

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7501621号
(P7501621)

(45)発行日 令和6年6月18日(2024.6.18)

(24)登録日 令和6年6月10日(2024.6.10)

(51)国際特許分類 F I
G 0 6 F 16/73 (2019.01) G 0 6 F 16/73

請求項の数 7 (全39頁)

(21)出願番号	特願2022-522434(P2022-522434)	(73)特許権者	000004237 日本電気株式会社 東京都港区芝五丁目7番1号
(86)(22)出願日	令和2年5月14日(2020.5.14)	(74)代理人	100110928 弁理士 速水 進治
(86)国際出願番号	PCT/JP2020/019255	(72)発明者	吉田 登 東京都港区芝五丁目7番1号 日本電気株式会社内
(87)国際公開番号	WO2021/229750	(72)発明者	潘 雅冬 東京都港区芝五丁目7番1号 日本電気株式会社内
(87)国際公開日	令和3年11月18日(2021.11.18)	(72)発明者	川合 諒 東京都港区芝五丁目7番1号 日本電気株式会社内
審査請求日	令和4年11月4日(2022.11.4)	(72)発明者	劉 健全

最終頁に続く

(54)【発明の名称】 画像選択装置、画像選択方法、およびプログラム

(57)【特許請求の範囲】

【請求項1】

複数の第1フレーム画像を含むクエリ動画を取得するクエリ取得手段と、
前記複数の第1フレーム画像から複数のクエリフレームを選択するクエリフレーム選択手段と、
前記複数のクエリフレームを用いて動画を選択する動画選択手段と、
を備え、
前記クエリフレーム選択手段は、第1の前記クエリフレームの次の前記クエリフレームとして前記第1のクエリフレームからの変化量が第3基準以上となった前記第1フレーム画像を選択する処理を繰り返す、

前記動画選択手段は、前記動画を選択する条件として、少なくとも以下の(1)及び(2)を用いる、画像選択装置。

(1)「前記クエリフレームに対する類似度が第1基準を満たす類似フレーム画像が存在する」という条件が少なくとも2つの前記クエリフレームについて満たされる。

(2)前記少なくとも2つの類似フレーム画像のそれぞれに対応する前記クエリフレームを前記少なくとも2つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも2つの前記クエリフレームの並び順に一致する。

【請求項2】

請求項1に記載の画像選択装置において、
前記クエリ動画は人物を含んでおり、

前記動画選択手段は、前記類似度として、前記人物の姿勢の類似度を用いる画像選択装置。

【請求項 3】

請求項 2 に記載の画像選択装置において、

前記クエリフレームが満たすべき条件の一つは、当該クエリフレームに含まれる前記人物に関する情報量が第 2 基準を満たしていることであり、

前記複数の第 1 フレーム画像のそれぞれは、前記人物の姿勢を示す情報として、当該人物の関節の位置を示す関節情報を含んでおり、

前記第 2 基準は、前記関節情報に含まれる関節の数が基準を満たしていることである画像選択装置。

10

【請求項 4】

請求項 1 ~ 3 のいずれか一項に記載の画像選択装置において、

前記クエリフレーム選択手段は、ユーザからの入力を用いて前記第 3 基準を設定する画像選択装置。

【請求項 5】

請求項 1 ~ 3 のいずれか一項に記載の画像選択装置において、

前記クエリフレーム選択手段は、前記複数の第 1 フレーム画像を用いて当該クエリフレームに用いる前記第 3 基準を設定する画像選択装置。

【請求項 6】

コンピュータが、

複数の第 1 フレーム画像を含むクエリ動画を取得する取得工程と、

前記複数の第 1 フレーム画像から複数のクエリフレームを選択するクエリフレーム選択工程と、

前記複数のクエリフレームを用いて動画を選択する動画選択工程と、

を行い、

前記クエリフレーム選択工程において、第 1 の前記クエリフレームの次の前記クエリフレームとして前記第 1 のクエリフレームからの変化量が第 3 基準以上となった前記第 1 フレーム画像を選択する処理を繰り返し、

前記動画選択工程において、前記動画を選択する条件として、少なくとも以下の (1) 及び (2) を用いる、画像選択方法。

20

(1) 「前記クエリフレームに対する類似度が第 1 基準を満たす類似フレーム画像が存在する」という条件が少なくとも 2 つの前記クエリフレームについて満たされる。

(2) 前記少なくとも 2 つの類似フレーム画像のそれぞれに対応する前記クエリフレームを前記少なくとも 2 つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも 2 つの前記クエリフレームの並び順に一致する。

30

【請求項 7】

コンピュータに、

複数の第 1 フレーム画像を含むクエリ動画を取得する取得機能と、

前記複数の第 1 フレーム画像から複数のクエリフレームを選択するクエリフレーム選択機能と、

前記複数のクエリフレームを用いて動画を選択する動画選択機能と、

前記クエリフレーム選択機能は、第 1 の前記クエリフレームの次の前記クエリフレームとして前記第 1 のクエリフレームからの変化量が第 3 基準以上となった前記第 1 フレーム画像を選択する処理を繰り返し、

前記動画選択機能は、前記動画を選択する条件として、少なくとも以下の (1) 及び (2) を用いるプログラム。

(1) 「前記クエリフレームに対する類似度が第 1 基準を満たす類似フレーム画像が存在する」という条件が少なくとも 2 つの前記クエリフレームについて満たされる。

(2) 前記少なくとも 2 つの類似フレーム画像のそれぞれに対応する前記クエリフレ

40

50

ムを前記少なくとも2つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも2つの前記クエリフレームの並び順に一致する。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、画像選択装置、画像選択方法、およびプログラムに関する。

【背景技術】

【0002】

近年、監視システム等において、監視カメラの画像から人物の姿勢や行動等の状態の検出や検索を行う技術が利用されている。関連する技術として、例えば、特許文献1及び2が知られている。特許文献1には、深さ映像に含まれる人物の頭や手足等のキーポイントに基づいて、類似する人物の姿勢を検索する技術が開示されている。特許文献2には、人物の姿勢と関連しないが、画像に付加された傾き等の姿勢情報を利用して類似画像を検索する技術が開示されている。なお、その他に、人物の骨格推定に関連する技術として、非特許文献1が知られている。

10

【0003】

一方、近年は動画をクエリとして利用し、このクエリに類似する動画を検索することも検討されている。例えば特許文献3には、クエリとなる参照映像を入力すると、登場人物の顔の数、並びに各登場人物の顔の位置、大きさ、及び向きを用いて、類似する映像を検索することが記載されている。

20

【先行技術文献】

【特許文献】

【0004】

【文献】特表2014-522035号公報

【文献】特開2006-260405号公報

【文献】国際公開第2006/025272号

【非特許文献】

【0005】

【文献】Zhe Cao, Tomas Simon, Shih-En Wei, Yaser Sheikh, "Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields", The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, P. 7291-7299

30

【発明の概要】

【発明が解決しようとする課題】

【0006】

動画をクエリとして、このクエリに類似する動画を選択する場合、選択の精度を上げることは難しい。本発明の目的の一つは、動画をクエリとして、このクエリに類似する動画を選択する場合において、選択の精度を上げることにある。

【課題を解決するための手段】

【0007】

本発明によれば、複数の第1フレーム画像を含むクエリ動画を取得するクエリ取得手段と、

40

前記複数の第1フレーム画像から複数のクエリフレームを選択するクエリフレーム選択手段と、

前記複数のクエリフレームを用いて動画を選択する動画選択手段と、
を備え、

前記動画選択手段は、前記動画を選択する条件として、少なくとも以下の(1)及び(2)を用いる、画像選択装置が提供される。

(1)「前記クエリフレームに対する類似度が第1基準を満たす類似フレーム画像が存在する」という条件が少なくとも2つの前記クエリフレームについて満たされる。

(2)前記少なくとも2つの類似フレーム画像のそれぞれに対応する前記クエリフレー

50

ムを前記少なくとも2つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも2つの前記クエリフレームの並び順に一致する。

【0008】

本発明によれば、コンピュータが、
 複数の第1フレーム画像を含むクエリ動画を取得する取得工程と、
 前記複数の第1フレーム画像から複数のクエリフレームを選択するクエリフレーム選択工程と、
 前記複数のクエリフレームを用いて動画を選択する動画選択工程と、
 を行い、

前記動画選択工程において、前記動画を選択する条件として、少なくとも以下の(1)及び(2)を用いる、画像選択方法が提供される。

(1)「前記クエリフレームに対する類似度が第1基準を満たす類似フレーム画像が存在する」という条件が少なくとも2つの前記クエリフレームについて満たされる。

(2)前記少なくとも2つの類似フレーム画像のそれぞれに対応する前記クエリフレームを前記少なくとも2つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも2つの前記クエリフレームの並び順に一致する。

【0009】

本発明によれば、コンピュータに、
 複数の第1フレーム画像を含むクエリ動画を取得する取得機能と、
 前記複数の第1フレーム画像から複数のクエリフレームを選択するクエリフレーム選択機能と、
 前記複数のクエリフレームを用いて動画を選択する動画選択機能と、
 を持たせ、

前記動画選択機能は、前記動画を選択する条件として、少なくとも以下の(1)及び(2)を用いるプログラムが提供される。

(1)「前記クエリフレームに対する類似度が第1基準を満たす類似フレーム画像が存在する」という条件が少なくとも2つの前記クエリフレームについて満たされる。

(2)前記少なくとも2つの類似フレーム画像のそれぞれに対応する前記クエリフレームを前記少なくとも2つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも2つの前記クエリフレームの並び順に一致する。

【発明の効果】

【0010】

本発明によれば、動画をクエリとして、このクエリに類似する動画を選択する場合において、選択の精度は上がる。

【図面の簡単な説明】

【0011】

上述した目的、およびその他の目的、特徴および利点は、以下に述べる好適な実施の形態、およびそれに付随する以下の図面によってさらに明らかになる。

【0012】

【図1】実施の形態に係る画像処理装置の概要を示す構成図である。

【図2】実施の形態1に係る画像処理装置の構成を示す構成図である。

【図3】実施の形態1に係る画像処理方法を示すフローチャートである。

【図4】実施の形態1に係る分類方法を示すフローチャートである。

【図5】実施の形態1に係る検索方法を示すフローチャートである。

【図6】実施の形態1に係る骨格構造の検出例を示す図である。

【図7】実施の形態1に係る人体モデルを示す図である。

【図8】実施の形態1に係る骨格構造の検出例を示す図である。

【図9】実施の形態1に係る骨格構造の検出例を示す図である。

【図10】実施の形態1に係る骨格構造の検出例を示す図である。

【図11】実施の形態1に係る分類方法の具体例を示すグラフである。

10

20

30

40

50

- 【図 1 2】実施の形態 1 に係る分類結果の表示例を示す図である。
- 【図 1 3】実施の形態 1 に係る検索方法を説明するための図である。
- 【図 1 4】実施の形態 1 に係る検索方法を説明するための図である。
- 【図 1 5】実施の形態 1 に係る検索方法を説明するための図である。
- 【図 1 6】実施の形態 1 に係る検索方法を説明するための図である。
- 【図 1 7】実施の形態 1 に係る検索結果の表示例を示す図である。
- 【図 1 8】実施の形態 2 に係る画像処理装置の構成を示す構成図である。
- 【図 1 9】実施の形態 2 に係る画像処理方法を示すフローチャートである。
- 【図 2 0】実施の形態 2 に係る身長画素数算出方法の具体例 1 を示すフローチャートである。
- 【図 2 1】実施の形態 2 に係る身長画素数算出方法の具体例 2 を示すフローチャートである。
- 【図 2 2】実施の形態 2 に係る身長画素数算出方法の具体例 3 を示すフローチャートである。
- 【図 2 3】実施の形態 2 に係る正規化方法を示すフローチャートである。
- 【図 2 4】実施の形態 2 に係る人体モデルを示す図である。
- 【図 2 5】実施の形態 2 に係る骨格構造の検出例を示す図である。
- 【図 2 6】実施の形態 2 に係る骨格構造の検出例を示す図である。
- 【図 2 7】実施の形態 2 に係る骨格構造の検出例を示す図である。
- 【図 2 8】実施の形態 2 に係る人体モデルを示す図である。
- 【図 2 9】実施の形態 2 に係る骨格構造の検出例を示す図である。
- 【図 3 0】実施の形態 2 に係る身長画素数算出方法を説明するためのヒストグラムである。
- 【図 3 1】実施の形態 2 に係る骨格構造の検出例を示す図である。
- 【図 3 2】実施の形態 2 に係る 3 次元人体モデルを示す図である。
- 【図 3 3】実施の形態 2 に係る身長画素数算出方法を説明するための図である。
- 【図 3 4】実施の形態 2 に係る身長画素数算出方法を説明するための図である。
- 【図 3 5】実施の形態 2 に係る身長画素数算出方法を説明するための図である。
- 【図 3 6】実施の形態 2 に係る正規化方法を説明するための図である。
- 【図 3 7】実施の形態 2 に係る正規化方法を説明するための図である。
- 【図 3 8】実施の形態 2 に係る正規化方法を説明するための図である。
- 【図 3 9】画像処理装置のハードウェア構成例を示す図である。
- 【図 4 0】検索方法 6 における検索部の機能構成を示す図である。
- 【図 4 1】動画選択部が類似フレームを選択する方法の一例を説明するための図である。
- 【図 4 2】検索部が行う処理の一例を示すフローチャートである。
- 【図 4 3】図 4 2 に示す処理を説明するための図である。
- 【図 4 4】図 4 2 の変形例を示すフローチャートである。
- 【図 4 5】クエリフレームの選択規則の第 1 例を説明するための図である。
- 【図 4 6】クエリフレームの選択規則の第 2 例を説明するための図である。
- 【図 4 7】第 3 基準を設定する方法を説明するための図である。

【発明を実施するための形態】

【0013】

以下、本発明の実施の形態について、図面を用いて説明する。尚、すべての図面において、同様な構成要素には同様の符号を付し、適宜説明を省略する。

【0014】

(実施の形態に至る検討)

近年、ディープラーニング等の機械学習を活用した画像認識技術が様々なシステムに活用されている。例えば、監視カメラの画像により監視を行う監視システムへの適用が進められている。監視システムに機械学習を活用することで、画像から人物の姿勢や行動等の状態をある程度把握することが可能とされつつある。

【0015】

10

20

30

40

50

しかしながら、このような関連する技術では、必ずしもオンデマンドにユーザが望む人物の状態を把握できない場合がある。例えば、ユーザが検索し把握したい人物の状態を事前に特定できている場合もあれば、未知の状態のように具体的に特定できていない場合もある。そうすると、場合によっては、ユーザが検索したい人物の状態を詳細に指定することができない。また、人物の体の一部が隠れているような場合には検索等を行うことができない。関連する技術では、特定の検索条件のみからしか人物の状態を検索できないため、所望の人物の状態を柔軟に検索や分類することが困難である。

【0016】

そこで、発明者らは、オンデマンドに画像からユーザ所望の人物の状態を認識するため、非特許文献1などの骨格推定技術を利用する方法を検討した。非特許文献1に開示されたOpenPose等のように、関連する骨格推定技術では、様々なパターンの正解付けされた画像データを学習することで、人物の骨格を推定する。以下の実施の形態では、このような骨格推定技術を活用することで、人物の状態を柔軟に認識することを可能とする。

10

【0017】

なお、OpenPose等の骨格推定技術により推定される骨格構造は、関節等の特徴的な点である「キーポイント」と、キーポイント間のリンクを示す「ボーン（ボーンリンク）」とから構成される。このため、以下の実施の形態では、骨格構造について「キーポイント」と「ボーン」という用語を用いて説明するが、特に限定されない限り、「キーポイント」は人物の「関節」に対応し、「ボーン」は人物の「骨」に対応している。そして「キーポイント」の位置は関節情報の一例になる。

20

【0018】

（実施の形態の概要）

図1は、実施の形態に係る画像処理装置10の概要を示している。図1に示すように、画像処理装置10は、骨格検出部11、特徴量算出部12、及び認識部13を備えている。骨格検出部11は、カメラ等から取得される2次元画像に基づいて、複数の人物の2次元骨格構造を検出する。特徴量算出部12は、骨格検出部11により検出された複数の2次元骨格構造の特徴量を算出する。認識部13は、特徴量算出部12により算出された複数の特徴量の類似度に基づいて、複数の人物の状態の認識処理を行う。認識処理は、人物の状態の分類処理や検索処理（選択処理）等である。このため、画像処理装置10は画像選択装置としても機能する。

30

【0019】

このように、実施の形態では、2次元画像から人物の2次元骨格構造を検出し、この2次元骨格構造から算出される特徴量に基づいて人物の状態の分類や検索等の認識処理を行うことで、所望の人物の状態を柔軟に認識することができる。

【0020】

（実施の形態1）以下、図面を参照して実施の形態1について説明する。図2は、本実施の形態に係る画像処理装置100の構成を示している。画像処理装置100は、カメラ200及びデータベース(DB)110とともに画像処理システム1を構成する。画像処理装置100を含む画像処理システム1は、画像から推定される人物の骨格構造に基づき、人物の姿勢や行動等の状態を分類及び検索するシステムである。なお、画像処理装置100も、画像選択装置としても機能する。

40

【0021】

カメラ200は、2次元の画像を生成する監視カメラ等の撮像部である。カメラ200は、所定の箇所に設置されて、設置個所から撮像領域における人物等を撮像する。カメラ200は、撮像した画像（映像）を画像処理装置100へ出力可能に直接接続、もしくはネットワーク等を介して接続されている。なお、カメラ200を画像処理装置100の内部に設けてもよい。

【0022】

データベース110は、画像処理装置100の処理に必要な情報（データ）や処理結果等を格納するデータベースである。データベース110は、画像取得部101が取得した

50

画像や、骨格構造検出部 102 の検出結果、機械学習用のデータ、特徴量算出部 103 が算出した特徴量、分類部 104 の分類結果、検索部 105 の検索結果等を記憶する。データベース 110 は、画像処理装置 100 と必要に応じてデータを入出力可能に直接接続、もしくはネットワーク等を介して接続されている。なお、データベース 110 をフラッシュメモリなどの不揮発性メモリやハードディスク装置等として、画像処理装置 100 の内部に設けてもよい。

【0023】

図 2 に示すように、画像処理装置 100 は、画像取得部 101、骨格構造検出部 102、特徴量算出部 103、分類部 104、検索部 105、入力部 106、及び表示部 107 を備えている。なお、各部（ブロック）の構成は一例であり、後述の方法（動作）が可能であれば、その他の各部で構成されてもよい。また、画像処理装置 100 は、例えば、プログラムを実行するパーソナルコンピュータやサーバ等のコンピュータ装置で実現されるが、1つの装置で実現してもよいし、ネットワーク上の複数の装置で実現してもよい。例えば、入力部 106 や表示部 107 等を外部の装置としてもよい。また、分類部 104 及び検索部 105 の両方を備えていてもよいし、いずれか一方のみを備えていてもよい。分類部 104 及び検索部 105 の両方、もしくは一方は、人物の状態の認識処理を行う認識部である。

10

【0024】

画像取得部 101 は、カメラ 200 が撮像した人物を含む 2 次元の画像を取得する。画像取得部 101 は、例えば、所定の監視期間にカメラ 200 が撮像した、人物を含む画像（複数の画像を含む映像）を取得する。なお、カメラ 200 からの取得に限らず、予め用意された人物を含む画像をデータベース 110 等から取得してもよい。

20

【0025】

骨格構造検出部 102 は、取得された 2 次元の画像に基づき、画像内の人物の 2 次元の骨格構造を検出する。骨格構造検出部 102 は、取得された画像の中で認識される全ての人物について、骨格構造を検出する。骨格構造検出部 102 は、機械学習を用いた骨格推定技術を用いて、認識される人物の関節等の特徴に基づき人物の骨格構造を検出する。骨格構造検出部 102 は、例えば、非特許文献 1 の Open Pose 等の骨格推定技術を用いる。

【0026】

特徴量算出部 103 は、検出された 2 次元の骨格構造の特徴量を算出し、算出した特徴量を、処理対象となった画像に紐づけてデータベース 110 に格納する。骨格構造の特徴量は、人物の骨格の特徴を示しており、人物の骨格に基づいて人物の状態を分類や検索するための要素となる。通常、この特徴量は、複数のパラメータ（例えば後述する分類要素）を含んでいる。そして特徴量は、骨格構造の全体の特徴量でもよいし、骨格構造の一部の特徴量でもよく、骨格構造の各部のように複数の特徴量を含んでもよい。特徴量の算出方法は、機械学習や正規化等の任意の方法でよく、正規化として最小値や最大値を求めてもよい。一例として、特徴量は、骨格構造を機械学習することで得られた特徴量や、骨格構造の頭部から足部までの画像上の大きさ等である。骨格構造の大きさは、画像上の骨格構造を含む骨格領域の上下方向の高さや面積等である。上下方向（高さ方向または縦方向）は、画像における上下の方向（Y 軸方向）であり、例えば、地面（基準面）に対し垂直な方向である。また、左右方向（横方向）は、画像における左右の方向（X 軸方向）であり、例えば、地面に対し平行な方向である。

30

40

【0027】

なお、ユーザが望む分類や検索を行うためには、分類や検索処理に対しロバスト性を有する特徴量を用いることが好ましい。例えば、ユーザが、人物の向きや体型に依存しない分類や検索を望む場合、人物の向きや体型にロバストな特徴量を使用してもよい。同じ姿勢で様々な方向に向いている人物の骨格や同じ姿勢で様々な体型の人物の骨格を学習することや、骨格の上下方向のみの特徴を抽出することで、人物の向きや体型に依存しない特徴量を得ることができる。

50

【 0 0 2 8 】

分類部 1 0 4 は、データベース 1 1 0 に格納された複数の骨格構造を、骨格構造の特徴量の類似度に基づいて分類する（クラスタリングする）。分類部 1 0 4 は、人物の状態の認識処理として、骨格構造の特徴量に基づいて複数の人物の状態を分類しているとも言える。類似度は、骨格構造の特徴量間の距離である。分類部 1 0 4 は、骨格構造の全体の特徴量の類似度により分類してもよいし、骨格構造の一部の特徴量の類似度により分類してもよく、骨格構造の第 1 の部分（例えば両手）及び第 2 の部分（例えば両足）の特徴量の類似度により分類してもよい。なお、各画像における人物の骨格構造の特徴量に基づいて人物の姿勢を分類してもよいし、時系列に連続する複数の画像における人物の骨格構造の特徴量の変化に基づいて人物の行動を分類してもよい。すなわち、分類部 1 0 4 は、骨格構造の特徴量に基づいて人物の姿勢や行動を含む人物の状態を分類できる。例えば、分類部 1 0 4 は、所定の監視期間に撮像された複数の画像における複数の骨格構造を分類対象とする。分類部 1 0 4 は、分類対象の特徴量間の類似度を求め、類似度の高い骨格構造が同じクラス（似た姿勢のグループ）となるように分類する。なお、検索と同様に、分類条件をユーザが指定できるようにしてもよい。分類部 1 0 4 は、骨格構造の分類結果をデータベース 1 1 0 に格納するとともに、表示部 1 0 7 に表示する。

10

【 0 0 2 9 】

検索部 1 0 5 は、データベース 1 1 0 に格納された複数の骨格構造の中から、検索クエリ（クエリ状態）の特徴量と類似度の高い骨格構造を検索する。検索部 1 0 5 は、人物の状態の認識処理として、骨格構造の特徴量に基づいて複数の人物の状態の中から、検索条件（クエリ状態）に該当する人物の状態を検索しているとも言える。分類と同様に、類似度は、骨格構造の特徴量間の距離である。検索部 1 0 5 は、骨格構造の全体の特徴量の類似度により検索してもよいし、骨格構造の一部の特徴量の類似度により検索してもよく、骨格構造の第 1 の部分（例えば両手）及び第 2 の部分（例えば両足）の特徴量の類似度により検索してもよい。なお、各画像における人物の骨格構造の特徴量に基づいて人物の姿勢を検索してもよいし、時系列に連続する複数の画像における人物の骨格構造の特徴量の変化に基づいて人物の行動を検索してもよい。すなわち、検索部 1 0 5 は、骨格構造の特徴量に基づいて人物の姿勢や行動を含む人物の状態を検索できる。例えば、検索部 1 0 5 は、分類対象と同様に、所定の監視期間に撮像された複数の画像における複数の骨格構造の特徴量を検索対象とする。また、分類部 1 0 4 が表示した分類結果の中からユーザが指定した骨格構造（姿勢）を検索クエリ（検索キー）とする。なお、分類結果に限らず、分類されていない複数の骨格構造の中から検索クエリを選択してもよいし、検索クエリとなる骨格構造をユーザが入力してもよい。検索部 1 0 5 は、検索対象の特徴量の中から、検索クエリの骨格構造の特徴量と類似度の高い特徴量を検索する。検索部 1 0 5 は、特徴量の検索結果をデータベース 1 1 0 に格納するとともに、表示部 1 0 7 に表示する。

20

30

【 0 0 3 0 】

入力部 1 0 6 は、画像処理装置 1 0 0 を操作するユーザから入力された情報を取得する入力インタフェースである。例えば、ユーザは、監視カメラの画像から不審な状態の人物を監視する監視者である。入力部 1 0 6 は、例えば、G U I（Graphical User Interface）であり、キーボードやマウス、タッチパネル等の入力装置から、ユーザの操作に応じた情報が入力される。例えば、入力部 1 0 6 は、分類部 1 0 4 により分類された骨格構造（姿勢）の中から、指定された人物の骨格構造を検索クエリとして受け付ける。

40

【 0 0 3 1 】

表示部 1 0 7 は、画像処理装置 1 0 0 の動作（処理）の結果等を表示する表示部であり、例えば、液晶ディスプレイや有機 E L（Electro Luminescence）ディスプレイ等のディスプレイ装置である。表示部 1 0 7 は、分類部 1 0 4 の分類結果や検索部 1 0 5 の検索結果を類似度等に応じて G U I に表示する。

【 0 0 3 2 】

図 3 9 は、画像処理装置 1 0 0 のハードウェア構成例を示す図である。画像処理装置 1 0 0 は、バス 1 0 1 0、プロセッサ 1 0 2 0、メモリ 1 0 3 0、ストレージデバイス 1 0

50

40、入出力インタフェース1050、及びネットワークインタフェース1060を有する。

【0033】

バス1010は、プロセッサ1020、メモリ1030、ストレージデバイス1040、入出力インタフェース1050、及びネットワークインタフェース1060が、相互にデータを送受信するためのデータ伝送路である。ただし、プロセッサ1020などを互いに接続する方法は、バス接続に限定されない。

【0034】

プロセッサ1020は、CPU (Central Processing Unit) やGPU (Graphics Processing Unit) などを実現されるプロセッサである。

10

【0035】

メモリ1030は、RAM (Random Access Memory) などを実現される主記憶装置である。

【0036】

ストレージデバイス1040は、HDD (Hard Disk Drive)、SSD (Solid State Drive)、メモ리카ード、又はROM (Read Only Memory) などを実現される補助記憶装置である。ストレージデバイス1040は画像処理装置100の各機能(例えば画像取得部101、骨格構造検出部102、特徴量算出部103、分類部104、検索部105、及び入力部106)を実現するプログラムモジュールを記憶している。プロセッサ1020がこれら各プログラムモジュールをメモリ1030上に読み込んで実行することで、そのプログラムモジュールに対応する各機能が実現される。また、ストレージデバイス1040はデータベース110としても機能することもある。

20

【0037】

入出力インタフェース1050は、画像処理装置100と各種入出力機器とを接続するためのインタフェースである。データベース110が画像処理装置100の外部に位置する場合、画像処理装置100は、入出力インタフェース1050を介してデータベース110と接続してもよい。

【0038】

ネットワークインタフェース1060は、画像処理装置100をネットワークに接続するためのインタフェースである。このネットワークは、例えばLAN (Local Area Network) やWAN (Wide Area Network) である。ネットワークインタフェース1060がネットワークに接続する方法は、無線接続であってもよいし、有線接続であってもよい。画像処理装置100は、ネットワークインタフェース1060を介してカメラ200と通信してもよい。データベース110が画像処理装置100の外部に位置する場合、画像処理装置100は、ネットワークインタフェース1060を介してデータベース110と接続してもよい。

30

【0039】

図3～図5は、本実施の形態に係る画像処理装置100の動作を示している。図3は、画像処理装置100における画像取得から検索処理までの流れを示し、図4は、図3の分類処理(S104)の流れを示し、図5は、図3の検索処理(S105)の流れを示している。

40

【0040】

図3に示すように、画像処理装置100は、カメラ200から画像を取得する(S101)。画像取得部101は、骨格構造から分類や検索を行うために人物を撮像した画像を取得し、取得した画像をデータベース110に格納する。画像取得部101は、例えば、所定の監視期間に撮像された複数の画像を取得し、複数の画像に含まれる全ての人物について以降の処理を行う。

【0041】

続いて、画像処理装置100は、取得した人物の画像に基づいて人物の骨格構造を検出する(S102)。図6は、骨格構造の検出例を示している。図6に示すように、監視力

50

メラ等から取得した画像には複数の人物が含まれており、画像に含まれる各人物について骨格構造を検出する。

【 0 0 4 2 】

図 7 は、このとき検出する人体モデル 3 0 0 の骨格構造を示しており、図 8 ~ 図 1 0 は、骨格構造の検出例を示している。骨格構造検出部 1 0 2 は、Open Pose 等の骨格推定技術を用いて、2次元の画像から図 7 のような人体モデル（2次元骨格モデル）3 0 0 の骨格構造を検出する。人体モデル 3 0 0 は、人物の関節等のキーポイントと、各キーポイントを結ぶボーンから構成された2次元モデルである。

【 0 0 4 3 】

骨格構造検出部 1 0 2 は、例えば、画像の中からキーポイントとなり得る特徴点を抽出し、キーポイントの画像を機械学習した情報を参照して、人物の各キーポイントを検出する。図 7 の例では、人物のキーポイントとして、頭 A 1、首 A 2、右肩 A 3 1、左肩 A 3 2、右肘 A 4 1、左肘 A 4 2、右手 A 5 1、左手 A 5 2、右腰 A 6 1、左腰 A 6 2、右膝 A 7 1、左膝 A 7 2、右足 A 8 1、左足 A 8 2 を検出する。さらに、これらのキーポイントを連結した人物の骨として、頭 A 1 と首 A 2 を結ぶボーン B 1、首 A 2 と右肩 A 3 1 及び左肩 A 3 2 をそれぞれ結ぶボーン B 2 1 及びボーン B 2 2、右肩 A 3 1 及び左肩 A 3 2 と右肘 A 4 1 及び左肘 A 4 2 をそれぞれ結ぶボーン B 3 1 及びボーン B 3 2、右肘 A 4 1 及び左肘 A 4 2 と右手 A 5 1 及び左手 A 5 2 をそれぞれ結ぶボーン B 4 1 及びボーン B 4 2、首 A 2 と右腰 A 6 1 及び左腰 A 6 2 をそれぞれ結ぶボーン B 5 1 及びボーン B 5 2、右腰 A 6 1 及び左腰 A 6 2 と右膝 A 7 1 及び左膝 A 7 2 をそれぞれ結ぶボーン B 6 1 及びボーン B 6 2、右膝 A 7 1 及び左膝 A 7 2 と右足 A 8 1 及び左足 A 8 2 をそれぞれ結ぶボーン B 7 1 及びボーン B 7 2 を検出する。骨格構造検出部 1 0 2 は、検出した人物の骨格構造をデータベース 1 1 0 に格納する。

【 0 0 4 4 】

図 8 は、直立した状態の人物を検出する例である。図 8 では、直立した人物が正面から撮像されており、正面から見たボーン B 1、ボーン B 5 1 及びボーン B 5 2、ボーン B 6 1 及びボーン B 6 2、ボーン B 7 1 及びボーン B 7 2 がそれぞれ重ならず検出され、右足のボーン B 6 1 及びボーン B 7 1 は左足のボーン B 6 2 及びボーン B 7 2 よりも多少折れ曲がっている。

【 0 0 4 5 】

図 9 は、しゃがみ込んでいる状態の人物を検出する例である。図 9 では、しゃがみ込んでいる人物が右側から撮像されており、右側から見たボーン B 1、ボーン B 5 1 及びボーン B 5 2、ボーン B 6 1 及びボーン B 6 2、ボーン B 7 1 及びボーン B 7 2 がそれぞれ検出され、右足のボーン B 6 1 及びボーン B 7 1 と左足のボーン B 6 2 及びボーン B 7 2 は大きく折れ曲がり、かつ、重なっている。

【 0 0 4 6 】

図 1 0 は、寝込んでいる状態の人物を検出する例である。図 1 0 では、寝込んでいる人物が左斜め前から撮像されており、左斜め前から見たボーン B 1、ボーン B 5 1 及びボーン B 5 2、ボーン B 6 1 及びボーン B 6 2、ボーン B 7 1 及びボーン B 7 2 がそれぞれ検出され、右足のボーン B 6 1 及びボーン B 7 1 と左足のボーン B 6 2 及びボーン B 7 2 は折れ曲がり、かつ、重なっている。

【 0 0 4 7 】

続いて、図 3 に示すように、画像処理装置 1 0 0 は、検出された骨格構造の特徴量を算出する（S 1 0 3）。例えば、骨格領域の高さや面積を特徴量とする場合、特徴量算出部 1 0 3 は、骨格構造を含む領域を抽出し、その領域の高さ（画素数）や面積（画素面積）を求める。骨格領域の高さや面積は、抽出される骨格領域の端部の座標や端部のキーポイントの座標から求められる。特徴量算出部 1 0 3 は、求めた骨格構造の特徴量をデータベース 1 1 0 に格納する。なお、この骨格構造の特徴量は、人物の姿勢を示す姿勢情報としても用いられる。

【 0 0 4 8 】

10

20

30

40

50

図 8 の例では、直立した人物の骨格構造から全てのボーンを含む骨格領域を抽出する。この場合、骨格領域の上端は頭部のキーポイント A 1、骨格領域の下端は左足のキーポイント A 8 2、骨格領域の左端は右肘のキーポイント A 4 1、骨格領域の右端は左手のキーポイント A 5 2 となる。このため、キーポイント A 1 とキーポイント A 8 2 の Y 座標の差分から骨格領域の高さを求める。また、キーポイント A 4 1 とキーポイント A 5 2 の X 座標の差分から骨格領域の幅を求め、骨格領域の高さと幅から面積を求める。

【 0 0 4 9 】

図 9 の例では、しゃがみ込んだ人物の骨格構造から全てのボーンを含む骨格領域を抽出する。この場合、骨格領域の上端は頭部のキーポイント A 1、骨格領域の下端は右足のキーポイント A 8 1、骨格領域の左端は右腰のキーポイント A 6 1、骨格領域の右端は右手のキーポイント A 5 1 となる。このため、キーポイント A 1 とキーポイント A 8 1 の Y 座標の差分から骨格領域の高さを求める。また、キーポイント A 6 1 とキーポイント A 5 1 の X 座標の差分から骨格領域の幅を求め、骨格領域の高さと幅から面積を求める。

10

【 0 0 5 0 】

図 10 の例では、画像の左右方向に寝込んだ人物の骨格構造から全てのボーンを含む骨格領域を抽出する。この場合、骨格領域の上端は左肩のキーポイント A 3 2、骨格領域の下端は左手のキーポイント A 5 2、骨格領域の左端は右手のキーポイント A 5 1、骨格領域の右端は左足のキーポイント A 8 2 となる。このため、キーポイント A 3 2 とキーポイント A 5 2 の Y 座標の差分から骨格領域の高さを求める。また、キーポイント A 5 1 とキーポイント A 8 2 の X 座標の差分から骨格領域の幅を求め、骨格領域の高さと幅から面積を求める。

20

【 0 0 5 1 】

続いて、図 3 に示すように、画像処理装置 100 は、分類処理を行う (S 104)。分類処理では、図 4 に示すように、分類部 104 は、算出された骨格構造の特徴量の類似度を算出し (S 111)、算出された類似度に基づいて骨格構造を分類する (S 112)。分類部 104 は、分類対象であるデータベース 110 に格納されている全ての骨格構造間の特徴量の類似度を求め、最も類似度が高い骨格構造 (姿勢) を同じクラスに分類する (クラスタリングする)。さらに、分類したクラス間での類似度を求めて分類し、所定の数のクラスとなるまで分類を繰り返す。図 11 は、骨格構造の特徴量の分類結果のイメージを示している。図 11 は、2次元の分類要素によるクラスタ分析のイメージであり、2つ分類要素は、例えば、骨格領域の高さと骨格領域の面積等である。図 11 では、分類の結果、複数の骨格構造の特徴量が3つのクラス C1 ~ C3 に分類されている。クラス C1 ~ C3 は、例えば、立っている姿勢、座っている姿勢、寝ている姿勢のように各姿勢に対応し、似ている姿勢ごとに骨格構造 (人物) が分類される。

30

【 0 0 5 2 】

本実施の形態では、人物の骨格構造の特徴量に基づいて分類することにより、多様な分類方法を用いることができる。なお、分類方法は、予め設定されていてもよいし、ユーザが任意に設定できるようにしてもよい。また、後述する検索方法と同じ方法により分類を行ってもよい。つまり、検索条件と同様の分類条件により分類してもよい。例えば、分類部 104 は、次の分類方法により分類を行う。いずれかの分類方法を用いてもよいし、任意に選択された分類方法を組み合わせてもよい。

40

【 0 0 5 3 】

(分類方法 1) 複数の階層による分類

__全身の骨格構造による分類や、上半身や下半身の骨格構造による分類、腕や脚の骨格構造による分類等を階層的に組み合わせて分類する。すなわち、骨格構造の第 1 の部分や第 2 の部分の特徴量に基づいて分類し、さらに、第 1 の部分や第 2 の部分の特徴量に重みづけを行って分類してもよい。

【 0 0 5 4 】

(分類方法 2) 時系列に沿った複数枚の画像による分類

__時系列に連続する複数の画像における骨格構造の特徴量に基づいて分類する。例えば、

50

時系列方向に特徴量を積み重ねて、累積値に基づいて分類してもよい。さらに、連続する複数の画像における骨格構造の特徴量の変化（変化量）に基づいて分類してもよい。

【0055】

（分類方法3）骨格構造の左右を無視した分類

人物の右側と左側が反対の骨格構造を同じ骨格構造として分類する。

【0056】

さらに、分類部104は、骨格構造の分類結果を表示する（S113）。分類部104は、データベース110から必要な骨格構造や人物の画像を取得し、分類結果として似ている姿勢（クラス）ごとに骨格構造及び人物を表示部107に表示する。図12は、姿勢を3つに分類した場合の表示例を示している。例えば、図12に示すように、表示ウィンドウW1に、姿勢ごとの姿勢領域WA1～WA3を表示し、姿勢領域WA1～WA3にそれぞれ該当する姿勢の骨格構造及び人物（イメージ）を表示する。姿勢領域WA1は、例えば立っている姿勢の表示領域であり、クラスC1に分類された、立っている姿勢に似た骨格構造及び人物を表示する。姿勢領域WA2は、例えば座っている姿勢の表示領域であり、クラスC2に分類された、座っている姿勢に似た骨格構造及び人物を表示する。姿勢領域WA3は、例えば寝ている姿勢の表示領域であり、クラスC3に分類された、寝ている姿勢に似た骨格構造及び人物を表示する。

10

【0057】

続いて、図3に示すように、画像処理装置100は、検索処理を行う（S105）。検索処理では、図5に示すように、検索部105は、検索条件の入力を受け付け（S121）、検索条件に基づいて骨格構造を検索する（S122）。検索部105は、入力部106から、ユーザの操作に応じて検索条件である検索クエリの入力を受け付ける。分類結果から検索クエリを入力する場合、例えば、図12の表示例では、ユーザは、表示ウィンドウW1に表示されている姿勢領域WA1～WA3の中から検索したい姿勢の骨格構造を指定（選択）する。そうすると、検索部105は、ユーザにより指定された骨格構造を検索クエリとして、検索対象であるデータベース110に格納されている全ての骨格構造の中から特徴量の類似度が高い骨格構造を検索する。検索部105は、検索クエリの骨格構造の特徴量と検索対象の骨格構造の特徴量との類似度を算出し、算出した類似度が所定の閾値よりも高い骨格構造を抽出する。検索クエリの骨格構造の特徴量は、予め算出された特徴量を使用してもよいし、検索時に求めた特徴量を使用してもよい。なお、検索クエリは、ユーザの操作に応じて骨格構造の各部を動かすことで入力してもよいし、ユーザがカメラの前で実演した姿勢を検索クエリとしてもよい。

20

30

【0058】

本実施の形態では、分類方法と同様に、人物の骨格構造の特徴量に基づいて検索することにより、多様な検索方法を用いることができる。なお、検索方法は、予め設定されていてもよいし、ユーザが任意に設定できるようにしてもよい。例えば、検索部105は、次の検索方法により検索を行う。いずれかの検索方法を用いてもよいし、任意に選択された検索方法を組み合わせてもよい。複数の検索方法（検索条件）を論理式（例えばAND（論理積）、OR（論理和）、NOT（否定））により組み合わせて検索してもよい。例えば、検索条件を「（右手を挙げている姿勢）AND（左足を挙げている姿勢）」として検索してもよい。

40

【0059】

（検索方法1）高さ方向の特徴量のみによる検索

人物の高さ方向の特徴量のみを用いて検索することで、人物の横方向の変化の影響を抑えることができ、人物の向きや人物の体型の変化に対しロバスト性が向上する。例えば、図13の骨格構造501～503のように、人物の向きや体型が異なる場合でも、高さ方向の特徴量は大きく変化しない。このため、骨格構造501～503では、検索時（分類時）に同じ姿勢であると判断することができる。

【0060】

（検索方法2）部分検索画像において人物の体の一部が隠れている場合、認識可能な部

50

分の情報のみを用いて検索する。例えば、図 1 4 の骨格構造 5 1 1 及び 5 1 2 のように、左足が隠れていることにより、左足のキーポイントが検出できない場合でも、検出されている他のキーポイントの特徴量を使用して検索できる。このため、骨格構造 5 1 1 及び 5 1 2 では、検索時（分類時）に同じ姿勢であると判断することができる。つまり、全てのキーポイントではなく、一部のキーポイントの特徴量を用いて、分類や検索を行うことができる。図 1 5 の骨格構造 5 2 1 及び 5 2 2 の例では、両足の向きが異なっているものの、上半身のキーポイント（A 1、A 2、A 3 1、A 3 2、A 4 1、A 4 2、A 5 1、A 5 2）の特徴量を検索クエリとすることで、同じ姿勢であると判断することができる。また、検索したい部分（特徴点）に対して、重みを付けて検索してもよいし、類似度判定の閾値を変化させてもよい。体の一部が隠れている場合、隠れた部分を無視して検索してもよいし、隠れた部分を加味して検索してもよい。隠れた部分も含めて検索することで、同じ部位が隠れているような姿勢を検索することができる。

10

【 0 0 6 1 】

（検索方法 3）骨格構造の左右を無視した検索

__人物の右側と左側が反対の骨格構造を同じ骨格構造として検索する。例えば、図 1 6 の骨格構造 5 3 1 及び 5 3 2 のように、右手を挙げている姿勢と、左手を挙げている姿勢を同じ姿勢として検索（分類）できる。図 1 6 の例では、骨格構造 5 3 1 と骨格構造 5 3 2 は、右手のキーポイント A 5 1、右肘のキーポイント A 4 1、左手のキーポイント A 5 2、左肘のキーポイント A 4 2 の位置が異なるものの、その他のキーポイントの位置は同じである。骨格構造 5 3 1 の右手のキーポイント A 5 1 及び右肘のキーポイント A 4 1 と骨格構造 5 3 2 の左手のキーポイント A 5 2 及び左肘のキーポイント A 4 2 のうち、一方の骨格構造のキーポイントを左右反転させると、他方の骨格構造のキーポイントと同じ位置となり、また、骨格構造 5 3 1 の左手のキーポイント A 5 2 及び左肘のキーポイント A 4 2 と骨格構造 5 3 2 の右手のキーポイント A 5 1 及び右肘のキーポイント A 4 1 のうち、一方の骨格構造のキーポイントを左右反転させると、他方の骨格構造のキーポイントと同じ位置となるため、同じ姿勢と判断する。

20

【 0 0 6 2 】

（検索方法 4）縦方向と横方向の特徴量による検索

__人物の縦方向（Y 軸方向）の特徴量のみで検索を行った後、得られた結果をさらに人物の横方向（X 軸方向）の特徴量を用いて検索する。

30

【 0 0 6 3 】

（検索方法 5）時系列に沿った複数枚の画像による検索

__時系列に連続する複数の画像における骨格構造の特徴量に基づいて検索する。例えば、時系列方向に特徴量を積み重ねて、累積値に基づいて検索してもよい。さらに、連続する複数の画像における骨格構造の特徴量の変化（変化量）に基づいて検索してもよい。

【 0 0 6 4 】

さらに、検索部 1 0 5 は、骨格構造の検索結果を表示する（S 1 2 3）。検索部 1 0 5 は、データベース 1 1 0 から必要な骨格構造や人物の画像を取得し、検索結果として得られた骨格構造及び人物を表示部 1 0 7 に表示する。例えば、検索クエリ（検索条件）が複数指定されている場合、検索クエリごとに検索結果を表示する。図 1 7 は、3 つの検索クエリ（姿勢）により検索した場合の表示例を示している。例えば、図 1 7 に示すように、表示ウィンドウ W 2 において、左端部に指定された検索クエリ Q 1 0、Q 2 0、Q 3 0 の骨格構造及び人物を表示し、検索クエリ Q 1 0、Q 2 0、Q 3 0 の右側に各検索クエリの検索結果 Q 1 1、Q 2 1、Q 3 1 の骨格構造及び人物を並べて表示する。

40

【 0 0 6 5 】

検索結果を検索クエリの隣から並べて表示する順番は、該当する骨格構造が見つかった順でもよいし、類似度が高い順でもよい。部分検索の部分（特徴点）に重みを付けて検索した場合に、重み付けて計算した類似度順に表示してもよい。ユーザが選択した部分（特徴点）のみから計算した類似度順に表示してもよい。また、検索結果の画像（フレーム）を中心に、時系列の前後の画像（フレーム）を一定時間分切り出して表示してもよい。

50

【 0 0 6 6 】

(検索方法 6) 本検索方法において、データベース 1 1 0 は動画を記憶している。動画を構成する複数のフレーム画像のそれぞれには、上記した処理が行われている。そして画像処理装置 1 0 0 はクエリとなる動画 (以下、クエリ動画と記載) に類似する動画をデータベース 1 1 0 から検索する。

【 0 0 6 7 】

図 4 0 は、本検索方法における検索部 1 0 5 の機能構成を示す図である。本図に示す検索部 1 0 5 は、クエリ取得部 6 1 0、クエリフレーム選択部 6 2 0、及び動画選択部 6 3 0 を備えている。

【 0 0 6 8 】

クエリ取得部 6 1 0 は、上記したクエリ動画を取得する。以下の説明において、クエリ動画に含まれるフレーム画像を第 1 フレーム画像と記載する。言い換えると、クエリ画像は複数の第 1 フレーム画像を含んでいる。

【 0 0 6 9 】

クエリフレーム選択部 6 2 0 は、複数の第 1 フレーム画像から複数のクエリフレームを選択する。クエリフレームは、検索対象となっている複数の動画からクエリ動画に類似する動画を選択する際に用いられる。クエリフレームの数は、2 以上であればよいが、多いほど好ましい。クエリフレームの選択基準については後述する。

【 0 0 7 0 】

なお、あるクエリフレーム画像と、その次のクエリフレーム画像の間に、他のフレーム画像が存在することがある。後述する選択基準のいずれにおいても、クエリフレームの間のフレーム画像の数 (又は時間) が基準以内であるのが好ましい。ここで用いられる基準は、例えば 1 フレーム以上 1 0 フレーム以下、又は 0 . 0 2 5 秒以上 1 秒以下である。

【 0 0 7 1 】

動画選択部 6 3 0 は、複数のクエリフレームを用いて動画を選択する。この際、動画選択部 6 3 0 は、動画を選択する条件として、以下の基準 (1) 及び (2) を用いる。

基準 (1) 「クエリフレームに対する類似度が第 1 基準を満たす類似フレーム画像が存在する」という条件が少なくとも 2 つのクエリフレームについて満たされる。

基準 (2) 少なくとも 2 つの類似フレーム画像のそれぞれに対応するクエリフレームを少なくとも 2 つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、クエリ動画における少なくとも 2 つのクエリフレームの並び順に一致する。

【 0 0 7 2 】

まず、上記 (1) について説明する。第 1 の動画と第 2 の動画が互いに類似する場合、第 1 の動画に含まれる少なくとも 2 つのフレーム画像は、それぞれ、第 2 の動画に含まれるいずれかのフレーム画像に類似しているはずである。上記 (1) は、これに対応している。

【 0 0 7 3 】

次に、上記 (2) について説明する。複数のフレーム画像が互いに類似していたとしても、これら複数のフレーム画像の出現順が一致していない場合、第 1 の動画と第 2 の動画は互いに異なると判断すべきである。基準 (2) は、複数のフレーム画像の出現順が一致していることを求めている。

【 0 0 7 4 】

上記した基準 (1) 及び (2) は、以下のようにも置き換えることができる。

基準 (1) ` 少なくとも第 1 のクエリフレームに対する類似度が基準を満たす第 1 の類似フレーム画像と、第 2 の前記クエリフレームに対する類似度が基準を満たす第 2 の類似フレーム画像とを含んでいる。

基準 (2) ` 第 1 の類似フレーム画像と第 2 の類似フレーム画像の順序が第 1 のクエリフレームと第 2 のクエリフレームの順序と同一である。

【 0 0 7 5 】

なお、第 1 の類似フレーム画像と第 2 の類似フレーム画像の間に、他のフレーム画像が

10

20

30

40

50

存在することがある。この場合、第1の類似フレームと第2の類似フレームの間のフレーム画像の数（又は時間）が基準以内であることを、基準（2）（又は基準（2）'）にさらに加えてもよい。ここで用いられる基準は、例えば1フレーム以上10フレーム以下、又は0.025秒以上1秒以下である。

【0076】

本検索方法が適用できる動画は、人物の姿勢を含む動画に限定されない。ただし、人物の姿勢を含む動画を検索する場合、フレーム画像間の類似度には、人物の姿勢の類似度を用いることができる。

【0077】

フレーム画像間の類似度として人物の姿勢の類似度を用いる場合、クエリフレームとして選択されるための条件は、「クエリフレームに含まれる人物に関する情報量が第2基準を満たしていること」を含むのが好ましい。言い換えると、画像解析によってその人物に関する情報が得られるが、クエリフレームは、この情報量がある程度多いことが必要である。第2基準の一例は、クエリフレームに含まれるキーポイント（関節）の数が基準を満たしていることである。この基準は、例えばキーポイントの全数に対して30%以上（例えば14個中5個）の値として設定される。ここで、手足に相当するキーポイントを少なくとも一つ以上含んでいるのが好ましい。なお、キーポイントのうち首、両肩、及び頭部に相当する部分は欠損しにくい。キーポイントの数が少ない場合、人物の姿勢に関する情報量が少なくなるため、動画の検索精度が低下してしまう。

【0078】

図41は、動画選択部630が類似フレームを選択する方法の一例を説明するための図である。本図に示す例において、人物の姿勢は、上記した骨格構造の特徴量（パラメータ）によって定義されている。そして動画選択部630は、この特徴量によって定義される空間（以下、特徴量空間と記載）における距離を用いて類似フレームを選択する。詳細には、動画選択部630は、特徴量空間においてクエリフレームからの距離が第1基準以内のフレーム画像を、類似フレーム画像として選択する。なお、動画選択部630は、第1基準をユーザからの入力に従って設定する。ただし動画選択部630は、第1基準として予め定められた値を用いてもよい。

【0079】

図42は、本検索方法に係る検索部105が行う処理の一例を示すフローチャートである。図43は、図42に示す処理を説明するための図である。

【0080】

まず、クエリ取得部610はクエリ動画を取得する（ステップS300）。一例として、クエリ取得部610は、ユーザが指定した動画をクエリ動画として取得する。ユーザは、データベース110に記憶されている動画からクエリ動画を選択してもよいし、外部の装置または記憶媒体からクエリ動画を検索部105に取得させてもよい。

【0081】

次いでクエリフレーム選択部620は、クエリ動画に含まれる複数の第1フレーム画像から、複数のクエリフレームを選択する。本図に示す例では、クエリフレーム選択部620は、ユーザの入力に従って複数のクエリフレームを選択する（ステップS302）。

【0082】

詳細には、クエリフレーム選択部620は、図43に示すように、クエリ動画に含まれる複数の第1フレーム画像（好ましくはすべての第1フレーム画像）を出現順に並べて表示部107に表示させる。するとユーザは、マウス等の入力デバイスを用いて、クエリフレームとすべき第1フレーム画像を選択する。

【0083】

次いで動画選択部630は、ステップS302で選択された複数のクエリフレームを用いて、クエリ動画に類似する動画を選択する（ステップS304）。この選択基準は、図40及び図41を用いて説明した通りである。なお、図43に示す図においても、選択された動画は上記した基準（1）及び（2）を満たしているとともに、基準（1）'及び（

10

20

30

40

50

2) も満たしている。

【0084】

なお、クエリ動画がデータベース110に記憶されていなかった場合、骨格構造検出部102及び特徴量算出部103は、クエリフレームを処理して骨格構造の特徴量を算出する。そして動画選択部630は、この骨格構造の特徴量を用いて、クエリ動画に類似する動画を選択する。

【0085】

その後、動画選択部630は、選択結果をデータベース110に記憶させる。ここで動画選択部630は、選択した動画そのものを記憶してもよいし、データベース110に既に記憶されている当該動画に、その画像がクエリ動画に類似することを示すフラグを紐づけてもよい(ステップS306)。

10

【0086】

図44は、図42の変形例を示すフローチャートである。本変形例は、クエリフレーム選択部620が予め定められた規則(以下、選択規則と記載)に従って複数のクエリ画像を選択する(ステップS303)点を除いて、図42に示した例と同様である。

【0087】

図45は、クエリフレームの選択規則の第1例を説明するための図である。本図に示す例において、クエリフレーム選択部620は、所定の間隔で複数の第1フレーム画像から複数のクエリフレームを選択する。ここでクエリフレーム選択部620は、最初のクエリフレームとして、 n 番目の第1フレーム画像を選択する。ここで n は、予め設定されている整数であってもよいし、ユーザ入力によって設定された整数であってもよい。また、所定の間隔は、例えば2フレーム以上10フレーム以下、又は0.05秒以上1秒以下である。

20

【0088】

図46は、クエリフレームの選択規則の第2例を説明するための図である。本図に示す例において、クエリフレーム選択部620は、第1のクエリフレームが特定された後、当該第1のクエリフレームからの変化量が第3基準以上となった第1フレーム画像を、第1のクエリフレームの次のクエリフレームとして選択する。そしてクエリフレーム選択部620は、この処理を繰り返す。これにより、隣り合うクエリフレームに含まれる人物の姿勢は、互いにある程度異なる。従って、クエリフレームが増加することを抑制しつつ、検索部105による動画の検索精度を高くすることができる。

30

【0089】

ここで、第3基準はユーザからの入力に従って設定されてもよいし、複数の第1クエリフレームを用いて設定されてもよい。

【0090】

図47は、クエリフレーム選択部620が複数の第1クエリフレームを用いて第3基準を設定する方法を説明するための図である。本図に示す例において、第3基準は、特徴量空間における距離である。そして具体例としては、以下の2つがある。

【0091】

第1の例において、クエリフレーム選択部620は、特徴量空間における2つのクエリフレームの最大距離を特定し、この最大距離を用いて第3基準を設定する。一例として、クエリフレーム選択部620は、最大距離を変数とした関数を用いて、第3基準を最大距離未満の値に設定する。一例として、クエリフレーム選択部620は、最大距離に1未満の係数を乗じたり、最大距離から所定の値を引くことにより、第3基準を設定する。

40

【0092】

第2の例において、クエリフレーム選択部620は、時間的に隣り合う2つの第1フレームの距離を統計処理した結果を用いて、第3基準を設定する。例えばクエリフレーム選択部620は、これらの距離の中央値、平均値、又は最頻値を第3基準としてもよいし、中央値、平均値、及び最頻値の少なくとも一つを変数とした関数を用いて第3基準を設定してもよい。

50

【 0 0 9 3 】

なお、図 4 6 及び図 4 7 に示す例においても、クエリフレーム選択部 6 2 0 は、図 4 5 に示した例と同様に、最初のクエリフレームとして n 番目の第 1 フレーム画像を選択する。

【 0 0 9 4 】

以上のように、本実施の形態では、2次元画像から人物の骨格構造を検出し、検出した骨格構造の特徴量に基づいて分類や検索を行うことを可能とした。これにより、類似度が高い似た姿勢ごとに分類することができ、また、検索クエリ（検索キー）と類似度が高い似た姿勢を検索することができる。画像から似ている姿勢を分類し表示することで、ユーザが姿勢等を指定することなく、画像中の人物の姿勢を把握することができる。分類結果の中からユーザが検索クエリの姿勢を指定できるため、予めユーザが検索したい姿勢を詳細に把握していない場合でも、所望の姿勢を検索することができる。例えば、人物の骨格構造の全体や一部等を条件として分類や検索を行うことができるため、柔軟な分類や検索が可能となる。

10

【 0 0 9 5 】

また、検索方法 6 によれば、クエリ動画に類似する動画を精度良く検索することができる。また、クエリ動画と検索対象となる動画との間で、フレームレートが異なったり再生時間が異なっていた場合でも、動画の検索を行うことができる。

【 0 0 9 6 】

（実施の形態 2）以下、図面を参照して実施の形態 2 について説明する。本実施の形態では、実施の形態 1 における特徴量算出の具体例について説明する。本実施の形態では、人物の身長を用いて正規化することで特徴量を求める。その他については、実施の形態 1 と同様である。

20

【 0 0 9 7 】

図 1 8 は、本実施の形態に係る画像処理装置 1 0 0 の構成を示している。図 1 8 に示すように、画像処理装置 1 0 0 は、実施の形態 1 の構成に加えて、さらに身長算出部 1 0 8 を備える。なお、特徴量算出部 1 0 3 と身長算出部 1 0 8 を一つの処理部としてもよい。

【 0 0 9 8 】

身長算出部（身長推定部）1 0 8 は、骨格構造検出部 1 0 2 により検出された 2 次元の骨格構造に基づき、2次元の画像内の人物の直立時の高さ（身長画素数という）を算出（推定）する。身長画素数は、2次元の画像における人物の身長（2次元画像空間上の人物の全身の長さ）であるとも言える。身長算出部 1 0 8 は、検出された骨格構造の各ボーンの長さ（2次元画像空間上の長さ）から身長画素数（ピクセル数）を求める。

30

【 0 0 9 9 】

以下の例では、身長画素数を求める方法として具体例 1 ~ 3 を用いる。なお、具体例 1 ~ 3 のいずれかの方法を用いてもよいし、任意に選択される複数の方法を組み合わせて用いてもよい。具体例 1 では、骨格構造の各ボーンのうち、頭部から足部までのボーンの長さを合計することで、身長画素数を求める。骨格構造検出部 1 0 2（骨格推定技術）が頭頂と足元を出力しない場合は、必要に応じて定数を乗じて補正することもできる。具体例 2 では、各ボーンの長さ（2次元画像空間上の身長）との関係を示す人体モデルを用いて、身長画素数を算出する。具体例 3 では、3次元人体モデルを2次元骨格構造にフィッティング（あてはめる）することで、身長画素数を算出する。

40

【 0 1 0 0 】

本実施の形態の特徴量算出部 1 0 3 は、算出された人物の身長画素数に基づいて、人物の骨格構造（骨格情報）を正規化する正規化部である。特徴量算出部 1 0 3 は、正規化した骨格構造の特徴量（正規化値）をデータベース 1 1 0 に格納する。特徴量算出部 1 0 3 は、骨格構造に含まれる各キーポイント（特徴点）の画像上での高さを、身長画素数で正規化する。本実施の形態では、例えば、高さ方向は、画像の 2 次元座標（X - Y 座標）空間における上下の方向（Y 軸方向）である。この場合、キーポイントの高さは、キーポイントの Y 座標の値（画素数）から求めることができる。あるいは、高さ方向は、実世界の 3 次元座標空間における地面（基準面）に対し垂直な鉛直軸の方向を、2次元座標空間に

50

投影した鉛直投影軸の方向（鉛直投影方向）でもよい。この場合、キーポイントの高さは、実世界における地面に対し垂直な軸を、カメラパラメータに基づいて2次元座標空間に投影した鉛直投影軸を求め、この鉛直投影軸に沿った値（画素数）から求めることができる。なお、カメラパラメータは、画像の撮像パラメータであり、例えば、カメラパラメータは、カメラ200の姿勢、位置、撮像角度、焦点距離等である。カメラ200により、予め長さや位置が分かっている物体を撮像し、その画像からカメラパラメータを求めることができる。撮像された画像の両端ではひずみが発生し、実世界の鉛直方向と画像の上下方向が合わない場合がある。これに対し、画像を撮影したカメラのパラメータを使用することで、実世界の鉛直方向が画像中でどの程度傾いているのかが分かる。このため、カメラパラメータに基づいて画像中に投影した鉛直投影軸に沿ったキーポイントの値を身長で正規化することで、実世界と画像のずれを考慮してキーポイントを特徴量化することができる。なお、左右方向（横方向）は、画像の2次元座標（X-Y座標）空間における左右の方向（X軸方向）であり、または、実世界の3次元座標空間における地面に対し平行な方向を、2次元座標空間に投影した方向である。

10

【0101】

図19～図23は、本実施の形態に係る画像処理装置100の動作を示している。図19は、画像処理装置100における画像取得から検索処理までの流れを示し、図20～図22は、図19の身長画素数算出処理（S201）の具体例1～3の流れを示し、図23は、図19の正規化処理（S202）の流れを示している。

【0102】

図19に示すように、本実施の形態では、実施の形態1における特徴量算出処理（S103）として、身長画素数算出処理（S201）及び正規化処理（S202）を行う。その他については実施の形態1と同様である。

20

【0103】

画像処理装置100は、画像取得（S101）及び骨格構造検出（S102）に続いて、検出された骨格構造に基づいて身長画素数算出処理を行う（S201）。この例では、図24に示すように、画像における直立時の人物の骨格構造の高さを身長画素数（ h ）とし、画像の人物の状態における骨格構造の各キーポイントの高さをキーポイント高さ（ y_i ）とする。以下、身長画素数算出処理の具体例1～3について説明する。

【0104】

<具体例1> 具体例1では、頭部から足部までのボーンの長さをを用いて身長画素数を求める。具体例1では、図20に示すように、身長算出部108は、各ボーンの長さを取得し（S211）、取得した各ボーンの長さを合計する（S212）。

30

【0105】

身長算出部108は、人物の頭部から足部の2次元の画像上のボーンの長さを取得し、身長画素数を求める。すなわち、骨格構造を検出した画像から、図24のボーンのうち、ボーンB1（長さ L_1 ）、ボーンB51（長さ L_{21} ）、ボーンB61（長さ L_{31} ）及びボーンB71（長さ L_{41} ）、もしくは、ボーンB1（長さ L_1 ）、ボーンB52（長さ L_{22} ）、ボーンB62（長さ L_{32} ）及びボーンB72（長さ L_{42} ）の各長さ（画素数）を取得する。各ボーンの長さは、2次元の画像における各キーポイントの座標から求めることができる。これらを合計した、 $L_1 + L_{21} + L_{31} + L_{41}$ 、もしくは、 $L_1 + L_{22} + L_{32} + L_{42}$ に補正定数を乗じた値を身長画素数（ h ）として算出する。両方の値を算出できる場合、例えば、長い方の値を身長画素数とする。すなわち、各ボーンは正面から撮像された場合が画像中で最も長くなり、カメラに対して奥行き方向に傾くと短く表示される。従って、長いボーンの方が正面から撮像されている可能性が高く、真実の値に近いと考えられる。このため、長い方の値を選択することが好ましい。

40

【0106】

図25の例では、ボーンB1、ボーンB51及びボーンB52、ボーンB61及びボーンB62、ボーンB71及びボーンB72がそれぞれ重ならず検出されている。これらのボーンの合計である、 $L_1 + L_{21} + L_{31} + L_{41}$ 、及び、 $L_1 + L_{22} + L_{32} +$

50

L 4 2 を求め、例えば、検出されたボーンの長さが長い左足側の $L 1 + L 2 2 + L 3 2 + L 4 2$ に補正定数を乗じた値を身長画素数とする。

【 0 1 0 7 】

図 2 6 の例では、ボーン B 1、ボーン B 5 1 及びボーン B 5 2、ボーン B 6 1 及びボーン B 6 2、ボーン B 7 1 及びボーン B 7 2 がそれぞれ検出され、右足のボーン B 6 1 及びボーン B 7 1 と左足のボーン B 6 2 及びボーン B 7 2 が重なっている。これらのボーンの合計である、 $L 1 + L 2 1 + L 3 1 + L 4 1$ 、及び、 $L 1 + L 2 2 + L 3 2 + L 4 2$ を求め、例えば、検出されたボーンの長さが長い右足側の $L 1 + L 2 1 + L 3 1 + L 4 1$ に補正定数を乗じた値を身長画素数とする。

【 0 1 0 8 】

図 2 7 の例では、ボーン B 1、ボーン B 5 1 及びボーン B 5 2、ボーン B 6 1 及びボーン B 6 2、ボーン B 7 1 及びボーン B 7 2 がそれぞれ検出され、右足のボーン B 6 1 及びボーン B 7 1 と左足のボーン B 6 2 及びボーン B 7 2 が重なっている。これらのボーンの合計である、 $L 1 + L 2 1 + L 3 1 + L 4 1$ 、及び、 $L 1 + L 2 2 + L 3 2 + L 4 2$ を求め、例えば、検出されたボーンの長さが長い左足側の $L 1 + L 2 2 + L 3 2 + L 4 2$ に補正定数を乗じた値を身長画素数とする。

【 0 1 0 9 】

具体例 1 では、頭から足までのボーンの長さを合計することで身長を求めることができるため、簡易な方法で身長画素数を求めることができる。また、機械学習を用いた骨格推定技術により、少なくとも頭から足までの骨格を検出できればよいため、しゃがみ込んでいる状態など、必ずしも人物の全体が画像に写っていない場合でも精度よく身長画素数を推定することができる。

【 0 1 1 0 】

< 具体例 2 > 具体例 2 では、2 次元骨格構造に含まれる骨の長さとの関係を示す 2 次元骨格モデルを用いて身長画素数を求める。

【 0 1 1 1 】

図 2 8 は、具体例 2 で用いる、2 次元画像空間上の各ボーンの長さとの関係を示す人体モデル (2 次元骨格モデル) 3 0 1 である。図 2 8 に示すように、平均的な人物の各ボーンの長さとの関係 (全身の長さに対する各ボーンの長さの割合) を、人体モデル 3 0 1 の各ボーンに対応付ける。例えば、頭のボーン B 1 の長さは全身の長さ $\times 0.2$ (20%) であり、右手のボーン B 4 1 の長さは全身の長さ $\times 0.15$ (15%) であり、右足のボーン B 7 1 の長さは全身の長さ $\times 0.25$ (25%) である。このような人体モデル 3 0 1 の情報をデータベース 1 1 0 に記憶しておくことで、各ボーンの長さから平均的な全身の長さを求めることができる。平均的な人物の人体モデルの他に、年代、性別、国籍等の人物の属性ごとに人体モデルを用意してもよい。これにより、人物の属性に応じて適切に全身の長さ (身長) を求めることができる。

【 0 1 1 2 】

具体例 2 では、図 2 1 に示すように、身長算出部 1 0 8 は、各ボーンの長さを取得する (S 2 2 1)。身長算出部 1 0 8 は、検出された骨格構造において、全てのボーンの長さ (2 次元画像空間上の長さ) を取得する。図 2 9 は、しゃがみ込んでいる状態の人物を右斜め後ろから撮像し、骨格構造を検出した例である。この例では、人物の顔や左側面が写っていないことから、頭のボーンと左腕及び左手のボーンが検出できていない。このため、検出されているボーン B 2 1、B 2 2、B 3 1、B 4 1、B 5 1、B 5 2、B 6 1、B 6 2、B 7 1、B 7 2 の各長さを取得する。

【 0 1 1 3 】

続いて、身長算出部 1 0 8 は、図 2 1 に示すように、人体モデルに基づき、各ボーンの長さから身長画素数を算出する (S 2 2 2)。身長算出部 1 0 8 は、図 2 8 のような、各ボーンと全身の長さとの関係を示す人体モデル 3 0 1 を参照し、各ボーンの長さから身長画素数を求める。例えば、右手のボーン B 4 1 の長さが全身の長さ $\times 0.15$ であるため、ボーン B 4 1 の長さ / 0.15 によりボーン B 4 1 に基づいた身長画素数を求める。ま

10

20

30

40

50

た、右足のボーン B 7 1 の長さが全身の長さ $\times 0.25$ であるため、ボーン B 7 1 の長さ / 0.25 によりボーン B 7 1 に基づいた身長画素数を求める。

【 0 1 1 4 】

このとき参照する人体モデルは、例えば、平均的な人物の人体モデルであるが、年代、性別、国籍等の人物の属性に応じて人体モデルを選択してもよい。例えば、撮像した画像に人物の顔が写っている場合、顔に基づいて人物の属性を識別し、識別した属性に対応する人体モデルを参照する。属性ごとの顔を機械学習した情報を参照し、画像の顔の特徴から人物の属性を認識することができる。また、画像から人物の属性が識別できない場合に、平均的な人物の人体モデルを用いてもよい。

【 0 1 1 5 】

また、ボーンの長さから算出した身長画素数をカメラパラメータにより補正してもよい。例えばカメラを高い位置において、人物を見下ろすように撮影した場合、二次元骨格構造において肩幅のボーン等の横の長さはカメラの俯角の影響を受けないが、首 - 腰のボーン等の縦の長さは、カメラの俯角が大きくなる程小さくなる。そうすると、肩幅のボーン等の横の長さから算出した身長画素数が実際より大きくなる傾向がある。そこで、カメラパラメータを活用すると、人物がどの程度の角度でカメラに見下ろされているかがわかるため、この俯角の情報を使って正面から撮影したような二次元骨格構造に補正することができる。これによって、より正確に身長画素数を算出できる。

【 0 1 1 6 】

続いて、身長算出部 1 0 8 は、図 2 1 に示すように、身長画素数の最適値を算出する (S 2 2 3)。身長算出部 1 0 8 は、ボーンごとに求めた身長画素数から身長画素数の最適値を算出する。例えば、図 3 0 に示すような、ボーンごとに求めた身長画素数のヒストグラムを生成し、その中で大きい身長画素数を選択する。つまり、複数のボーンに基づいて求められた複数の身長画素数の中で他よりも長い身長画素数を選択する。例えば、上位 3 0 % を有効な値とし、図 3 0 ではボーン B 7 1、B 6 1、B 5 1 による身長画素数を選択する。選択した身長画素数の平均を最適値として求めてもよいし、最も大きい身長画素数を最適値としてもよい。2 次元画像のボーンの長さから身長を求めるため、ボーンを正面から撮像できていない場合、すなわち、ボーンがカメラから見て奥行き方向に傾いて撮像された場合、ボーンの長さが正面から撮像した場合よりも短くなる。そうすると、身長画素数が大きい値は、身長画素数が小さい値よりも、正面から撮像された可能性が高く、より尤もらしい値となることから、より大きい値を最適値とする。

【 0 1 1 7 】

具体例 2 では、2 次元画像空間上のボーンと全身の長さとの関係を示す人体モデルを用いて、検出した骨格構造のボーンに基づき身長画素数を求めるため、頭から足までの全ての骨格が得られない場合でも、一部のボーンから身長画素数を求めることができる。特に、複数のボーンから求められた値のうち、より大きい値を採用することで、精度よく身長画素数を推定することができる。

【 0 1 1 8 】

< 具体例 3 > 具体例 3 では、2 次元骨格構造を 3 次元人体モデル (3 次元骨格モデル) にフィッティングさせて、フィッティングした 3 次元人体モデルの身長画素数を用いて全身の骨格ベクトルを求める。

【 0 1 1 9 】

具体例 3 では、図 2 2 に示すように、身長算出部 1 0 8 は、まず、カメラ 2 0 0 の撮像した画像に基づき、カメラパラメータを算出する (S 2 3 1)。身長算出部 1 0 8 は、カメラ 2 0 0 が撮像した複数の画像の中から、予め長さが分かっている物体を抽出し、抽出した物体の大きさ (画素数) からカメラパラメータを求める。なお、カメラパラメータを予め求めておき、求めておいたカメラパラメータを必要に応じて取得してもよい。

【 0 1 2 0 】

続いて、身長算出部 1 0 8 は、3 次元人体モデルの配置及び高さを調整する (S 2 3 2)。身長算出部 1 0 8 は、検出された 2 次元骨格構造に対し、身長画素数算出用の 3 次元

10

20

30

40

50

人体モデルを用意し、カメラパラメータに基づいて、同じ2次元画像内に配置する。具体的には、カメラパラメータと、2次元骨格構造から、「実世界におけるカメラと人物の相対的な位置関係」を特定する。例えば、仮にカメラの位置を座標(0, 0, 0)としたときに、人物が立っている(または座っている)位置の座標(x, y, z)を特定する。そして、特定した人物と同じ位置(x, y, z)に3次元人体モデルを配置して撮像した場合の画像を想定することで、2次元骨格構造と3次元人体モデルを重ね合わせる。

【0121】

図31は、しゃがみ込んでいる人物を左斜め前から撮像し、2次元骨格構造401を検出した例である。2次元骨格構造401は、2次元の座標情報を有する。なお、全てのボーンを検出していることが好ましいが、一部のボーンが検出されていなくてもよい。この2次元骨格構造401に対し、図32のような、3次元人体モデル402を用意する。3次元人体モデル(3次元骨格モデル)402は、3次元の座標情報を有し、2次元骨格構造401と同じ形状の骨格のモデルである。そして、図33のように、検出した2次元骨格構造401に対し、用意した3次元人体モデル402を配置し重ね合わせる。また、重ね合わせるとともに、3次元人体モデル402の高さを2次元骨格構造401に合うように調整する。

10

【0122】

なお、このとき用意する3次元人体モデル402は、図33のように、2次元骨格構造401の姿勢に近い状態のモデルでもよいし、直立した状態のモデルでもよい。例えば、機械学習を用いて2次元画像から3次元空間の姿勢を推定する技術を用いて、推定した姿勢の3次元人体モデル402を生成してもよい。2次元画像の関節と3次元空間の関節の情報を学習することで、2次元画像から3次元の姿勢を推定することができる。

20

【0123】

続いて、身長算出部108は、図22に示すように、3次元人体モデルを2次元骨格構造にフィッティングする(S233)。身長算出部108は、図34のように、3次元人体モデル402を2次元骨格構造401に重ね合わせた状態で、3次元人体モデル402と2次元骨格構造401の姿勢が一致するように、3次元人体モデル402を変形させる。すなわち、3次元人体モデル402の身長、体の向き、関節の角度を調整し、2次元骨格構造401との差異がなくなるように最適化する。例えば、3次元人体モデル402の関節を、人の可動範囲で回転させていき、また、3次元人体モデル402の全体を回転させたり、全体のサイズを調整する。なお、3次元人体モデルと2次元骨格構造のフィッティング(あてはめ)は、2次元空間(2次元座標)上で行う。すなわち、2次元空間に3次元人体モデルを写像し、変形させた3次元人体モデルが2次元空間(画像)でどのように変化するかを考慮して、3次元人体モデルを2次元骨格構造に最適化する。

30

【0124】

続いて、身長算出部108は、図22に示すように、フィッティングさせた3次元人体モデルの身長画素数を算出する(S234)。身長算出部108は、図35のように、3次元人体モデル402と2次元骨格構造401の差異がなくなり、姿勢が一致すると、その状態の3次元人体モデル402の身長画素数を求める。最適化された3次元人体モデル402を直立させた状態として、カメラパラメータに基づき、2次元空間上の全身の長さを求める。例えば、3次元人体モデル402を直立させた場合の頭から足までのボーンの長さ(画素数)により身長画素数を算出する。具体例1と同様に、3次元人体モデル402の頭部から足部までのボーンの長さを合計してもよい。

40

【0125】

具体例3では、カメラパラメータに基づいて3次元人体モデルを2次元骨格構造にフィッティングさせて、その3次元人体モデルに基づいて身長画素数を求めることで、全てのボーンが正面に写っていない場合、すなわち、全てのボーンが斜めに映っているため誤差が大きい場合でも、精度よく身長画素数を推定することができる。

【0126】

<正規化処理> 図19に示すように、画像処理装置100は、身長画素数算出処理に続

50

いて、正規化処理（S202）を行う。正規化処理では、図23に示すように、特徴量算出部103は、キーポイント高さを算出する（S241）。特徴量算出部103は、検出された骨格構造に含まれる全てのキーポイントのキーポイント高さ（画素数）を算出する。キーポイント高さは、骨格構造の最下端（例えばいずれかの足のキーポイント）からそのキーポイントまでの高さ方向の長さ（画素数）である。ここでは、一例として、キーポイント高さを、画像におけるキーポイントのY座標から求める。なお、上記のように、キーポイント高さは、カメラパラメータに基づいた鉛直投影軸に沿った方向の長さから求めてもよい。例えば、図24の例で、首のキーポイントA2の高さ（ y_i ）は、キーポイントA2のY座標から右足のキーポイントA81または左足のキーポイントA82のY座標を引いた値である。

10

【0127】

続いて、特徴量算出部103は、正規化のための基準点を特定する（S242）。基準点は、キーポイントの相対的な高さを表すための基準となる点である。基準点は、予め設定されていてもよいし、ユーザが選択できるようにしてもよい。基準点は、骨格構造の中心もしくは中心よりも高い（画像の上下方向における上である）ことが好ましく、例えば、首のキーポイントの座標を基準点とする。なお、首に限らず頭やその他のキーポイントの座標を基準点としてもよい。キーポイントに限らず、任意の座標（例えば骨格構造の中心座標等）を基準点としてもよい。

【0128】

続いて、特徴量算出部103は、キーポイント高さ（ y_i ）を身長画素数で正規化する（S243）。特徴量算出部103は、各キーポイントのキーポイント高さ、基準点、身長画素数を用いて、各キーポイントを正規化する。具体的には、特徴量算出部103は、基準点に対するキーポイントの相対的な高さを身長画素数により正規化する。ここでは、高さ方向のみに着目する例として、Y座標のみを抽出し、また、基準点を首のキーポイントとして正規化を行う。具体的には、基準点（首のキーポイント）のY座標を（ y_c ）として、次の式（1）を用いて、特徴量（正規化値）を求める。なお、カメラパラメータに基づいた鉛直投影軸を用いる場合は、（ y_i ）及び（ y_c ）を鉛直投影軸に沿った方向の値に変換する。

20

【数1】

$$f_i = (y_i - y_c) / h \quad \dots (1)$$

30

【0129】

例えば、キーポイントが18個の場合、各キーポイントの18点の座標（ x_0, y_0 ）、（ x_1, y_1 ）、 \dots （ x_{17}, y_{17} ）を、上記式（1）を用いて、次のように18次元の特徴量に変換する。

【数2】

$$\begin{aligned} f_0 &= (y_0 - y_c) / h \\ f_1 &= (y_1 - y_c) / h \\ &\vdots \\ f_{17} &= (y_{17} - y_c) / h \end{aligned} \quad \dots (2)$$

40

【0130】

図36は、特徴量算出部103が求めた各キーポイントの特徴量の例を示している。この例では、首のキーポイントA2を基準点とするため、キーポイントA2の特徴量は0.0となり、首と同じ高さの右肩のキーポイントA31及び左肩のキーポイントA32の特

50

徴量も 0.0 である。首よりも高い頭のキーポイント A 1 の特徴量は - 0.2 である。首よりも低い右手のキーポイント A 5 1 及び左手のキーポイント A 5 2 の特徴量は 0.4 であり、右足のキーポイント A 8 1 及び左足のキーポイント A 8 2 の特徴量は 0.9 である。この状態から人物が左手を挙げると、図 3 7 のように左手が基準点よりも高くなるため、左手のキーポイント A 5 2 の特徴量は - 0.4 となる。一方で、Y 軸の座標のみを用いて正規化を行っているため、図 3 8 のように、図 3 6 に比べて骨格構造の幅が変わっても特徴量は変わらない。すなわち、本実施の形態の特徴量（正規化値）は、骨格構造（キーポイント）の高さ方向（Y 方向）の特徴を示しており、骨格構造の横方向（X 方向）の変化に影響を受けない。

【0131】

以上のように、本実施の形態では、2 次元画像から人物の骨格構造を検出し、検出した骨格構造から求めた身長画素数（2 次元画像空間上の直立時の高さ）を用いて、骨格構造の各キーポイントを正規化する。この正規化された特徴量を用いることで、分類や検索等を行った場合のロバスト性を向上することができる。すなわち、本実施の形態の特徴量は、上記のように人物の横方向の変化に影響を受けないため、人物の向きや人物の体型の変化に対しロバスト性が高い。

【0132】

さらに、本実施の形態では、Open Pose 等の骨格推定技術を用いて人物の骨格構造を検出することで実現できるため、人物の姿勢等を学習する学習データを用意する必要がない。また、骨格構造のキーポイントを正規化し、データベースに格納しておくことで、人物の姿勢等の分類や検索が可能となるため、未知な姿勢に対しても分類や検索を行うことができる。また、骨格構造のキーポイントを正規化することで、明確でわかりやすい特徴量を得ることができるため、機械学習のようにブラックボックス型のアルゴリズムと異なり、処理結果に対するユーザの納得性が高い。

【0133】

以上、図面を参照して本発明の実施形態について述べたが、これらは本発明の例示であり、上記以外の様々な構成を採用することもできる。

【0134】

また、上述の説明で用いた複数のフローチャートでは、複数の工程（処理）が順番に記載されているが、各実施形態で実行される工程の実行順序は、その記載の順番に制限されない。各実施形態では、図示される工程の順番を内容的に支障のない範囲で変更することができる。また、上述の各実施形態は、内容が相反しない範囲で組み合わせることができる。

【0135】

上記の実施形態の一部または全部は、以下の付記のようにも記載されうるが、以下に限られない。

1. 複数の第 1 フレーム画像を含むクエリ動画を取得するクエリ取得手段と、
前記複数の第 1 フレーム画像から複数のクエリフレームを選択するクエリフレーム選択手段と、

前記複数のクエリフレームを用いて動画を選択する動画選択手段と、
を備え、

前記動画選択手段は、前記動画を選択する条件として、少なくとも以下の（1）及び（2）を用いる、画像選択装置。

（1）「前記クエリフレームに対する類似度が第 1 基準を満たす類似フレーム画像が存在する」という条件が少なくとも 2 つの前記クエリフレームについて満たされる。

（2）前記少なくとも 2 つの類似フレーム画像のそれぞれに対応する前記クエリフレームを前記少なくとも 2 つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも 2 つの前記クエリフレームの並び順に一致する。

2. 上記 1 に記載の画像選択装置において、

前記クエリ動画は人物を含んでおり、

10

20

30

40

50

前記動画選択手段は、前記類似度として、前記人物の姿勢の類似度を用いる画像選択装置。

3．上記2に記載の画像選択装置において、

前記クエリフレームが満たすべき条件の一つは、当該クエリフレームに含まれる前記人物に関する情報量が第2基準を満たしていることである、画像選択装置。

4．上記3に記載の画像選択装置において、

前記複数の第1フレーム画像のそれぞれは、前記人物の姿勢を示す情報として、当該人物の関節の位置を示す関節情報を含んでおり、

前記第2基準は、前記関節情報に含まれる関節の数が基準を満たしていることである画像選択装置。

10

5．上記1～4のいずれか一項に記載の画像選択装置において、

前記クエリフレーム選択手段は、ユーザからの入力に従って少なくとも一つの前記クエリフレームを選択する画像選択装置。

6．上記1～4のいずれか一項に記載の画像選択装置において、

前記クエリフレーム選択手段は、第1の前記クエリフレームの次の前記クエリフレームとして前記第1のクエリフレームからの変化量が第3基準以上となった前記第1フレーム画像を選択する処理を、繰り返す画像選択装置。

7．上記6に記載の画像選択装置において、

前記クエリフレーム選択手段は、ユーザからの入力を用いて前記第3基準を設定する画像選択装置。

20

8．上記6に記載の画像選択装置において、

前記クエリフレーム選択手段は、前記複数の第1フレーム画像を用いて当該クエリフレームに用いる前記第3基準を設定する画像選択装置。

9．上記1～4のいずれか一項に記載の画像選択装置において、

前記クエリフレーム選択手段は、所定の間隔で前記複数の第1フレーム画像から複数のクエリフレームを選択する画像選択装置。

10．上記6～9のいずれか一項に記載の画像選択装置において、

前記クエリフレーム選択手段は、最初の前記クエリフレームとして、n番目の前記第1フレーム画像を選択する画像選択装置。

ここで、nは予め又はユーザ入力により設定された整数である。

30

11．コンピュータが、

複数の第1フレーム画像を含むクエリ動画を取得する取得処理と、

前記複数の第1フレーム画像から複数のクエリフレームを選択するクエリフレーム選択処理と、

前記複数のクエリフレームを用いて動画を選択する動画選択処理と、
を行い、

前記動画選択処理において、前記動画を選択する条件として、少なくとも以下の(1)及び(2)を用いる、画像選択方法。

(1)「前記クエリフレームに対する類似度が第1基準を満たす類似フレーム画像が存在する」という条件が少なくとも2つの前記クエリフレームについて満たされる。

40

(2)前記少なくとも2つの類似フレーム画像のそれぞれに対応する前記クエリフレームを前記少なくとも2つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも2つの前記クエリフレームの並び順に一致する。

12．上記11に記載の画像選択方法において、

前記クエリ動画は人物を含んでおり、

前記動画選択処理において、前記コンピュータは、前記類似度として、前記人物の姿勢の類似度を用いる画像選択方法。

13．上記12に記載の画像選択方法において、

前記クエリフレームが満たすべき条件の一つは、当該クエリフレームに含まれる前記人物に関する情報量が第2基準を満たしていることである、画像選択方法。

50

14. 上記13に記載の画像選択方法において、

前記複数の第1フレーム画像のそれぞれは、前記人物の姿勢を示す情報として、当該人物の関節の位置を示す関節情報を含んでおり、

前記第2基準は、前記関節情報に含まれる関節の数が基準を満たしていることである画像選択方法。

15. 上記11～14のいずれか一項に記載の画像選択方法において、

前記クエリフレーム選択処理において、前記コンピュータは、ユーザからの入力に従って少なくとも一つの前記クエリフレームを選択する画像選択方法。

16. 上記11～14のいずれか一項に記載の画像選択方法において、

前記クエリフレーム選択処理において、前記コンピュータは、第1の前記クエリフレームの次の前記クエリフレームとして前記第1のクエリフレームからの変化量が第3基準以上となった前記第1フレーム画像を選択する処理を、繰り返す画像選択方法。

10

17. 上記16に記載の画像選択方法において、

前記クエリフレーム選択処理において、前記コンピュータは、ユーザからの入力を用いて前記第3基準を設定する画像選択方法。

18. 上記16に記載の画像選択方法において、

前記クエリフレーム選択処理において、前記コンピュータは、前記複数の第1フレーム画像を用いて当該クエリフレームに用いる前記第3基準を設定する画像選択方法。

19. 上記11～14のいずれか一項に記載の画像選択方法において、

前記クエリフレーム選択処理において、前記コンピュータは、所定の間隔で前記複数の第1フレーム画像から複数のクエリフレームを選択する画像選択方法。

20

20. 上記16～19のいずれか一項に記載の画像選択方法において、

前記クエリフレーム選択処理において、前記コンピュータは、最初の前記クエリフレームとして、 n 番目の前記第1フレーム画像を選択する画像選択方法。

ここで、 n は予め又はユーザ入力により設定された整数である。

21. コンピュータに、

複数の第1フレーム画像を含むクエリ動画を取得する取得機能と、

前記複数の第1フレーム画像から複数のクエリフレームを選択するクエリフレーム選択機能と、

30

前記複数のクエリフレームを用いて動画を選択する動画選択機能と、
を持たせ、

前記動画選択機能は、前記動画を選択する条件として、少なくとも以下の(1)及び(2)を用いるプログラム。

(1)「前記クエリフレームに対する類似度が第1基準を満たす類似フレーム画像が存在する」という条件が少なくとも2つの前記クエリフレームについて満たされる。

(2)前記少なくとも2つの類似フレーム画像のそれぞれに対応する前記クエリフレームを前記少なくとも2つの類似フレーム画像と同じ順序で並べたときに、当該並び順が、前記クエリ動画における前記少なくとも2つの前記クエリフレームの並び順に一致する。

22. 上記21に記載のプログラムにおいて、

40

前記クエリ動画は人物を含んでおり、

前記動画選択機能は、前記類似度として、前記人物の姿勢の類似度を用いるプログラム。

23. 上記22に記載のプログラムにおいて、

前記クエリフレームが満たすべき条件の一つは、当該クエリフレームに含まれる前記人物に関する情報量が第2基準を満たしていることである、プログラム。

24. 上記23に記載のプログラムにおいて、

前記複数の第1フレーム画像のそれぞれは、前記人物の姿勢を示す情報として、当該人物の関節の位置を示す関節情報を含んでおり、

前記第2基準は、前記関節情報に含まれる関節の数が基準を満たしていることであるプログラム。

50

25. 上記21～24のいずれか一項に記載のプログラムにおいて、
前記クエリフレーム選択機能は、ユーザからの入力に従って少なくとも一つの前記クエリフレームを選択するプログラム。

26. 上記21～24のいずれか一項に記載のプログラムにおいて、
前記クエリフレーム選択機能は、第1の前記クエリフレームの次の前記クエリフレームとして前記第1のクエリフレームからの変化量が第3基準以上となった前記第1フレーム画像を選択する処理を、繰り返すプログラム。

27. 上記26に記載のプログラムにおいて、
前記クエリフレーム選択機能は、ユーザからの入力を用いて前記第3基準を設定するプログラム。

10

28. 上記26に記載のプログラムにおいて、
前記クエリフレーム選択機能は、前記複数の第1フレーム画像を用いて当該クエリフレームに用いる前記第3基準を設定するプログラム。

29. 上記21～24のいずれか一項に記載のプログラムにおいて、
前記クエリフレーム選択機能は、所定の間隔で前記複数の第1フレーム画像から複数のクエリフレームを選択するプログラム。

30. 上記26～29のいずれか一項に記載のプログラムにおいて、
前記クエリフレーム選択機能は、最初の前記クエリフレームとして、n番目の前記第1フレーム画像を選択するプログラム。

ここで、nは予め又はユーザ入力により設定された整数である。

20

【符号の説明】

【0136】

- 1 画像処理システム
- 10 画像処理装置（画像選択装置）
- 11 骨格検出部
- 12 特徴量算出部
- 13 認識部
- 100 画像処理装置（画像選択装置）
- 101 画像取得部
- 102 骨格構造検出部
- 103 特徴量算出部
- 104 分類部
- 105 検索部
- 106 入力部
- 107 表示部
- 108 身長算出部
- 110 データベース
- 200 カメラ
- 300、301 人体モデル
- 401 2次元骨格構造
- 610 クエリ取得部
- 620 クエリフレーム選択部
- 630 動画選択部

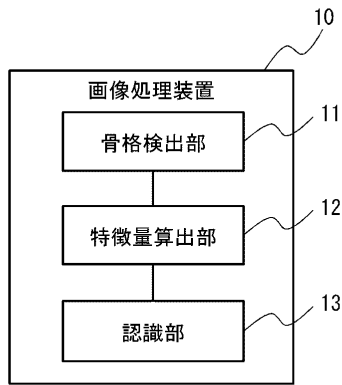
30

40

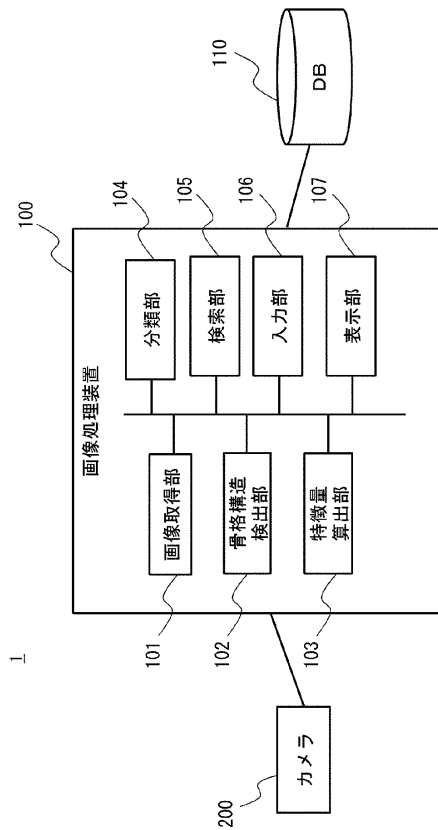
50

【図面】

【図 1】



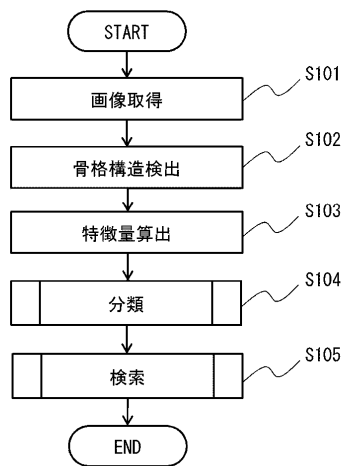
【図 2】



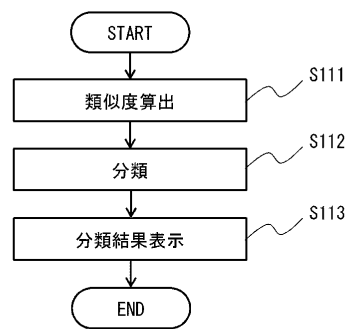
10

20

【図 3】



【図 4】

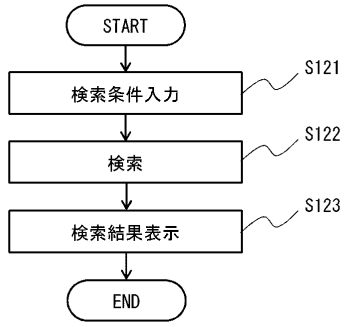


30

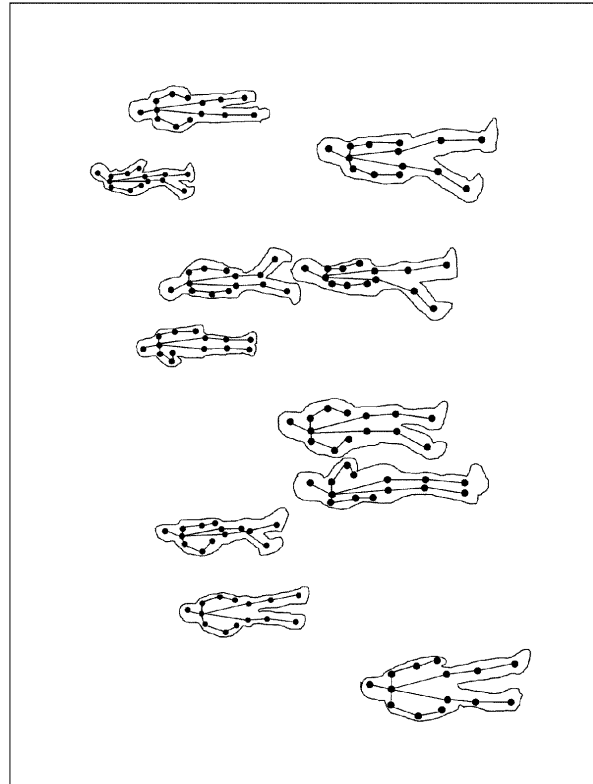
40

50

【 図 5 】



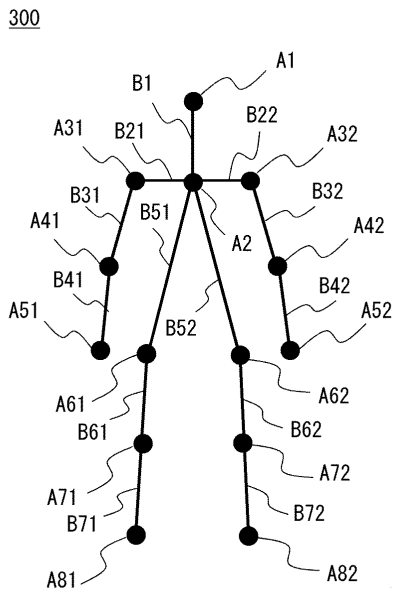
【 図 6 】



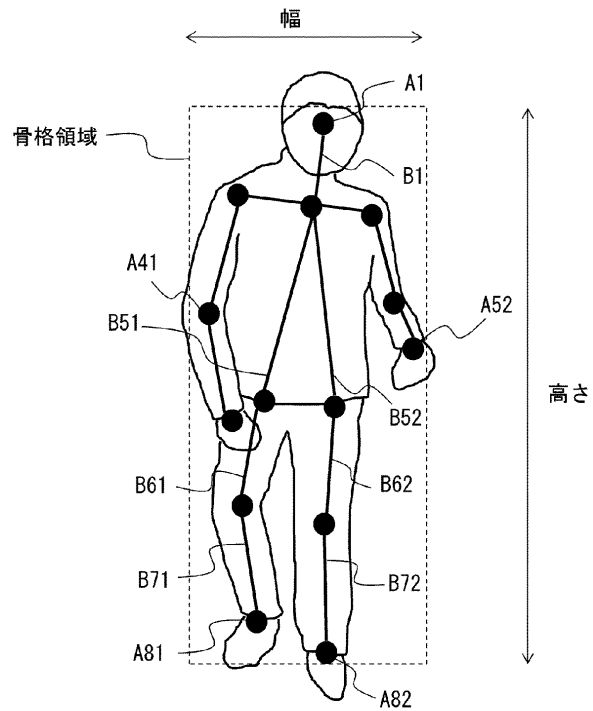
10

20

【 図 7 】



【 図 8 】

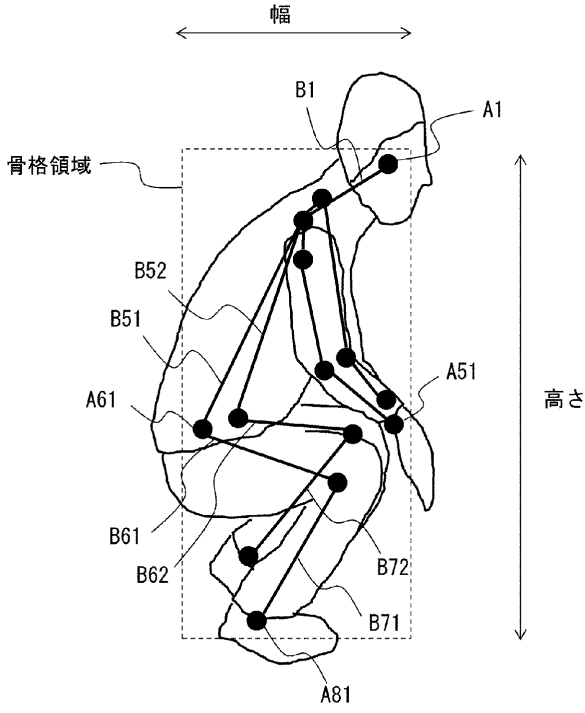


30

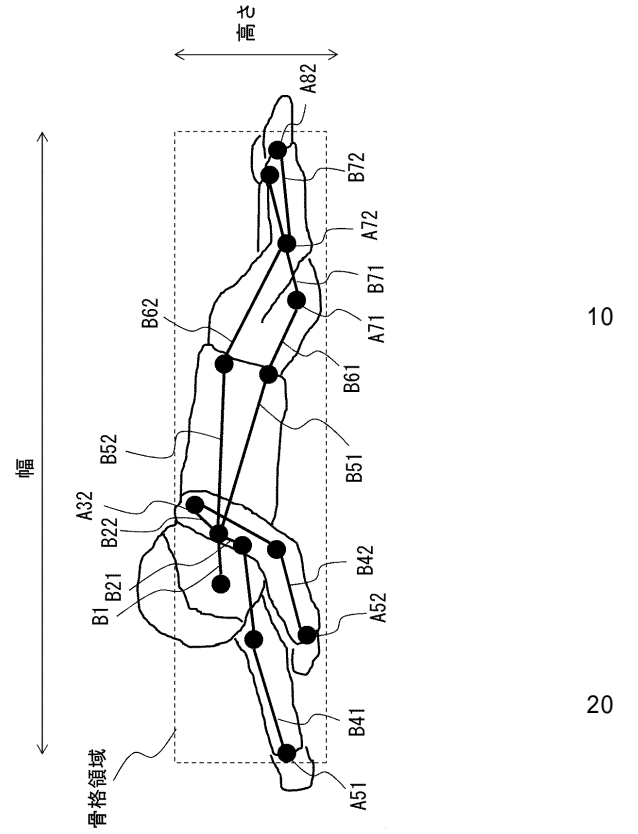
40

50

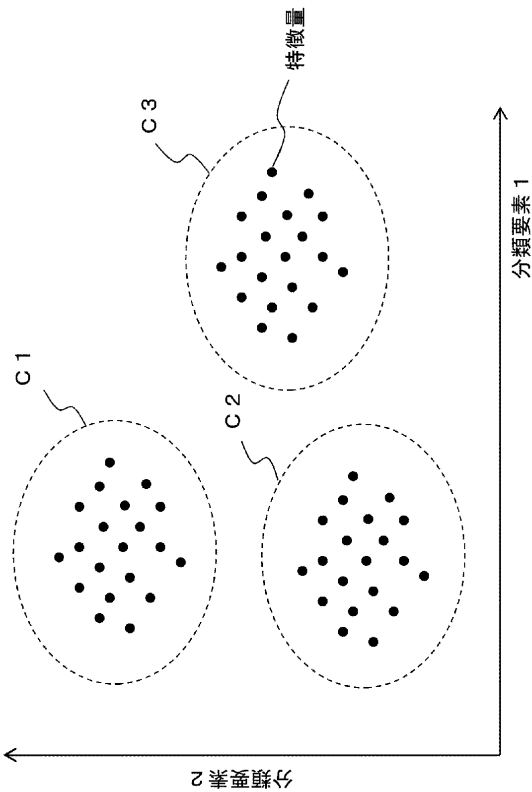
【 図 9 】



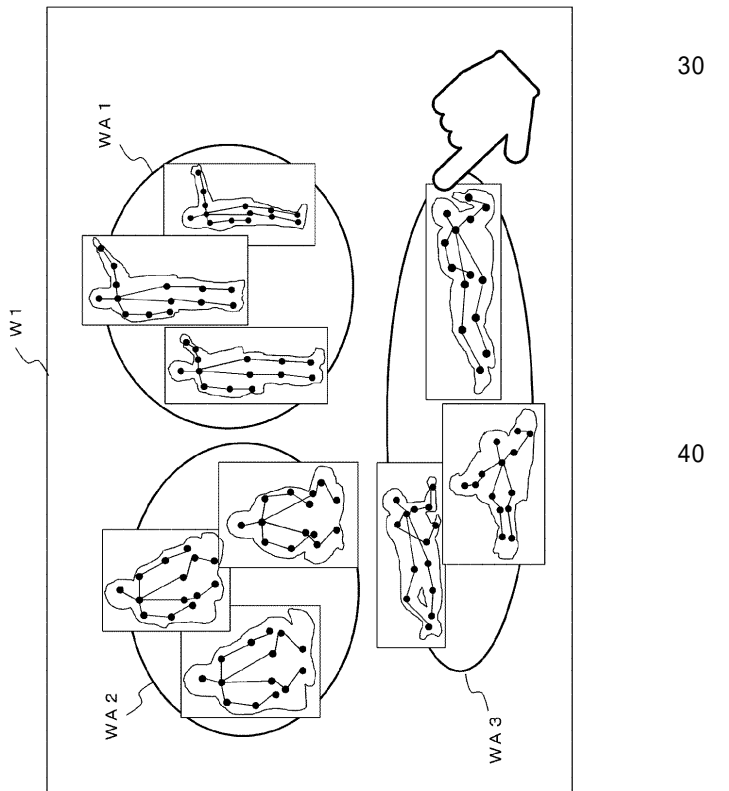
【 図 10 】



【 図 11 】



【 図 12 】



10

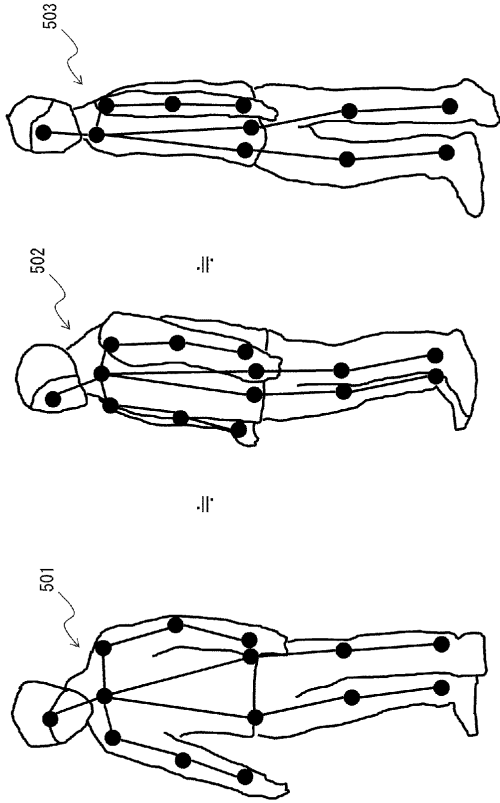
20

30

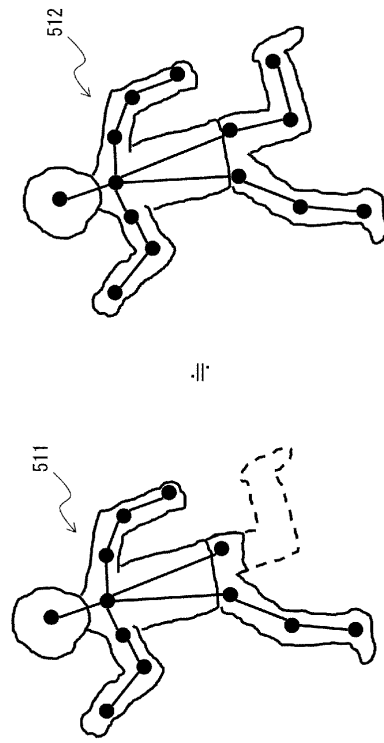
40

50

【図 13】



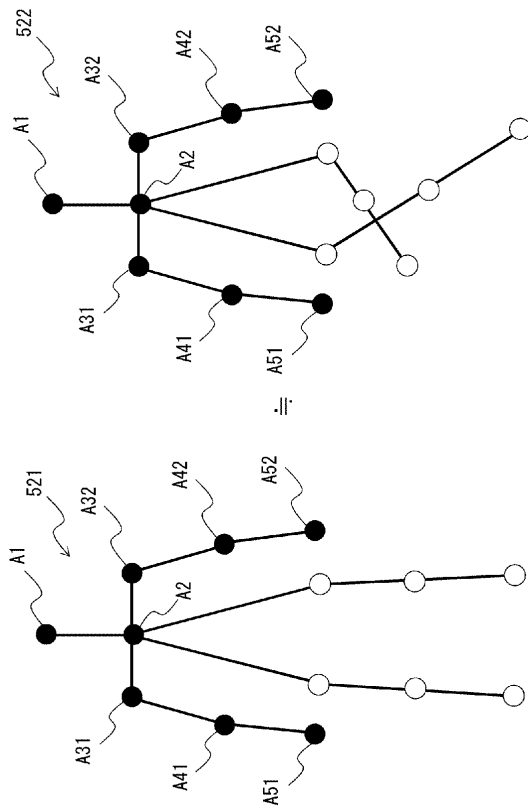
【図 14】



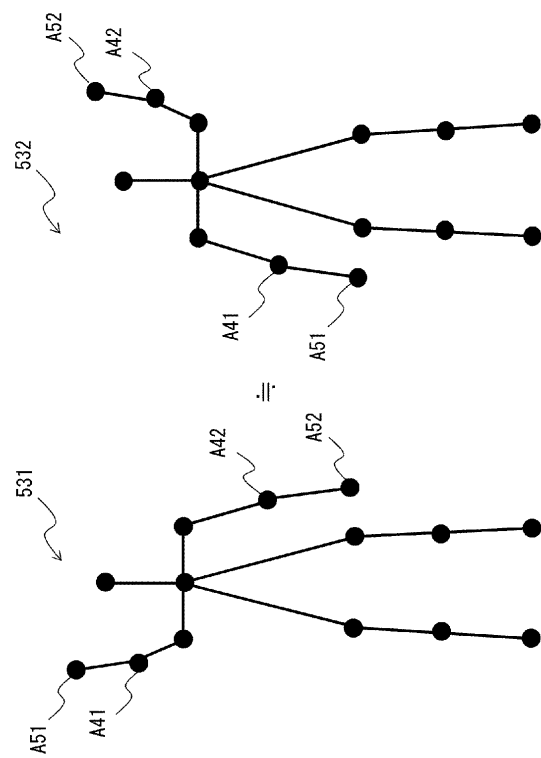
10

20

【図 15】



【図 16】

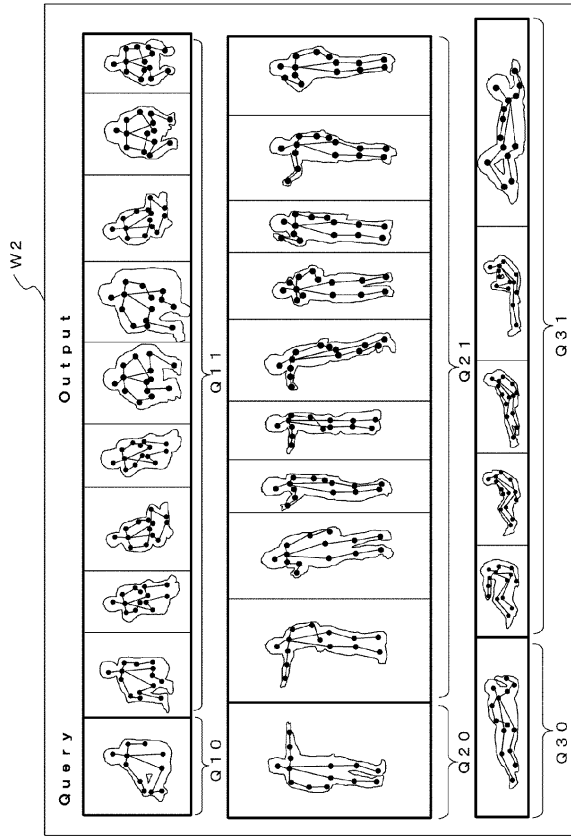


30

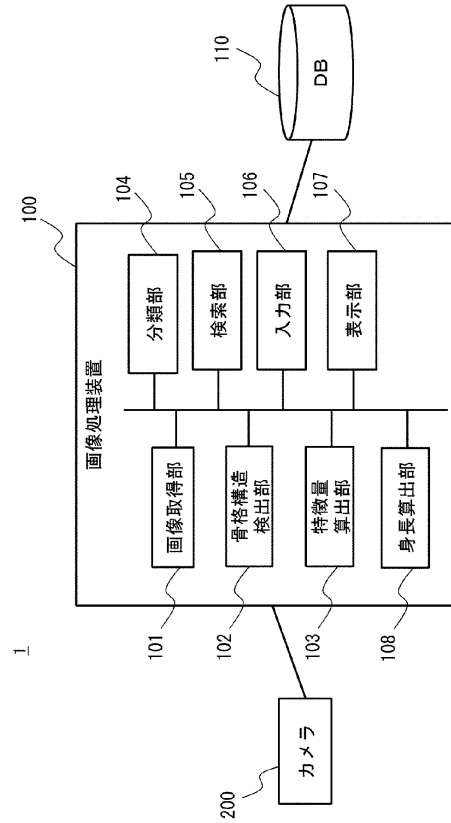
40

50

【図 17】



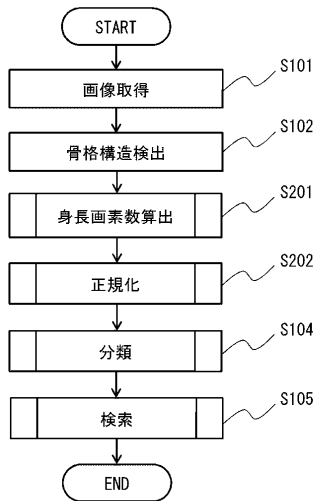
【図 18】



10

20

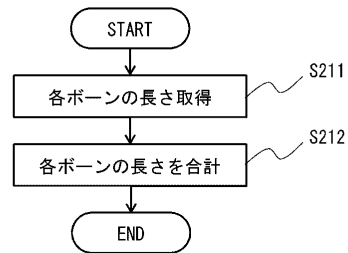
【図 19】



30

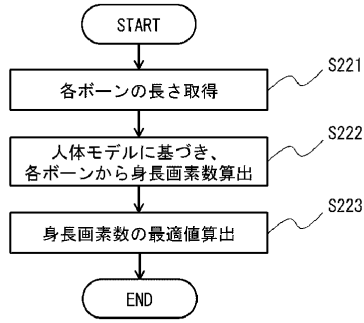
40

【図 20】

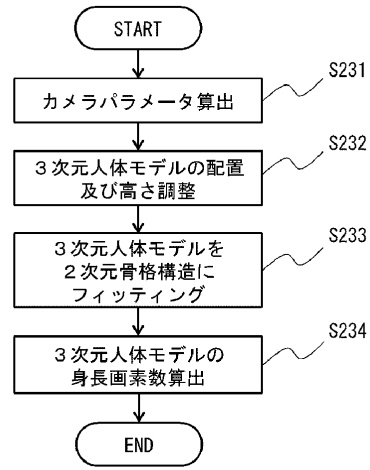


50

【 図 2 1 】

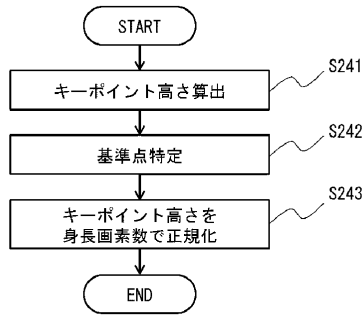


【 図 2 2 】

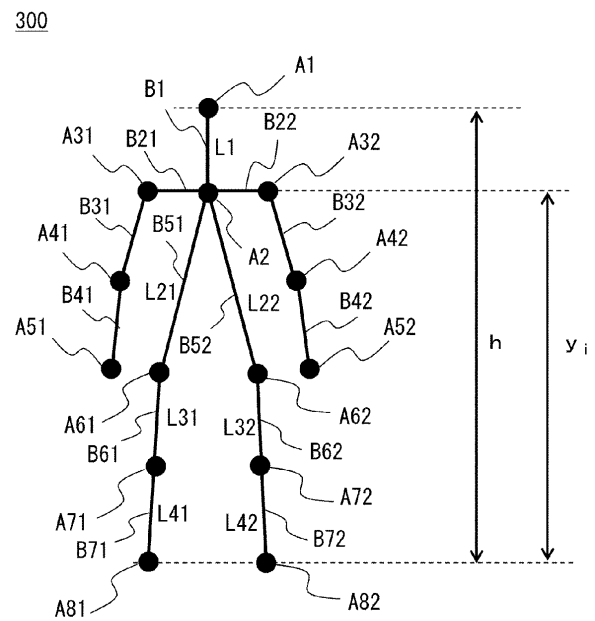


10

【 図 2 3 】



【 図 2 4 】



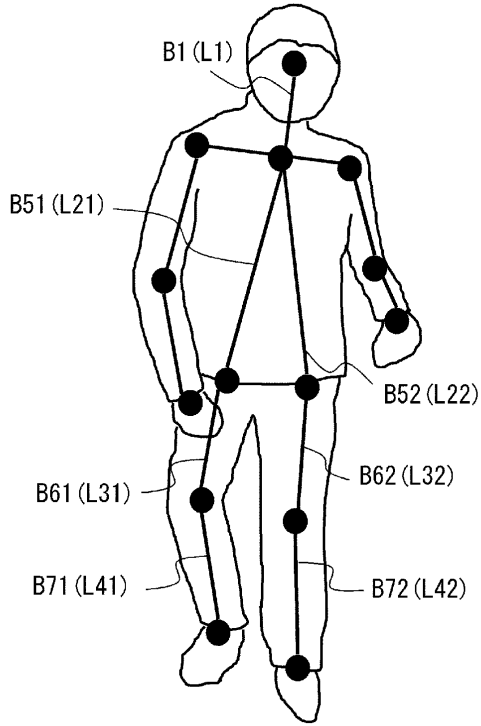
20

30

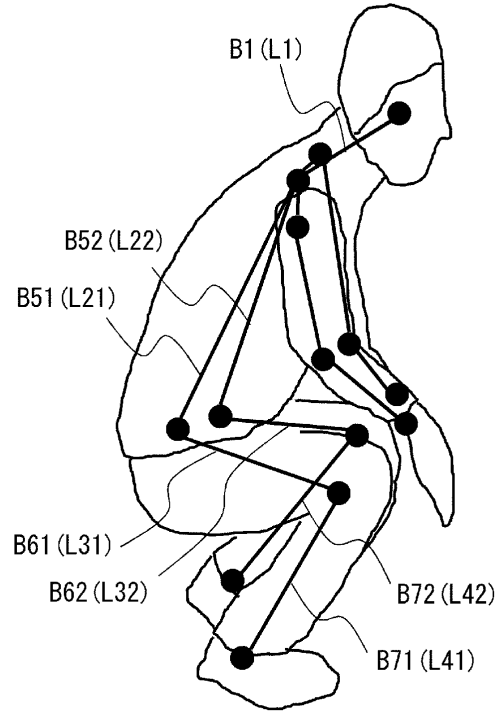
40

50

【 図 2 5 】



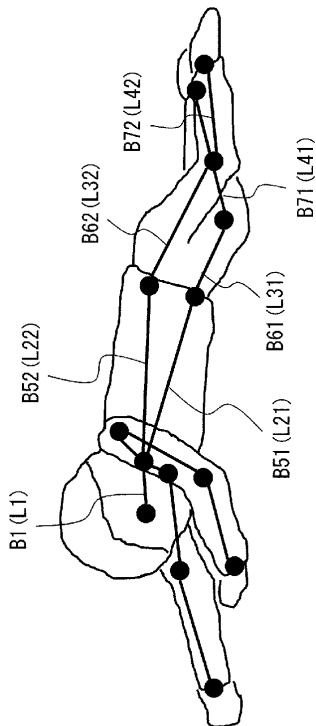
【 図 2 6 】



10

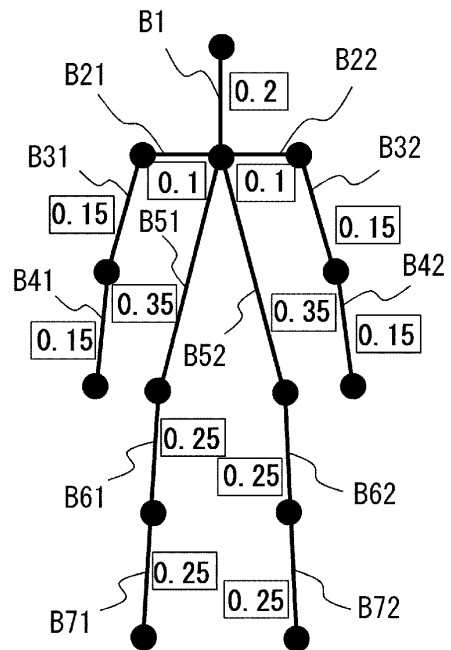
20

【 図 2 7 】



【 図 2 8 】

301

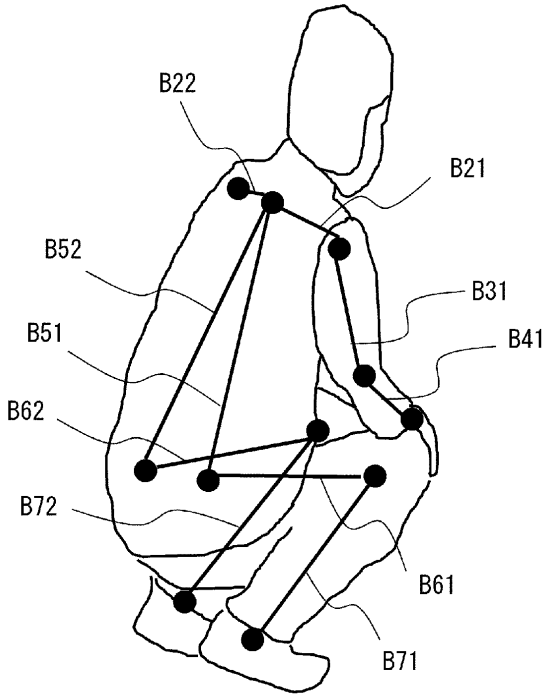


30

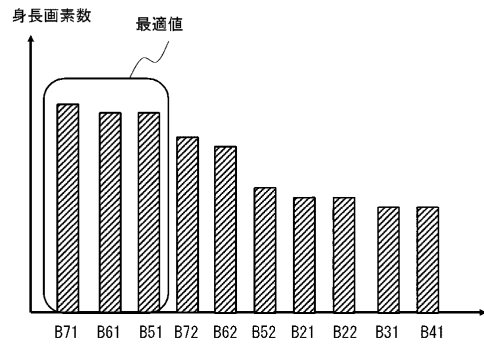
40

50

【 図 2 9 】



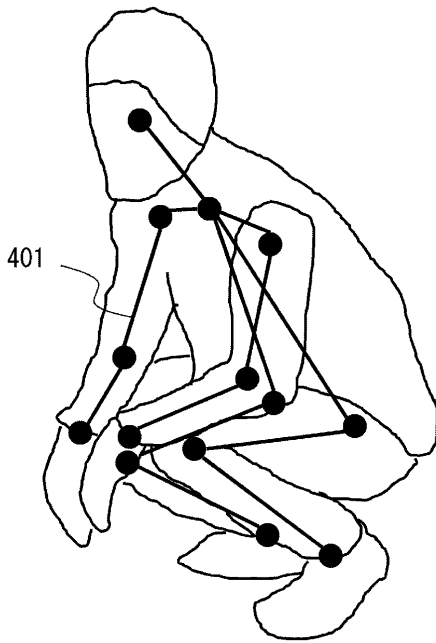
【 図 3 0 】



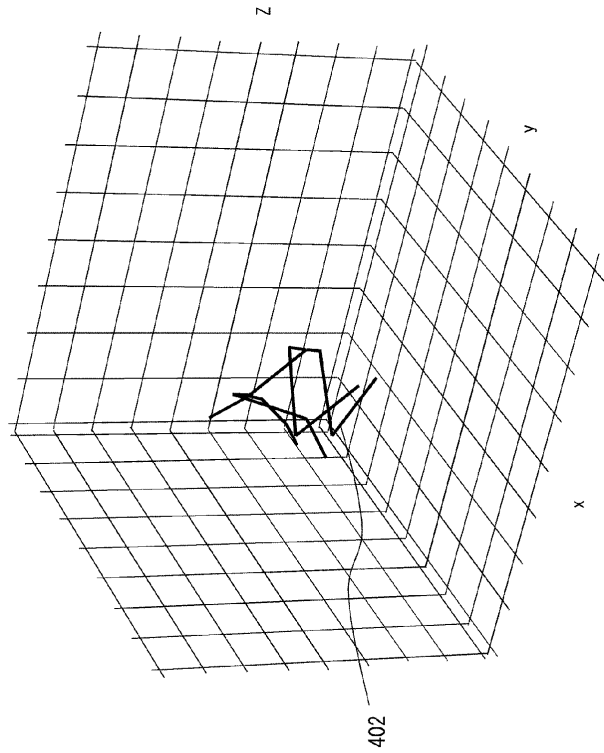
10

20

【 図 3 1 】



【 図 3 2 】

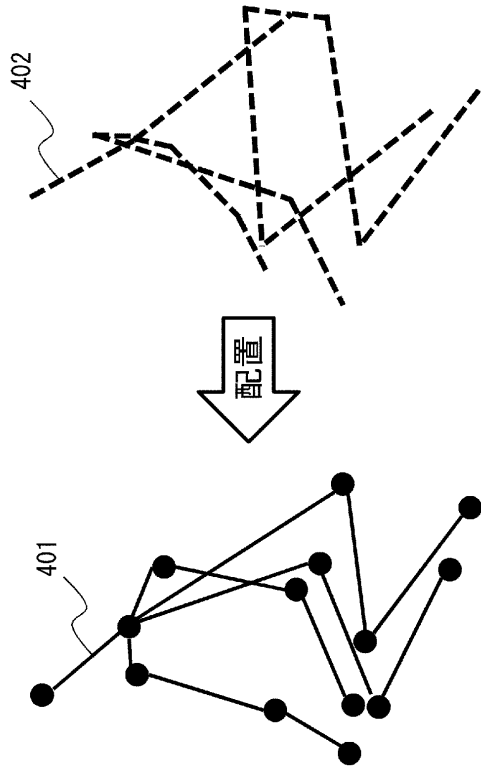


30

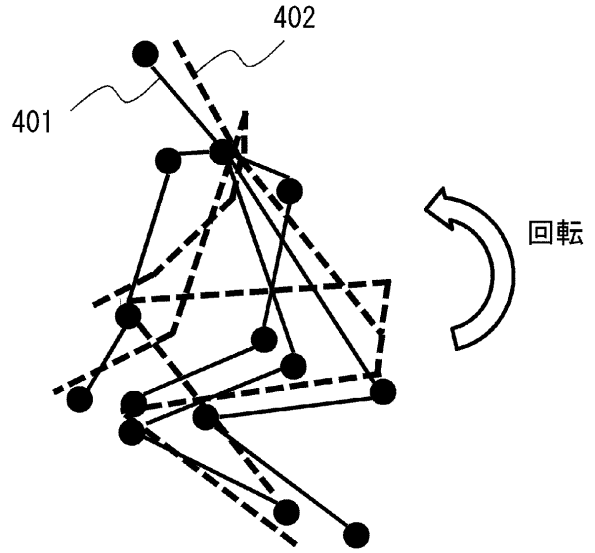
40

50

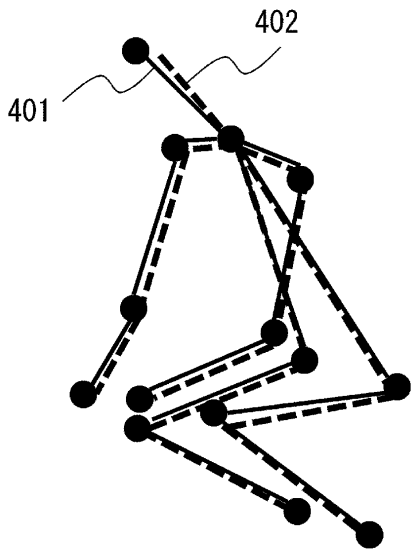
【図 3 3】



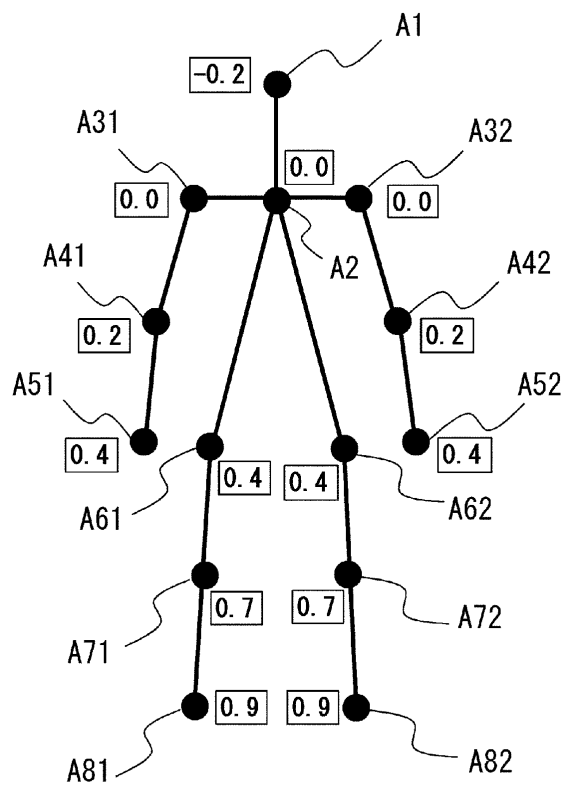
【図 3 4】



【図 3 5】



【図 3 6】



10

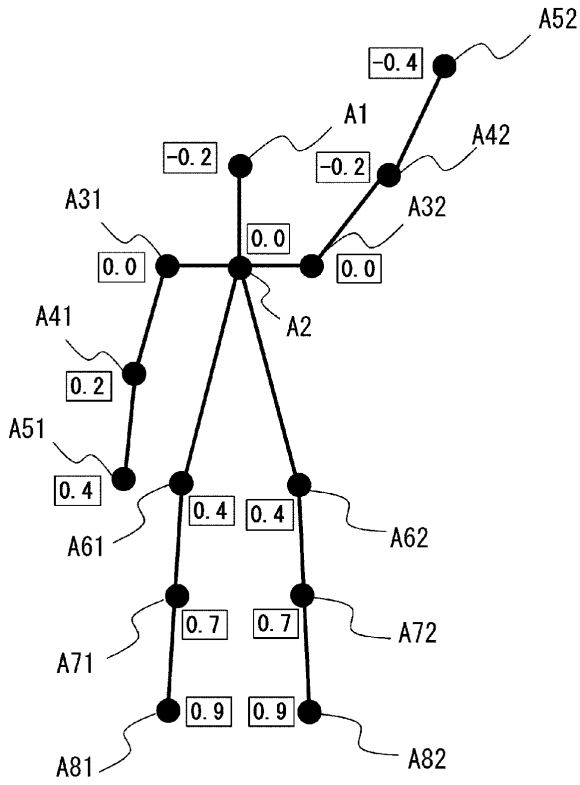
20

30

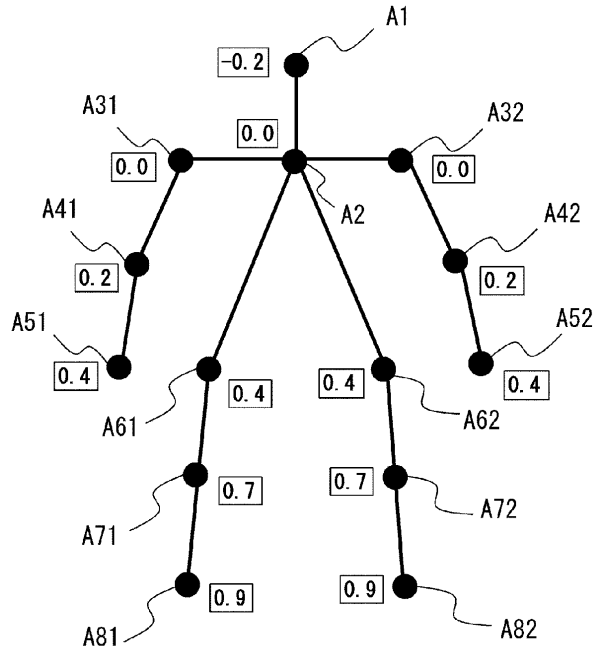
40

50

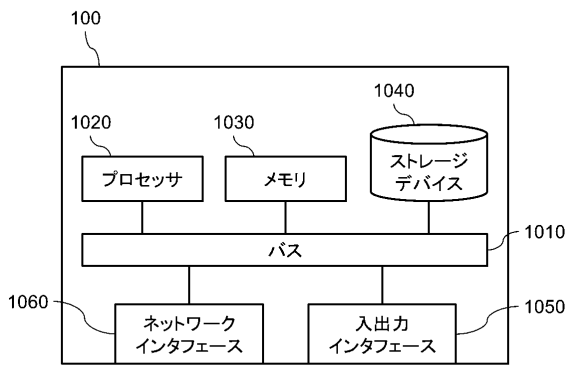
【図 3 7】



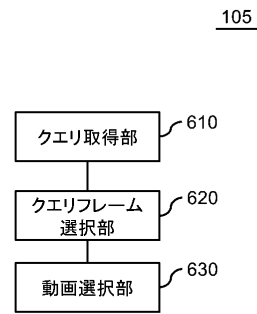
【図 3 8】



【図 3 9】



【図 4 0】



10

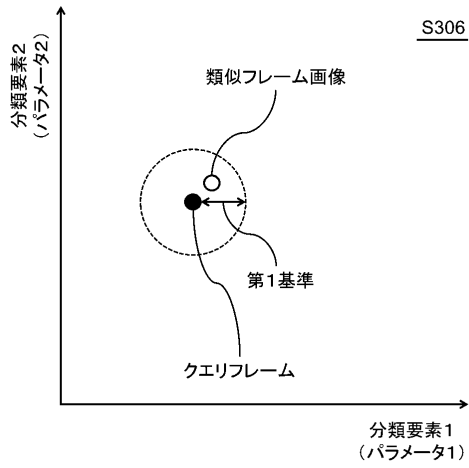
20

30

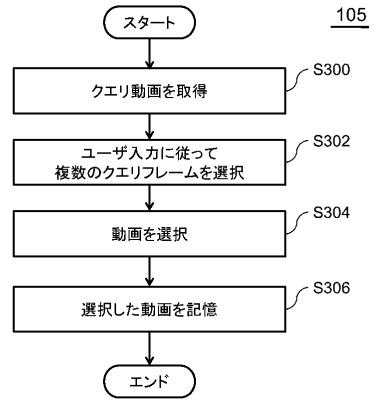
40

50

【図 4 1】

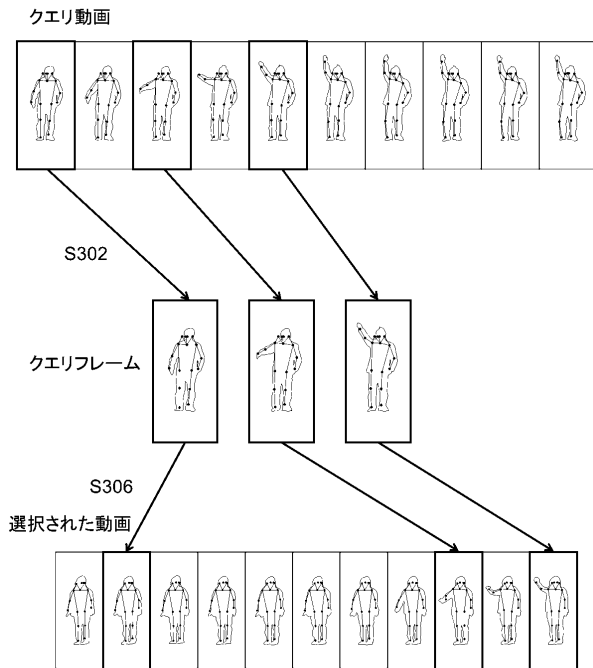


【図 4 2】

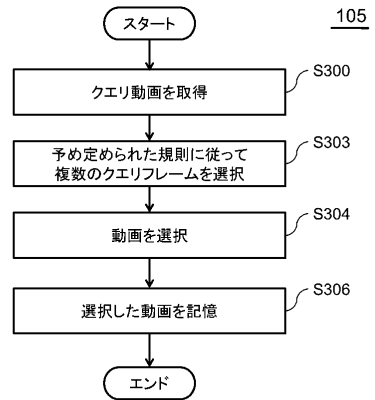


10

【図 4 3】



【図 4 4】



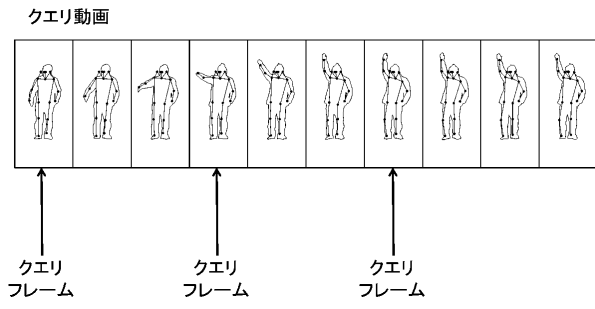
20

30

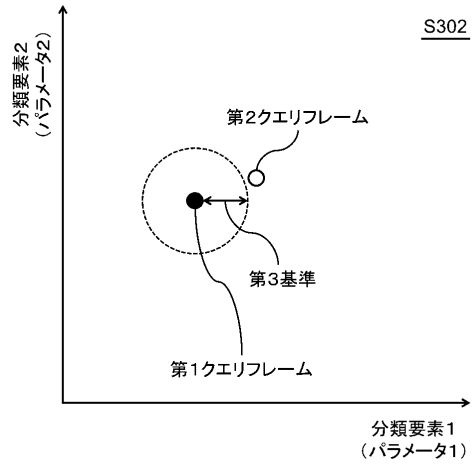
40

50

【 図 4 5 】

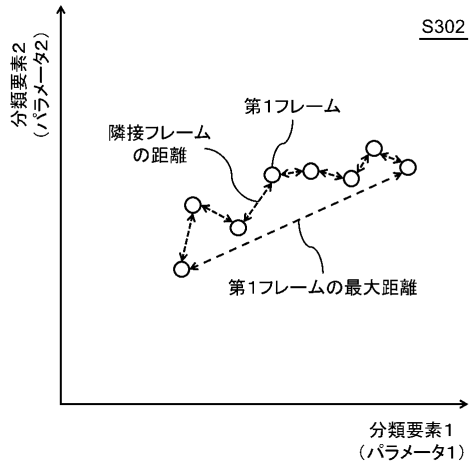


【 図 4 6 】



10

【 図 4 7 】



20

30

40

50

フロントページの続き

東京都港区芝五丁目7番1号 日本電気株式会社内

(72)発明者 西村 祥治

東京都港区芝五丁目7番1号 日本電気株式会社内

審査官 松尾 真人

- (56)参考文献 特開2017-117302(JP,A)
特開2019-091138(JP,A)
特開2011-118790(JP,A)
中国特許出願公開第110728209(CN,A)
米国特許出願公開第2016/0110453(US,A1)
特開2001-134589(JP,A)
国際公開第2006/025272(WO,A1)

(58)調査した分野 (Int.Cl., DB名)

G06F 16/00 - 16/958
G06T 1/00 - 1/40
G06T 3/00 - 5/50
G06T 11/60 - 13/80
G06T 19/00 - 19/20
G06Q 10/00 - 99/00