何回目の処理であるかを示す変数 cnt を、 $cnt = 1$ に初期化する。

【0222】

ステップS113において、姿勢パラメータ推定部67は、候補対応特徴点ペア群からランダムに所定数のペア(例えば、4ペア)を選択し、上述した式(14)乃至式(16)を用いて、3次元アフィン変換パラメータを計算する。

【0223】

ステップS114において、姿勢パラメータ推定部67は、ステップS113において計算された3次元アフィン変換パラメータをパラメータ空間に投射し、NN法により3次

50

元アフィン変換パラメータ空間をクラスタリングする。なお、 $cnt = 1$ のときは、ステップS113において計算された3次元アフィン変換パラメータをセントロイドとする1つ目のクラスタを生成する。

【0224】

ステップS115において、姿勢パラメータ推定部67は、例えば、最多メンバ数が所定の閾値（例えば15）を超え、かつ、最多メンバ数と2番目に多いメンバ数との差が所定の閾値（例えば3）を超える場合、または、処理の繰り返し回数が、所定の閾値（例えば5000回）を超える場合などの、繰り返し終了条件を満たすか否かを判断する。

【0225】

ステップS115において、繰り返し終了条件を満たしていないと判断された場合、ステップS116において、姿勢パラメータ推定部67は、 $cnt = cnt + 1$ として、処理はステップS113に戻り、それ以降の処理が繰り返される。

10

【0226】

ステップS115において、繰り返し終了条件を満たしていると判断された場合、ステップS117において、姿勢パラメータ推定部67は、 $i =$ 登録モデルアクション数であるか否かを判断する。

【0227】

ステップS117において、 $i =$ 登録モデルアクション数でないと判断された場合、ステップS118において、姿勢パラメータ推定部67は、 $i = i + 1$ とし、処理は、ステップS112に戻り、それ以降の処理が繰り返される。

20

【0228】

ステップS117において、 $i =$ 登録モデルアクション数であると判断された場合、ステップS119において、姿勢パラメータ推定部67は、クラスタリングの結果に基づいて、全てのモデルアクションについて、それぞれ、最多メンバ数を持つクラスタ Z_{max} のメンバである対応特徴点ペア群を、それぞれのモデルアクションに対する認識対応特徴点ペア群とし、クラスタ Z_{max} のセントロイド C_{max} を認識姿勢パラメータとして、認識結果生成部に出力し、処理は、図7のステップS46に戻り、ステップS47に進む。

【0229】

このような処理により、それぞれのモデルアクションに対する認識対応特徴点と、認識姿勢パラメータとを求めることができる。

30

【0230】

以上説明したように、認識処理装置11は、登録したアクションを入力動画から検出する処理を実行するものであって、1学習サンプルから認識が可能であるため、統計学習を用いた認識処理を行う場合とは異なり、大量の学習用データを用意する必要がない。

【0231】

また、認識処理装置11においては、ユーザが認識対象となるモデルアクションを容易に新規登録することができるので、事前に学習されているアクションのみが認識可能な統計学習を用いた手法や、登録可能なアクションにより認識アルゴリズムがことなるアクションおよびボディパーツを限定した手法とは異なり、認識対象となるモデルアクションの数を容易に増やすことができる。

40

【0232】

また、認識処理装置11によって実行される認識処理は、入力動画の部分隠れやカメラ視点の変化に対して頑強である。例えば、ボディパーツ同定を行う手法や、動き領域の形状または重心の移動を用いる方法は、画像変形に弱いため、認識処理装置11によって実行される認識処理は、これらの手法と比較して、有利である。認識処理装置11によって実行される認識処理においては、例えば、手を振る動作が認識されるべきモデルアクションである場合、入力動画において、認識対象の人物が、立った姿勢で手を振っていても、据わった姿勢で手を振っていても、寝転んだ姿勢で手を振っていても、柔軟に認識することが可能となる。

【0233】

50

なお、ここでは、認識処理装置 11 は 1 つの装置であるものとして説明したが、特徴抽出処理部 21 および認識処理部 22 は、同様の機能を有するそれぞれ 1 つの装置として構成されていても良い。

【0234】

また、特徴抽出処理と、認識処理は連続して行われなくても良く、特徴抽出処理部 21 および認識処理部 22 が、それぞれ、異なる 1 つの装置として構成され、乖離して設置されていても良いことはいうまでもない。換言すれば、特徴抽出処理部 21 に対応する装置により生成される特徴量の情報が辞書登録部 61 に記憶されている認識処理部 22 に対応する装置は、特徴抽出処理部 21 に対応する装置と乖離した場所に設置されても、単独で、入力動画のアクションを認識する処理を行うことができる。

10

【0235】

上述した一連の処理は、ハードウェアにより実行させることもできるし、ソフトウェアにより実行させることもできる。そのソフトウェアは、そのソフトウェアを構成するプログラムが、専用のハードウェアに組み込まれているコンピュータ、または、各種のプログラムをインストールすることで、各種の機能を実行することが可能な、例えば汎用のパーソナルコンピュータなどに、記録媒体からインストールされる。この場合、上述した処理は、図 10 に示されるようなパーソナルコンピュータ 500 により実行される。

【0236】

図 10 において、CPU (Central Processing Unit) 501 は、ROM (Read Only Memory) 502 に記憶されているプログラム、または、記憶部 508 から RAM (Random Access Memory) 503 にロードされたプログラムに従って各種の処理を実行する。RAM 503 にはまた、CPU 501 が各種の処理を実行する上において必要なデータなどが適宜記憶される。

20

【0237】

CPU 501、ROM 502、および RAM 503 は、内部バス 504 を介して相互に接続されている。この内部バス 504 にはまた、入出力インターフェース 505 も接続されている。

【0238】

入出力インターフェース 505 には、キーボード、マウスなどの操作入力部、または、カメラなどの撮像装置などよりなる入力部 506、CRT、LCD などよりなるディスプレイ、スピーカなどよりなる出力部 507、ハードディスクなどより構成される記憶部 508、並びに、モデム、ターミナルアダプタなどより構成される通信部 509 が接続されている。通信部 509 は、電話回線や CATV を含む各種のネットワークを介しての通信処理を行う。

30

【0239】

入出力インターフェース 505 にはまた、必要に応じてドライブ 510 が接続され、磁気ディスク、光ディスク、光磁気ディスク、あるいは半導体メモリなどによりなるリムーバブルメディア 521 が適宜装着され、それから読み出されたコンピュータプログラムが、必要に応じて記憶部 508 にインストールされる。

【0240】

一連の処理をソフトウェアにより実行させる場合には、そのソフトウェアを構成するプログラムが、ネットワークや記録媒体からインストールされる。

40

【0241】

この記録媒体は、図 10 に示されるように、コンピュータとは別に、ユーザにプログラムを提供するために配布される、プログラムが記録されているリムーバブルメディア 521 よりなるパッケージメディアにより構成されるだけでなく、装置本体に予め組み込まれた状態でユーザに提供される、プログラムが記録されている ROM 502 や記憶部 508 が含まれるハードディスクなどで構成される。

【0242】

上述した一連の処理を実行するソフトウェアを実行するパーソナルコンピュータ 500 において、上述した認識処理は、例えば、次のようなアプリケーションに適用可能である

50

。

【0243】

パーソナルコンピュータ500において、内部（辞書登録部61に対応して、モデルアクションを、認識処理時に参照可能な形で保存する記憶部508）に記憶されているジェスチャを、パーソナルコンピュータ500が実行するソフトウェアを操作するためのコマンドと予め関係付けておく。ジェスチャは、上述したモデル動画に含まれるアクションとして、ユーザが固有に登録可能である。例えば、手の上下ジェスチャを、ブラウザやワードプロセッサ等の実行ウィンドウのスクロールを指令するコマンドに対応付けて登録することができる。

【0244】

10

そして、パーソナルコンピュータ500において、出力部507のディスプレイ付近に設置された入力部506のカメラ（入力動画バッファ部62に対応する）により、パーソナルコンピュータ500を使用するユーザを撮像し、これを、入力動画とする。そして、前処理実行部63乃至認識結果生成部68に対応する機能を実現することができるCPU501は、ユーザが固有に登録し、内部（記憶部508）に記憶されているジェスチャとの認識処理を実行する。

【0245】

例えば、ブラウザやワードプロセッサ等の対応するソフトウェアの起動中は、逐次認識を行うようにする。そして、例えば、手の上下など、登録されたジェスチャが検出された時には、それに対応したコマンド処理が実行される。

20

【0246】

なお、記憶部508にジェスチャを記憶する処理は、予め、パーソナルコンピュータ500において実行されていても良いし、他の装置において実行され、得られたモデルアクションが、記憶部508に、認識処理時に参照可能な形で保存されるものとしてもよい。

【0247】

また、本明細書において、記録媒体に記録されるプログラムを記述するステップは、記載された順序に沿って時系列的に行われる処理はもちろん、必ずしも時系列的に処理されなくとも、並列的あるいは個別に実行される処理をも含むものである。

【0248】

なお、本明細書において、システムとは、複数の装置により構成される装置全体を表すものである。

30

【0249】

なお、本発明の実施の形態は、上述した実施の形態に限定されるものではなく、本発明の要旨を逸脱しない範囲において種々の変更が可能である。

【図面の簡単な説明】

【0250】

【図1】認識処理装置の構成を示すブロック図である。

【図2】画像平面 $x-y$ に対して、時間 t を空間的奥行き方向の次元とした3次元の画像シーケンスについて説明する図である。

【図3】背景の分離について説明するための図である。

40

【図4】背景の分離について説明するための図である。

【図5】背景の分離について説明するための図である。

【図6】特徴抽出処理について説明するためのフローチャートである。

【図7】認識処理について説明するためのフローチャートである。

【図8】特徴量比較処理について説明するためのフローチャートである。

【図9】姿勢パラメータ推定処理について説明するためのフローチャートである。

【図10】パーソナルコンピュータの構成を示すブロック図である。

【符号の説明】

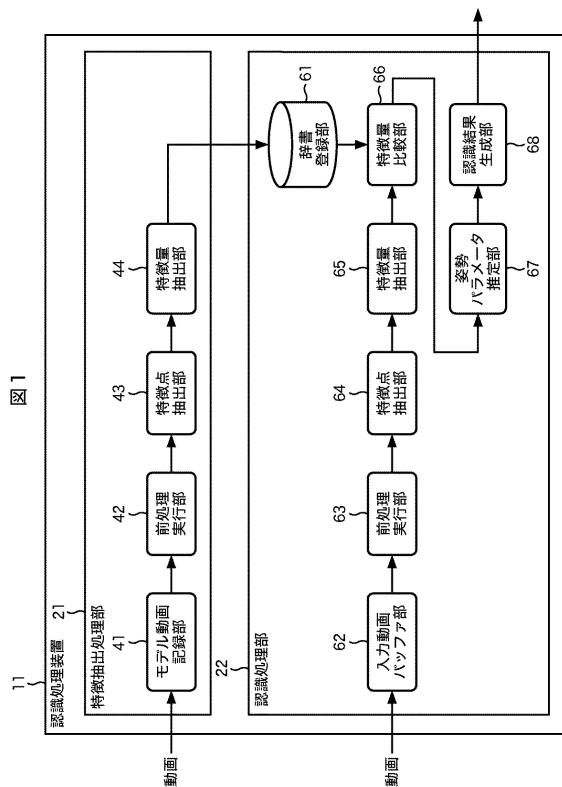
【0251】

11 認識処理装置, 21 特徴抽出処理部, 22 認識処理部, 41 モデル

50

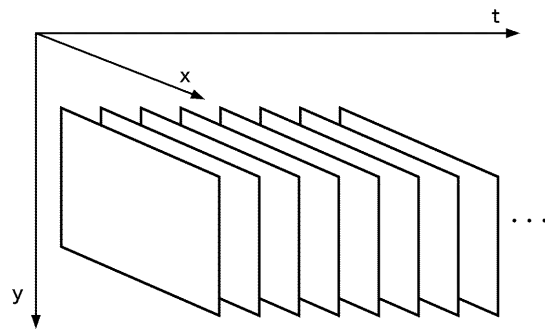
動画記録部， 4 2 前処理実行部， 4 3 特徴点抽出部， 4 4 特徴量抽出部，
 6 1 辞書登録部， 6 2 入力動画バッファ部， 6 3 前処理実行部， 6 4 特徴
 点抽出部， 6 5 特徴量抽出部， 6 7 姿勢パラメータ推定部， 6 8 認識結果生
 成部

【図 1】

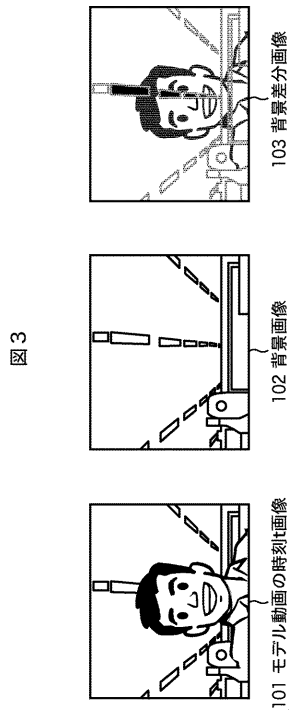


【図 2】

図 2

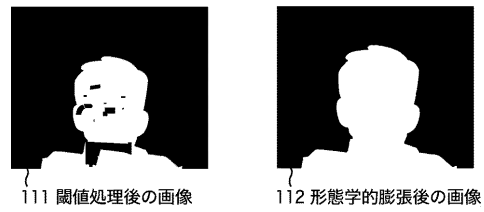


【図 3】



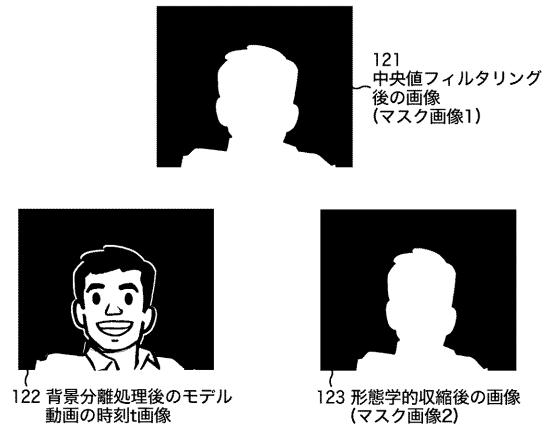
【図 4】

図 4



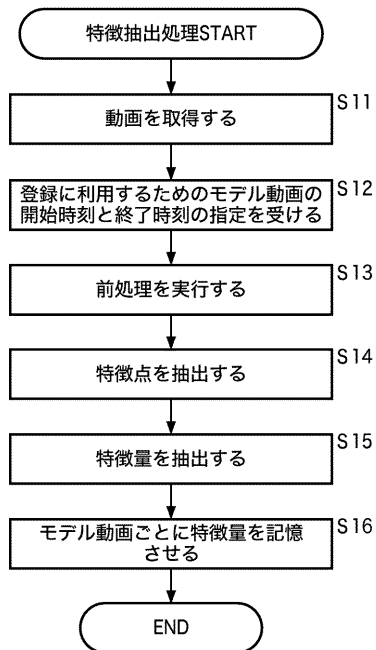
【図 5】

図 5



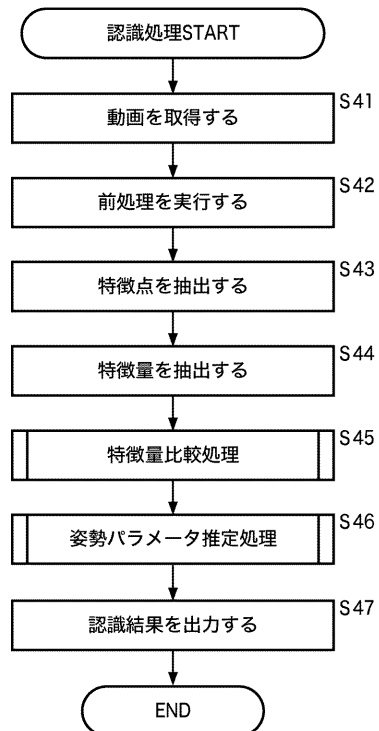
【図 6】

図 6



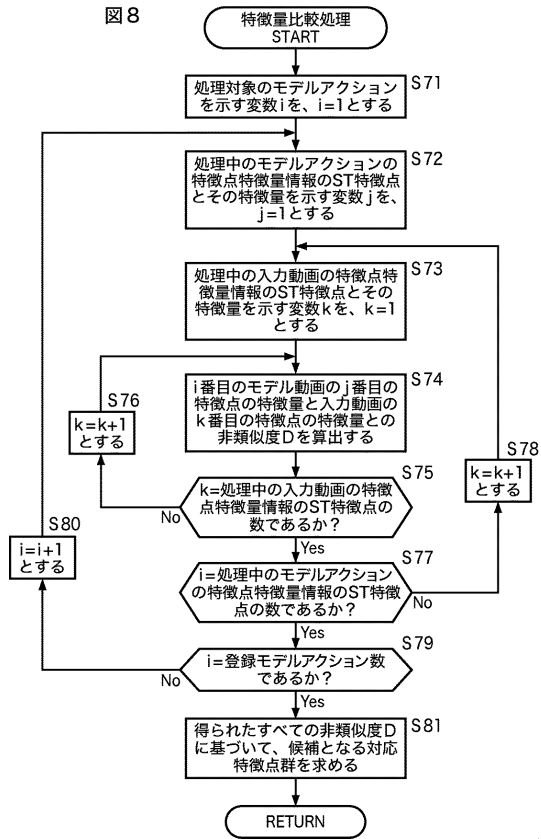
【図 7】

図 7



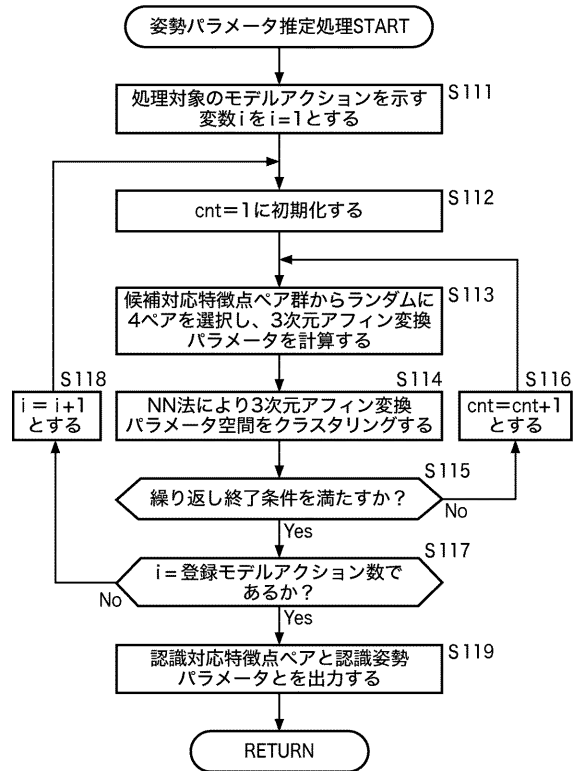
【図 8】

図 8



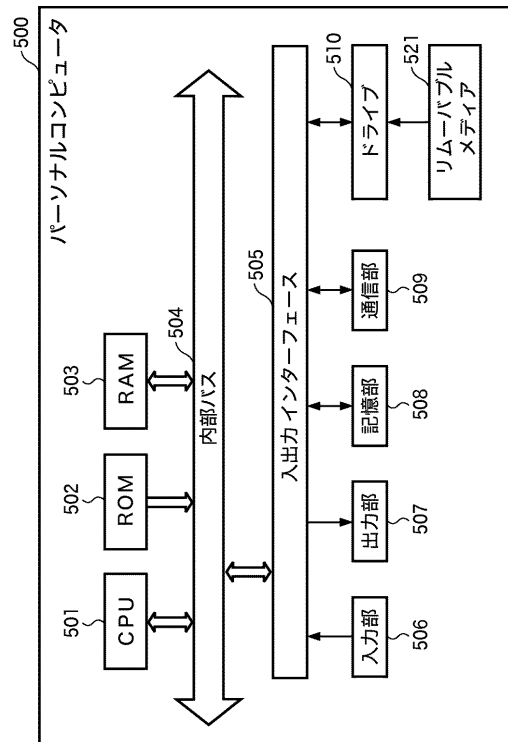
【図 9】

図 9



【図 10】

図 10



フロントページの続き

(56)参考文献 特開平 8 - 2 1 2 3 2 7 (J P , A)

国際公開第 2 0 0 5 / 0 2 0 1 5 2 (W O , A 1)

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 T 7 / 2 0