

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7228569号  
(P7228569)

(45)発行日 令和5年2月24日(2023.2.24)

(24)登録日 令和5年2月15日(2023.2.15)

(51)国際特許分類

F I

G 0 6 T 7/20 (2017.01) G 0 6 T 7/20 3 0 0 Z

G 0 6 T 7/00 (2017.01) G 0 6 T 7/00 3 5 0 C

請求項の数 13 (全69頁)

(21)出願番号	特願2020-507669(P2020-507669)	(73)特許権者	518394684
(86)(22)出願日	平成30年7月26日(2018.7.26)		スタンダード コグニション コーポレー
(65)公表番号	特表2020-530170(P2020-530170		ション
	A)		アメリカ合衆国 カリフォルニア州 9 4
(43)公表日	令和2年10月15日(2020.10.15)		1 0 3 サンフランシスコ, 7 番フロア
(86)国際出願番号	PCT/US2018/043933		ミッション ストリート 9 6 5
(87)国際公開番号	WO2019/032304	(74)代理人	100114476
(87)国際公開日	平成31年2月14日(2019.2.14)		弁理士 政木 良文
審査請求日	令和3年7月9日(2021.7.9)	(72)発明者	フィッシャー, ジョーダン
(31)優先権主張番号	62/542,077		アメリカ合衆国 カリフォルニア州 9 4
(32)優先日	平成29年8月7日(2017.8.7)		1 0 7 サンフランシスコ, アパートメ
(33)優先権主張国・地域又は機関			ント 7 1 6 キング ストリート 2 5 0
	米国(US)	(72)発明者	ヴァルドマン, デイヴィッド
(31)優先権主張番号	15/847,796		アメリカ合衆国 カリフォルニア州 9 4
(32)優先日	平成29年12月19日(2017.12.19)		1 0 3 サンフランシスコ, ナンバー 1
	最終頁に続く		最終頁に続く

(54)【発明の名称】 画像認識を用いた被写体識別及び追跡

(57)【特許請求の範囲】

【請求項 1】

実空間のエリア内の多関節被写体を追跡するシステムであって、  
複数のカメラ内のカメラが前記実空間内の対応する視野のそれぞれの画像シーケンスを生成し、前記複数のカメラにおいて各カメラの前記視野が少なくとも1つの他のカメラの前記視野と重なる、前記複数のカメラと、  
前記複数のカメラと結合された処理システムと、を備えてなり、  
前記処理システムが、  
前記複数のカメラから前記画像シーケンスを受信し、画像を処理し、関節タイプ、特定の画像の時間、及び前記特定の画像内の要素の座標によって、前記特定の画像の要素を分類する関節データ構造の配列を、前記特定の画像毎に、生成する画像認識エンジン、  
重なる視野を有するカメラからの画像シーケンス内の画像に対応する前記関節データ構造の配列を受信し、様々な画像シーケンス内の画像に対応する前記関節データ構造の配列内の前記要素の前記座標を、実空間内の座標を有する候補関節に変換するように構成された追跡エンジン、及び、  
実空間内の座標を有する候補関節のセットを前記実空間内の多関節被写体として識別するロジック、を備え、  
前記画像認識エンジンが、画像を処理して、  
前記特定の画像の要素毎に信頼度配列をそれぞれ生成し、且つ、  
前記信頼度配列に基づいて、前記特定の画像の特定の要素の前記関節データ構造の関節

10

20

タイプを選択し、

前記特定の画像の特定の要素についての信頼度配列が、複数の関節タイプの各関節タイプについての対応する信頼値を含み、関節タイプ毎に、前記特定の要素が前記各関節タイプである信頼度を示すものであることを特徴とするシステム。

【請求項 2】

前記処理システムが、前記識別された候補関節のセットを処理して、前記画像シーケンス内の前記多関節被写体の手の画像を含む有界ボックスを指定するように構成された有界ボックス生成器を、更に備える請求項 1 に記載のシステム。

【請求項 3】

多関節被写体として識別された前記候補関節のセットを記憶するロジックを備え、  
候補関節のセットを識別する前記ロジックが、特定の時間に撮影された画像において識別された候補関節が先行する画像において多関節被写体として識別された前記候補関節のセットのうちの 1 つのメンバーに対応するかどうかを判定するロジックを含む請求項 1 または 2 に記載のシステム。

10

【請求項 4】

実空間のエリア内の多関節被写体を追跡する方法であって、  
各カメラの視野が少なくとも 1 つの他のカメラの視野と重なり合う複数のカメラを使用して、前記実空間内の対応する視野のそれぞれの画像シーケンスを生成すること、  
関節タイプ、特定の画像の時間、及び前記特定の画像内の要素の座標によって、前記特定の画像の要素を分類する関節データ構造の配列を、前記特定の画像毎に、生成するために、前記画像シーケンス内の画像を処理すること、  
様々な画像シーケンス内の画像に対応する前記関節データ構造の配列内の前記要素の前記座標を、前記実空間内の座標を有する候補関節に変換すること、及び、  
実空間内の座標を有する候補関節のセットを前記実空間内の多関節被写体として識別すること、を備え、

20

前記画像を処理することが、  
前記特定の画像の要素毎に信頼度配列をそれぞれ生成すること、及び、  
前記信頼度配列に基づいて、前記特定の画像の特定の要素の前記関節データ構造の関節タイプを選択することを含み、  
前記特定の画像の特定の要素についての信頼度配列が、複数の関節タイプの各関節タイプについての対応する信頼値を含み、関節タイプ毎に、前記特定の要素が前記各関節タイプである信頼度を示すものであることを特徴とする方法。

30

【請求項 5】

前記画像を処理することが、畳み込みニューラル・ネットワークを使用することを含む請求項 4 に記載の方法。

【請求項 6】

前記候補関節のセットを識別することが、前記識別された候補関節のセットを処理して、前記画像シーケンス内の前記多関節被写体の手の画像を含む有界ボックスを指定することを含む請求項 4 または 5 に記載の方法。

【請求項 7】

多関節被写体として識別された前記候補関節のセットを記憶することを含み、  
前記候補関節のセットを識別することが、特定の時間に撮影された画像において識別された候補関節が先行する画像において多関節被写体として識別された前記候補関節のセットのうちの 1 つのメンバーに対応するかどうかを判定することを含む請求項 4 ~ 6 のいずれか 1 項に記載の方法。

40

【請求項 8】

前記画像シーケンスが同期されている請求項 4 ~ 7 のいずれか 1 項に記載の方法。

【請求項 9】

前記複数のカメラが、前記実空間内のエリアのそれぞれの部分を包含する視野を有し、その上に配置されたカメラを備え、

50

多関節被写体として識別された候補関節のセットのメンバーの実空間内の座標は、前記多関節被写体の前記エリア内の位置を識別する請求項 4 ~ 8 のいずれか 1 項に記載の方法。

【請求項 10】

前記実空間のエリア内の複数の多関節被写体の位置を追跡することを含む請求項 4 ~ 9 のいずれか 1 項に記載の方法。

【請求項 11】

前記複数の多関節被写体内の多関節被写体が前記実空間のエリアを離れるときを特定することを含む請求項 10 に記載の方法。

【請求項 12】

特定の多関節被写体として識別された候補関節のセットのメンバーである複数の候補関節の前記実空間のエリア内の座標を追跡することを含む請求項 4 ~ 11 のいずれか 1 項に記載の方法。

10

【請求項 13】

非一時的なコンピュータ可読記憶媒体であって、

請求項 4 ~ 12 のいずれか 1 項に係る実空間のエリア内の多関節被写体を追跡する方法のためのコンピュータ命令が格納されていることを特徴とする非一時的なコンピュータ可読記憶媒体。

【発明の詳細な説明】

【著作権通知】

【0001】

20

本特許書類の開示の一部は、著作権保護の対象となる資料を含んでいる。著作権所有者は、特許文献や特許開示を誰でも特許庁の特許ファイルや記録に記載されている通りにファクシミリで複写することに異議はないが、それ以外はあらゆる著作権を保有するものとする。

【技術分野】

【0002】

本発明は、レジレス・チェックアウトに使用可能なシステム及びその構成要素に関する。

【背景技術】

【0003】

画像処理における困難な問題は、大きな空間上に配置された複数のカメラからの画像が被写体の行為を識別し追跡するために使用される場合に生じる。

30

【0004】

ショッピングストア内の人々のような実空間のエリア内の被写体の行為を追跡することは、多くの技術的課題を提示する。例えば、複数の顧客がショッピングストア内の棚と棚の間の通路及びオープンスペースを移動するショッピングストアに配備される当該画像処理システムを考える。顧客は棚から商品を取り、それらをそれぞれのショッピングカートまたはバスケットに置く。顧客は、商品を望まない場合には商品を棚に置くこともできる。

【0005】

顧客がこれらの行為を実行している間、顧客の異なる部分、及び、棚の異なる部分、または店舗の在庫を保持する他の陳列構成は、他の顧客、棚、及び製品陳列などの存在のために、異なるカメラからの画像において、塞がれることになる。また、いつでも店舗内に多くの顧客がいる可能性があり、個人及びその行為を経時的に識別し追跡することが困難になる。

40

【0006】

より効果的かつ自動的に、大きなスペース内の被写体の取る及び置くという行為を識別及び追跡し、レジレス・チェックアウトなどの機能を含む、被写体とその環境との複雑な交流をサポートする他の処理を実行することができるシステムを提供することが望ましい。

【発明の概要】

【0007】

システム及びシステムを操作する方法は、画像処理を使用して、実空間のエリア内の人

50

物などの被写体による変化、及び被写体とその環境とその他の複雑な交流を追跡するために提供される。画像処理による変化を追跡するこの機能は、処理されるべき画像データのタイプ、画像データの如何なる処理を実行すべきか、及び、如何にして画像データから高い信頼性で行為を決定するかに関連して、コンピュータ工学の複雑な問題を提示する。本明細書に記載のシステムは、実空間の頭上に配置されたカメラからの画像のみを使用してこれらの機能を実行することができ、その結果、所与の設定での展開のために、店舗の棚及びフロアスペースにセンサなどを改装する必要がない。

【 0 0 0 8 】

在庫陳列構造の上方に配置された複数のカメラを使用して、各カメラの視野が複数のカメラ内の少なくとも1つの他のカメラの視野と重なる、実空間内の対応する視野内に在庫陳列構造の画像のそれぞれのシーケンスを生成することを備えた在庫陳列構造を含む、実空間のエリア内の被写体による在庫商品の置くこと及び取ることを追跡するシステム及び方法が提供される。これらの画像シーケンスを使用して、在庫陳列構造上の在庫商品に関連する画像シーケンスにおける意味的に重要な変化を識別し、意味的に重要な変化を画像シーケンス内に表される被写体に関連付けることによって、在庫商品を置くこと及び取ることを検出するシステム及び方法が説明される。

10

【 0 0 0 9 】

実空間のエリア内の被写体による在庫商品を置くこと及び取ることを追跡するためのシステム及び方法が提供され、このシステム及び方法は在庫陳列構造の上方に配置された複数のカメラを使用して、実空間内の対応する視野内に在庫陳列構造の画像のそれぞれのシーケンスを生成することを含み、各カメラの視野は、複数のカメラ内の少なくとも1つの他のカメラの視野と重複する。これらの画像シーケンスを使用して、画像シーケンス内の前景データを処理することにより、被写体のジェスチャ及び該ジェスチャに関連する在庫商品を識別することによって、在庫商品を置くこと及び取ることを検出するシステム及び方法が説明される。

20

【 0 0 1 0 】

また、前景処理と背景処理とを同じ画像シーケンスで組み合わせるシステム及び方法が説明される。この組み合わせられたアプローチでは、提供されるシステム及び方法が、画像シーケンス内の前景データを処理することにより、被写体のジェスチャ及び該ジェスチャに関連付けられた在庫商品を識別することによって、在庫商品を置くこと及び取ることを検出するためにこれらの画像のシーケンスを使用することと、画像シーケンス内の背景データを処理することによって、在庫陳列構造上の在庫商品に関連する画像シーケンス内の意味的に重要な変化を識別することによって、在庫商品を置くこと及び取ることを検出するためにこれらの画像シーケンスを使用することと、意味的に重要な変化を画像シーケンス内に表される被写体に関連付けることとを含む。

30

【 0 0 1 1 】

本明細書で説明される実施形態では、システムが複数のカメラを使用して、実空間内の対応する視野のそれぞれの画像シーケンスを生成する。各カメラの視野は、複数のカメラのうちの少なくとも1つの他のカメラの視野と重なる。このシステムは、複数のカメラから対応する画像シーケンスを受信する、被写体画像認識エンジンを含む第1の画像プロセッサを含む。第1の画像プロセッサは、画像を処理して、対応する画像シーケンス内の画像に表される被写体を識別する。システムは、複数のカメラから対応する画像シーケンスを受信する、背景画像認識エンジンを含む第2の画像プロセッサを更に含む。第2の画像プロセッサは、識別された被写体をマスクしてマスクされた画像を生成し、マスクされた画像を処理して、対応する画像シーケンス内の画像に表される背景変化を識別し且つ分類する。

40

【 0 0 1 2 】

一実施形態では、背景画像認識エンジンは、畳み込みニューラル・ネットワークを含む。システムは、識別された背景変化を識別された被写体に関連付けるロジックを含む。

【 0 0 1 3 】

50

一実施形態では、第2の画像プロセッサが、対応する画像シーケンスの背景画像を格納する背景画像格納装置を含む。第2の画像プロセッサは、識別された被写体を表す前景画像データを背景画像データで置き換えるために、画像シーケンス内の画像を処理するマスクロジックを更に含む。背景画像データは、マスクされた画像を提供するために、対応する画像シーケンスの背景画像から収集される。

【0014】

一実施形態では、マスクロジックが画像シーケンス内のN個のマスクされた画像のセットを組み合わせて、各カメラのファクタ化画像のシーケンスを生成する。第2の画像プロセッサは、ファクタ化画像のシーケンスを処理することによって、背景変化を識別し且つ分類する。

【0015】

一実施形態では、第2の画像プロセッサが、対応する画像シーケンスのための変化データ構造を生成するロジックを含む。変化データ構造は、識別された背景変化のマスクされた画像内の座標、識別された背景変化の在庫商品被写体の識別子、及び識別された背景変化の分類を含む。第2の画像プロセッサは更に、重なり合う視野を有するカメラのセットからの変化データ構造を処理して、実空間内での識別された背景変化の位置を見つけるための調整ロジックを含む。

【0016】

一実施形態では、変化データ構造における識別された背景変化の分類が識別された在庫商品が背景画像に対して追加されたか除去されたかを示す。

【0017】

別の実施形態では、変化データ構造における識別された背景変化の分類が識別された在庫商品が背景画像に対して追加されたか除去されたかを示す。システムは、背景変化を識別された被写体に関連付けるためのロジックを更に含む。最後に、システムは、識別された被写体による在庫商品を取ることに、及び識別された被写体による在庫陳列構造上に在庫商品を置くことの検出を行うロジックを含む。

【0018】

別の実施形態では、システムは、背景変化を識別された被写体に関連付けるロジックを含む。システムは、識別された被写体による在庫商品を取ることに、及び識別された被写体による在庫陳列構造上に在庫商品を置くことの検出を行うロジックを更に含む。

【0019】

システムは、複数のカメラから対応する画像シーケンスを受信する前景画像認識エンジンを含む、本明細書で説明する第3の画像プロセッサを含むことができる。第3の画像プロセッサは画像を処理して、対応する画像シーケンス内の画像に表される前景変化を識別し且つ分類する。

【0020】

システム及びシステムを操作する方法が、実空間において、人などの多関節被写体を追跡するために提供される。このシステムは、複数のカメラを使用して、実空間における対応する視野のそれぞれの画像シーケンスを生成する。各カメラの視野は、複数のカメラのうちの少なくとも1つの他のカメラの視野と重なる。このシステムは、画像シーケンス内の画像を処理して、各画像に対応する関節データ構造の配列を生成する。特定の画像に対応する関節データ構造の配列は、関節タイプ、特定の画像の時間、及び特定の画像内の要素の座標によって、特定の画像の要素を分類する。次に、システムは、様々な画像シーケンスに対応する関節データ構造の配列内の要素の座標を、実空間内の座標を有する候補関節に変換する。最後に、システムは候補関節のコンステレーションを識別し、コンステレーションは、実空間内の座標を有する候補関節のそれぞれのセットを、実空間における多関節被写体として含む。

【0021】

一実施形態では、画像認識エンジンが畳み込みニューラル・ネットワークを含む。画像認識エンジンによる画像の処理は、画像の要素に対する信頼度配列を生成することを含む

10

20

30

40

50

。画像の特定の要素についての信頼度配列は、特定の要素についての複数の関節タイプについての信頼値を含む。信頼度配列は、信頼度配列に基づいて、特定の要素の関節データ構造の関節タイプを選択するために使用される。

【 0 0 2 2 】

多関節被写体を追跡するためのシステムの一実施形態では、候補関節のセットを識別することは、候補関節のセットを多関節被写体として識別するために、実空間における被写体の関節間の物理的關係に基づいてヒューリスティック関数を適用することを含む。この処理は、多関節被写体として識別された関節のセットを記憶することを含む。候補関節のセットを識別することは、特定の時間に撮影された画像において識別された候補関節が先行する画像において多関節被写体として識別された候補関節のセットのうちの1つのメンバーに対応するかどうかを判定することを含む。

10

【 0 0 2 3 】

一実施形態では、複数のカメラによって取得された画像シーケンスのそれぞれの画像が、空間を通る被写体の移動の時間スケール上の単一の時点で実空間を表すように、画像シーケンスが同期される。

【 0 0 2 4 】

多関節被写体として識別された候補関節のセットのメンバーの実空間内の座標は、多関節被写体のエリア内の位置を識別する。いくつかの実施形態では、処理が、実空間のエリア内の複数の多関節被写体の位置の同時追跡を含む。いくつかの実施形態では、処理が、複数の多関節被写体内の多関節被写体が実空間のエリアを離れるときを特定することを含む。いくつかの実施形態では、処理が、多関節被写体が所与の時点で向いている方向を判定することを含む。本明細書で説明される実施形態では、システムが、複数のカメラを使用して、実空間内の対応する視野のそれぞれの画像シーケンスを生成する。各カメラの視野は、複数のカメラのうちの少なくとも1つの他のカメラの視野と重なる。システムは、複数のカメラから受け取った画像シーケンス内の画像を処理して、画像内に表された被写体を識別し、識別された被写体の分類を生成する。最後に、システムは、識別された被写体による在庫商品を取ることと、識別された被写体による棚に在庫商品を置くことを検出するために、画像シーケンス内の画像のセットに対する識別された被写体の分類を処理する。

20

【 0 0 2 5 】

一実施形態では、分類が、識別された被写体が在庫商品を保持しているかどうかを識別する。分類はまた、識別された被写体の手が棚の近くにあるかどうか、または識別された被写体の手が識別された被写体の近くにあるかどうかを識別する。手が識別された被写体の近くにあるかどうかの分類は、識別された被写体の手が識別された被写体に関連付けられたバスケットの近くにあり、識別された被写体の身体の近くにあるかどうかを含むことができる。

30

【 0 0 2 6 】

複数の画像内の被写体の手の分類を時系列で生成するために、視野内の被写体の手を表す画像を処理できる技術が記載されている。画像シーケンスからの手の分類は、いくつかの実施形態では、被写体による行為を識別するために、畳み込みニューラル・ネットワークを使用して処理することができる。行為は、本明細書に記載された実施形態に記載されているように、在庫商品を置くこと及び取ること、或いは、手の画像を処理することによって解読可能な他のタイプの行為であり得る。

40

【 0 0 2 7 】

画像を処理して視野内の被写体を識別し、被写体の関節の位置を見つける技術が記載されている。被写体の関節の位置は、被写体の手を含む対応する画像内の有界ボックスを識別するために、本明細書で説明するように処理することができる。有界ボックス内のデータは、対応する画像内の被写体の手の処理された分類とすることができる。画像シーケンスからこのようにして生成された識別された被写体からの手の分類は、被写体による行為を識別するために処理することができる。

50

## 【 0 0 2 8 】

前景と背景の画像認識エンジンのような複数の画像認識エンジンを含むシステムにおいて、該システムは、識別された被写体による在庫商品を取ることに、識別された被写体による在庫陳列構造上に在庫商品を置くことの第1の検出セット、及び、識別された被写体による在庫商品を取ることに識別された被写体による在庫陳列構造上に在庫商品を置くことの第2の検出セットを、作成できる。第1及び第2の検出セットを処理するための選択ロジックを使用して、ログ・データ構造を生成することができる。ログ・データ構造は、識別された被写体に関する在庫商品のリストを含む。

## 【 0 0 2 9 】

本明細書で説明する実施形態では、複数のカメラにおいて、カメラからの画像シーケンスが同期される。1つの好ましい実施態様では、同じカメラ及び同じ画像シーケンスが前景及び背景イメージプロセッサの両方によって使用される。その結果、同じ入力データを用いて、在庫商品を置くこと及び取ることの冗長な検出が行われ、結果として得られるデータにおいて高い信頼性と高い精度を可能にする。

10

## 【 0 0 3 0 】

本明細書で説明される1つの技術では、システムが、画像シーケンスで表されるジェスチャに関連付けられた被写体のジェスチャ及び在庫商品を識別することによって、在庫商品を置くこと及び取ることを検出するロジックを備える。これは、本明細書に記載されるように、被写体画像認識エンジンと協調して前景画像認識エンジンを使用して行うことができる。

20

## 【 0 0 3 1 】

本明細書で説明される別の技術では、システムは、棚のような在庫陳列構造上の在庫商品の意味的に重要な変化を経時的に識別し、意味的に重要な変化を画像シーケンスで表される被写体に関連付けることによって、在庫商品を置くこと及び取ることを検出するロジックを備える。これは、本明細書で説明するように、背景画像認識エンジンを被写体画像認識エンジンと協調させて使用して行うことができる。

## 【 0 0 3 2 】

本明細書で説明するテクノロジーを適用するシステムでは、ジェスチャ分析と意味的差異分析の両方を組み合わせて、カメラの配列からの同期画像の同じシーケンスに対して実行することができる。

30

## 【 0 0 3 3 】

コンピュータ・システムによって実行することができる方法及びコンピュータ・プログラム製品も、本明細書に記載されている。

## 【 0 0 3 4 】

本発明の他の実施態様及び利点は、以下の図面、詳細な説明、及び特許請求の範囲を検討することによって理解することができる。

## 【図面の簡単な説明】

## 【 0 0 3 5 】

【図1】追跡エンジンが画像認識エンジンによって生成された関節データを使用して被写体を追跡するシステムのアーキテクチャレベル概略図を示す。

40

## 【 0 0 3 6 】

【図2】カメラ配置を示すショッピングストアの通路の側面図である。

## 【 0 0 3 7 】

【図3】カメラ配置を示すショッピングストアにおける図2の通路の上面図である。

## 【 0 0 3 8 】

【図4】図1の画像認識エンジンをホストするように構成されたカメラ及びコンピュータ・ハードウェア構成である。

## 【 0 0 3 9 】

【図5】図1の画像認識エンジンにおける関節の識別を示す畳み込みニューラル・ネットワークを示す。

50

【 0 0 4 0 】

【 図 6 】 関節情報を記憶するための例示的なデータ構造を示す。

【 0 0 4 1 】

【 図 7 】 グローバル・メトリック計算器を有する図 1 の追跡エンジンを示す。

【 0 0 4 2 】

【 図 8 】 関連する関節の情報を含む被写体を記憶するための例示的なデータ構造を示す。

【 0 0 4 3 】

【 図 9 】 図 1 のシステムによって被写体を追跡するための処理ステップを示すフローチャートである。

【 0 0 4 4 】

【 図 1 0 】 図 9 のカメラ校正ステップのより詳細な処理ステップを示すフローチャートである。

【 0 0 4 5 】

【 図 1 1 】 図 9 のビデオ処理ステップのより詳細な処理ステップを示すフローチャートである。

【 0 0 4 6 】

【 図 1 2 A 】 図 9 のシーン処理のためのより詳細な処理ステップの第 1 の部分を示すフローチャートである。

【 0 0 4 7 】

【 図 1 2 B 】 図 9 のシーン処理のためのより詳細な処理ステップの第 2 の部分を示すフローチャートである。

【 0 0 4 8 】

【 図 1 3 】 図 1 のシステムの実施形態が使用される環境の図である。

【 0 0 4 9 】

【 図 1 4 】 図 1 のシステムの一実施形態におけるビデオ処理及びシーン処理の図である。

【 0 0 5 0 】

【 図 1 5 A 】 実空間において被写体毎にショッピングカート・データ構造を生成するための関節 CNN、What CNN、及び When CNN を含む複数の畳み込みニューラル・ネットワーク ( CNN ) を有するパイプラインを示す概略図である。

【 0 0 5 1 】

【 図 1 5 B 】 複数のカメラからの複数の画像チャンネルと、被写体及びそれらのそれぞれのショッピングカート・データ構造のための調整ロジックとを示す。

【 0 0 5 2 】

【 図 1 6 】 実空間内の被写体を識別して更新する処理ステップを示すフローチャートである。

【 0 0 5 3 】

【 図 1 7 】 在庫商品を識別するために被写体の手関節を処理するための処理ステップを示すフローチャートである。

【 0 0 5 4 】

【 図 1 8 】 被写体毎のショッピングカート・データ構造を作成するための、手関節毎の在庫商品の時系列分析のための処理ステップを示すフローチャートである。

【 0 0 5 5 】

【 図 1 9 】 図 1 5 A のシステムの実施形態における What CNN モデルの図である。

【 0 0 5 6 】

【 図 2 0 】 図 1 5 A のシステムの一実施形態における When CNN モデルの図である。

【 0 0 5 7 】

【 図 2 1 】 畳み込み層の次元を識別する What CNN モデルの例示的なアーキテクチャを示す。

【 0 0 5 8 】

【 図 2 2 】 手画像の分類のための What CNN モデルの実施形態の高レベルブロック図

10

20

30

40

50



を示す。

【 0 0 5 9 】

【図 2 3】図 2 2 に示される W h a t C N N モデルの高レベルブロック図の第 1 のブロックの詳細を示す。

【 0 0 6 0 】

【図 2 4】図 2 2 に提示された例示的 W h a t C N N モデルにおける全結合層における演算子を提示する。

【 0 0 6 1 】

【図 2 5】W h a t C N N モデルのためのトレーニング・データセットの一部として記憶される画像ファイルの例示的なファイル名である。

【 0 0 6 2 】

【図 2 6】背景意味的差分抽出を使用する第 1 の検出と、前景領域提案を使用する冗長検出との間で選択ロジックが選択する、実空間のエリア内の被写体による変化を追跡するためのシステムの高レベルアーキテクチャである。

【 0 0 6 3 】

【図 2 7】図 2 6 のシステムを実施するサブシステムの構成要素を示す。

【 0 0 6 4 】

【図 2 8 A】在庫イベントを決定し、ショッピングカート・データ構造を生成するための詳細な処理ステップの第 1 の部分を示すフローチャートである。

【 0 0 6 5 】

【図 2 8 B】在庫イベントを決定し、ショッピングカート・データ構造を生成するための詳細な処理ステップの第 2 の部分を示すフローチャートである。

【発明を実施するための形態】

【 0 0 6 6 】

以下の説明は、当業者が本発明を作成し使用することを可能にするために提示され、特定の用途及びその要件に即して提供される。開示された実施態様に対する様々な修正は、当業者には容易に明らかであり、本明細書で定義される一般原則は、本発明の精神及び範囲から逸脱することなく、他の実施態様及び用途に適用され得る。従って、本発明は、示された実施態様に限定されることを意図するものではなく、本明細書に開示された原理及び特徴と一致する最も広い範囲が与えられるべきである。

[ システム概要 ]

【 0 0 6 7 】

図 1 ~ 図 2 8 A / 2 8 B を参照して、対象技術のシステム及び様々な実施態様を説明する。システム及び処理は、本実施態様によるシステムのアーキテクチャレベル概略図である図 1 を参照して説明される。図 1 は、アーキテクチャ図であるため、説明の明確性を向上させるために、特定の詳細は省略されている。

【 0 0 6 8 】

図 1 の説明は、以下のように編成される。最初に、システムの要素を説明し、次にそれらの相互接続を説明する。次に、システムにおける要素の使用についてより詳細に説明する。

【 0 0 6 9 】

図 1 は、システム 1 0 0 のブロック図レベルの説明図を提供する。本システム 1 0 0 は、カメラ 1 1 4、ネットワーク・ノードがホスティングする画像認識エンジン 1 1 2 a、1 1 2 b 及び 1 1 2 n、ネットワーク上の 1 つまたは複数のネットワーク・ノードに配置される追跡エンジン 1 1 0、較正器 1 2 0、被写体データベース 1 4 0、トレーニング・データベース 1 5 0、関節ヒューリスティックス用、置く及び取るヒューリスティックス用、及び、後述する複数の画像認識エンジンの出力を調整し、結合するための他のヒューリスティックス用のヒューリスティックス・データベース 1 6 0、較正データベース 1 7 0、及び、1 つまたは複数の通信ネットワーク 1 8 1 を含む。ネットワーク・ノードは、1 つの画像認識エンジンのみ、または本明細書で説明されるように、複数の画像認識エンジ

10

20

30

40

50

ンをホストすることができる。システムはまた、在庫データベース及び他のサポートデータを含むことができる。

【 0 0 7 0 】

本明細書で使用されるように、ネットワーク・ノードは、ネットワークに接続され、通信チャネルを介して他のネットワーク・ノードとの間で情報を送信、受信、または転送することができる、アドレス可能なハードウェア・デバイスまたは仮想デバイスである。ハードウェア・ネットワーク・ノードとして配置することができる電子デバイスの例には、あらゆる種類のコンピュータ、ワークステーション、ラップトップ・コンピュータ、ハンドヘルド・コンピュータ、及びスマートフォンが含まれる。ネットワーク・ノードは、クラウドベースのサーバ・システムで実施することができる。ネットワーク・ノードとして構成された複数の仮想デバイスを、単一の物理デバイスを使用して実施することができる。

10

【 0 0 7 1 】

明確性のために、画像認識エンジンをホストする3つのネットワーク・ノードのみがシステム 1 0 0 に示されている。しかしながら、画像認識エンジンをホストする任意の数のネットワーク・ノードを、ネットワーク 1 8 1 を介して追跡エンジン 1 1 0 に接続することができる。また、本明細書で説明する画像認識エンジン、追跡エンジン、及び他の処理エンジンは、分散アーキテクチャ内の複数のネットワーク・ノードを使用して実行することができる。

【 0 0 7 2 】

次に、システム 1 0 0 の要素の相互接続について説明する。ネットワーク 1 8 1 は、画像認識エンジン 1 1 2 a、1 1 2 b、及び 1 1 2 n をそれぞれホストするネットワーク・ノード 1 0 1 a、1 0 1 b、及び 1 0 1 c、追跡エンジン 1 1 0 をホストするネットワーク・ノード 1 0 2、較正器 1 2 0、被写体データベース 1 4 0、トレーニング・データベース 1 5 0、関節ヒューリスティックス・データベース 1 6 0、及び較正データベース 1 7 0 を結合する。カメラ 1 1 4 は、画像認識エンジン 1 1 2 a、1 1 2 b、及び 1 1 2 n をホストするネットワーク・ノードを介して追跡エンジン 1 1 0 に接続される。一実施形態では、カメラ 1 1 4 がショッピングストア（スーパーマーケットなど）に設置され、重なり合う視野を有するカメラ 1 1 4 のセット（2 つ以上）が各通路の上に配置されて、店舗内の実空間の画像を取得する。図 1 では、2 つのカメラが通路 1 1 6 a の上に配置され、2 つのカメラが通路 1 1 6 b の上に配置され、3 つのカメラが通路 1 1 6 n の上に配置されている。カメラ 1 1 4 は、重なり合う視野を有する通路上に設置される。斯かる実施形態では、カメラは、ショッピングストアの通路内を移動する顧客がいつの時点でも 2 つ以上のカメラの視野内に存在することを目標として構成される。

20

30

【 0 0 7 3 】

カメラ 1 1 4 は互いに時間的に同期させることができ、その結果、画像は、同時にまたは時間的に近く、かつ同じ画像キャプチャレートで取得される。カメラ 1 1 4 は、画像認識エンジン 1 1 2 a ~ 1 1 2 n をホストするネットワーク・ノードに、所定のレートでそれぞれの継続的な画像ストリームを送ることができる。同時にまたは時間的に近くに、実空間のエリアをカバーする全てのカメラにおいて取得された画像は、同期された画像が実空間において固定された位置を有する被写体の異なる光景を表すものとして処理エンジンにおいて識別され得るという意味で、同期している。例えば、一実施形態では、カメラが、3 0 フレーム / 秒 ( f p s ) のレートで、画像認識エンジン 1 1 2 a ~ 1 1 2 n をホストするそれぞれのネットワーク・ノードに画像フレームを送信する。各フレームは、画像データと共に、タイムスタンプ、カメラの識別情報（「カメラ I D」と略される）、及びフレーム識別情報（「フレーム I D」と略される）を有する。

40

【 0 0 7 4 】

通路上に設置されたカメラは、それぞれの画像認識エンジンに接続される。例えば、図 1 において、通路 1 1 6 a 上に設置された 2 つのカメラは、画像認識エンジン 1 1 2 a をホストするネットワーク・ノード 1 0 1 a に接続される。同様に、通路 1 1 6 b 上に設置された 2 つのカメラは、画像認識エンジン 1 1 2 b をホストするネットワーク・ノード 1

50

01bに接続される。ネットワーク・ノード101a～101n内でホストされる各画像認識エンジン112a～112nは、図示の例ではそれぞれ1つのカメラから受信した画像フレームを別々に処理する。

【0075】

一実施形態では、各画像認識エンジン112a、112b、及び112nは、畳み込みニューラル・ネットワーク(CNNと略す)などの深層学習アルゴリズムとして実装される。斯かる実施形態では、CNNがトレーニング・データベース150を使用してトレーニングされる。本明細書で説明される実施形態では、実空間内の被写体の画像認識が、画像内で認識可能な関節を識別しグループ化することに基づいており、関節のグループは個々の被写体に帰属することができる。この関節ベースの分析のために、トレーニング・データベース150は、被写体のための異なるタイプの関節の各々に対して膨大な画像を収集している。ショッピングストアの例示的な実施形態では、被写体は、棚の間の通路を移動する顧客である。例示的な実施形態では、CNNのトレーニング中に、システム100は「トレーニング・システム」と呼ばれる。トレーニング・データベース150を使用してCNNをトレーニングした後、CNNは、プロダクション・モードに切り替えられ、ショッピングストア内の顧客の画像をリアルタイムで処理する。例示的な実施形態では、プロダクション中に、システム100はランタイム・システムと呼ばれる(推論システムとも呼ばれる)。それぞれの画像認識装置のCNNは、それぞれの画像ストリーム中の画像に対して関節データ構造の配列を生成する。本明細書に記載される実施形態では、関節データ構造の配列が、各処理された画像に対して生成されることで、各画像認識エンジン112a～112nが、関節データ構造の配列の出力ストリームを生成する。重なり合う視野を有するカメラからの関節データ構造のこれらの配列は、関節のグループを形成し、斯かる関節のグループを被写体として識別するために、更に処理される。

【0076】

カメラ114は、CNNをプロダクション・モードに切り替える前に較正される。キャリブレーション120はカメラを較正し、較正データを較正データベース170に格納する。

【0077】

追跡エンジン110は、ネットワーク・ノード102上でホストされ、画像認識エンジン112a～112nから被写体の関節データ構造の配列の継続的なストリームを受信する。追跡エンジン110は、関節データ構造の配列を処理し、様々なシーケンスの画像に対応する関節データ構造の配列内の要素の座標を、実空間内の座標を有する候補関節に変換する。同期画像の各セットについて、実空間全体にわたって識別された候補関節の組み合わせは、類推目的のために、候補関節の銀河に似ていると考えることができる。後続の各時点において、銀河が経時的に変化するように、候補関節の動きが記録される。追跡エンジン110の出力は、被写体データベース140に格納される。

【0078】

追跡エンジン110は、実空間内の座標を有する候補関節のグループまたはセットを、実空間内の被写体として識別するロジックを使用する。類推目的のために、候補点の各セットは、各時点における候補関節の星座(コンステレーション)に似ている。候補関節のコンステレーションは、時間とともに移動することができる。

【0079】

候補関節のセットを識別するロジックは、実空間における被写体の関節間の物理的關係に基づくヒューリスティック関数を含む。これらのヒューリスティック関数は、候補関節のセットを被写体として識別するために使用される。ヒューリスティック関数はヒューリスティックス・データベース160に格納される。追跡エンジン110の出力は、被写体データベース140に格納される。従って、候補関節のセットは、他の個々の候補関節とヒューリスティックス・パラメータに従った関係を有する個々の候補関節、及び、個々の被写体として識別された、または識別することができる所与のセット内の候補関節のサブセットを含む。

【0080】

10

20

30

40

50

ネットワーク 181 を通る実際の通信経路は、公衆ネットワーク及び/またはプライベート・ネットワーク上のポイント・ツー・ポイントとすることができる。通信は、プライベート・ネットワーク、VPN、MPLS 回路、またはインターネットなどの様々なネットワーク 181 を介して行うことができ、適切なアプリケーション・プログラミング・インターフェース (API) 及びデータ交換フォーマット、例えば、REST (Representational State Transfer)、JSON (JavaScript (商標) Object Notation)、XML (Extensible Markup Language)、SOAP (Simple Object Access Protocol)、JMS (Java (商標) Message Service)、及び/または Java プラットフォーム・モジュール・システムなどを使用することができる。すべての通信は、暗号化することができる。通信は、一般に、EDGE、3G、4G LTE、Wi-Fi、及び WiMAX などのプロトコルを介して、LAN (ローカル・エリア・ネットワーク)、WAN (ワイド・エリア・ネットワーク)、電話ネットワーク (公衆交換電話網 (PSTN))、セッション開始プロトコル (SIP)、無線ネットワーク、ポイント・ツー・ポイント・ネットワーク、星型ネットワーク、トークンリング型ネットワーク、ハブ型ネットワーク、インターネット (モバイルインターネットを含む) などのネットワーク上で行われる。更に、ユーザ名/パスワード、オープン許可 (OAuth)、Kerberos、SecureID、デジタル証明書などの様々な承認及び認証技術を使用して、通信を保護することができる。

10

#### 【0081】

本明細書に開示される技術は、データベースシステム、マルチテナント環境、または、Oracle (商標) と互換性のあるデータベース実施態様、IBM DB2 Enterprise Server (商標) と互換性のあるリレーショナル・データベース実施態様、MySQL (商標) または PostgreSQL (商標) と互換性のあるリレーショナル・データベース実施態様または Microsoft SQL Server (商標) と互換性のあるリレーショナル・データベース実施態様等のリレーショナル・データベース実施態様、または、Vampire (商標) と互換性のある非リレーショナル・データベース実施態様、Apache Cassandra (商標) と互換性のある非リレーショナル・データベース実施態様、BigTable (商標) と互換性のある非リレーショナル・データベース実施態様、または HBase (商標) または DynamoDB (商標) と互換性のある非リレーショナル・データベース実施態様、等の NoSQL (商標) の非リレーショナル・データベース実施態様を含む何かのコンピュータ実装システムという状況下で実施され得る。更に、開示された技術は、MapReduce (商標)、バルク同期プログラミング、MPI プリミティブ等の様々なプログラミングモデル、または、Apache Storm (商標)、Apache Spark (商標)、Apache Kafka (商標)、Apache Flink (商標)、Truviso (商標)、Amazon Elasticsearch Service (商標)、Amazon Web Services (AWS) (商標)、IBM Info Sphere (商標)、Borealis (商標)、及び Yahoo! S4 (商標) 等の様々なスケーラブルなバッチ及びストリーム管理システムを使用して実施され得る。

20

30

#### 【カメラ配置】

40

#### 【0082】

カメラ 114 は、3次元 (3D と略される) 実空間において多関節存在物 (または被写体) を追跡するように配置される。ショッピングストアの例示的な実施形態では、実空間は、販売用の商品が棚に積み重ねられるショッピングストアのエリアを含むことができる。実空間内の点は、(x, y, z) 座標系で表すことができる。システムが適用される実空間のエリア内の各点は、2つ以上のカメラ 114 の視野によってカバーされる。

#### 【0083】

ショッピングストアでは、棚及び他の在庫陳列構造は、ショッピングストアの側壁に沿って、または通路を形成する列に、または2つの構成の組合せでなど、様々な方法で配置することができる。図2は、通路 116a の一端から見た、通路 116a を形成する棚の

50

配置を示す。2つのカメラ、カメラA 206及びカメラB 208は、棚のような在庫陳列構造の上のショッピングストアの天井230及びフロア220から所定の距離で通路116aの上に配置される。カメラ114は、実空間内の在庫陳列構造及びフロアエリアのそれぞれの部分を包含する視野を有し、その上に配置されたカメラを備える。被写体として識別された候補関節のセットのメンバーの実空間内の座標は、被写体のフロアエリア内の位置を識別する。ショッピングストアの例示的な実施形態では、実空間は、在庫にアクセスできるショッピングストア内のフロア220のすべてを含むことができる。カメラ114は、フロア220及び棚のエリアが少なくとも2つのカメラによって見えるように配置され、配向される。カメラ114はまた、棚202及び204の少なくとも一部と、棚202及び204の前のフロアスペースとを覆う。カメラの角度は急峻な視点、真っ直ぐな視点及び角度の付いた視点の両方を有するように選択され、これにより、顧客のより完全な身体画像が得られる。一実施形態では、カメラ114が、ショッピングストア全体を通して、8フィート高さ以上で構成される。図13に、斯かる実施形態の説明図を示す。

10

#### 【0084】

図2では、カメラ206及び208が重なり合う視野を有し、それぞれ重なり合う視野216及び218で棚A 202と棚B 204との間の空間をカバーする実空間内の位置は、実空間座標系の $(x, y, z)$ 点として表される。「 $x$ 」及び「 $y$ 」は、ショッピングストアのフロア220とすることができる2次元(2D)平面上の位置を表し、値「 $z$ 」は、1つの構成ではフロア220における2Dプレーン上の点の高さである。

#### 【0085】

20

図3は、図2の上から見た通路116aを示し、通路116a上のカメラ206及び208の位置の例示的な配置を更に示す。カメラ206及び208は、通路116aの両端の近くに配置される。カメラA 206は棚A 202から所定の距離に配置され、カメラB 208は棚B 204から所定の距離に配置される。3つ以上のカメラが通路上に配置される別の実施形態では、カメラは互いに等しい距離に配置される。このような実施形態では、2つのカメラが両端の近くに配置され、第3のカメラが通路の中央に配置される。多数の異なるカメラ配置が可能であることが理解される。

#### [カメラ校正]

#### 【0086】

カメラ校正器120は2つのタイプの校正、即ち、内部及び外部校正を実行する。内部校正では、カメラ114の内部パラメータが校正される。内部カメラパラメータの例には、焦点距離、主点、スキュー、魚眼係数などが含まれる。内部カメラ校正のための様々な手法を使用することができる。斯かる手法の1つは、Zhangによって、IEEE Transactions on Pattern Analysis and Machine Intelligence、Volume 22、No. 11、November 2000に発行された「A flexible new technique for camera calibration」に示されている。

30

#### 【0087】

外部校正では、外部カメラパラメータが、2D画像データを実空間の3D座標に変換するためのマッピング・パラメータを生成するために校正される。一実施形態では、人物などの1つの被写体が実空間に導入される。被写体は、各カメラ114の視野を通過する経路上で実空間を移動する。実空間内の任意の所与の点において、被写体は、3Dシーンを形成する少なくとも2つのカメラの視野内に存在する。しかしながら、2つのカメラは、それぞれの2次元(2D)画像平面において同じ3Dシーンの異なるビューを有する。被写体の左手首などの3Dシーン内の特徴は、それぞれの2D画像平面内の異なる位置にある2つのカメラによって見られる。

40

#### 【0088】

点対応は、所与のシーンについて重複する視野を有する全てのカメラ・ペアの間で確立される。各カメラは同じ3Dシーンの異なる視野を有するので、点対応は3Dシーンにおける同じ点の投影を表す2つのピクセル位置(重なり合う視野を有する各カメラからの1

50

つの位置)である。外部較正のために、画像認識エンジン 112a ~ 112n の結果を使用して、各 3D シーンについて多くの点对応が識別される。画像認識エンジンは関節の位置を、それぞれのカメラ 114 の 2D 画像平面内のピクセルの (x, y) 座標、例えば、行及び列番号として識別する。一実施形態では、関節は、被写体の 19 の異なるタイプの関節のうちの 1 つである。被写体が異なるカメラの視野を通して移動するとき、追跡エンジン 110 は、較正に使用される被写体の 19 の異なるタイプの関節の各 (x, y) 座標を、画像毎にカメラ 114 から受け取る。

#### 【0089】

例えば、カメラ A からの画像と、カメラ B からの画像との両方が同じ時点に、重なり合う視野で撮影された場合を考える。カメラ A からの画像には、カメラ B からの同期画像のピクセルに対応するピクセルがあり、カメラ A とカメラ B の両方の視野内の或る物体または表面の特定の点があり、その点が両方の画像フレームのピクセルに取り込まれていると考える。外部カメラ較正では、多数のそのような点が識別され、対応点と呼ばれる。較正中にカメラ A 及びカメラ B の視野内に 1 つの被写体があるので、この被写体の主要な関節、例えば左手首の中心が識別される。これらの主要な関節がカメラ A 及びカメラ B の両方からの画像フレーム内に見える場合、これら是对応点を表すと仮定される。この処理は、多くの画像フレームについて繰り返され、重なり合う視野を有する全てのカメラ・ペアについて対応点の大きな集合を構築する。一実施形態では、画像が 30 FPS (フレーム/秒) 以上のレートで、フル RGB (赤、緑、及び青) カラーで 720 ピクセルの解像度で、すべてのカメラからストリーミングされる。これらの画像は、一次元配列 (フラット配列とも呼ばれる) の形態である。

#### 【0090】

被写体について上記で収集された多数の画像を使用して、重なり合う視野を有するカメラ間の対応点を決定することができる。重なり合う視野を有する 2 つのカメラ A 及び B を考える。カメラ A、B のカメラ中心と 3D シーンの関節位置 (特徴点ともいう) を通る平面を「エピポーラ平面」と呼び、エピポーラ平面とカメラ A、B の 2D 画像平面との交差箇所を「エピポーラ線」と定義する。これらの対応点が与えられると、カメラ A からの対応点を、カメラ B の画像フレーム内の対応点と交差することが保証されるカメラ B の視野内のエピポーラ線に正確にマッピングすることができる変換が決定される。被写体について上記で収集された画像フレームを使用して、変換が生成される。この変換は非線形であることが当技術分野で知られている。更に、一般形態では、投影された空間へ及び投影された空間から移動する非線形座標変換と同様に、それぞれのカメラのレンズの半径方向の歪み補正が必要であることが知られている。外部カメラ較正では、理想的な非線形変換への近似が非線形最適化問題を解くことによって決定される。この非線形最適化機能は、重なり合う視野を有するカメラ 114 の画像を処理する様々な画像認識エンジン 112a ~ 112n の出力 (関節データ構造の配列) 内の同じ関節を識別するために、追跡エンジン 110 によって使用される。内部カメラ較正及び外部カメラ較正の結果は、較正データベース 170 に格納される。

#### 【0091】

実空間におけるカメラ 114 の画像内の点の相対位置を決定するための様々な手法を使用することができる。例えば、Longuet-Higgins が、「A computer algorithm for reconstructing a scene from two projections」(Nature、第 293 巻、1981 年 9 月 10 日) を公表している。本論文では、2 つの投影間の空間的關係が未知であるとき、遠近投影の相関ペアからシーンの 3 次元構造を計算することが提示されている。Longuet-Higgins の論文は、実空間での各カメラの他のカメラに対する位置を決定する手法を提示する。更に、その手法は、実空間における被写体の三角測量を可能にし、重なり合う視野を有するカメラ 114 からの画像を使用して z 座標の値 (フロアからの高さ) を識別する。実空間の任意の点、例えば、実空間の一角の棚の端を、実空間の (x, y, z) 座標系上の (0, 0, 0) 点とする。

10

20

30

40

50

## 【 0 0 9 2 】

本技術の一実施形態では、外部較正のパラメータが2つのデータ構造に格納される。第1のデータ構造は、固有パラメータを格納する。固有パラメータは、3D座標から2D画像座標への射影変換を表す。第1のデータ構造は以下に示すように、カメラ毎の固有パラメータを含む。データ値はすべて浮動小数点数値である。このデータ構造は、「K」及び歪み係数として表される3×3固有行列を格納する。歪み係数は、6つの半径方向歪み係数と2つの接線方向歪み係数とを含む。半径方向の歪みは、光線がその光学的中心よりも、レンズの縁部の近傍でより大きく屈曲するときに生じる。接線方向の歪みは、レンズと像平面が平行でないときに生じる。以下のデータ構造は、第1のカメラのみの値を示す。同様のデータが全てのカメラ114に対して記憶される。

10

```
{
  1: {
    K: [[x, x, x], [x, x, x], [x, x, x]],
    distortion_coefficients: [x, x, x, x, x, x, x, x]
  },
  .....
}
```

## 【 0 0 9 3 】

第2のデータ構造はカメラ・ペア毎に、3×3基本行列(F)、3×3必須行列(E)、3×4投影行列(P)、3×3回転行列(R)、及び3×1平行移動ベクトル(t)を記憶する。このデータは、1つのカメラの基準フレーム内の点を別のカメラの基準フレームに変換するために使用される。カメラの各ペアについて、1つのカメラから別のカメラへフロア220の平面をマッピングするために、8つのホモグラフィ係数も記憶される。基本行列は、同じシーンの2つの画像間の関係であり、シーンからの点の投影が両方の画像において起こり得る場所を制約する。必須行列は、カメラが較正されている状態での、同じシーンの2つの画像間の関係でもある。投影行列は、3D実空間から部分空間へのベクトル空間投影を与える。回転行列は、ユークリッド空間における回転を実行するために使用される。平行移動ベクトル「t」は、図形または空間のすべての点を所与の方向に同じ距離だけ移動させる幾何学的変形を表す。ホモグラフィ・フロア係数は、重なり合う視野を有するカメラによって見られるフロア220上の被写体の特徴の画像を結合するために使用される。第2のデータ構造を以下に示す。同様のデータが、全てのカメラ・ペアについて記憶される。前述のように、xは浮動小数点数値を表す。

20

30

```
{
  1: {
    2: {
      F: [[x, x, x], [x, x, x], [x, x, x]],
      E: [[x, x, x], [x, x, x], [x, x, x]],
      P: [[x, x, x, x], [x, x, x, x], [x, x, x, x]],
      R: [[x, x, x], [x, x, x], [x, x, x]],
      t: [x, x, x],
      homography_floor_coefficients: [x, x, x, x, x, x, x, x]
    }
  },
  .....
}
```

40

## [ ネットワーク構成 ]

## 【 0 0 9 4 】

図4は、画像認識エンジンをホストするネットワークのアーキテクチャ400を示す。システムは、図示する実施形態では、複数のネットワーク・ノード101a～101nを含む。該実施形態では、ネットワーク・ノードは、処理プラットフォームとも呼ばれる。

50

処理プラットフォーム 101a ~ 101n 及びカメラ 412、414、416、418 は、ネットワーク 481 に接続される。

【0095】

図4は、ネットワークに接続された複数のカメラ412、414、416、418を示す。多数のカメラを特定のシステムに配備することができる。一実施形態では、カメラ412 ~ 418が、それぞれイーサネット（登録商標）ベースのコネクタ422、424、426、及び428を使用してネットワーク481に接続される。該実施形態では、イーサネット（登録商標）ベースのコネクタがギガビットイーサネット（登録商標）とも呼ばれる1ギガビット/秒のデータ転送速度を有する。他の実施形態では、カメラ114が、ギガビットイーサネット（登録商標）よりも高速または低速のデータ転送速度を有することができる他のタイプのネットワーク接続を使用してネットワークに接続されると理解される。また、代替の実施形態では、1組のカメラを各処理プラットフォームに直接接続することができる、処理プラットフォームをネットワークに結合することができる。

10

【0096】

記憶サブシステム430は、本発明の特定の実施形態の機能を提供する基本的なプログラミング及びデータ構成を記憶する。例えば、複数の画像認識エンジンの機能を実施する様々なモジュールを記憶サブシステム430に格納することができる。記憶サブシステム430は、非一時的なデータ記憶媒体を備えるコンピュータ可読メモリの一例であり、コンピュータによって実行可能なメモリに記憶されたコンピュータ命令を有し、本明細書で説明されるデータ処理機能及び画像処理機能のすべてまたは任意の組合せを実行し、これには、実空間の変化を識別し、被写体を追跡し、本明細書で説明されるような処理によって実空間のエリア内において在庫商品を置くこと及び取ることを検出するためのロジックが含まれる。他の例では、コンピュータ命令は、1つまたは複数のコンピュータ可読非一時的データ記憶媒体を備えるポータブルメモリを含む他のタイプのメモリに格納することができる。

20

【0097】

これらのソフトウェアモジュールは、一般に、プロセッサ・サブシステム450によって実行される。ホスト・メモリ・サブシステム432は、通常、プログラム実行中に命令及びデータを記憶するためのメイン・ランダム・アクセス・メモリ（RAM）434と、固定命令が記憶される読取り専用メモリ（ROM）436とを含むいくつかのメモリを含む。一実施形態では、RAM 434がプラットフォーム101aに接続されたカメラ114からのビデオストリームを格納するためのバッファとして使用される。

30

【0098】

ファイル記憶サブシステム440は、プログラム・ファイル及びデータ・ファイルのための永続的記憶を提供する。例示的な一実施形態では、記憶サブシステム440が番号442で識別されるRAID0（独立ディスクの冗長配列）構成内に4つの120ギガバイト（GB）ソリッド・ステート・ディスク（SSD）を有する。CNNが被写体の関節を識別するために使用される例示的な実施形態では、RAID0 442が訓練データを記憶するために使用される。訓練中、RAM 434にないトレーニング・データはRAID0 442から読み出される。同様に、画像がトレーニングのために記録されているとき、RAM 434にないデータはRAID0 442に記憶される。例示的な実施形態では、ハードディスク・ドライブ（HDD）446が10テラバイトのストレージである。これは、RAID0 442ストレージよりもアクセス速度が遅い。ソリッド・ステート・ディスク（SSD）444は、画像認識エンジン112aのためのオペレーティング・システム及び関連ファイルを格納する。

40

【0099】

例示的な構成では、3つのカメラ412、414、及び416が処理プラットフォーム101aに接続される。各カメラは、カメラによって送られた画像を処理するために、専用グラフィックス処理ユニットGPU1 462、GPU2 464、及びGPU3 466を有する。1つの処理プラットフォームにつき、3つより少ないまたは多いカメラを接続

50



することできると理解される。従って、各カメラがカメラから受信した画像フレームを処理するための専用GPUを有するように、より少ないまたはより多いGPUがネットワーク・ノード内に構成される。プロセッサ・サブシステム450、記憶サブシステム430、及びGPU462、464、466は、バス・サブシステム454を使用して通信する。  
【0100】

ネットワーク・インターフェース・サブシステム、ユーザ・インターフェース出力デバイス、及びユーザ・インターフェース入力デバイスなどのいくつかの周辺デバイスも、処理プラットフォーム101aの一部を形成するバス・サブシステム454に接続される。これらのサブシステム及びデバイスは説明の明確性を改善するために、図4には意図的に示されていない。バス・サブシステム454は単一のバスとして概略的に示されているが、バス・サブシステムの代わりの実施形態では複数のバスを使用することができる。

10

【0101】

一実施形態では、カメラ412が、1288×964の解像度、30FPSのフレームレート、及び1.3メガピクセル/イメージで、300mm～無限大の作動距離を有する可変焦点レンズ、98.2°～23.8°の1/3インチセンサによる視野を有するChameleon3 1.3MP Color USB3 Vision(Sony ICX445)を使用して実装することができる。

[畳み込みニューラル・ネットワーク]

【0102】

処理プラットフォーム内の画像認識エンジンは、所定のレートで継続的な画像ストリームを受信する。一実施形態では、画像認識エンジンが畳み込みニューラル・ネットワーク(CNNと略す)を含む。

20

【0103】

図5は、符号500で示されるCNNによる画像フレームの処理を示す。入力画像510は、行列状に配置された画像ピクセルからなる行列である。一実施形態では、入力画像510が1280ピクセルの幅、720ピクセルの高さ、及びRGBとも呼ばれる3チャンネルの赤、青、及び緑を有する。チャンネルは、互いに積み重ねられた3つの1280×720の2次元画像として想像することができる。従って、入力画像は図5に示すように、1280×720×3の寸法を有する。

【0104】

2×2フィルタ520は、入力画像510と畳み込まれる。この実施形態では、フィルタが入力と畳み込まれるとき、パディングは適用されない。これに続いて、非線形関数が畳み込み画像に適用される。本実施形態では、正規化線形ユニット(ReLU)活性化を用いる。非線形関数の他の例には、シグモイド、双曲正接(tanh)、及びリーキーReLUなどのReLUの変形が含まれる。探索は、ハイパー・パラメータ値を見つけるために実行される。ハイパー・パラメータは、C<sub>1</sub>、C<sub>2</sub>、・・・、C<sub>N</sub>であり、C<sub>N</sub>は、畳み込み層「N」に対するチャンネル数を意味する。N及びCの典型的な値を図5に示す。N=25で表されるCNNには25層がある。Cの値は、層1～25の各畳み込み層におけるチャンネルの数である。他の実施形態では、残留接続、スクイズ励起モジュール、及び複数の解像度などの追加の特徴がCNN500に追加される。

30

40

【0105】

画像分類に使用される典型的なCNNでは、画像が畳み込み層を介して処理されるにつれて、画像のサイズ(幅及び高さ)が低減される。これは、入力画像のクラスを予測することを目的とするので、特徴識別に役立つ。しかし、図示の実施形態では、画像フレーム内の関節(特徴とも呼ばれる)を識別するだけでなく、実空間内の座標にマッピングできるように画像内のその位置を識別することも目標とするので、入力画像のサイズ(すなわち、画像の幅及び高さ)は縮小されない。従って、図5に示すように、この例では、CNNの畳み込み層を介して処理が進行することにつれて、画像の幅及び高さの寸法は変化しないままである。

【0106】

50

一実施形態では、CNN500が画像の各要素における被写体の19個の可能な関節のうちの1つを識別する。可能な関節は、足関節と非足関節の2つのカテゴリに分類することができる。関節分類の19番目のタイプは、被写体の全ての非関節特徴（すなわち、関節として分類されない画像の要素）に対するものである。

足関節：

足首関節（左右）

非足関節：

首

鼻

眼（左右）

耳（左右）

肩（左右）

肘（左右）

手首（左右）

尻（左右）

膝（左右）

非関節

【0107】

以上のように、本説明の目的のための「関節」は、実空間における被写体の追跡可能な特徴である。関節は、被写体の生理学的関節、または眼もしくは鼻などの他の特徴に対応し得る。

【0108】

入力画像のストリーム上の第1の分析セットは、実空間内の被写体の追跡可能な特徴を識別する。一実施形態では、これは「関節分析」と呼ばれる。このような実施形態では、関節分析に使用されるCNNは「関節CNN」と呼ばれる。一実施形態では、関節分析は、対応するカメラから受信される毎秒30フレームにわたって毎秒30回実行される。分析は時間的に同期され、すなわち、実空間における全ての被写体の関節を識別するために、1/30秒で、全てのカメラ114からの画像が、対応する関節CNNにおいて分析される。複数のカメラからの一時点での画像の分析の結果は、「スナップショット」として記憶される。

【0109】

スナップショットは、システムによってカバーされる実空間のエリア内の候補関節のコンステレーションを表す、ある時点の全てのカメラ114の画像からの関節データ構造の配列を含む辞書形式であり得る。一実施形態では、スナップショットは被写体データベース140に格納される。

【0110】

このCNNの例では、ソフトマックス関数が畳み込み層530の最終層内の画像のすべての要素に適用される。ソフトマックス関数は、任意の実数値のK次元ベクトルを、合計で1になる範囲[0, 1]の実数値のK次元ベクトルに変換する。一実施形態では、画像の要素は単一のピクセルである。ソフトマックス関数は、ピクセル毎の任意の実数値の19次元配列（19次元ベクトルとも呼ばれる）を、合計で1になる[0, 1]の実数値の19次元信頼度配列に変換する。画像フレーム内のピクセルの19次元は、被写体の19タイプの関節に更に対応するCNNの最終層内の19個のチャンネルに対応する。

【0111】

多数の画素は、その画像に対するソースカメラの視野内の被写体の数に応じて、1つの画像内の19タイプの関節の各々の1つとして分類することができる。

【0112】

画像認識エンジン112a~112nは、画像を処理して、画像の要素に対する信頼度配列を生成する。画像の特定の要素についての信頼度配列は、その特定の要素についての複数の関節タイプについての信頼値を含む。画像認識エンジン112a~112nの各々

10

20

30

40

50

は、それぞれ、画像毎に信頼度配列の出力行列 5 4 0 を生成する。最後に、各画像認識エンジンは、画像当たりの信頼度配列の各出力行列 5 4 0 に対応する関節データ構造の配列を生成する。特定の画像に対応する関節データ構造の配列は、関節タイプ、特定の画像の時間、及び特定の画像内の要素の座標によって、特定の画像の要素を分類する。信頼度配列の値に基づいて、各イメージ内の特定の要素の関節データ構造の関節タイプが選択される。

#### 【 0 1 1 3 】

被写体の各関節は、ヒートマップとして出力行列 5 4 0 に分布していると考えることができる。ヒートマップは、各関節タイプについて最高値（ピーク）を有する画像要素を示すように分解することができる。理想的には、特定の関節タイプの高い値を有する所与の画素について、所与の画素からの範囲外の周囲の画素はその関節タイプについてより低い値を有し、その結果、その関節タイプを有する特定の関節の位置を画像空間座標において識別することができる。それに対応して、その画像要素に対する信頼度配列はその関節に対して最も高い信頼値を有し、残りの 1 8 種類の関節に対してより低い信頼値を有する。

#### 【 0 1 1 4 】

一実施形態では、各カメラ 1 1 4 からの画像のバッチがそれぞれの画像認識エンジンによって処理される。例えば、6 つの連続的にタイムスタンプされた画像は、キャッシュ・コヒーレンスを有効に利用するためにバッチで連続的に処理される。CNN 5 0 0 の 1 つの層に対するパラメータは、メモリにロードされ、6 つの画像フレームのバッチに適用される。次に、次の層のパラメータがメモリにロードされ、6 つの画像のバッチに適用される。これは、CNN 5 0 0 内のすべての畳み込み層 5 3 0 について繰り返される。キャッシュ・コヒーレンスは処理時間を短縮し、画像認識エンジンの性能を改善する。

#### 【 0 1 1 5 】

3 次元（3 D）畳み込みと呼ばれる 1 つの斯かる実施形態では、CNN 5 0 0 の性能の更なる改善がバッチ内の画像フレームにわたって情報を共有することによって達成される。これは、関節のより正確な識別に役立ち、誤検知を減少させる。例えば、所与のバッチ内の複数の画像フレームにわたってピクセル値が変化しない画像フレーム内の特徴は、シェルフなどの静的物体である可能性が高い。所与のバッチ内の画像フレームにわたる同じピクセルの値の変化は、このピクセルが関節である可能性が高いことを示す。従って、CNN 5 0 0 はそのピクセルによって識別された関節を正確に識別するために、そのピクセルの処理により焦点を当てることができる。

#### [ 関節データ構造 ]

#### 【 0 1 1 6 】

CNN 5 0 0 の出力は、カメラ当たりの各画像に対する信頼度配列の行列である。信頼度配列の行列は、関節データ構造の配列に変換される。図 6 に示すような関節データ構造 6 0 0 は、各関節の情報を記憶するために使用される。関節データ構造 6 0 0 は、画像が受信されるカメラの 2 D 画像空間内の特定の画像内の要素の x 位置及び y 位置を識別する。関節番号は、識別された関節のタイプを識別する。例えば、一実施形態では、値は 1 ~ 1 9 の範囲である。値 1 は関節が左足首であることを示し、値 2 は関節が右足首であることを示し、以下同様である。関節のタイプは、出力行列 5 4 0 内のその要素に対する信頼度配列を使用して選択される。例えば、一実施形態では、左足首関節に対応する値がその画像要素の信頼度配列において最も高い場合、関節番号の値は「1」である。

#### 【 0 1 1 7 】

信頼度数は、その関節を予測する際の CNN 5 0 0 の信頼度の程度を示す。信頼度数の値が高ければ、CNN は自身の予想に確信していることになる。関節データ構造を一意に識別するために、関節データ構造に整数 ID が割り当てられる。上記マッピングに続いて、画像毎の信頼度配列の出力行列 5 4 0 は、画像毎の関節データ構造の配列に変換される。

#### 【 0 1 1 8 】

画像認識エンジン 1 1 2 a ~ 1 1 2 n はカメラ 1 1 4 から画像のシーケンスを受信し、画像を処理して、上述のように関節データ構造の対応する配列を生成する。特定の画像の

10

20

30

40

50

関節データ構造の配列は、関節タイプ、特定の画像の時間、及び特定の画像内の要素の座標によって、特定の画像の要素を分類する。一実施形態では画像認識エンジン 1 1 2 a ~ 1 1 2 n が畳み込みニューラル・ネットワーク C N N 5 0 0 であり、関節タイプは被写体の 1 9 種類の関節のうちの 1 つ、特定の画像の時間は特定の画像についてソースカメラ 1 1 4 によって生成された画像のタイムスタンプであり、座標 ( x , y ) は 2 D 画像平面上の要素の位置を特定する。

#### 【 0 1 1 9 】

一実施形態では、関節分析が、各入力画像に対して、k 最近傍、ガウス混合、様々な画像形態変換、及び関節 C N N の組み合わせを実行することを含む。この結果は、各時点において画像数をビットマスクにマッピングするリング・バッファ内にビットマスクの形式で格納することができる関節データ構造の配列を含む。

#### [ 追跡エンジン ]

#### 【 0 1 2 0 】

追跡エンジン 1 1 0 は、重なり合う視野を有するカメラからの画像のシーケンス内の画像に対応する、画像認識エンジン 1 1 2 a ~ 1 1 2 n によって生成された関節データ構造の配列を受信するように構成される。画像当たりの関節データ構造の配列は、図 7 に示すように、画像認識エンジン 1 1 2 a ~ 1 1 2 n によってネットワーク 1 8 1 を介して追跡エンジン 1 1 0 に送られる。追跡エンジン 1 1 0 は、様々なシーケンスの画像に対応する関節データ構造の配列内の要素の座標を、実空間内の座標を有する候補関節に変換する。追跡エンジン 1 1 0 は、実空間における座標 ( 関節のコンステレーション ) を有する候補関節のセットを、実空間における被写体として識別するためのロジックを備える。一実施形態では、追跡エンジン 1 1 0 が、所与の時点におけるすべてのカメラについて、画像認識エンジンからの関節データ構造の配列を蓄積し、候補関節のコンステレーションを識別するために使用されるように、この情報を辞書として被写体データベース 1 4 0 に格納する。辞書は、キー値ペアの形式で編成することができ、ここで、キーはカメラ I D であり、値はカメラからの関節データ構造の配列である。斯かる実施形態では、この辞書が候補関節を決定し、関節を被写体に割り当てるために、ヒューリスティックス・ベースの分析で使用される。斯かる実施形態では、追跡エンジン 1 1 0 の高レベル入力、処理、及び出力が表 1 に示されている。

表 1 : 例示的な実施形態における追跡エンジン 1 1 0 からの入力、処理、及び出力。

入力	処理	出力
画像毎、及び、関節データ構造毎の関節データ構造の配列 - 個有 I D - 信頼度数 - 関節番号 - 画像空間内の ( x , y ) の位置	- 関節辞書の作成 - 重なり合う視野を有するカメラの視野における関節の位置を候補関節に再投影	- ある時点での実空間内の被写体のリスト

#### [ 関節の候補関節へのグループ化 ]

#### 【 0 1 2 1 】

追跡エンジン 1 1 0 は 2 つの次元、すなわち、時間及び空間に沿った関節データ構造の配列を受け取る。時間次元に沿って、追跡エンジンは、カメラ当たり画像識別エンジン 1 1 2 a ~ 1 1 2 n によって処理された関節データ構造のタイムスタンプ付き配列を連続的

に受け取る。関節データ構造は、重なり合う視野を有するカメラからの画像において、ある期間にわたる同じ被写体の同じ関節の複数のインスタンスを含む。特定の画像内の要素の $(x, y)$ 座標は、通常、特定の関節が属する被写体の動きのために、関節データ構造の連続的にタイムスタンプされた配列において異なっている。例えば、左手首関節として分類された20個の画素は、特定のカメラからの多くの連続的にタイムスタンプされた画像に現れることができ、各左手首関節は、画像毎に変化していること或いは変化しないことができる実空間内の位置を有する。その結果、多くの連続的にタイムスタンプされた関節データ構造の配列内の20個の左手首関節データ構造600は、経時的に実空間内の同じ20個の関節を表すことができる。

#### 【0122】

重なり合う視野を有する複数のカメラは実空間内の各位置をカバーするので、任意の所与の時点に、カメラ114のうちの2つ以上の画像に同じ関節が現れる可能性がある。カメラ114は時間的に同期され、従って、追跡エンジン110は任意の所与の時点に、重なり合う視野を有する複数のカメラから、特定の関節の関節データ構造を受信する。これは空間次元であり、2つの次元、すなわち、時間及び空間のうちの第2の次元であり、追跡エンジン110は、空間次元に沿って関節データ構造の配列内のデータを受け取る。

#### 【0123】

追跡エンジン110は、ヒューリスティックス・データベース160に格納されたヒューリスティックスの最初の組を使用して、関節データ構造の配列から関節データ構造の候補を識別する。目標は、ある期間にわたってグローバル・メトリックを最小化することである。グローバル・メトリック計算器702は、グローバル距離を計算する。グローバル・メトリックは、以下に説明する複数の値の合計である。直観的には、追跡エンジン110によって時間次元と空間次元に沿って受信された関節データ構造の配列における関節がそれぞれの被写体に正しく割り当てられる場合、グローバル・メトリックの値は最小である。例えば、顧客が通路内を移動するショッピングストアの実施形態を考える。顧客Aの左手首が顧客Bに誤って割り当てられた場合、グローバル・メトリックの値は増加する。従って、各顧客に対する各関節のグローバル・メトリックを最小化することは、最適化問題である。この問題を解決する1つの選択肢は、関節の全ての可能な接続を試みることである。しかしながら、これは、顧客の数が増加することにつれて、扱いにくくなる可能性がある。

#### 【0124】

この問題を解決するための第2のアプローチは、ヒューリスティックスを使用して、単一の被写体に対する候補関節のセットのメンバーとして識別される関節の可能な組み合わせを低減することである。例えば、関節の相対位置の既知の生理学的特性のために、左手首関節は被写体の他の関節から空間的に遠く離れた被写体に属することができない。同様に、画像から画像への位置の変化が小さい左手首関節は、被写体が非常に高速で動くことが期待されないため、時間的に遠く離れた画像から同じ位置に同じ関節を有する被写体に属する可能性が低い。これらの初期ヒューリスティックスは、特定の被写体として分類され得る候補関節のコンステレーションのための時間及び空間における境界を構築するために使用される。特定の時間及び空間境界内の関節データ構造内の関節は、実空間内に存在する被写体としての候補関節のセットに割り当てるための「候補関節」と見なされる。これらの関節候補は、ある期間（時間次元）にわたる同じカメラからの多数の画像からの関節データ構造の配列において、重なり合う視野（空間次元）を有する様々なカメラにわたって識別された関節を含む。

#### 〔足関節〕

#### 【0125】

関節は、関節のリストで上述したように、関節をコンステレーションに、足関節及び非足関節にグループ化するための手順を目的として分割することができる。本実施例における左及び右足首関節タイプは、この手順の目的として足関節と考えられる。追跡エンジン110は、足関節を使用して、特定の被写体の候補関節のセットの識別を開始することが

10

20

30

40

50

できる。ショッピングストアの実施形態では、顧客の足が図2に示すようにフロア220上にある。カメラ114のフロア220までの距離は既知である。従って、重なり合う視野を有するカメラの画像に対応するデータ関節データ構造の配列からの足関節の関節データ構造を組み合わせる場合、追跡エンジン110は、既知の深さ（ $z$ 軸に沿った距離）を仮定することができる。足関節の深さの値はゼロ、すなわち、実空間の（ $x, y, z$ ）座標系において（ $x, y, 0$ ）である。この情報を使用して、画像追跡エンジン110は候補足関節を識別するために、重なり合う視野を有するカメラからの足関節の関節データ構造を組み合わせるために、ホモグラフィック・マッピングを適用する。このマッピングを使用して、画像空間における（ $x, y$ ）座標における関節の位置が、実空間における（ $x, y, z$ ）座標における位置に変換され、候補足関節が得られる。この処理は、それぞれの関節データ構造を使用して候補左足関節及び候補右足関節を識別するために別々に実行される。

10

#### 【0126】

これに続いて、追跡エンジン110は、候補左足関節及び候補右足関節を組み合わせ（候補関節のセットにそれらを割り当て）、被写体を作成することができる。候補関節の銀河からの他の関節は、作成された被写体の関節タイプの幾つかまたは全てのコンステレーションを構築するために、被写体にリンクすることができる。

#### 【0127】

左側候補足関節が1つしかなく、右側候補足関節が1つしかない場合、特定の時点で特定の空間に1つの被写体しか存在しないことを意味する。追跡エンジン110は、関節のセットに属する候補左足関節及び候補右足関節を有する新しい被写体を作成する。被写体は、被写体データベース140に保存される。複数の候補左足関節及び候補右足関節がある場合、グローバル・メトリック計算器702はグローバル・メトリックの値が最小化されるように、各候補左足関節を各候補右足関節に結合して被写体を作成することを試みる。

20

#### 〔非足関節〕

#### 【0128】

特定の時間及び空間境界内の関節データ構造の配列から候補非足関節を識別するために、追跡エンジン110は任意の所与のカメラAから、重なり合う視野を有するその隣接するカメラBへの非線形変換（基本行列とも呼ばれる）を使用する。非線形変換は、単一の多関節被写体を使用して計算され、上述のように較正データベース170に格納される。例えば、重なり合う視野を有する2つのカメラA及びBについて、候補非足関節は、以下のように識別される。カメラAからの画像フレーム内の要素に対応する関節の配列データ構造内の非足関節はカメラBからの同期画像フレーム内のエピソード線にマッピングされる。カメラAの特定の画像の関節データ構造の配列内の関節データ構造によって識別される関節（マシンビジョン文献では特徴とも呼ばれる）が、カメラBの画像内に現れる場合、対応するエピソード線上に現れる。例えば、カメラAからの関節データ構造内の関節が左手関節である場合、カメラBの画像内のエピソード線上の左手関節はカメラBの視点から見て同じ左手関節を表す。カメラA及びBの画像内のこれら2つの点が実空間内の3D場面内の同じ点の投影であり、「共役ペア」と呼ばれる。

30

#### 【0129】

Nature、Volume 293、1981年9月10日号に、「A computer algorithm for reconstructing a scene from two projections」という表題の論文に掲載されたLonguet-Higginsによる手法などのマシンビジョン技術は、実空間におけるフロア220からの関節の高さを決定するために、対応点の共役ペアに適用される。上記の方法を適用するには、重なり合う視野を有するカメラ間の所定のマッピングが必要である。そのデータは、上述のカメラ114の較正中に決定された非線形関数として較正データベース170に格納される。

40

#### 【0130】

追跡エンジン110は、重なり合う視野を有するカメラからの画像シーケンス内の画像

50

に対応する関節データ構造の配列を受信し、様々なシーケンス内の画像に対応する関節データ構造の配列内の要素の座標を、実空間内の座標を有する候補非足関節に変換する。識別された候補非足関節は、グローバル・メトリック計算器 702 を使用して、実空間内の座標を有する被写体の集合にグループ化される。グローバル・メトリック計算器 702 は、グローバル・メトリック値を計算し、非足関節の異なる組み合わせをチェックすることによって値を最小化することを試みる。一実施形態では、グローバル・メトリックが 4 つのカテゴリに編成されたヒューリスティックスの合計である。候補関節のセットを識別するロジックは、候補関節のセットを被写体として識別するために、実空間における被写体の関節間の物理的関係に基づくヒューリスティック関数を含む。関節間の物理的関係の例は、以下に記載されるようなヒューリスティックスにおいて考慮される。

10

[ 第 1 カテゴリのヒューリスティックス ]

【 0 1 3 1 】

第 1 カテゴリのヒューリスティックスは、同じまたは異なる時点における同じカメラ視野内の 2 つの提案された被写体関節位置間の類似性を確認するためのメトリックを含む。一実施形態ではこれらのメトリックは浮動小数点値であり、より高い値は関節の 2 つのリストが同じ被写体に属する可能性が高いことを意味する。ショッピングストアの例示的な実施形態を考えると、メトリックは、時間次元に沿った 1 つの画像から次の画像までの、1 つのカメラ内の顧客の同じ関節間の距離を決定する。カメラ 412 の視野内の顧客 A が与えられると、メトリックの第 1 のセットは、カメラ 412 からの 1 つの画像からカメラ 412 からの次の画像までの人物 A の関節の各々の間の距離を決定する。メトリックは、カメラ 114 からの画像当たりの関節データ構造の配列において関節データ構造 600 に適用される。

20

【 0 1 3 2 】

一実施形態では、第 1 カテゴリのヒューリスティックスにおける 2 つの例示的なメトリックを以下に列挙する：

1. フロア上の 2 人の被写体の左足首関節とフロア上の 2 人の被写体の右足首関節との間の合計されたユークリッド 2 D 座標距離の逆数（特定のカメラからの特定の画像の x、y 座標値を使用する）。

2. 画像フレーム内の被写体の非足関節のすべてのペアの間のユークリッド 2 D 座標距離の逆数の合計。

30

[ 第 2 カテゴリのヒューリスティックス ]

【 0 1 3 3 】

第 2 カテゴリのヒューリスティックスは、同じ時点で複数のカメラの視野から 2 つの提案された被写体関節位置間の類似性を確認するためのメトリックを含む。一実施形態では、これらのメトリックは浮動小数点値であり、より高い値は関節の 2 つのリストが同じ被写体に属する可能性が高いことを意味する。ショッピングストアの例示的な実施形態を考えると、第 2 のセットのメトリックは、同じ時点で（重なり合う視野を有する）2 つ以上のカメラからの画像フレーム内の顧客の同じ関節間の距離を決定する。

【 0 1 3 4 】

一実施形態では、第 2 カテゴリのヒューリスティックスにおける 2 つの例示的なメトリックを以下に列挙する：

40

【 0 1 3 5 】

フロア上の 2 人の被写体の左足首関節とフロア上の 2 人の被写体の右足首関節との間のユークリッド 2 D 座標距離の逆数（特定のカメラからの特定の画像の x、y 座標値を使用する）を合計した。第 1 の被写体の足首関節位置は、第 2 の被写体がホモグラフィック・マッピングを通して見えるカメラに投影される。

【 0 1 3 6 】

線と点との間のユークリッド 2 D 座標距離の逆数の関節の全てのペアの和であり、ここで、線は、視野内に第 1 の被写体を有する第 1 のカメラから、視野内に第 2 の被写体を有する第 2 のカメラまでの画像の関節のエピポーラ線であり、点は、第 2 のカメラからの画

50

像内の第 2 の被写体の関節である。

[ 第 3 カテゴリのヒューリスティックス ]

【 0 1 3 7 】

第 3 カテゴリのヒューリスティックスは、同じカメラビュー内の提案された被写体関節位置のすべての関節間の類似性を同じ時点に確認するためのメトリックを含む。ショッピングストアの例示的な実施形態を考えると、このカテゴリのメトリックは、1つのカメラからの1つのフレームにおける顧客の関節間の距離を決定する。

[ 第 4 カテゴリのヒューリスティックス ]

【 0 1 3 8 】

第 4 カテゴリのヒューリスティックスは、提案された被写体関節位置間の相違を確認するためのメトリックを含む。一実施形態では、これらのメトリックは浮動小数点値である。より高い値は、関節の2つのリストが同じ被写体ではない可能性がより高いことを意味する。一実施形態では、このカテゴリにおける2つの例示的なメトリックが以下を含む：

- 1 . 2 人の提案された被写体の頸部関節間の距離。
- 2 . 2 人の被写体間の関節のペア間の距離の合計。

【 0 1 3 9 】

一実施形態では、経験的に決定され得る様々な閾値が以下に記載されるように、上記に列挙されたメトリックに適用される：

- 1 . メトリック値が、関節が既知の被写体に属すると考えるのに十分に小さい場合を判定するための閾値。
- 2 . メトリック類似性スコアが良好すぎる状態で関節が属する可能性がある潜在的な候補被写体が多すぎる場合を判定するための閾値。
- 3 . 関節の集合が、経時的に、以前は実空間には存在しなかった新しい被写体と見なされるのに十分に高いメトリック類似性を有する場合を判定するための閾値。
- 4 . 被写体が既に実空間にいない場合を判定するための閾値。
- 5 . 追跡エンジン 1 1 0 が、間違っ 2 つの被写体を混同した場合を判定するための閾値。

【 0 1 4 0 】

追跡エンジン 1 1 0 は、被写体として識別された関節のセットを記憶するロジックを含む。候補関節のセットを識別するロジックは、特定の時間に撮影された画像において識別された候補関節が先行する画像において被写体として識別された候補関節のセットのうちの1つのメンバーに対応するかどうかを判定するロジックを含む。一実施形態では、追跡エンジン 1 1 0 が被写体の現在の関節位置を、同じ被写体の以前に記録された関節位置と、定期的に比較する。この比較により、追跡エンジン 1 1 0 は、実空間内の被写体の関節位置を更新することができる。更に、これを使用して、追跡エンジン 1 1 0 は誤検知（すなわち、誤って識別された被写体）を識別し、実空間に既に存在しない被写体を除去する。

【 0 1 4 1 】

ショッピングストアの実施形態の例を考えると、追跡エンジン 1 1 0 はより早い時点に、顧客（被写体）を生成したが、ある時間の後、追跡エンジン 1 1 0 はその特定の顧客に対して現在の関節位置を有していない。それは、顧客が誤って生成されたことを意味する。追跡エンジン 1 1 0 は、誤って生成された被写体を被写体データベース 1 4 0 から削除する。一実施形態では、追跡エンジン 1 1 0 はまた、上述の処理を用いて、実空間から積極的に識別された被写体を除去する。ショッピングストアの例を考えると、顧客がショッピングストアを離れると、追跡エンジン 1 1 0 は、被写体データベース 1 4 0 から対応する顧客レコードを削除する。斯かる一実施形態では、追跡エンジン 1 1 0 が「顧客が店を出た」ことを示すために、被写体データベース 1 4 0 内のこの顧客レコードを更新する。

【 0 1 4 2 】

一実施形態では、追跡エンジン 1 1 0 が足ヒューリスティックスと非足ヒューリスティックスを同時に適用することによって、被写体を識別しようと試みる。これにより、被写体の連結関節の「アイランド」が生成される。追跡エンジン 1 1 0 が、時間次元と空間次

10

20

30

40

50



元に沿って関節データ構造の配列を更に処理すると、アイランドの大きさが増加する。最終的に、関節のアイランドは被写体を形成する関節の他のアイランドと融合し、そして、被写体データベース 140 に格納される。一実施形態では、追跡エンジン 110 が所定の期間、未割り当ての関節の記録を維持する。この間、追跡エンジンは、未割り当ての関節を既存の被写体に割り当てるか、またはこれらの未割り当ての関節から新しい多関節存在物を作成しようと試みる。追跡エンジン 110 は、所定の期間の後、未割り当ての関節を破棄する。他の実施形態では、被写体を識別し追跡するために、上述の列挙したものとは異なるヒューリスティックスが使用されることを理解されたい。

#### 【0143】

一実施形態では、追跡エンジン 110 をホストするノード 102 に接続されたユーザ・インターフェース出力デバイスが、実空間内の各被写体の位置を表示する。斯かる一実施形態では、出力デバイスの表示が、被写体の新しい位置でもって、定期的にリフレッシュされる。

#### [被写体データ構造]

#### 【0144】

被写体の関節は、上述のメトリックを使用して互いに接続される。その際、追跡エンジン 110 は新しい被写体を生成し、それぞれの関節位置を更新することによって既存の被写体の位置を更新する。図 8 は、被写体を格納するための被写体データ構造 800 を示す。該データ構造 800 は、被写体関連データをキー値辞書として格納する。キーはフレーム番号であり、値は別のキー値辞書であり、ここでは、キーはカメラ ID であり、値は（被写体の）18 個の関節と実空間内のそれらの位置のリストである。被写体データは、被写体データベース 140 に格納される。新しい被写体毎に、被写体データベース 140 内の被写体のデータにアクセスするために使用される固有識別子も割り当てられる。

#### 【0145】

一実施形態では、システムが被写体の関節を識別し、被写体の骨格を作成する。骨格は、実空間に投影され、実空間における被写体の位置及び向きを示す。これは、マシビジョンの分野では「姿勢推定」とも呼ばれる。一実施形態では、システムがグラフィカル・ユーザ・インターフェース (GUI) 上に実空間内の被写体の向き及び位置を表示する。一実施形態では、画像分析は匿名であり、すなわち、関節分析によって作成された被写体に割り当てられた固有識別子は実空間内の任意の特定被写体の詳細な個人識別情報（名前、電子メールアドレス、郵送先住所、クレジットカード番号、銀行口座番号、運転免許証番号など）を識別しない。

#### [被写体追跡の処理フロー]

#### 【0146】

本明細書では、ロジックを示す幾つかのフローチャートを説明する。ロジックは、プロセッサによってアクセス可能かつ実行可能なメモリに格納されたコンピュータ・プログラムを使用してプログラムされ、上述のように構成されたプロセッサを使用して、及び他の構成では、フィールドプログラマブル集積回路を含む専用ロジックハードウェアによって、及び専用ロジックハードウェアとコンピュータ・プログラムとの組合せによって実装され得る。本明細書のすべてのフローチャートでは、達成される機能に影響を及ぼすことなく、ステップの多くを組み合わせること、並列に実行すること、または異なる順序で実行することができることが理解されよう。幾つか場合では、読者が理解するように、ステップの再編は、特定の他の変更が同様に行われる場合にのみ、同じ結果を達成する。他の場合には、読者が理解するように、ステップの再編は特定の条件が満たされる場合にのみ、同じ結果を達成する。更に、本明細書のフローチャートは実施形態の理解に関連するステップのみを示し、他の機能を達成するための多数の追加のステップが、示されたステップの前、後、及びそれらの間で実行され得ることが理解されるのであろう。

#### 【0147】

図 9 は、被写体を追跡するための処理ステップを示すフローチャートである。処理はステップ 902 で開始する。実空間のエリア内に視野を有するカメラ 114 は、ステップ 9

10

20

30

40

50

04の進行中に較正される。ビデオ処理は、ステップ906において、画像認識エンジン112a~112nによって実行される。一実施形態では、ビデオ処理がそれぞれのカメラから受信された画像フレームのバッチを処理するために、カメラ毎に実行される。それぞれの画像認識エンジン112a~112nからのすべてのビデオ処理の出力は、ステップ908で追跡エンジン110によって実行されるシーン処理への入力として与えられる。シーン処理は新しい被写体を識別し、既存の被写体の共同位置を更新する。ステップ910では、処理すべき画像フレームがまだあるかどうかチェックされる。更に画像フレームがある場合、処理はステップ906に進み、なければ、処理はステップ914で終了する。

#### 【0148】

処理ステップ904「実空間でカメラを較正する」のより詳細な処理ステップが、図10のフローチャートに示されている。較正処理は、ステップ1002で、実空間の(x, y, z)座標に対する(0, 0, 0)点を識別することによって開始する。ステップ1004において、視野内に位置(0, 0, 0)を有する第1のカメラが較正される。カメラ較正の更なる詳細は、本出願において以前に提示されている。ステップ1006において、第1のカメラと重なり合う視野を有する次のカメラが較正される。ステップ1008では、較正すべきカメラがまだあるかどうかチェックされる。この処理は、全てのカメラ114が較正されるまで、ステップ1006で繰り返される。

#### 【0149】

次の処理ステップ1010では、被写体の実空間に導入され、重なり合う視野を有するカメラ間の対応点の共役ペアが識別される。この処理のいくつかの詳細は上述されている。この処理は、ステップ1012で、重なり合うカメラのすべてのペアについて繰り返される。カメラがこれ以上存在しない場合、処理は終了する(ステップ1014)。

#### 【0150】

図11のフローチャートは、「ビデオ処理」ステップ906のより詳細なステップを示す。ステップ1102では、カメラ当たりk個の連続してタイムスタンプされた画像が更なる処理のためのバッチとして選択される。一実施形態では、kの値=6で、画像認識エンジン112a~112nをそれぞれホストするネットワーク・ノード101a~101nにおけるビデオ処理のための利用可能なメモリに基づいて計算される。次のステップ1104では、画像のサイズが適切な寸法に設定される。一実施形態では、画像が1280ピクセルの幅、720ピクセルの高さ、及び3つのチャンネルRGB(赤色、緑色、及び青色を表す)を有する。ステップ1106では、複数のトレーニングされた畳み込みニューラル・ネットワーク(CNN)が画像を処理し、画像当たりの関節データ構造の配列を生成する。CNNの出力は、画像当たりの関節データ構造の配列である(ステップ1108)。この出力は、ステップ1110でシーン処理に送られる。

#### 【0151】

図12Aは、図9「シーン処理」ステップ908のより詳細なステップの第1の部分を示すフローチャートである。シーン処理は、ステップ1202において、複数のビデオ処理からの出力を結合する。ステップ1204では、関節データ構造が足関節または非足関節を識別するかどうかチェックされる。関節データ構造が足関節のものである場合、ステップ1206において、重なり合う視野を有するカメラからの画像に対応する関節データ構造を結合するために、ホモグラフィック・マッピングが適用される。この処理は、候補足関節(左足関節及び右足関節)を識別する。ステップ1208で、ステップ1206で識別された候補足関節にヒューリスティクスを適用して、候補足関節のセットを被写体として識別する。ステップ1210において、候補足関節のセットが既存の被写体に属するかどうかチェックされる。属さない場合には、ステップ1212で、新しい被写体を作成される。属する場合は、ステップ1214で、既存の被写体が更新される。

#### 【0152】

図12Bは、「シーン処理」ステップ908のより詳細なステップの第2の部分を示すフローチャートである。ステップ1240では、重なり合う視野を有するカメラからの画

10

20

30

40

50

像シーケンス内の画像に対応する関節データ構造の複数の配列から、非足関節のデータ構造が組み合わされる。これは、第1のカメラからの第1の画像からの対応点を、重なり合う視野を有する第2のカメラからの第2の画像にマッピングすることによって実行される。この処理の幾つかの詳細は上述されている。ヒューリスティックスは、ステップ1242において、候補非足関節に適用される。ステップ1246では、候補非足関節が既存の被写体に属するかどうか判定される。属する場合、ステップ1248において、既存の被写体が更新される。属さない場合は、ステップ1250において、所定の時間の後に、候補非足関節を既存の被写体と一致させるために、候補非足関節が再び処理される。ステップ1252において、非足関節が既存の被写体に属するかどうかチェックされる。属するのであれば、ステップ1256で被写体が更新される。属さない場合は、ステップ1254で関節は破棄される。

10

#### 【0153】

例示的な実施形態では、新しい被写体を識別し、被写体を追跡し、被写体（実空間を離れたか、または間違っ生成された）を削除する処理はランタイムシステム（推論システムとも呼ばれる）によって実行される「存在物結束アルゴリズム」の一部として実装される。存在物は、上記の被写体と呼ばれる関節のコンステレーションである。存在物結束アルゴリズムは、実空間内の存在物を識別し、実空間内の関節の位置を更新して、存在物の移動を追跡する。

#### 【0154】

図14は、ビデオ処理1411及びシーン処理1415を示す。図示の実施形態では、4つのビデオ処理が示されており、それぞれが、1または複数のカメラ114からの画像を処理する。ビデオ処理は、上述のように画像を処理し、フレーム毎に関節を識別する。一実施形態では、それぞれのビデオ処理が、フレーム当たりの関節毎に、2D座標、信頼度数、関節番号、及び固有IDを識別する。すべてのビデオ処理の出力1452は、入力1453としてシーン処理1415に与えられる。一実施形態では、シーン処理が、キーがカメラIDであり、値が関節の配列である時点毎の関節キー値辞書を作成する。関節は、重なり合う視野を有するカメラの視点に再投影される。再投影された関節はキー値辞書として記憶され、後述するように、各カメラ内の各画像について前景被写体マスクを生成するために使用することができる。この辞書のキーは、関節IDとカメラIDの組み合わせである。辞書内の値は、対象のカメラの視野に再投影された関節の2D座標である。

20

30

#### 【0155】

シーン処理1415は、ある時点での実空間内のすべての被写体のリストを含む出力1457を生成する。リストは、被写体毎にキー値辞書を含む。キーは、被写体の固有識別子であり、値は、キーをフレーム番号とし、値をカメラ?被写体関節キー値辞書とする別のキー値辞書である。カメラ-被写体関節キー値辞書は、キーがカメラ識別子であり、値が関節のリストである被写体毎の辞書である。

[被写体毎に在庫商品を識別し追跡するための画像分析]

#### 【0156】

図15～図25を参照して、実空間のエリア内の被写体による在庫商品を置くこと及び取ることを追跡するシステム及び様々な実施態様について説明する。システム及び処理は、一実施態様によるシステムのアーキテクチャレベル概略図である図15Aを参照して説明される。図15Aはアーキテクチャ図であるため、説明の明確性を向上させるために、特定の詳細は省略される。

40

[マルチCNNパイプラインのアーキテクチャ]

#### 【0157】

図15Aは、カメラ114から受信した画像フレームを処理して、実空間内の各被写体についてショッピングカート・データ構造を生成する畳み込みニューラル・ネットワークのパイプライン（マルチCNNパイプラインとも呼ばれる）の高レベルアーキテクチャである。本明細書に記載のシステムは、多関節被写体を識別し追跡するための、上述のカメラ毎の画像認識エンジンを含む。個人毎に1つの「関節」のみが認識され追跡される例、

50

または空間及び時間にわたる他の特徴または他のタイプの画像データが、処理されている実空間内の被写体を認識し追跡するために利用される例を含む、代替の画像認識エンジンを使用することができる。

【 0 1 5 8 】

マルチCNNパイプラインは、カメラ毎に並列に作動し、各カメラからの画像を、カメラ毎に循環バッファ1502を介して画像認識エンジン112a~112nに移動させる。一実施形態では、システムが第1の画像プロセッサ・サブシステム2602、第2の画像プロセッサ・サブシステム2604、及び第3の画像プロセッサ・サブシステム2606の3つのサブシステムから構成される。一実施形態では、第1の画像プロセッサ・サブシステム2602が、畳み込みニューラル・ネットワーク(CNN)として実装され、関節CNN112a~112nと呼ばれる画像認識エンジン112a~112nを含む。図1に関連して説明したように、カメラ114は互いに時間的に同期させることができ、その結果、画像は同時に、または時間的に近く、かつ同じ画像キャプチャレートで取得される。同時にまたは時間的に近い実空間のエリアをカバーする全てのカメラにおいて取得された画像は、同期された画像が実空間において固定された位置を有する被写体のある時点での様々な光景を表すものとして処理エンジンにおいて識別されることができるという意味で同期される。

10

【 0 1 5 9 】

一実施形態では、カメラ114がショッピングストア(スーパーマーケットなど)に設置され、重なり合う視野を有するカメラのセット(2つ以上)が各通路の上に配置されて、店舗内の実空間の画像を取得する。実空間にはN台のカメラがあるが、簡略化のために、図17Aではカメラ(i)として1台のカメラしか示されておらず、iの値は1からNまでの範囲である。各カメラは、それぞれの視野に対応する実空間の画像シーケンスを生成する。

20

【 0 1 6 0 】

一実施形態では、各カメラからの画像シーケンスに対応する画像フレームが每秒30フレーム(fps)のレートでそれぞれの画像認識エンジン112a~112nに送られる。各画像フレームは画像データと共に、タイムスタンプ、カメラの識別情報(「カメラID」と略される)、及びフレーム識別情報(「フレームID」と略される)を有する。画像フレームは、カメラ114毎に循環バッファ1502(リング・バッファとも呼ばれる)に格納される。循環バッファ1502は、それぞれのカメラ114からの連続的にタイムスタンプされた画像フレームのセットを格納する。

30

【 0 1 6 1 】

関節CNNはカメラ当たりの画像フレームのシーケンスを処理し、それぞれの視野に存在する各被写体の18個の異なるタイプの関節を識別する。重なり合う視野を有するカメラに対応する関節CNN112a~112nの出力は、各カメラの2D画像座標から実空間の3D座標に関節の位置をマッピングするために組み合わせられる。jが1~xに等しい被写体(j)毎の関節データ構造800は、実空間における被写体(j)の関節の位置を識別する。被写体データ構造800の詳細を図8に示す。1つの例示的な実施形態では、関節データ構造800が各被写体の関節の2レベルのキー値辞書である。第1のキーはフレーム番号であり、値は、キーがカメラIDであり、値が被写体に割り当てられた関節のリストである第2のキー値辞書である。

40

【 0 1 6 2 】

関節データ構造800によって識別される被写体と、カメラ当たりの画像フレームのシーケンスからの対応する画像フレームとを含むデータセットは、第3の画像プロセッサ・サブシステム2606内の有界ボックス生成器1504への入力として与えられる。第3の画像プロセッサ・サブシステムは、前景画像認識エンジンを更に備える。一実施形態では、前景画像認識エンジンが、例えば、前景における意味的に重要な物体(すなわち、買物客、その手及び在庫商品)が、各カメラからの画像において経時的に、在庫商品を置くこと及び取ることに関連するときに、当該物体を認識する。図15Aに示される例示的な

50

実施態様では、前景画像認識エンジンがWhatCNN1506及びWhenCNN1508として実装される。有界ボックス生成器1504は、データセットを処理して、画像シーケンス内の画像内の識別された被写体の手の画像を含む有界ボックスを指定するロジックを実装する。有界ボックス生成器1504は、それぞれのソース画像フレームに対応する関節データ構造800内の手関節の位置を使用して、カメラ毎に各ソース画像フレーム内の手関節の位置を識別する。被写体データ構造内の関節の座標が3D実空間座標内の関節の位置を示す一実施形態では、有界ボックス生成器が、関節位置を3D実空間座標からそれぞれのソース画像の画像フレーム内の2D座標にマッピングする。

#### 【0163】

有界ボックス生成器1504は、カメラ114毎に循環バッファ内の画像フレーム内の手関節のための有界ボックスを作成する。一実施形態では有界ボックスが、画像フレームの128ピクセル(幅)×128ピクセル(高さ)部分であり、手関節は有界ボックスの中心に位置する。他の実施形態では、有界ボックスのサイズが64ピクセル×64ピクセルまたは32ピクセル×32ピクセルである。カメラからの画像フレーム内のm個の被写体について、最大2m個の手関節、従って2m個の有界ボックスが存在し得る。しかしながら、実際には、他の被写体または他の物体による遮蔽のために、2mより少ない手が画像フレーム内で見える。1つの例示的な実施形態では、被写体の手の位置が肘関節及び手首関節の位置から推測される。例えば、被写体の右手の位置は、右肘の位置(p1として識別される)及び右手首の位置(p2として識別される)を用いて、外挿量 $\times (p2 - p1) + p2$ として外挿される。ここで外挿量は0.4である。別の実施形態では、関節CNN112a~112nが左手画像及び右手画像を使用してトレーニングされる。従って、斯かる実施形態では、関節CNN112a~112nがカメラ当たりの画像フレーム内の手関節の位置を直接識別する。画像フレーム当たりの手の位置は、識別された手関節当たりの有界ボックスを生成するために有界ボックス生成器1504によって使用される。

#### 【0164】

WhatCNN1506は、識別された被写体の手の分類を生成するために、画像内の指定された有界ボックスを処理するようにトレーニングされた畳み込みニューラル・ネットワークである。1つの訓練されたWhatCNN1506は、1つのカメラからの画像フレームを処理する。ショッピングストアの例示的な実施形態では、各画像フレーム内の各手関節について、WhatCNN1506は手関節が空であるかどうかを識別する。WhatCNN1506はまた、手関節内の在庫商品のSKU(在庫管理単位)番号、手関節内の商品を示す信頼値が非SKU商品(すなわち、ショッピングストア在庫に属さない)、及び画像フレーム内の手関節位置の状況を識別する。

#### 【0165】

すべてのカメラ114のWhatCNNモデル1506の出力は、所定の時間帯の間、単一のWhenCNNモデル1508によって処理される。ショッピングストアの例では、WhenCNN1508が被写体の両手について時系列分析を実行して、被写体が棚から店舗在庫商品を取るか、または店舗在庫商品を棚に置くかを識別する。ショッピングカート・データ構造1510(在庫商品のリストを含むログ・データ構造とも呼ばれる)は、被写体に関連するショッピングカート(またはバスケット)内の店舗在庫商品の記録を保持するために、被写体毎に作成される。

#### 【0166】

第2の画像プロセッサ・サブシステム2604は、関節データ構造800によって識別される被写体と、第3の画像プロセッサへの入力として与えられるカメラ当たりの画像フレームのシーケンスからの対応する画像フレームとを含む同じデータセットを受信する。サブシステム2604は、背景画像認識エンジンを含み、背景(すなわち、棚のような在庫陳列構造)における意味的に重要な差異を、例えば、当該差異が、各カメラからの画像において経時的に、在庫商品を置くこと及び取ることに関連するときに認識する。選択ロジック・コンポーネント(図15Aには図示せず)は信頼度スコアを使用して、第2の画像プロセッサまたは第3の画像プロセッサのいずれかからの出力を選択し、ショッピング

10

20

30

40

50

カート・データ構造 1 5 1 0 を生成する。

【 0 1 6 7 】

図 1 5 B は、複数の W h a t C N N モデルの結果を結合し、それを単一の W h e n C N N モデルへの入力として与える調整ロジック・モジュール 1 5 2 2 を示す。上述したように、重なり合う視野を有する 2 つ以上のカメラは、実空間における被写体の画像を取得する。単一の被写体の関節は、それぞれの画像チャネル 1 5 2 0 内の複数のカメラの画像フレーム内に現れることができる。別個の W h a t C N N モデルは、被写体の手（手関節によって表される）における在庫商品の S K U を識別する。調整ロジック・モジュール 1 5 2 2 は、W h a t C N N モデルの出力を結合して、W h e n C N N モデルのための単一の統合入力とする。W h e n C N N モデル 1 5 0 8 は被写体のショッピングカートを生

10

【 0 1 6 8 】

図 1 5 A のマルチ C N N パイプラインを含むシステムの詳細な実施態様は、図 1 6 、 1 7 、 及び 1 8 に提示される。ショッピングストアの例では、システムが、実空間のエリア内の被写体による在庫商品を置くこと及び取ることを追跡する。実空間のエリアは、図 2 及び図 3 に示すように通路に設置された棚に配置された在庫商品を有するショッピングストアである。在庫商品を含む棚は、様々な異なる配置で構成され得ることを理解されたい。例えば、棚はそれらの背面がショッピングストアの側壁に当接し、前面が実空間の開放エリアに面した状態で一列に配置することができる。実空間において重なり合う視野を有する複数のカメラ 1 1 4 は、それらの対応する視野の画像シーケンスを生成する。図 2 及び図 3 に示すように、1 つのカメラの視野は、少なくとも 1 つの他のカメラの視野と重なる。

20

[ 関節 C N N - 被写体の識別と更新 ]

【 0 1 6 9 】

図 1 6 は、関節 C N N 1 1 2 a ~ 1 1 2 n が実空間内の被写体を識別するために実行する処理ステップのフローチャートである。ショッピングストアの例では、被写体は、棚と他のオープンスペースとの間の通路内で店舗内を移動する顧客である。処理はステップ 1 6 0 2 で開始する。上述したように、カメラは、被写体を識別するためにカメラからの画像シーケンスが処理される前に較正されることに留意されたい。カメラ較正の詳細は、上述されている。重なり合う視野を有するカメラ 1 1 4 は、被写体が存在する実空間の画像を取得する（ステップ 1 6 0 4）。一実施形態では、カメラは同期された画像シーケンスを生成するように構成される。各カメラの画像シーケンスは、カメラ毎にそれぞれの循環バッファ 1 5 0 2 に保存される。循環バッファ（リング・バッファとも呼ばれる）は、スライドする時間帯に画像のシーケンスを格納する。一実施形態では、循環バッファが対応するカメラからの画像フレームを格納する（1 1 0）。別の実施形態では、各循環バッファ 1 5 0 2 が 3 . 5 秒間、画像フレームを格納する。他の実施形態では、画像フレーム（または期間）の数が上記の列挙した例示的な値よりも大きくても小さくてもよいことを理解されたい。

30

【 0 1 7 0 】

関節 C N N 1 1 2 a ~ 1 1 2 n は、対応するカメラ 1 1 4 から画像フレームのシーケンスを受信する（ステップ 1 6 0 6）。各関節 C N N は対応するカメラからの画像のバッチを複数の畳み込みネットワーク層を介して処理し、対応するカメラからの画像フレーム内の被写体の関節を識別する。例示的な畳み込みニューラル・ネットワークによる画像のアーキテクチャ及び処理を図 5 に示す。カメラ 1 1 4 は重なり合う視野を有するので、被写体の関節は、2 つ以上の関節 C N N によって識別される。関節 C N N によって生成される関節データ構造 6 0 0 の 2 次元（2 D）座標は、実空間の 3 次元（3 D）座標にマッピングされ、実空間における関節位置を識別する。このマッピングの詳細は、追跡エンジン 1 1 0 が様々な画像シーケンス内の画像に対応する関節データ構造の配列内の要素の座標を、実空間内の座標を有する候補関節に変換する、図 7 の説明において提示される。

40

【 0 1 7 1 】

50

被写体の関節は上述のように、関節をコンステレーションにグループ化するために、2つのカテゴリ（足関節及び非足関節）に編成される。本実施例における左及び右足首関節タイプは、この手順の目的として足関節と考えられる。ステップ1608で、ヒューリスティックスを適用して、候補左足関節及び候補右足関節を候補関節のセットに割り当てて、被写体を作成する。これに続いて、ステップ1610において、新たに識別された被写体が既に実空間に存在するかどうか判定される。存在していない場合には、ステップ1614で、新しい被写体が生成され、存在している場合は、ステップ1612で既存の被写体が更新される。

#### 【0172】

候補関節の銀河からの他の関節は、作成された被写体の関節タイプのいくつかまたはすべてのコンステレーションを構築するために、被写体にリンクすることができる。ステップ1616において、ヒューリスティックスが非足関節に適用され、それらが識別された被写体に割り当てられる。グローバル・メトリック計算器702はグローバル・メトリック値を計算し、非足関節の異なる組み合わせをチェックすることによって値を最小化することを試みる。一実施形態では、グローバル・メトリックは上述のように4つのカテゴリに編成されたヒューリスティックスの合計である。

#### 【0173】

候補関節のセットを識別するロジックは、候補関節のセットを被写体として識別するために、実空間における被写体の関節間の物理的関係に基づくヒューリスティック関数を含む。ステップ1618において、既存の被写体は、対応する非足関節を使用して更新される。処理する画像がまだある場合（ステップ1620）、ステップ1606～1618が繰り返され、なければ、処理はステップ1622で終了する。第1のデータセットは、上述の処理の終わりに生成される。第1のデータセットは、被写体と、実空間における識別された被写体の位置とを識別する。一実施形態では、第1のデータセットが図15Aに関連して、被写体毎の関節データ構造800として上述される。

#### [WhatCNN - 手関節の分類]

#### 【0174】

図17は、実空間で特定された被写体の手の中の在庫商品を特定する処理ステップを示すフローチャートである。ショッピングストアの例では、被写体はショッピングストア内の顧客である。顧客が通路及びオープンスペースを移動すると、顧客は棚に貯蔵された在庫商品を取り上げ、その商品をショッピングカートまたはバスケット内に置く。画像認識エンジンは、複数のカメラから受け取った画像シーケンス内の画像セット内の被写体を識別する。このシステムは、識別された被写体によって在庫商品を取ることと、識別された被写体によって棚に在庫商品を置くことを検出するために、識別された被写体を含む画像シーケンス内の画像のセットを処理するロジックを含む。

#### 【0175】

一実施形態では、画像のセットを処理するロジックが、識別された被写体に対して、識別された被写体の画像の分類を生成するために画像を処理するロジックを含む。分類は、識別された被写体が在庫商品を保持しているかどうかを含む。分類は、棚との相対的な識別された被写体の手の位置を示す第1の近似度分類を含む。分類は、識別された被写体の身体との相対的な識別された被写体の手の位置を示す第2の近似度分類を含む。分類は、識別された被写体に関連するバスケットとの相対的な識別された被写体の手の位置を示す第3の近似度分類を更に含む。最後に、分類は、可能性のある在庫商品の識別子を含む。

#### 【0176】

別の実施形態では、画像のセットを処理するロジックが、識別された被写体について、識別された被写体の画像のセット内の画像内の手を表すデータの有界ボックスを識別するロジックを含む。有界ボックス内のデータは、識別された被写体の有界ボックス内のデータの分類を生成するために処理される。斯かる実施形態では、分類は識別された被写体が在庫商品を保持しているかどうかを含む。分類は、棚との相対的な識別された被写体の手の位置を示す第1の近似度分類を含む。分類は、識別された被写体の身体との相対的な識

10

20

30

40

50

別された被写体の手の位置を示す第2の近似度分類を含む。分類は、識別された被写体に関連するバスケットとの相対的な識別された被写体の手の位置を示す第3の近似度分類を含む。最後に、分類は、可能性のある在庫商品の識別子を含む。

#### 【0177】

処理はステップ1702で開始する。ステップ1704では、画像フレーム内の被写体の手（手関節によって表される）の位置が識別される。有界ボックス生成器1504は、図18で説明したように、関節CNN112a～112nによって生成された第1のデータセット内で識別された関節位置を使用して、各カメラからフレーム当たりの被写体の手の位置を識別する。これに続いて、ステップ1706で、有界ボックス生成器1504は、第1のデータセットを処理して、画像シーケンス内の画像内の識別された多関節被写体の手の画像を含む有界ボックスを指定する。有界ボックス生成器の詳細は、図15Aの議論において上述されている。

#### 【0178】

第2の画像認識エンジンは複数のカメラから画像シーケンスを受け取り、画像内の指定された有界ボックスを処理して、識別された被写体の手の分類を生成する（ステップ1708）。一実施形態では、手の画像に基づいて被写体を分類するために使用される画像認識エンジンのそれぞれは、WhatCNN1506と呼ばれるトレーニングされた畳み込みニューラル・ネットワークを備える。WhatCNNは、図15Aに関連して上述したように、マルチCNNパイプラインに配置される。一実施形態では、WhatCNNへの入力が多次元配列 $B \times W \times H \times C$ （ $B \times W \times H \times C$ テンソルとも呼ばれる）である。「B」はWhatCNNによって処理される画像のバッチ内の画像フレームの数を示すバッチサイズであり、「W」及び「H」は有界ボックスの幅及び高さをピクセルで示し、「C」は、チャンネルの数である。一実施形態では、バッチ内に30個の画像があり（ $B = 30$ ）、それで、有界ボックスのサイズは32ピクセル（幅） $\times$  32ピクセル（高さ）である。赤、緑、青、前景マスク、前腕マスク、及び上腕マスクをそれぞれ表す6つのチャンネルが存在し得る。前景マスク、前腕マスク、及び上腕マスクは、この例ではWhatCNNのための追加の任意的な入力データソースであり、CNNは、これをRGB画像データ内の情報を分類する処理に含めることができる。前景マスクは、例えば、ガウス・アルゴリズムの混合を使用して生成することができる。前腕マスクは、関節データ構造内の情報を使用して生成される状況を提供する、手首と肘との間の線とすることができる。同様に、上腕マスクは、関節データ構造内の情報を使用して生成される肘と肩との間の線とすることができる。他の実施形態では、B、W、H、及びCパラメータの異なる値を使用することができる。例えば、別の実施形態では、有界ボックスのサイズはより大きく、例えば、64ピクセル（幅） $\times$  64ピクセル（高さ）または128ピクセル（幅） $\times$  128ピクセル（高さ）である。

#### 【0179】

各WhatCNN1506は、画像のバッチを処理して、識別された被写体の手の分類を生成する。分類は、識別された主題が在庫商品を保持しているかどうかを含む。分類は、置くこと及び取ることを検出するために使用可能な、棚及び被写体に対する相対的な手の位置を示す1または複数の分類を含む。この例では、第1の近似度分類が棚との相対的な識別された被写体の手の位置を示す。分類は、この例では、識別された被写体の身体との相対的な識別された被写体の手の位置を示す第2の近似度分類を含み、その場合に、被写体は買い物中に在庫商品を保持することができる。この例における分類は、識別された被写体に関連するバスケットとの相対的な識別された被写体の手の位置を示す第3の近似度分類を更に含み、この状況における「バスケット」は、買い物中に在庫商品を保持するために被写体によって使用されるバッグ、バスケット、カート、または他の物体である。最後に、分類は、可能性のある在庫商品の識別子を含む。WhatCNN1506の最終レイヤは、未加工の予測値であるロジットを生成する。ロジットは浮動小数点値として表され、以下に説明するように、分類結果を生成するために更に処理される。一実施形態では、WhatCNNモデルの出力が多次元配列 $B \times L$ （ $B \times L$ テンソルとも呼ばれる）を



含む。「B」はバッチサイズであり、「 $L = N + 5$ 」は画像フレーム当たりのロジット出力数であり、「N」は、ショッピングストアで販売される「N」個の固有在庫商品を表すSKUの数である。

#### 【0180】

1フレーム当たりの出力「L」は、WhatCNN1506からの生の活性化である。ロジット「L」がステップ1710で処理され、在庫商品及び状況を識別する。最初の「N」個のロジットは被写体が「N」個の在庫商品の1つを保持していることの信頼度を表す。ロジット「L」が以下に説明する追加の5つのロジットを含む。第1のロジットは、被写体の手の中にある商品の画像が店舗SKU商品（非SKU商品とも呼ばれる）の1つでないという信頼度を表す。第2のロジットは、被写体が商品を保持しているか否かの信頼度を示す。大きな正の値は、WhatCNNモデルが、被写体が商品を保持しているという高いレベルの信頼度を有することを示す。大きな負の値は、モデルが、被写体が商品を保持していないことを確信していることを示す。第2のロジットのゼロに近い値は、WhatCNNモデルが、被写体が商品を保持しているか否かを予測することに確信がないことを示す。

10

#### 【0181】

次の3つのロジットは第1、第2、及び第3の近似度分類を表す。第1の近似度分類は、棚との相対的な識別された被写体の手の位置を示し、第2の近似度分類は、識別された被写体の身体との相対的な識別された被写体の手の位置を示し、第3の近似度分類は、識別された被写体に関連するバスケットとの相対的な識別された被写体の手の位置を示す。従って、3つのロジットは手の位置の状況を表し、1つのロジットはそれぞれ、手の状況が棚の近く、バスケット（またはショッピングカート）の近く、または被写体の身体の近くにあるという信頼度を示す。一実施形態では、WhatCNNが棚の近く、バスケット（またはショッピングカート）の近く、及び被写体の身体の近くの3つの状況で手の画像を含むトレーニング・データセットを使用してトレーニングされる。別の実施形態では、「近似度分類」パラメータが手の状況を分類するためにシステムによって使用される。斯かる実施形態では、システムが状況を分類するために、棚、バスケット（またはショッピングカート）、及び被写体の身体までの識別された被写体の手の距離を決定する。

20

#### 【0182】

WhatCNNの出力は上述したように、N個のSKUロジット、1個の非SKUロジット、1個の保持ロジット、及び3個の状況ロジットから構成される「L」個のロジットである。SKUロジット（最初のNロジット）及び非SKUロジット（Nロジットに続く最初のロジット）は、softmax関数によって処理される。図5を参照して上述したように、softmax関数は、任意の実数値のK次元ベクトルを、合計で1になる範囲[0, 1]の実数値のK次元ベクトルに変換する。softmax関数は、N+1個の商品にわたる商品の確率分布を計算する。出力値は0と1の間であり、すべての確率の合計は1に等しい。（複数クラス分類のための）softmax関数は、各クラスの確率を返す。最高の確率を有するクラスは、予測クラス（目標クラスとも呼ばれる）である。

30

#### 【0183】

保持ロジットは、シグモイド関数によって処理される。シグモイド関数は入力として実数値をとり、0～1の範囲の出力値を生成する。シグモイド関数の出力は、手が空であるか、商品を保持しているかを識別する。3つの状況ロジットは、手関節位置の状況を識別するためにsoftmax関数によって処理される。ステップ1712では、処理すべき画像がまだあるかどうかチェックされる。処理すべき画像がまだあれば、ステップ1704～1710が繰り返され、なければ、処理はステップ1714で終了する。

40

[WhenCNN - 商品を置くこと及び取ることを識別するための時系列分析]

#### 【0184】

一実施形態では、システムが被写体の前景画像処理に基づいて、識別された被写体による置くこと及び取ることを検出するために、被写体の分類にわたって時系列分析を実行するロジックを実装する。時系列分析は、被写体のジェスチャと、画像シーケンスで表され

50

るジェスチャに関連する在庫商品とを識別する。

#### 【 0 1 8 5 】

マルチCNNパイプラインにおけるWhatCNN1506の出力は、WhenCNN1508への入力として与えられ、WhenCNN1508は、識別された被写体による置くこと及び取ることを検出するために、これらの入力を処理する。最後に、システムは、検出された置くこと及び取ることに応答して、識別された各被写体に対して在庫商品のリストを含むログ・データ構造を生成するロジックを含む。ショッピングストアの例では、ログ・データ構造は、被写体毎のショッピングカート・データ構造1510とも呼ばれる。

#### 【 0 1 8 6 】

図18は、被写体毎にショッピングカート・データ構造を生成するためのロジックを実施する処理を示す。処理はステップ1802で開始する。WhenCNN1508への入力は、ステップ1804で準備される。WhenCNNへのインプットは多次元配列 $B \times C \times T \times Cams$ であり、ここで、 $B$ はバッチサイズであり、 $C$ はチャンネルの数であり、 $T$ は時間帯の間考慮されるフレームの数であり、 $Cams$ はカメラ114の数である。一実施形態では、バッチサイズ「 $B$ 」は64であり、「 $T$ 」の値は110画像フレームまたは3.5秒の時間内の画像フレームの数である。

#### 【 0 1 8 7 】

画像フレーム毎に識別された各被写体に対して、カメラ毎に、手関節毎に10ロジット（両手に対して20ロジット）のリストが生成される。保持ロジット及び状況ロジットは、上述のようにWhatCNN1506によって生成される「 $L$ 」ロジットの一部である。

[

```

    holding,                # 1 logit
    context,                # 3 logits
    slice_dot(sku, log_sku), # 1 logit
    slice_dot(sku, log_other_sku), # 1 logit
    slice_dot(sku, roll(log_sku, -30)), # 1 logit
    slice_dot(sku, roll(log_sku, 30)), # 1 logit
    slice_dot(sku, roll(log_other_sku, -30)), # 1 logit
    slice_dot(sku, roll(log_other_sku, 30)) # 1 logit

```

]

#### 【 0 1 8 8 】

上記のデータ構造は、画像フレーム内の手ごとに生成され、同じ被写体の他方の手に関するデータも含む。例えば、データが被写体の左手関節に対するものである場合、右手に対する対応する値は「他の」ロジットとして含まれる。5番目のロジット（ $\log\_sku$ と呼ばれる上記リストの項目番号3）は、上述の「 $L$ 」ロジットにおけるSKUロジットのログである。6番目のロジットが他の手に対するSKUロジットのログである。「 $roll$ 」関数が現在のフレームの前後で同じ情報を生成する。例えば、第7のロジット（ $roll(\log\_sku, -30)$ と呼ばれる）は、現在のフレームより30フレーム早いSKUロジットのログである。8番目のロジットは手のSKUロジットのログであり、現在のフレームより30フレーム遅い。リスト内の第9及び第10のデータ値は、現在のフレームよりも30フレーム前及び30フレーム後の他方の手についての類似データである。他方の手についての同様のデータ構造も生成され、その結果、カメラ当たり画像フレーム当たり被写体当たり合計20ロジットとなる。従って、WhenCNNへの入力におけるチャンネル数は20である（すなわち、多次元配列 $B \times C \times T \times Cams$ において $C = 20$ ）。

#### 【 0 1 8 9 】

各カメラからの画像フレームのバッチ（例えば、 $B = 64$ ）内のすべての画像フレームについて、画像フレーム内で識別される、被写体当たり20個の手ロジットの同様のデータ構造が生成される。時間帯（ $T = 3.5$ 秒または110画像フレーム）を使用して、被写体の手関節に対して画像フレームのシーケンス内の前方及び後方画像フレームを探索する

10

20

30

40

50

。ステップ1806では、フレーム当たりの被写体当たり20個の手ロジットがマルチCNNパイプラインから統合される。一実施形態では、画像フレームのバッチ(64)が、前方及び後方探索のための追加の画像フレームを両側に有する、画像フレーム110のより大きなウィンドウの中央に配置された画像フレームのより小さなウィンドウとして想像することができる。WhenCNN1508への入力 $B \times C \times T \times Cams$ は、全てのカメラ114(「Cams」と呼ばれる)からの画像フレームのバッチ「B」で識別された被写体の両手に対する20個のロジットから構成される。統合された入力は、WhenCNNモデル1508と呼ばれる単一のトレーニングされた畳み込みニューラル・ネットワークに与えられる。

#### 【0190】

WhenCNNモデルの出力は3つのロジットで構成され、識別された被写体の3つの可能な行為、すなわち棚から在庫商品を取ること、在庫商品を棚に置くこと、及び行為を行わないことに対する信頼度を表す。3つの出力ロジットは、実行される行為を予測するためにsoftmax関数によって処理される。3つの分類ロジットは各被写体に対して一定の間隔で生成され、結果はタイムスタンプと共に個人毎に記憶される。一実施形態では、3つのロジットが被写体当たり20フレーム毎に生成される。斯かる実施形態では、カメラ当たり20画像フレーム毎の間隔で、110画像フレームのウィンドウが現在の画像フレームの周りに形成される。

#### 【0191】

ある期間にわたる被写体当たりのこれら3つのロジットの時系列分析が実行されて(ステップ1808)、真のイベント及びそれらの発生時間に対応するジェスチャが識別される。この目的のために、非最大抑制(NMS)アルゴリズムが使用される。1つのイベント(すなわち、被写体による商品を置くことまたは取ること)がWhenCNN1508によって複数回(同じカメラ及び複数のカメラの両方から)検出されると、NMSは、被写体に対する余分なイベントを除去する。NMSは、2つの主要なタスク、すなわち、余分な検出にペナルティを課す「マッチングロス」と、より良好な検出が手近に存在するかどうかを知るための近隣の「ジョイント処理」とを含む再スコアリング技術である。

#### 【0192】

各被写体に対する取ること及び置くことの真のイベントは、真のイベントを有する画像フレームの前の30画像フレームに対するSKUロジットの平均を計算することによって更に処理される。最後に、最大値の引数(arg maxまたはargmaxと略す)を使用して、最大値を決定する。argmax値によって分類された在庫商品は、棚に置かれたまたは棚から取られた在庫商品を識別するために使用される。在庫商品は、ステップ1810で、それぞれの被写体のSKU(ショッピングカートまたはバスケットとも呼ばれる)のログに追加される。分類データが更にある場合(ステップ1812でチェックされる)、処理ステップ1804~1810年が繰り返される。ある期間にわたって、この処理の結果、各被写体のショッピングカートまたはバスケットが更新される。処理はステップ1814で終了する。

#### [シーン処理とビデオ処理を伴うWhatCNN]

#### 【0193】

図19は、シーン処理1415及びビデオ処理1411からのデータがWhatCNNモデル1506に入力として与えられ、手の画像分類を生成するシステムの実施形態を示す。各ビデオ処理の出力は、別個のWhatCNNモデルに与えられることに留意されたい。シーン処理1415からの出力は関節辞書である。この辞書ではキーは固有関節識別子であり、値は関節が関連付けられる固有被写体識別子である。関節に関連する被写体がない場合、それは辞書に含まれない。各ビデオ処理1411はシーン・処理から関節辞書を受け取り、フレーム番号を返された辞書にマッピングするリング・バッファにそれを格納する。返されたキー値辞書を使用して、ビデオ処理は、識別された被写体に関連付けられた手の近くにある各時点における画像のサブセットを選択する。手の関節の周りの画像フレームのこれらの部分は、領域提案と呼ぶことができる。

10

20

30

40

50

## 【 0 1 9 4 】

ショッピングストアの事例では、領域提案が 1 または複数のカメラからの手の位置のフレームイメージであり、被写体は対応する視野にある。領域提案は、システム内のすべてのカメラによって生成される。これには、空の手だけでなく、ショッピングストア在庫商品及びショッピングストア在庫に属さない商品を持ち運ぶ手も含まれる。ビデオ処理は、時点毎に手の関節を含む画像フレームの部分を選択する。前景マスクの同様のスライスが生成される。上記（手関節の画像部分、前景マスク）を関節辞書（各手関節が属する被写体を示す）に連結して多次元配列を作成する。ビデオ処理からのこの出力は、WhatCNNモデルへの入力として与えられる。

## 【 0 1 9 5 】

WhatCNNモデルの分類結果は、領域提案データ構造（ビデオ処理によって生成される）に格納される。ある時点での全ての領域は、その後、シーン処理への入力として戻される。シーン処理は結果をキー値辞書に格納する。但し、キーは被写体識別子であり、値はキー値辞書であり、但し、キーはカメラ識別子であり、値は領域のロジットである。次に、この集約されたデータ構造は、フレーム番号を時点毎に集約された構造にマッピングするリング・バッファに格納される。

[ シーン処理とビデオ処理を伴うWhenCNN ]

## 【 0 1 9 6 】

図 20 は、WhenCNN 1508 が、図 19 で説明したように、ビデオ処理毎にWhatCNNモデルによって実行される手画像分類に続くシーン処理から出力を受け取るシステムの実施形態を示す。ある期間、例えば、1 秒間の領域提案データ構造が、シーン処理への入力として与えられる。カメラが毎秒 30 フレームの速度で画像を撮影している一実施形態では、入力が 30 の期間と、対応する領域提案とを含む。シーン処理は、30 個の領域提案（手当たり）を、在庫商品SKUを表す単一の整数に縮小する。シーン処理の出力は、キーが被写体識別子であり、値がSKU整数であるキー値辞書である。

## 【 0 1 9 7 】

WhenCNNモデル 1508 は、時系列分析を実行して、この辞書の経時変化を判定する。この結果、棚から取り出され、ショッピングストアの棚に置かれた商品が識別される。WhenCNNモデルの出力は、キーが被写体識別子であり、値がWhenCNNによって生成されたロジットであるキー値辞書である。一実施形態では、1 組のヒューリスティックス 2002 を使用して、被写体毎のショッピングカート・データ構造 1510 を決定する。ヒューリスティックスは、WhenCNNの出力、それぞれの関節データ構造によって示される被写体の関節位置、及びプラノグラムに適用される。プラノグラムは、棚上の在庫商品の予め計算されたマップである。ヒューリスティックス 2002 は、在庫商品が棚に置かれているか棚から取られているか、在庫商品がショッピングカート（またはバスケット）に置かれているか、またはショッピングカート（またはバスケット）から取られているか、または在庫商品が識別された被写体の身体に近いかを、取ることまたは置くことの夫々に対して判定する。

[ What - CNNモデルのアーキテクチャ例 ]

## 【 0 1 9 8 】

図 21 は、WhatCNNモデル 1506 の例示的なアーキテクチャを示す。この例示的なアーキテクチャでは、合計 26 の畳み込み層がある。それぞれの幅（ピクセル単位）、高さ（ピクセル単位）、及びチャネル数に関する異なる層の次元も提示される。第 1 の畳み込み層 2113 は入力 2111 を受け取り、64 ピクセルの幅、64 ピクセルの高さ、及び 64 チャネル（ $64 \times 64 \times 64$  と記載）を有する。WhatCNNへの入力の詳細は、上述されている。矢印の方向は、1 つの層から次の層へのデータの流れを示す。第 2 の畳み込み層 2115 は、 $32 \times 32 \times 64$  の次元を有する。続いて第 2 の層があり、それぞれ  $32 \times 32 \times 64$  の次元を有する 8 つの畳み込み層（ボックス 2117 に示される）がある。2 つの層 2119 及び 2121 のみが、例示の目的のためにボックス 2117 に示されている。この後に、 $16 \times 16 \times 128$  の次元の別の 8 つの畳み込み層 212

10

20

30

40

50

3が続く。このような2つの畳み込み層2 1 2 5及び2 1 2 7が図2 1に示されている。最後に、最後の8つの畳み込み層2 1 2 9は、それぞれ $8 \times 8 \times 256$ の次元数を有する。2つの畳み込み層2 1 3 1及び2 1 3 3が、説明のためにボックス2 1 2 9に示されている。

#### 【0 1 9 9】

N + 5個の出力を生成する最後の畳み込み層2 1 3 3からの256個の入力を有する1つの全結合層2 1 3 5がある。上述したように、「N」は、ショッピングストアで販売される「N」個の固有在庫商品を表すSKUの数である。5つの追加のロジットは、画像内の商品が非SKU商品であるという信頼性を表す第1のロジットと、被写体が商品を持しているかどうかの信頼度を表す第2のロジットとを含む。次の3つのロジットは上述したように、第1、第2及び第3の近似度分類を表す。WhatCNNの最終出力は2 1 3 7に示されている。例示的なアーキテクチャは、バッチ正規化(BN)を使用する。畳み込みニューラル・ネットワーク(CNN)における各層の分布はトレーニング中に変化し、層別に変化する。これは、最適化アルゴリズムの収束速度を低下させる。バッチ正規化(Ioffe及びSzegedyの2015年の論文)は、この問題を克服するための技術である。ReLU(正規化線形ユニット)活性化は、softmaxが使用される最終出力を除いて、各層の非線形性のために使用される。

#### 【0 2 0 0】

図2 2、図2 3、及び図2 4は、WhatCNN1506の実施態様の様々な部分の図式的な視覚化である。これらの図は、TensorBoard(商標)によって生成されたWhatCNNモデルの図式的視覚化から編集された図である。TensorBoard(商標)は、深層学習モデル、例えば、畳み込みニューラル・ネットワークを検査し、理解するための一連の視覚化ツールである。

#### 【0 2 0 1】

図2 2は、片手(「片手」モデル2 2 1 0)を検出する畳み込みニューラル・ネットワークモデルの高レベルアーキテクチャを示す。WhatCNNモデル1506は、それぞれ左手及び右手を検出するための2つの当該畳み込みニューラル・ネットワークを備える。図示の実施形態では、アーキテクチャが、ブロック0 2 2 1 6、ブロック1 2 2 1 8、ブロック2 2 2 2 0、及びブロック3 2 2 2 2と呼ばれる4つのブロックを含む。ブロックはより高レベルの抽象化であり、畳み込み層を表す複数のノードを含む。ブロックは1つのブロックからの出力が次のブロックに入力されるように、下から上への順序で配置される。このアーキテクチャは、プーリング層2 2 1 4及び畳み込み層2 2 1 2も含む。ブロック間では、異なる非線形性を使用することができる。図示の実施形態では、上述のようにReLU非線形性が使用される。

#### 【0 2 0 2】

図示の実施形態では、片手モデル2 2 1 0への入力がWhatCNN1506の説明において上記で定義した $B \times W \times H \times C$ テンソルである。「B」はバッチサイズであり、「W」及び「H」は入力画像の幅及び高さを示し、「C」はチャンネル数である。片手モデル2 2 1 0の出力は、第2の片手モデルと結合され、全結合ネットワークに転送される。

#### 【0 2 0 3】

トレーニング中、片手モデル2 2 1 0の出力は、グラウンドトゥールズと比較される。出力とグラウンドトゥールズとの間で計算された予測誤差は、畳み込み層の重みを更新するために使用される。図示の実施形態では、WhatCNN1506をトレーニングするために確率的勾配降下法(SGD)が使用される。

#### 【0 2 0 4】

図2 3は、図2 2の片手畳み込みニューラル・ネットワークモデルのブロック0 2 2 1 6の更なる詳細を示す。これは、ボックス2 3 1 0内のconv0、conv1 2 3 1 8、conv2 2 3 2 0、及びconv3 2 3 2 2とラベル付けされた4つの畳み込み層を含む。畳み込み層conv0の更なる詳細は、ボックス2 3 1 0内に提示されている。入力は、畳み込み層2 3 1 2によって処理される。畳み込み層の出力は、バッチ正規化層

10

20

30

40

50

2 3 1 4によって処理される。ReLU非線形性2 3 1 6は、バッチ正規化層2 3 1 4の出力に適用される。畳み込み層conv 0の出力は、次の層conv 1 2 3 1 8に転送される。最終的な畳み込み層conv 3の出力は、加算演算2 3 2 4を介して処理される。この演算は、層conv 3 2 3 2 2からの出力を、スキップ接続2 3 2 6を介して到来する修正されていない入力に合計する。Heらの論文「深層残余ネットワークにおけるアイデンティティ・マッピング」(2016年7月25日に<https://arxiv.org/pdf/1603.05027.pdf>で公開)では、順方向信号及び逆方向信号が1つのブロックから任意の他のブロックに直接的に伝播することができることが示されている。信号は、畳み込みニューラル・ネットワークを通して変化せずに伝播する。この技術は、深い畳み込みニューラル・ネットワークのトレーニング及び試験性能を改善する。

10

#### 【0205】

図21で説明したように、WhatCNNの畳み込み層の出力は、全結合層によって処理される。2つの片手モデル2 2 1 0の出力は結合され、入力として全結合層に転送される。図24は、全結合層(FC)2 4 1 0の例示的な実施態様である。FC層への入力は、再整形演算子2 4 1 2によって処理される。再整形演算子は、テンソルを次の層2 4 2 0に転送する前にテンソルの形状を変更する。再整形は、畳み込み層からの出力を平坦化すること、すなわち、多次元行列からの出力を1次元行列またはベクトルに再整形することを含む。再構築演算子2 4 1 2の出力はMatMul 2 4 2 2と表示される行列乗算演算子にパスされ、MatMul 2 4 2 2からの出力はxw\_\_plus\_\_b 2 4 2 4と表示される行列加算演算子に転送される。入力「x」毎に、演算子2 4 2 4は入力に行列「w」及びベクトル「b」を乗算して出力を生成する。「w」が入力「x」に関連するトレーニング可能なパラメータであり、「b」がバイアスまたはインターセプトと呼ばれる別のトレーニング可能なパラメータである。全結合層2 4 1 0からの出力2 4 2 6が、WhatCNN 1 5 0 6の説明において上述したように、B x Lテンソルである。「B」はバッチサイズであり、「L = N + 5」は画像フレーム当たりの出力ロジット数である。「N」がショッピングストアで販売するための「N」個の固有在庫商品を表すSKUの数である。

20

#### [WhatCNNモデルのトレーニング]

#### 【0206】

様々な状況における空の手のみならず、様々な状況における様々な在庫商品を保持する手の画像のトレーニング・データセットが作成される。これを達成するために、人間の行為者が、試験環境の様々な場所で、多数の異なる方法で、各々の固有のSKU在庫商品を保持する。彼らの手の状況は、行為者の身体に近いこと、店舗の棚に近いこと、及び行為者のショッピングカートまたはバスケットに近いことに及ぶ。行為者は、空の手でも上記の行為を行う。この手順は、左手及び右手の両方について完了する。複数の行為者が、実際のショッピングストアで起こる自然な閉塞をシミュレーションするために、同じテスト環境でこれらの行為を同時に実行する。

30

#### 【0207】

カメラ1 1 4は、上記行為を実行する行為者の画像を撮影する。一実施形態では、20台のカメラがこの処理で使用される。関節CNN 1 1 2 a ~ 1 1 2 n及び追跡エンジン1 1 0は、関節を識別するために画像を処理する。有界ボックス生成器1 5 0 4は、プロダクションまたは推論に類似した手領域の有界ボックスを作成する。WhatCNN 1 5 0 6を介してこれらの手領域を分類する代わりに、画像は記憶ディスクに保存される。保存された画像は、精査され、ラベル付けされる。画像には、在庫商品SKU、状況、及び手が何かを保持しているか否かという3つのラベルが割り当てられる。この処理は、多数の画像(数百万枚までの画像)に対して行われる。

40

#### 【0208】

画像ファイルは、データ収集シーンに従って編成される。画像ファイルの命名規則は、画像のコンテンツ及び状況を識別する。図25は、一実施形態における画像ファイル名を示す図である。数表示2 5 0 2によって参照されるファイル名の第1の部分は、データ収

50

集シーンを識別し、画像のタイムスタンプも含む。ファイル名の第2の部分2504は、ソースカメラを識別する。図25に示す例では、「カメラ4」で撮影されている。ファイル名の第3の部分2506は、ソースカメラからのフレーム番号を識別する。図示の例では、ファイル名が、それがカメラ4からの94,600番目の画像フレームであることを示す。ファイル名の第4の部分2508は、この手領域画像が取得されるソース画像フレーム内のx座標領域及びy座標領域の範囲を識別する。図示の例では、領域がピクセル117から370までのx座標値と、ピクセル370から498までのy座標値との間で定義される。ファイル名の第5の部分2510は、シーン内の行為者の個人IDを識別する。図示の例では、シーン内の人物がID「3」を有する。最後に、ファイル名の第6の部分2512は、画像内で識別された在庫商品のSKU数(商品=68)を識別する。

10

#### 【0209】

WhatCNN1506のトレーニング・モードでは、順方向パスのみが実行されるプロダクション・モードとは対照的に、順方向パスと逆方向伝播が実行される。トレーニング中、WhatCNNは、順方向パスにおいて識別された被写体の手の分類を生成する。WhatCNNの出力は、グラントゥルースと比較される。逆伝播では、1または複数のコスト関数の勾配が計算される。次いで、勾配は、畳み込みニューラル・ネットワーク(CNN)及び全結合(FC)ニューラル・ネットワークに伝播され、その結果、予測誤差が低減され、出力がグラントゥルースに近づく。一実施形態では、WhatCNN1506をトレーニングするために、確率的勾配降下法(SGD)が使用される。

#### 【0210】

20

一実施形態では、64個の画像がトレーニング・データからランダムに選択され、増強される。画像増強の目的はトレーニング・データを多様化し、モデルの性能を向上させることである。画像増強は、画像のランダムフリッピング、ランダム回転、ランダム色相シフト、ランダムガウスノイズ、ランダムコントラスト変化、及びランダムクロッピングを含む。増強の量はハイパー・パラメータであり、ハイパー・パラメータ探索によって調整される。増強された画像は、トレーニング中にWhatCNN1506によって分類される。分類はグラントゥルースと比較され、WhatCNN1506の係数または重みは、勾配損失関数を計算し、勾配に学習レートを乗算することによって更新される。上記処理は、エポックを形成するために何度も(例えば、約1000回)繰り返される。50から200のエポックが実行される。各エポックの間、学習速度は、余弦アニーリングスケジュールに従ってわずかに減少する。

30

#### [WhenCNNモデルのトレーニング]

#### 【0211】

WhenCNN1508のトレーニングは、予測誤差を低減するために逆伝播を使用する、上述のWhatCNN1506のトレーニングと同様である。行為者は、トレーニング環境において様々な行為を実行する。例示的な実施形態では、トレーニングは、在庫商品が貯蔵された棚を有するショッピングストアで実行される。行為者によって実行される行為の例には、棚から在庫商品を取り出すこと、在庫商品を棚に置いて戻すこと、在庫商品をショッピングカート(またはバスケット)に置くこと、ショッピングカートから在庫商品を取り戻すこと、商品を左手と右手との間で交換すること、在庫商品を行為者のスックに入れることが含まれる。スックとは、左手及び右手以外の在庫商品を保持することができる行為者の身体上の位置を指す。スックの幾つかの例は、在庫商品を、前腕と上腕との間で挟み込むこと、前腕と胸との間で挟み込むこと、首と肩との間で挟み込むことが含まれる。

40

#### 【0212】

カメラ114は、トレーニング中に上述した全ての行為のビデオを記録する。ビデオは精査され、全ての画像フレームはタイムスタンプ及び実行された行為を示すラベルが付される。これらのラベルは、それぞれの画像フレームに対する行為ラベルと呼ばれる。画像フレームはプロダクションまたは推論のために、上述したように、WhatCNN1506までのマルチCNNパイプラインを介して処理される。次に、関連付けられた行為ラベ

50

ルに沿ったWhatCNNの出力を、グラウンドトゥースとして作用する行為ラベルとともに使用して、WhenCNN1508をトレーニングする。WhatCNN1506のトレーニングについて上述したように、余弦アニーリングスケジュールを有する確率的勾配降下法(SGD)がトレーニングのために使用される。

#### 【0213】

画像増強(WhatCNNのトレーニングに使用される)に加えて、時間増強は、WhenCNNのトレーニング中の画像フレームにも適用される。幾つかの例は、ミラーリング、ガウスノイズの追加、左手及び右手に関連するロジットの交換、時間の短縮、画像フレームをドロップすることによる時系列の短縮、フレームを複製することによる時系列の延長、及びWhenCNNのための入力を生成する基礎となるモデルにおけるスポッティ性をシミュレーションするための時系列におけるデータポイントのドロップを含む。ミラーリングは時系列及びそれぞれのラベルを反転させることを含み、例えば、置く行為は、反転されると取る行為になる。

[背景画像処理を使用した在庫イベントの予測]

#### 【0214】

図26～図28Bを用いて、実空間のエリアにおける被写体による変化を追跡するシステム及び各種実施態様について説明する。

[システム・アーキテクチャ]

#### 【0215】

図26は、本実施態様に係るシステムの高レベル概略図である。図26はアーキテクチャ図であるため、説明の明確性を向上させるために、特定の詳細は省略されている。

#### 【0216】

図26に示すシステムは、複数のカメラ114から画像フレームを受信する。上述のように、一実施形態では、カメラ114が、画像が同時に、または時間的に近く、かつ同じ画像キャプチャレートで取得されるように、互いに時間的に同期させることができる。同時にまたは時間的に近い実空間のエリアをカバーする全てのカメラにおいて取得された画像は、同期された画像が実空間において固定された位置を有する被写体のある時点での様々な光景を表すものとして処理エンジンにおいて識別されることができるという意味で同期される。

#### 【0217】

一実施形態では、カメラ114がショッピングストア(スーパーマーケットなど)に設置され、重なり合う視野を有するカメラのセット(2つ以上)が各通路の上に配置されて、店舗内の実空間の画像を取得する。実空間には「n」台のカメラがある。各カメラは、それぞれの視野に対応する実空間の画像シーケンスを生成する。

#### 【0218】

被写体識別サブシステム2602(第1の画像プロセッサとも呼ばれる)は、カメラ114から受け取った画像フレームを処理して、実空間内の被写体を識別し追跡する。第1の画像プロセッサは、被写体画像認識エンジンを含む。被写体画像認識エンジンは、複数のカメラから対応する画像シーケンスを受け取り、画像を処理して、対応する画像シーケンス内の画像に表される被写体を識別する。一実施形態では、システムが多関節被写体を識別し追跡するための、上述したようなカメラ毎の画像認識エンジンを含む。個人毎に1つの「関節」のみが認識され追跡される例、または空間及び時間にわたる他の特徴または他のタイプの画像データが、処理されている実空間内の被写体を認識し追跡するために利用される例を含む、代替の画像認識エンジンを使用することができる。

#### 【0219】

「意味的差分抽出」サブシステム2604(第2の画像プロセッサとも呼ばれる)は背景画像認識エンジンを含み、複数のカメラから対応する画像シーケンスを受信し、例えば、背景(すなわち棚のような在庫陳列構造)内の意味的に重要な差異が、各カメラからの画像において経時的に在庫商品を置くこと及び取ることに関連するときに、当該差異を認識する。第2の画像プロセッサは、被写体識別サブシステム2602の出力と、カメラ1

10

20

30

40

50



14からの画像フレームとを入力として受け取る。第2の画像プロセッサは、前景内の識別された被写体をマスクして、マスクされた画像を生成する。マスクされた画像は、前景被写体に対応する有界ボックスを背景画像データに置き換えることによって生成される。これに続いて、背景画像認識エンジンはマスクされた画像を処理して、対応する画像シーケンス内の画像に表される背景変化を識別し且つ分類する。一実施形態では、背景画像認識エンジンが畳み込みニューラル・ネットワークを含む。

#### 【0220】

最後に、第2の画像プロセッサは、識別された背景変化を処理して、識別された被写体による在庫商品を取ることに、識別された被写体による在庫陳列構造上に在庫商品を置くことの第1の検出セットを行う。第1の検出セットは、在庫商品を置くこと及び取ることの背景検出とも呼ばれる。ショッピングストアの例では、第1の検出が店舗の顧客または従業員によって棚から取られた、または棚に置かれた在庫商品を識別する。意味的差分抽出サブシステムは、識別された背景変化を識別された被写体に関連付けるロジックを含む。

#### 【0221】

領域提案サブシステム2606（第3の画像プロセッサとも呼ばれる）は前景画像認識エンジンを含み、複数のカメラ114から対応する画像シーケンスを受信し、例えば、前景（すなわち、買物客、買物客の手、及び在庫商品）内の意味的に重要な物体が、各カメラからの画像において経時的に、在庫商品を置くこと及び取ることに関連するときに、当該物体を認識する。サブシステム2606はまた、被写体識別サブシステム2602の出力を受信する。第3の画像プロセッサは、カメラ114からの画像シーケンスを処理して、対応する画像シーケンス内の画像に表される前景変化を識別し且つ分類する。第3の画像プロセッサは、識別された前景変化を処理して、識別された被写体による在庫商品を取ることに、識別された被写体による在庫陳列構造上に在庫商品を置くことの第2の検出セットを行う。第2の検出セットは、在庫商品を置くこと及び取ることの前景検出とも呼ばれる。ショッピングストアの例では、第2の検出セットが、在庫商品を取ることに、店舗の顧客及び従業員による在庫陳列構造上に在庫商品を置くこととを識別する。

#### 【0222】

図26に記載されるシステムは、第1及び第2の検出セットを処理して、識別された被写体についての在庫商品のリストを含むログ・データ構造を生成するための選択ロジック・コンポーネント2608を含む。実空間内の置くこと及び取ることのために、選択ロジック2608は、意味的差分抽出サブシステム2604または領域提案サブシステム2606の何れかからの出力を選択する。一実施形態では、選択ロジック2608が、第1の検出セットについて意味的差分抽出サブシステムによって生成された信頼度スコアと、第2の検出セットについて領域提案サブシステムによって生成された信頼度スコアとを使用して、選択を行う。特定の検出に対するより高い信頼度スコアを有するサブシステムの出力が選択され、識別された前景被写体に関連付けられた在庫商品のリストを含むログ・データ構造1510（ショッピングカート・データ構造とも呼ばれる）を生成するために使用される。

#### 【サブシステム・コンポーネント】

#### 【0223】

図27は、実空間のエリア内の被写体による変化を追跡するためのシステムを実施するサブシステム・コンポーネントを示す。システムは、実空間における対応する視野のそれぞれの画像シーケンスを生成する複数のカメラ114を備える。各カメラの視野は上述したように、複数のカメラのうちの少なくとも1つの他のカメラの視野と重なる。一実施形態では、複数のカメラ114によって生成された画像に対応する画像フレームのシーケンスがカメラ114毎に循環バッファ1502（リング・バッファとも呼ばれる）に格納される。各画像フレームは、画像データと共に、タイムスタンプ、カメラの識別情報（「カメラID」と略される）、及びフレーム識別情報（「フレームID」と略される）を有する。循環バッファ1502は、それぞれのカメラ114からの連続的にタイムスタンプされた画像フレームのセットを格納する。一実施形態では、カメラ114が同期された画像

シーケンスを生成するように構成される。

【 0 2 2 4 】

1つの好ましい実施態様では、同じカメラ及び同じ画像シーケンスが前景及び背景画像プロセッサの両方によって使用される。その結果、同じ入力データを用いて、在庫商品を置くこと及び取ることの冗長な検出が行われ、結果として得られるデータにおいて高い信頼度と高い精度を可能にする。

【 0 2 2 5 】

被写体識別サブシステム 2 6 0 2 (第1の画像プロセッサとも呼ばれる)は、複数のカメラ 1 1 4 から対応する画像シーケンスを受信する被写体画像認識エンジンを含む。被写体画像認識エンジンは、画像を処理して、対応する画像シーケンス内の画像に表される被写体を識別する。一実施形態では、被写体画像認識エンジンが関節 CNN 1 1 2 a ~ 1 1 2 n と呼ばれる畳み込みニューラル・ネットワーク (CNN) として実装される。重なり合う視野を有するカメラに対応する関節 CNN 1 1 2 a ~ 1 1 2 n の出力は、各カメラの 2 D 画像座標から実空間の 3 D 座標に関節の位置をマッピングするために組み合わせられる。j が 1 ~ x に等しい被写体 (j) 毎の関節データ構造 8 0 0 は、各画像について実空間及び 2 D 空間における被写体 (j) の関節の位置を識別する。被写体データ構造 8 0 0 の幾つかの詳細を図 8 に示す。

【 0 2 2 6 】

背景画像格納装置 2 7 0 4 は、意味的差分抽出サブシステム 2 6 0 4 において、カメラ 1 1 4 からの対応する画像シーケンスのためのマスクされた画像 (前景被写体がマスクによって除去された背景画像とも呼ばれる) を記憶する。背景画像格納装置 2 7 0 4 は、背景バッファとも呼ばれる。一実施形態では、マスクされた画像のサイズが循環バッファ 1 5 0 2 内の画像フレームのサイズと同じである。一実施形態では、マスクされた画像が、カメラ当たりの画像フレームのシーケンス内の各画像フレームに対応する背景画像格納装置 2 7 0 4 に格納される。

【 0 2 2 7 】

意味的差分抽出サブシステム 2 6 0 4 (または第2の画像プロセッサ) は、カメラからの対応する画像シーケンス内の画像に表される前景被写体のマスクを生成するマスク生成器 2 7 2 4 を含む。一実施形態では、1つのマスク生成器がカメラ毎に画像シーケンスを処理する。ショッピングストアの例では、前景被写体が、販売用の商品を含む背景棚の前の顧客または店舗の従業員である。

【 0 2 2 8 】

一実施形態では、関節データ構造 8 0 0 及び循環バッファ 1 5 0 2 からの画像フレームがマスク生成器 2 7 2 4 への入力として与えられる。関節データ構造は、各画像フレームにおける前景被写体の位置を識別する。マスク生成器 2 7 2 4 は、画像フレーム内で識別された前景被写体毎に有界ボックスを生成する。斯かる実施形態では、マスク生成器 2 7 2 4 が、2 D 画像フレーム内の関節位置の x 座標及び y 座標の値を使用して、有界ボックスの4つの境界を決定する。x の最小値 (被写体の関節のすべての x 値からの) は、被写体の有界ボックスの左側垂直境界を定義する。y の最小値 (被写体に対する関節の全ての y 値からの) は、有界ボックスの下側水平境界を定義する。同様に、x 座標及び y 座標の最大値は、有界ボックスの右側垂直境界及び上側水平境界を識別する。第2の実施形態では、マスク生成器 2 7 2 4 が畳み込みニューラル・ネットワークベースの人物検出及び位置特定アルゴリズムを使用して、前景被写体の有界ボックスを生成する。斯かる実施形態では、マスク生成器 2 7 2 4 が前景被写体のための有界ボックスを生成するために関節データ構造 8 0 0 を使用しない。

【 0 2 2 9 】

意味的差分抽出サブシステム 2 6 0 4 (または第2の画像プロセッサ) は、識別された被写体を表す前景画像データを、対応する画像シーケンスに対する背景画像からの背景画像データで置き換えるための、画像シーケンス内の画像を処理するマスクロジックを含み、処理用の新しい背景画像となるマスクされた画像を提供する。循環バッファがカメラ 1

10

20

30

40

50

14 から画像フレームを受け取ると、マスクロジックは、画像マスクによって定義された前景画像データを背景画像データで置き換えるために、画像シーケンス内の画像を処理する。背景画像データは、対応するマスクされた画像を生成するために、対応する画像シーケンスの背景画像から取得される。

#### 【0230】

ショッピングストアの例を考える。最初に時間  $t = 0$  において、店舗内に顧客がいない場合、背景画像格納装置 2704 内の背景画像は、カメラ当たりの画像シーケンス内の対応する画像フレームと同じである。次に、時間  $t = 1$  において、顧客が棚の前を移動して棚内の商品を購入する場合を考える。マスク生成器 2724 は、顧客の有界ボックスを作成し、それをマスクロジック・コンポーネント 2702 に送る。マスクロジック・コンポーネント 2702 は、有界ボックス内の  $t = 1$  における画像フレーム内のピクセルを、 $t = 0$  における背景画像フレーム内の対応するピクセルで置き換える。この結果、循環バッファ 1502 内の  $t = 1$  における画像フレームに対応する  $t = 1$  におけるマスクされた画像が得られる。マスクされた画像は、 $t = 0$  で背景画像フレームからのピクセルによって置き換えられる前景被写体（または顧客）のピクセルを含まない。 $t = 1$  におけるマスクされた画像は、背景画像格納装置 2704 に格納され、対応するカメラからの画像シーケンス内の  $t = 2$  における次の画像フレームに対する背景画像として作用する。

#### 【0231】

一実施形態では、マスクロジック・コンポーネント 2702 が、ピクセルによる平均化または加算などによって、画像シーケンス内の  $N$  個のマスクされた画像のセットを組み合わせ、各カメラのファクタ化画像のシーケンスを生成する。斯かる実施形態では、第2の画像プロセッサが、ファクタ化画像のシーケンスを処理することによって背景変化を識別し且つ分類する。ファクタ化画像は、例えば、カメラ当たりのマスクされた画像シーケンスにおける  $N$  個のマスクされた画像内のピクセルの平均値をとることによって生成することができる。一実施形態では、 $N$  の値がカメラ 114 のフレームレートに等しく、例えば、フレームレートが 30 FPS (フレーム/秒) である場合、 $N$  の値は 30 である。斯かる実施形態では、1 秒の期間に対してマスクされた画像がファクタ化画像を生成するために組み合わせられる。ピクセル値の平均をとることにより、実空間のエリアでのセンサノイズ及び明度変化によるピクセル変動が最小限に抑えられる。

#### 【0232】

第2の画像プロセッサはファクタ化画像のシーケンスを処理することによって、背景変化を識別し且つ分類する。ファクタ化画像のシーケンス内のファクタ化画像は、ビットマスク計算器 2710 によって、同じカメラに対する先行するファクタ化画像と比較される。ファクタ化画像 2706 のペアは、2つのファクタ化画像の対応するピクセルの変化を識別するビットマスクを生成するために、ビットマスク計算器 2710 への入力として与えられる。ビットマスクは、対応するピクセル（現在及び前のファクタ化画像）の RGB（赤、緑及び青チャネル）値間の差が「差閾値」よりも大きいピクセル位置に 1 を有する。差閾値の値は調整可能である。一実施形態では、差閾値の値は 0.1 に設定される。

#### 【0233】

ビットマスクと、カメラ当たりのファクタ化画像のシーケンスからのファクタ化画像のペア（現在及び前）は、背景画像認識エンジンへの入力として与えられる。一実施形態では、背景画像認識エンジンが畳み込みニューラル・ネットワークを含み、変化 CNN 2714a ~ 2714n と呼ばれる。単一の変化 CNN は、カメラ毎にファクタ化画像のシーケンスを処理する。別の実施形態では、対応する画像シーケンスからのマスクされた画像は結合されない。ビットマスクは、マスクされた画像のペアから計算される。この実施形態では、マスクされた画像とビットマスクのペアが次に、変化 CNN への入力として与えられる。

#### 【0234】

この例での変化 CNN モデルへの入力は、ファクタ化画像毎に 3 つの画像チャンネル（赤、緑、青）とビットマスクの 1 つのチャンネルを含む 7 チャンネルから構成されている

10

20

30

40

50

。変化CNNは、複数の畳み込み層と、1または複数の全結合（FC）層とを含む。一実施形態では、変化CNNが、図5に示す関節CNN 112a ~ 112nと同じ数の畳み込み層及びFC層を含む。

#### 【0235】

背景画像認識エンジン（変化CNN 2714a - 2714n）は、ファクタ化画像の変化を識別し且つ分類し、対応する画像シーケンスに対して変化データ構造を生成する。変化データ構造は、識別された背景変化のマスクされた画像内の座標、識別された背景変化の在庫商品被写体の識別子、及び識別された背景変化の分類を含む。変化データ構造における識別された背景変化の分類は、識別された在庫商品が背景画像に対して追加されたか除去されたかを分類する。

10

#### 【0236】

複数の商品が1つまたは複数の被写体によって同時に棚上で取られ、または置かれ得るので、変化CNNは出力位置毎に数「B」の重複有界ボックス予測を生成する。有界ボックス予測はファクタ化画像の変化に対応する。ショッピングストアが固有のSKUによって識別される数「C」の固有の在庫商品を有すると考える。変化CNNは、変化の在庫商品被写体のSKUを予測する。最後に、変化CNNは識別された商品が棚から取られるか、または棚に置かれるかを示す、出力内のすべての位置（ピクセル）についての変化（または在庫イベントタイプ）を識別する。変化CNNからの出力の上記3つの部分は式「 $5 \times B + C + 1$ 」によって記述される。各有界ボックス「B」予測が5つの数字を含むので、「B」は5で乗算される。これらの5つの数字は、有界ボックスの中心の「x」及び「y」座標、有界ボックスの幅及び高さを表す。5番目の数字は有界ボックスの予測のための変化CNNモデルの信頼度スコアを表す。「B」は変化CNNモデルの性能を改善するために調整可能なハイパー・パラメータである。一実施形態では、「B」の値が4に等しい。変化CNNからの出力の幅及び高さ（ピクセル単位）がそれぞれ、W及びHによって表され则认为する。変化CNNの出力は「 $W \times H \times (5 \times B + C + 1)$ 」として表される。有界ボックス出力モデルは、論文「YOLO9000: Better, Faster, Stronger」（2016年12月25日発行）においてRedmon及びFarhadiによって提案された物体検出システムに基づく。この論文は<https://arxiv.org/pdf/1612.08242.pdf>で入手可能である。

20

#### 【0237】

重なり合う視野を有するカメラからの画像シーケンスに対応する変化CNN 2714a ~ 2714nの出力は、調整ロジック・コンポーネント2718によって結合される。調整ロジック・コンポーネントは、重なり合う視野を有するカメラのセットからの変化データ構造を処理して、実空間内での識別された背景変化の位置を確認する。調整ロジック・コンポーネント2718は重なり合う視野を有する複数のカメラから、同じSKU及び同じ在庫イベントタイプ（取るまたは置く）を有する在庫商品を表す有界ボックスを選択する。次いで、選択された有界ボックスは3D実空間における在庫商品の位置を識別するために、上述の三角測量技法を使用して3D実空間において三角測量される。実空間における棚の位置は、3D実空間における在庫商品の三角測量された位置と比較される。誤検知予測は廃棄される。例えば、有界ボックスの三角測量された位置が実空間内の棚の位置にマッピングされない場合、出力は破棄される。棚にマップする3D実空間内の有界ボックスの三角測量された位置は、在庫イベントの真の予測と考えられる。

30

40

#### 【0238】

一実施形態では、第2の画像プロセッサによって生成された変化データ構造における識別された背景変化の分類が、識別された在庫商品が背景画像に対して追加されたか除去されたかを分類する。別の実施形態では、変化データ構造における識別された背景変化の分類が、識別された在庫商品が背景画像に対して追加されたか除去されたかを示し、システムは背景変化を識別された被写体に関連付けるロジックを含む。システムは、識別された被写体による在庫商品を取ることに、識別された被写体による在庫陳列構造上に在庫商品を置くことの検出を行う。

50

## 【 0 2 3 9 】

ログ生成器 2 7 2 0 は、変化の真の予測によって識別された変化を、変化の位置付近の識別された被写体に関連付けるためのロジックを実施する。関節識別エンジンを利用して被写体を識別する実施形態では、ログ生成器 2 7 2 0 が関節データ構造 8 0 0 を使用して 3 D 実空間内の被写体の手関節の位置を決定する。手関節位置が、変化時の変化の位置までの閾値距離内にある被写体が識別される。ログ生成器は、変化を識別された被写体に関連付ける。

## 【 0 2 4 0 】

一実施形態では、上述のように、N 個のマスクされた画像が組み合わされてファクタ化画像が生成され、次いで、ファクタ化画像が変化 C N N への入力として与えられる。N はカメラ 1 1 4 のフレームレート（フレーム / 秒）に等しいと考える。従って、斯かる実施形態では、1 秒の期間中の被写体の手の位置を変化の位置と比較して、変化を識別された被写体に関連付ける。2 つ以上の被写体の手関節位置が変化の位置までの閾値距離内にある場合、被写体との変化の関連付けは、前景画像処理サブシステム 2 6 0 6 の出力に対して保留される。

## 【 0 2 4 1 】

前景画像処理（領域提案）サブシステム 2 6 0 6（第 3 の画像プロセッサとも呼ばれる）は、複数のカメラからの画像シーケンスから画像を受信する前景画像認識エンジンを含む。第 3 の画像プロセッサは、対応する画像シーケンス内の画像に表される前景変化を識別し且つ分類するロジックを含む。領域提案サブシステム 2 6 0 6 は、識別された被写体による在庫商品を取ることと、識別された被写体による在庫陳列構造上に在庫商品を置くことの第 2 の検出セットを生成する。図 2 7 に示すように、サブシステム 2 6 0 6 は、有界ボックス生成器 1 5 0 4、What CNN 1 5 0 6、及び When CNN 1 5 0 8 を含む。循環バッファ 1 5 0 2 からのカメラ当たりの関節データ構造 8 0 0 及び画像フレームは、有界ボックス生成器 1 5 0 4 への入力として与えられる。有界ボックス生成器 1 5 0 4、What CNN 1 5 0 6、及び When CNN 1 5 0 8 の詳細は、以前に提示されている。

## 【 0 2 4 2 】

図 2 7 に記載されたシステムは、識別された被写体に対する在庫商品のリストを含むログ・データ構造を生成するために、第 1 及び第 2 の検出セットを処理する選択ロジックを含む。識別された被写体による在庫商品を取ることと、識別された被写体による在庫陳列構造上に在庫商品を置くことの第 1 の検出セットは、ログ生成器 2 7 2 0 によって生成される。第 1 の検出セットは、上述したように、第 2 の画像プロセッサの出力及び関節データ構造 8 0 0 を使用して決定される。識別された被写体による在庫商品を取ることと、識別された被写体による在庫陳列構造上に在庫商品を置くことが、第 3 の画像処理装置の出力を用いて決定される。各真の在庫イベント（取るまたは置く）について、選択ロジック・コントローラ 2 6 0 8 は、第 2 の画像プロセッサ（意味的差分抽出サブシステム 2 6 0 4）または第 3 の画像プロセッサ（領域提案サブシステム 2 6 0 6）の何れかからの出力を選択する。一実施形態では、選択ロジックが、その在庫イベントの予測のために、より高い信頼度スコアを有する画像プロセッサからの出力を選択する。

[ 背景画像意味的差分抽出の処理フロー ]

## 【 0 2 4 3 】

図 2 8 A 及び図 2 8 B は、実空間のエリア内の被写体による変化を追跡するために意味的差分抽出サブシステム 2 6 0 4 によって実行される詳細なステップを示す。ショッピングストアの例では、被写体が棚と他の空きスペースとの間の通路内で店舗内を移動する顧客及び店舗の従業員である。処理はステップ 2 8 0 2 で開始する。上述のように、カメラ 1 1 4 は、被写体を識別するためにカメラからの画像シーケンスが処理される前に較正される。カメラ較正の詳細は、上述されている。重なり合う視野を有するカメラ 1 1 4 は、被写体が存在する実空間の画像を取得する。一実施形態では、カメラが毎秒 N フレームの速度で同期された画像シーケンスを生成するように構成される。各カメラの画像シーケン

10

20

30

40

50

スは、ステップ 2804 において、カメラ毎にそれぞれの循環バッファ 1502 に格納される。循環バッファ（リング・バッファとも呼ばれる）は、スライドする時間帯に画像シーケンスを格納する。背景画像格納装置 2704 は、前景被写体のないカメラ当たりの画像フレームのシーケンス内の初期画像フレームで初期化される（ステップ 2806）。

#### 【0244】

被写体が棚の前を移動することにつれて、被写体当たりの有界ボックスが上述のように、それらの対応する関節データ構造 800 を使用して生成される（ステップ 2808）。ステップ 2810 では、画像フレーム当たりの有界ボックス内のピクセルを、背景画像格納装置 2704 からの背景画像からの同じ位置のピクセルで置き換えることによって、マスクされた画像が作成される。カメラ毎の画像のシーケンス内の各画像に対応するマスクされた画像は、背景画像格納装置 2704 に格納される。i 番目のマスクされた画像は、カメラ当たりの画像フレームのシーケンス内の次の (i + 1) 画像フレーム内のピクセルを置換するための背景画像として使用される。

#### 【0245】

ステップ 2812 において、N 個のマスクされた画像が組み合わされて、ファクタ化画像が生成される。ステップ 2814 では、ファクタ化画像のペアのピクセル値を比較することによって、差異ヒートマップが生成される。一実施形態では 2 つのファクタ化画像 (f i 1 及び f i 2) の 2 D 空間内の位置 (x, y) におけるピクセル間の差は以下の式 1 に示すように計算される：

$$\sqrt{((f i 1[x, y][red] - f i 2[x, y][red])^2 + (f i 1[x, y][green] - f i 2[x, y][green])^2 + (f i 1[x, y][blue] - f i 2[x, y][blue])^2)} \quad (1)$$

#### 【0246】

2 D 空間内の同じ x 及び y 位置におけるピクセル間の差は式に示されるように、赤、緑及び青 (R G B) チャネルのそれぞれの強度値を使用して決定される。上記の式は、2 つのファクタ化画像における対応するピクセル間の差（ユークリッドノルムとも呼ばれる）の大きさを与える。

#### 【0247】

差異ヒートマップは、実空間のエリアにおけるセンサノイズ及び明度変化によるノイズを含み得る。図 28B では、ステップ 2816 で、差異ヒートマップのためのビットマスクが生成される。意味的に重要な変化は、ビットマスク内の 1 のクラスタによって識別される。これらのクラスタは、棚から取られた、または棚に置かれた在庫商品を識別する変化に対応する。しかしながら、差異ヒートマップのノイズは、ビットマスクにランダムな 1 を導入する可能性がある。更に、複数の変化（複数の商品が棚から取り出されるか、または棚に置かれる）は、1 の重なり合うクラスタを導入し得る。処理フローの次のステップ (2818) では、画像形態操作がビットマスクに適用される。画像形態操作はノイズ（望ましくない 1）を除去し、また、1 の重なり合うクラスタを分離しようと試みる。この結果、意味的に重要な変更に対応する 1 のクラスタを含む、よりクリーンなビットマスクが得られる。

#### 【0248】

形態的操作には 2 つの入力が与えられる。第 1 の入力ビットマスクであり、第 2 の入力は構造化要素またはカーネルと呼ばれる。2 つの基本的な形態的操作は、「収縮」及び「膨張」である。カーネルは、様々なサイズの矩形行列に配置された 1 からなる。異なる形状（例えば、円形、楕円形、または十字形）のカーネルは、行列内の特定の位置に 0 を加えることによって生成される。異なる形状のカーネルがビットマスクをクリーニングする際に所望の結果を達成するために、画像形態操作に使用される。収縮操作では、カーネルはビットマスク上をスライド（または移動）する。カーネルの下すべてのピクセルが 1 である場合、ビットマスク内のピクセル（1 または 0 の何れか）は 1 と見なされる。そ

れ以外では、それは収縮される（0に変化する）。収縮操作は、ビットマスク内の孤立した1を除去するのに有用である。しかしながら、収縮はまた、エッジを収縮することによって1のクラスタを収縮させる。

#### 【0249】

膨張操作は、収縮とは逆である。この操作では、カーネルがビットマスク上をスライドするとき、カーネルの下少なくとも1つのピクセルの値が1である場合、カーネルによってオーバーラップされたビットマスクエリア内のすべてのピクセルの値が1に変更される。1のサイズクラスタを増大させるために、収縮後にビットマスクに膨張が適用される。ノイズが収縮において除去されるので、膨張は、ビットマスクにランダムノイズを導入しない。よりクリーンなビットマスクを達成するために、収縮操作と膨張操作との組み合わせが適用される。例えば、コンピュータ・プログラム・コードの以下の行は、ビットマスクに1の3×3フィルタを適用して、「オープン」操作を実行し、この「オープン」操作では、収縮操作と、それに続く膨張操作とを適用して、ノイズを除去し、上述のようにビットマスク内の1のクラスタのサイズを復元する。上記のコンピュータ・プログラム・コードは、リアルタイム・コンピュータ・ビジョン・アプリケーション用のプログラミング機能のOpenCV(オープンソース・コンピュータ・ビジョン)ライブラリを使用する。ライブラリは、<https://opencv.org/>で入手できる。

```
_bit_mask = cv2.morphologyEx(bit_mask, cv2.MORPH_OPEN, self.kernel_3x3,
dst=_bit_mask)
```

#### 【0250】

「クローズ」操作は、膨張操作に続いて収縮操作を適用する。これは、1のクラスタの内側の小さな穴を閉じるのに有用である。以下のプログラム・コードは、30×30十字形(クロス・シェープ)フィルタを使用してビットマスクにクローズ操作を適用する。

```
_bit_mask = cv2.morphologyEx(bit_mask, cv2.MORPH_CLOSE, self.kernel_30x
30_cross, dst=_bit_mask)
```

#### 【0251】

ビットマスク及び2つのファクタ化画像(前後)は、カメラ毎に畳み込みニューラル・ネットワーク(上記の変化CNNと呼ばれる)への入力として与えられる。変化CNNの出力は、変化データ構造である。ステップ2822では、重なり合う視野を有する変化CNNからの出力が前述の三角測量技法を使用して結合される。3D実空間における変化の位置は、棚の位置と一致する。在庫イベントの位置が棚上の位置にマップされる場合、変化は真のイベントと見なされる(ステップ2824)。マップされない場合は、変化は誤検知であり、廃棄される。真のイベントは、前景被写体に関連付けられる。ステップ2826において、前景被写体が識別される。一実施形態では、関節データ構造800が変化の閾値距離内の手関節の位置を決定するために使用される。ステップ2828で前景被写体が識別された場合、ステップ2830で、その変化が識別された被写体に関連付けられる。例えば、変化の閾値距離内の複数の被写体の手の関節位置のために、ステップ2828で前景被写体が識別されない場合は、次に、ステップ2832において、領域提案サブシステムによる変化の冗長検出が選択される。処理はステップ2834で終了する。

#### [変化CNNのトレーニング]

#### 【0252】

7つのチャンネル入力のトレーニング・データセットが、変化CNNをトレーニングするために作成される。顧客として行為する1または複数の被写体が、ショッピングストアで買い物をするふりをすることによって、取る及び置く行為を実行する。被写体が通路を移動し、棚から在庫商品を取り、在庫商品を棚に置き戻す。取る行為及び置く行為を実行する行為者の画像は、循環バッファ1502に収集される。画像は上述のように、ファクタ化画像を生成するために処理される。2つのファクタ化画像間の変化を視覚的に識別するために、ファクタ化画像2706のペア及びビットマスク計算器2710によって出力された対応するビットマスクが手動で精査される。変化を有するファクタ化画像については、有界ボックスが変化の周りに手動で描かれる。これは、ビットマスクの変化に対応する

1のクラスタを含む最小の有界ボックスである。変化における在庫商品のSKU数が識別され、有界ボックスと共に画像のラベルに含まれる。在庫商品の取ることまたは置くことを識別するイベントタイプも、有界ボックスのラベルに含まれる。従って、各有界ボックスのラベルは、ファクタ化画像上のその位置、商品のSKU、及びイベントタイプを識別する。ファクタ化画像は、2つ以上の有界ボックスを有することができる。上記の処理は、トレーニング・データセット内の全ての収集されたファクタ化画像における全ての変化について繰り返される。1対のファクタ化画像はビットマスクと共に、変化CNNへの7チャンネル入力を形成する。

#### 【0253】

変化CNNのトレーニング中に、順方向パス及び逆方向伝播が実行される。順方向パスでは、変化CNNが、トレーニング・データセット内の画像の対応するシーケンス内のファクタ化画像内に表される背景変化を識別し且つ分類する。変化CNNは、識別された背景変化を処理し、識別された被写体による在庫商品を取ることと、識別された被写体による在庫陳列構造上に在庫商品を置くことの第1の検出セットを生成する。逆方向伝播の間、変化CNNの出力はトレーニング・データセットのラベルに示されるように、グラウンドトゥールースと比較される。1または複数のコスト関数に対する勾配が計算される。次いで、勾配は、畳み込みニューラル・ネットワーク(CNN)及び全結合(FC)ニューラル・ネットワークに伝播され、その結果、予測誤差が低減され、出力がグラウンドトゥールースに近づく。一実施形態では、ソフトマックス関数及びクロスエントロピー損失関数が、出力のクラス予測部分に対する変化CNNのトレーニングのために使用される。出力のクラス予測部分は、在庫商品のSKU識別子及びイベントタイプ、すなわち取ることまたは置くことを含む。

#### 【0254】

第2の損失関数は、有界ボックスの予測のために変化CNNをトレーニングするために使用される。この損失関数は、予測されたボックスとグラウンドトゥールース・ボックスとの間の共通集合/和集合(IOU)を計算する。変化CNNによって予測された有界ボックスと真の有界ボックスラベルとの共通集合の面積が、同じ有界ボックスの和集合の面積によって割り算される。IOUの値は、予測ボックスとグラウンドトゥールース・ボックスとの間のオーバーラップが大きい場合に高い。2つ以上の予測された有界ボックスがグラウンドトゥールース有界ボックスとオーバーラップする場合、最も高いIOU値を有するものが選択され、損失関数が計算される。損失関数の詳細は、2016年5月9日に発行されたRedmonらの論文「You Only Look Once: Unified, Real-Time Object Detection」に記載されている。この論文は<https://arxiv.org/pdf/1506.02640.pdf>で入手可能である。

#### [ 特定の実施態様 ]

#### 【0255】

様々な実施形態において、上述の実空間のエリア内において被写体による在庫商品を置くこと及び取ることを追跡するためのシステムは、また、以下の特徴の1または複数を含む。

#### [ 1. 領域提案 ]

#### 【0256】

領域提案は、人物をカバーする全ての異なるカメラからの手の位置のフレーム画像である。領域提案は、システム内のすべてのカメラによって生成される。それは、店舗の商品を持っている手だけでなく、空の手も含む。

#### [ 1.1 WhatCNNモデル ]

#### 【0257】

領域提案は、深層学習アルゴリズムを使用して画像分類への入力として使用することができる。この分類エンジンは、「WhatCNN」モデルと呼ばれる。それは、インハンド分類モデルである。それは手の中の物を分類する。インハンド画像分類は、物体の部分が手によって遮蔽されていても、動作することができる。より小さい商品は、手で90%

10

20

30

40

50



まで遮閉することができる。WhatCNNモデルによる画像分析のための領域は、計算コストが高いため、幾つかの実施形態では意図的に小さく保たれる。各カメラは、専用のGPUを有することができる。これは、すべてのカメラからのすべての手の画像について、すべてのフレームについて実行される。WhatCNNモデルによる上記の画像分析に加えて、信頼度重みもその画像（1つのカメラ、1つの時点）に割り当てられる。分類アルゴリズムは、在庫管理単位（SKU）のリスト全体にわたるロジットを出力して、 $n$ 個の商品について店舗の製品及びサービス識別コードリストを生成し、空の手（ $n + 1$ ）について1つの追加を生成する。

【0258】

シーン処理は、キー値辞書を各ビデオに送ることによって、その結果を各ビデオ処理に送り返す。ここで、キーは固有関節IDであり、値は、関節が関連付けられた固有の個人IDである。関節に関連する人物が見つからなかった場合、それは辞書に含まれない。

【0259】

各ビデオ処理はシーン処理からキー値辞書を受け取り、フレーム番号を返された辞書にマッピングするリング・バッファにそれを格納する。

【0260】

返されたキー値辞書を使用して、ビデオは、既知の人々に関連付けられた手の近くにある時点毎の画像のサブセットを選択する。これらの領域は、numpyスライスである。また、前景マスクと関節CNNの生の出力特徴配列の周りに同様のスライスを取る。これらの組み合わせられた領域は一緒に連結されて単一の多次元numpy配列になり、領域が関連付けられている個人IDと、その領域が人物からのどの手から来たかと同様に、numpy配列を保持するデータ構造に格納される。

【0261】

次に、全ての提案された領域がFIFOキューに供給される。このキューは領域を取り込み、それらのnumpy配列をGPU上のメモリにプッシュする。

【0262】

配列がGPUに到着すると、それらは、WhatCNNと呼ばれる、分類専用のCNNに供給される。このCNNの出力は、 $N + 1$ の大きさのフロートの平坦な配列であり、ここで、 $N$ は店舗内の固有のSKUの数であり、最終クラスは、空クラス、すなわち空の手を表す。この配列のフロートは、ロジットと呼ばれる。

【0263】

WhatCNNの結果は、領域データ構造に記憶される。

【0264】

その後、ある時点での全ての領域が、各ビデオ処理からシーン処理に送り返される。

【0265】

シーン処理は、ある時点ですべてのビデオからすべての領域を受け取り、その結果をキー値辞書に格納し、そこでは、キーは個人IDであり、値はキー値辞書であり、そこでは、キーはカメラIDであり、値は領域のロジットである。

【0266】

次に、この集約されたデータ構造は、時点毎にフレーム番号を集約された構造にマッピングするリング・バッファに格納される。

[ 1.2 WhenCNNモデル ]

【0267】

WhatCNNモデルによって処理された様々なカメラからの画像は、ある期間にわたって結合される（ある期間にわたって複数のカメラ）。このモデルへの追加の入力は、複数のカメラから三角測量された3D空間における手の位置である。このアルゴリズムへのもう1つの入力は、店舗のプラノグラムからの手の距離である。いくつかの実施形態では、プラノグラムを使用して、手が特定の商品を含む棚（例えば、チェリオボックス）に近いかどうかを識別することができる。このアルゴリズムへの別の入力は、店舗上の足の位置である。

10

20

30

40

50

## 【 0 2 6 8 】

S K Uを使用する物体分類に加えて、第 2 の分類モデルは、時系列分析を使用して、物体が棚から取り出されたか、または棚上に置かれたかを判定する。画像は、以前の画像フレームにおいて手の中にあった物体が棚に戻されたか、または、棚から取り出されたか否かの判定を行うために、ある期間にわたって分析される。

## 【 0 2 6 9 】

1 秒間 ( 3 0 フレーム / 秒 ) 及び 3 台のカメラについて、システムは、信頼度の付加された同じ手に対して 9 0 の分類出力を有する。この組み合わせられた画像分析は、手の中の物体を正確に識別する確率を劇的に増加させる。時間分析は、個々のフレームの幾つかの非常に低い信頼度レベル出力にもかかわらず、出力の品質を改善する。このステップは例えば、8 0 % の精度から 9 5 % の精度までの出力信頼度を取ることができる。

10

## 【 0 2 7 0 】

このモデルはまた、棚モデルからの出力をその入力として含み、この人物がどの物体を選択したかを識別する。

## 【 0 2 7 1 】

シーン処理は、3 0 以上の集約された構造が蓄積されるのを待ち、少なくとも 1 秒のリアルタイムを表し、次いで、集約された構造を、個人 I D と手のペア毎に単一の整数に縮小するために、更なる分析を実行し、ここで、整数は、店舗内の S K U を表す固有 I D である。一時点において、この情報はキーが個人 I D と手のペアであり、値が S K U 整数であるキー値辞書に記憶される。この辞書は、その時点においてフレーム番号を各辞書にマッピングするリング・バッファに経時的に格納される。

20

## 【 0 2 7 2 】

次に、人が何かを何時取ったか、取られた物が何かを識別するために、この辞書が時間の経過と共にどのように変化するかを見て、追加の分析を実行することができる。このモデル ( W h e n C N N ) は、何かが取られたか？、何かが置かれたか？というブール式の各質問に対するロジットのみならず、S K U ロジットも発する。

## 【 0 2 7 3 】

W h e n C N N の出力は、フレーム番号を、キーが個人 I D であり、値が W h e n C N N によって発せられる拡張ロジットであるキー値辞書にマッピングするリング・バッファに格納される。

30

## 【 0 2 7 4 】

次いで、ヒューリスティックスの更なる集合が、W h e n C N N 及び格納された人々の関節位置の両方の格納された結果、並びに、店舗の棚上の商品の事前に計算されたマップ上で実行される。このヒューリスティックスの集合は、取ること及び置くことの結果、商品がどこに追加されるか、またはどこから除去されるかを決定する。各取ること / 置くことについて、ヒューリスティックスは、取ることまたは置くことが、棚からまたは棚へ、バスケットからまたはバスケットへ、または、人からまたは人へ、であったかどうかを決定する。出力は、S K U の指数における配列の値が個人の有する S K U の数である配列として格納された、個人毎の在庫である。

## 【 0 2 7 5 】

40

買物客が店舗の出口に近づくと、システムは、物品明細リストを買物客の電話に送ることができる。次に、電話はユーザの物品明細を表示し、格納されたクレジットカード情報に課金するための確認を求める。ユーザが了承すると、そのクレジットカードに課金される。システム内で知られているクレジットカードを有していない場合、クレジットカード情報を提供するように要求される。

## 【 0 2 7 6 】

代替的に、買物客は店舗内キオスクに接近することもできる。システムは買物客が何時キオスクの近くにいるかを識別し、その買物客の物品明細を表示するためにキオスクにメッセージを送信する。キオスクは、買物客に物品明細の料金を受け入れるように求める。買物客が了承した場合、買物客は、自分のクレジットカードを通すか、または支払いのた

50

めに現金を投入することができる。図 16 は、領域提案のための W h e n C N N モデルの図を示す。

[ 2. 置き間違えられた商品 ]

【 0 2 7 7 】

この特徴では、置き間違えられた商品を、それらが人によってでたらめな棚に置き戻されたときに識別する。これにより、プラノグラムに対する足及び手の位置が不正確になるので、物体識別に問題が生じる。従って、システムは、経時的に修正されたプラノグラムを構築する。以前の時系列分析に基づいて、システムは、人が商品を棚に戻したかどうかを判定することができる。次に、物体がその棚位置から取り出されると、システムは、その手の位置に少なくとも 1 つの置き間違えられた商品があることを知る。それに対応して、アルゴリズムは、人がその棚から置き間違えられた商品を取り上げることができるというある程度の信頼度を有する。置き間違えられた商品が棚から取り出される場合、システムはその商品をその位置から取り去り、従って、棚は、もはやその商品を有していない。システムはまた、アプリを介して置き間違えられた商品について店員に通知し、店員がその商品をその正しい棚に移動させることができるようにすることができる。

10

[ 3. 意味的差分抽出 (棚モデル) ]

【 0 2 7 8 】

背景画像処理のための代替技術は、棚上の商品 (取り除かれた、または、置かれた商品) に対する変化を識別するための背景減算アルゴリズムを含む。これは、ピクセルレベルでの変化に基づく。棚の前に人がいる場合、人の存在によるピクセル変化を考慮しないようにアルゴリズムは停止する。背景減算はノイズの多い処理である。従って、クロスカメラ分析が行われる。シェルフに「意味的に重要な」変化があることを十分な数のカメラが認める場合、システムは、棚のその部分に変化があることを記録する。

20

【 0 2 7 9 】

次のステップは、その変化が「置く」または「取る」変化であるかどうかを確認することである。このために、第 2 の分類モデルの時系列分析が使用される。棚のその特定の部分に対する領域提案が生成され、深層学習アルゴリズムを通過する。これは、物体が手の中で遮蔽されないで、手の中の画像分析よりも容易である。第 4 の入力、3 つの典型的な R G B 入力に加えてアルゴリズムに与えられる。第 4 のチャンネルは背景情報である。棚または意味的差分抽出の出力は、第 2 の分類モデル (時系列分析モデル) に再び入力される。

30

【 0 2 8 0 】

このアプローチにおける意味的差分抽出は、以下のステップを含む：

1. カメラからの画像は、同じカメラからの以前の画像と比較される。
2. 2 つの画像間の対応する各ピクセルは、R G B 空間におけるユークリッド距離を介して比較される。
3. ある閾値を超える距離がマーキングされ、その結果、マーキングされたばかりのピクセルの新しい画像が得られる。
4. マーキングされた画像からノイズを除去するために、画像形態フィルタの集合が使用される。
5. 次に、マーキングされたピクセルの大きな集合を探索し、それらの周りに有界ボックスを形成する。
6. 次に、各有界ボックスについて、2 つの画像中の元のピクセルを見て、2 つの画像スナップショットを得る。
7. 次に、これらの 2 つの画像スナップショットは、画像領域が取られている商品を表すか、または置かれている商品を表すか、及び商品が何であることを分類するようにトレーニングされた C N N にプッシュされる。

40

[ 3. 店舗監査 ]

【 0 2 8 1 】

各棚の在庫は、システムによって維持される。それは、商品が顧客によって取り出され

50

ると更新される。任意の時点で、システムは、店舗在庫の監査報告書を生成することができる。

[ 4 . 手の中の複数の商品 ]

【 0 2 8 2 】

複数の商品に対して異なる画像が使用される。手の中の2つの商品は、1つの商品と比較して別に扱われる。幾つかのアルゴリズムは、1つの商品のみを予測することができるが、複数の商品を予測することはできない。従って、CNNは、商品の「2つの」量のためのアルゴリズムが手の中の単一の商品とは別個に実行され得るようにトレーニングされる。

[ 5 . データ収集システム ]

【 0 2 8 3 】

所定のショッピングスクリプトが、画像の良質なデータを収集するために使用される。これらの画像は、アルゴリズムのトレーニングに使用される。

[ 5 . 1 ショッピングスクリプト ]

【 0 2 8 4 】

データ収集は、以下のステップを含む：

1 . スクリプトが自動的に生成され、どのような行為を取るべきかを人間の行為者に伝える。

2 . これらの行為は、商品Xを取る、商品Xの置く、商品XをY秒間保持することを含む行為の集合から任意に抽出される。

3 . これらの行為を実行しながら、行為者は所与の行為を持続させながら、可能な限り多くの方法で、自分自身を移動させ、向きを合わせる。

4 . 一連の行為の間、カメラの集合体は、多くの視点から行為者を記録する。

5 . 行為者がスクリプトを終了した後、カメラビデオは一緒に束ねられ、元のスクリプトと共に保存される。

6 . スクリプトは、行為者のビデオでトレーニングする機械学習モデル（CNNなど）への入力ラベルとして機能する。

[ 6 . 製品ライン ]

【 0 2 8 5 】

本システム及びその一部は、以下のアプリでサポートされるレジレス・チェックアウトに使用できる。

[ 6 . 1 店舗アプリ ]

【 0 2 8 6 】

店舗アプリは、幾つかの主要な機能を有しており、データ分析視覚化を提供し、損失防止をサポートし、人々が店舗内のどこにいるか、及びどの商品を収集したかを小売業者に示すことによって顧客を支援するプラットフォームを提供する。従業員に対する許可レベル及びアプリアクセスは、小売業者の裁量で指示することができる。

[ 6 . 1 . 1 標準分析 ]

【 0 2 8 7 】

データは、プラットフォームによって収集され、様々な方法で 사용할 ことができる。

【 0 2 8 8 】

1 . デリバティブデータは、店舗、店舗が提供するショッピング体験、及び、製品、環境、及び他の人々との顧客の交流に関する様々な種類の分析を実行するために使用される。

a . データは、店舗と顧客の交流の分析を実行するために、格納され背景で使用される。店舗アプリは、このデータの視覚化の一部を小売業者に表示する。他のデータは、データポイントが要求されるときに格納され、照会される。

2 . ヒートマップ：

プラットフォームは、小売業者のフロアプラン、棚レイアウト、及び他の店舗環境を、様々な種類の活動のレベルを示すオーバーレイとともに視覚化する。

1 . 例：

10

20

30

40

50

1. 人は通り過ぎるが、どんな製品も扱っていない場所のマップ。
2. 製品と接触するときに、人がフロアのどこに立っているかのマップ。
3. 置き間違えられた商品：

プラットフォームは、店舗のすべてのSKUを追跡する。商品が間違った場所に置かれると、プラットフォームはその商品がどこにあるかを知り、ログを構築する。ある閾値で、または直ちに、店舗の従業員は、置き間違えられた商品に対し注意喚起され得る。或いは、スタッフは、店舗アプリ内の置き間違えられた商品マップにアクセスすることができる。都合の良いときに、スタッフが置き間違えられた商品を迅速に見つけ出し、修正することができる。

#### [ 6 . 1 . 2 標準アシスト ]

- ・ 店舗アプリは店舗のフロアプランを表示する。
- ・ それは、店舗内の各人物を表す図画を表示する。
- ・ タッチ、クリック、または他の手段を介して図画が選択されると、店舗の従業員に対して適切な情報が表示される。例えば、ショッピングカートの商品（収集した商品）がリストに表示される。

- ・ プラットフォームが、特定の商品（単数または複数）に対する、及び個人が所有（ショッピングカート）している期間に対する、所定の閾値より低い信頼度レベルを有する場合、その図画（現在はドット）は差を示す。アプリは色の変化を使用する。緑は高い信頼度を示し、黄色／オレンジは低い信頼度を示す。

- ・ 店舗アプリを所持した店舗従業員には、より低い信頼度を通知することができる。店舗従業員は、顧客のショッピングカートが正確であることを確認することができる。

- ・ 店舗アプリを介して、小売業者の従業員は、顧客のショッピングカート商品を調整（追加または削除）することができる。

#### [ 6 . 1 . 3 標準LP ]

- ・ 買物客が買物客アプリを使用している場合、買物客は単に店舗を出て、課金される。しかし、買物客が買物客アプリを使用していない場合、買物客は、彼らのショッピングカート内の商品に対する支払い用のゲストアプリを使用する必要があるだろう。

- ・ 買物客が、店舗から出る途中でゲストアプリを迂回する場合、買物客の図画は、買物客が店舗を出る前にアプローチしなければならないことを示す。このアプリでは、色を赤色に変更する。スタッフはまた、潜在的な損失の通知を受ける。

- ・ 店舗アプリを介して、小売業者の従業員は、顧客のショッピングカート商品を調整（追加または削除）することができる。

#### [ 6 . 2 非店舗アプリ ]

##### 【 0 2 8 9 】

以下の分析の特徴は、プラットフォームの追加機能を表す。

#### [ 6 . 2 . 1 標準分析 ]

##### 1 . 製品交流：

次のような製品交流の細かな内訳：

- a. 各製品の交流時間対転化率。
- b. A / B 比較（色、スタイル等）。ディスプレイ上のより小さい製品の一部は、色、風味などのような複数の選択肢を有する。
  - ・ バラの金は銀よりも取扱いが多いか？
  - ・ 青い缶は赤い缶よりも多くの交流を招くのだろうか

##### 2 . 方向性インプレッション：

位置ベースのインプレッションと、買物客の注視がどこにあるかの差を知る。もし買物客が15フィート離れた製品を20秒間見ているならば、インプレッションは、彼らがどこにいるかが重要ではなく、彼らがどこを見ているかが重要である。

##### 3 . 顧客認識：

リピート買物客とそれに付随する電子メールアドレス（小売業者によって様々な方法で収集された）及び買物プロフィールを記憶する。

#### 4．グループダイナミックス：

買物客が、他の誰かが製品と接触するのを何時見ているかを判定する。

- ・ その後、その人が製品と接触するかどうかを答える。
- ・ その人たちは一緒に店舗に入ったのか、或いは、他人同士だろうか？
- ・ 個人或いは集団が、店舗でより多くの時間を費やしているか？

#### 5．顧客タッチバック：

顧客に対象情報、店舗後体験の提供。この特徴は、特定の慣行及び方針に応じて、それぞれの小売業者とわずかに異なった実施態様を有することができる。この特徴を採用するためには、小売業者からの統合及び／または開発が必要となる場合がある。

・ 買物客は、関心のある製品に関する通知を受領したいかどうかを尋ねられる。そのステップは、電子メールを収集する店舗の方法と統合されてもよい。

・ 店舗を出た後、顧客は、店舗で時間を費やした製品を伴う電子メールを受け取ることができる。持続時間、接触、及び視界（方向インプレッション）に対する交流閾値が決定される。閾値が満たされると、製品はそれを顧客のリストに送り、店舗を出た後すぐに顧客に送る。

#### 【0290】

追加的に、または代替的に、買物客に、販売中の製品または他の特別な情報を提供した後の期間に電子メールを送ることができる。これらの製品は興味を表明した商品であるが、購入しなかった商品である。

#### [ 6．3 ゲストアプリ ]

#### 【0291】

買物客アプリは、店舗を出るときに自動的に人々をチェックアウトする。しかしながら、プラットフォームは、買物客が店舗を使用するために買物客アプリを有することも使用することも必要としない。

#### 【0292】

買物客／個人が買物客アプリを持っていないか、または使用していないとき、買物客はキオスク（iPad（登録商標）／タブレットまたは他の画面）まで歩いて行くか、または予めインストールされたセルフ・チェックアウト・マシンまで歩いて行く。プラットフォームと一体化されたディスプレイは、顧客のショッピングカートを自動的に表示する。

#### 【0293】

買物客は、何が表示されているかを見直す機会を有する。買物客がディスプレイ上の情報に同意する場合、買物客は、現金をマシンに投入するか（その機能がハードウェア（例えば、セルフ・チェックアウト・マシン）に装備されている場合）、または、買物客のクレジットカードまたはデビットカードを通すことができる。そして、店舗を出ることができる。

#### 【0294】

買物客が、ディスプレイに同意しない場合に、タッチスクリーン、ボタン、または他の手段を介して、異議を申し立てることを選択することで、店員に通知される（店舗アプリの「店舗アシスト」を参照）。

#### [ 6．4 買物客アプリ ]

#### 【0295】

アプリ、買物客アプリを使用することにより、顧客は商品と共に店舗を出ることができ、自動的に課金され、デジタルレシートが与えられる。買物客は、店舗のショッピングエリア内にいる間は常に自分のアプリを開かなければならない。プラットフォームは、買物客のデバイスに表示される固有画像を認識する。プラットフォームは、それらを買物客のアカウントに結びつけ（顧客関連付け）、買物客がアプリを開いたままにしているかどうかにかかわらず、店舗のショッピングエリア内での時間中、誰がいるかを覚えておくことができる。

10

20

30

40

50

## 【 0 2 9 6 】

買物客が商品を集めると、買物客アプリは、買物客のショッピングカートに商品を表示する。買物客が望む場合、買物客は取り出した（すなわち、ショッピングカートに追加された）各商品に関する製品情報を見ることができる。製品情報は、店舗のシステムに格納されるか、またはプラットフォームに追加される。製品販売を提供すること、または価格を表示することなど、その情報を更新する能力は、小売業者が要求／購入または開発することができるオプションである。

## 【 0 2 9 7 】

買物客が商品を下に置くと、バックエンド及び買物客アプリ上のショッピングカートから商品が取り除かれる。

## 【 0 2 9 8 】

買物客アプリが開かれ、顧客関連付けが完了した後に閉じられると、プラットフォームは、買物客のショッピングカートを維持し、買物客が店舗を出ると、それらに正しく課金する。

## 【 0 2 9 9 】

買物客アプリはまた、開発ロードマップに関するマッピング情報を有する。それは、顧客が、捜している商品をタイプ入力することによって情報を要求する場合、店舗内の商品をどこで見つけるべきかを顧客に伝えることができる。後日、買物客のショッピングリスト（手動で、または他のインテリジェントシステムを介してアプリに入力された）を取得し、すべての所望の商品を収集する店舗を通る最速ルートを表示する。「袋詰め傾向」などの他のフィルタを追加することができる。袋詰め傾向フィルタにより、買物客は最も速いルートをたどらず、最初に頑丈な商品を収集し、その後、より壊れやすい商品を収集することができる。

## [ 7 . 顧客のタイプ ]

## 【 0 3 0 0 】

メンバー顧客： 最初のタイプの顧客が、アプリを使用してシステムにログインする。顧客は画面で促され、クリックすると、システムはそれをその顧客の内部IDにリンクする。顧客がアカウントを有する場合、顧客が店舗から出るときにアカウントに自動的に課金される。これは、会員制店舗である。

## 【 0 3 0 1 】

ゲスト顧客： すべての店舗が会員権を持っているわけではない。或いは、顧客がスマートフォンやクレジットカードを持っていないこともある。このタイプの顧客はキオスクまで歩いていこう。キオスクは顧客が有する商品を表示し、顧客にお金を入れるように依頼する。キオスクは、顧客が購入した全ての商品について既に知っている。このタイプの顧客の場合、システムは顧客がショッピングカート内の商品に対して支払っていないかどうかを識別し、顧客がそこに到着する前に、チェッカーに未支払い商品について知らせるようにドアのチェッカーに促すことができる。システムはまた、支払いが行われていないか、システムがその商品について低い信頼度を有する1つの商品に対してプロンプトを表示することもできる。これは、予測経路探索と呼ばれる。

## 【 0 3 0 2 】

システムは、信頼度レベルに基づいて、店舗内を歩いている顧客にカラーコード（緑色及び黄色）を割り当てる。緑色で色分けされた顧客は、システムにログインされているか、またはシステムがそれら顧客について高い信頼度を有している顧客である。黄色の色分けされた顧客は、高い信頼度で予測されない1または複数の商品を有している顧客である。店員は黄色の点を見て、それらをクリックして、問題商品を特定し、顧客まで歩いて行き、問題を修正することができる。

## [ 8 . 分析 ]

## 【 0 3 0 3 】

顧客が特定の棚の前でどれだけの時間を費やしたかといった、顧客に関する多くの分析情報が収集される。更に、システムは、顧客が見ている場所（システム上のインプレッシ

10

20

30

40

50

ョン)と、顧客が取り出して棚に戻した商品とを追跡する。このような分析は現在、電子商取引で利用可能であるが、小売店では利用可能ではない。

#### [ 9. 機能モジュール ]

##### 【 0 3 0 4 】

以下は、機能モジュールのリストである：

- 1 . 同期カメラを使用して、店舗内の画像の配列を取得するシステム。
- 2 . 画像中の関節を識別し、個々の人物の関節のセットを識別するシステム。
- 3 . 関節セットを使用して新しい人物を作成するシステム。
- 4 . 関節セットを使用してゴースト人物を削除するシステム。
- 5 . 関節セットを追跡することによって、経時的に個々の人物を追跡するシステム。
- 6 . 店舗内にいる各人に対して手の中の商品のSKU数を示す領域提案を生成するシステム (WhatCNN)。

10

7 . 手の中の商品が棚上において取り出されたか、または置かれたかを示す領域提案のための取ること/置くこと分析を実行するシステム (WhenCNN)。

8 . 領域提案及び取ること/置くこと分析を用いて、1人当たりの在庫配列を生成するためのシステム (ヒューリスティクスと人物の保存された関節位置と店舗の棚上の事前に計算された商品のマップを組み合わせたWhenCNNの出力)。

9 . 棚上の置き間違えられた商品の位置を識別し、追跡し、更新するシステム。

10 . ピクセルベースの分析を使用して、棚上の商品に対する変化 (取る/置く) を追跡するシステム。

20

11 . 店舗の在庫監査を実施するシステム。

12 . 手の中の複数の商品を識別するシステム。

13 . ショッピングスクリプトを用いて店舗から商品画像データを収集するシステム。

14 . 会員顧客からのチェックアウトを実行し、集金を行うシステム。

15 . ゲスト顧客からのチェックアウトを実行し、集金を行うシステム。

16 . カート内の未払商品を特定し、損失防止を行うシステム。

17 . 顧客のカート内で誤って識別された商品を店員が識別するのに支援するために、例えばカラーコードを使用して顧客を追跡するシステム。

18 . 位置ベースのインプレッション、方向性インプレッション、A/B分析、顧客認識、グループダイナミクス等を含む顧客ショッピング分析を生成するシステム。

30

19 . ショッピング分析を使用して目標顧客タッチバックを生成するシステム。

20 . 様々な活動を視覚化するために店舗のヒートマップオーバーレイを生成するシステム。

##### 【 0 3 0 5 】

本明細書に記載されている技術は、レジレス・チェックアウトをサポートすることができる。店舗に行く。商品を取る。去る。

##### 【 0 3 0 6 】

レジレス・チェックアウトは、純粋なマシンビジョンと深層学習に基づくシステムである。買物客は、列に並ばず、より早くより簡単に欲しいものを得る。RFIDタグは不要。店舗のバックエンドシステムに対する変更は不要。第三者の販売時点在庫管理システムと統合することができる。

40

各ビデオフィールドのリアルタイム30FPS分析。

最先端の構内GPUクラスタ。

買物客と彼らが交流する商品を認識する。

例示的な実施形態では、インターネットに依存しない。

マシンビジョン技術のギャップを初めて解決するために、独自のカスタムアルゴリズムを含む複数の最先端の深層学習モデル。

##### 【 0 3 0 7 】

技術と機能には以下が含まれる：

- 1 . スタンダード・コグニションの機械学習パイプラインは、以下を解決する：

50



- a. 人物検出。
- b. 存在物追跡。
- c. マルチカメラ人物一致。
- d. 手検出。
- e. 商品分類。
- f. 商品所有権決定。

**【 0 3 0 8 】**

これらの技術を組み合わせると、以下のことができる：

- 1. 買い物体験を通じて、すべての人々をリアルタイムで追跡する。
- 2. 買物客が手に持っているもの、どこに立っているか、どんな商品を戻すかを知る。
- 3. 買物客が、どの方向にどれだけ長く向いているのかを知る。
- 4. 置き間違えられた商品を認識し、24 / 7 の目視商品化監査を実施する。

10

**【 0 3 0 9 】**

買物客が手に持っているものとバスケットに持っているものを正確に検出することができる。

店舗の学習：

**【 0 3 1 0 】**

特定の店舗や商品についてトレーニングされたカスタムニューラル・ネットワーク。トレーニング・データは、全ての店舗位置にわたって再利用可能である。

標準配備：

20

**【 0 3 1 1 】**

天井カメラは、店舗の全エリアを二重にカバーするように設置しなければならない。典型的な通路には2～6台のカメラが必要である。

**【 0 3 1 2 】**

構内GPUクラスタは、バックオフィス内の1つまたは2つのサーバックに収容できる。

**【 0 3 1 3 】**

例示的なシステムは、販売時点在庫管理システムと統合することができ、またはそれらを含むことができる。

**【 0 3 1 4 】**

30

同期カメラを使用して店舗内の画像配列を取得する第1のシステム、方法、コンピュータ・プログラム製品。

**【 0 3 1 5 】**

画像内の関節、及び個々の人物の関節のセットを識別する第2のシステム、方法、及びコンピュータ・プログラム製品。

**【 0 3 1 6 】**

関節のセットを使用して新しい人物を作成する第3のシステム、方法、及びコンピュータ・プログラム製品。

**【 0 3 1 7 】**

関節のセットを使用してゴースト人物を削除する第4のシステム、方法、及びコンピュータ・プログラム製品。

40

**【 0 3 1 8 】**

関節のセットを追跡することにより経時的に個々の人物を追跡する第5のシステム、方法、及びコンピュータ・プログラム製品。

**【 0 3 1 9 】**

手の中の商品のSKU数を示す、店舗内にいる各人物のための領域提案を生成する第6のシステム、方法、及びコンピュータ・プログラム製品(WhatCNN)。

**【 0 3 2 0 】**

手の中の商品が棚上に取り出されたか、または置かれたかを示す領域提案のために取る／置く分析を実行する第7のシステム、方法、及びコンピュータ・プログラム製品(Wh

50

e n C N N )。

【 0 3 2 1 】

領域提案と取る / 置く分析 (例えば、ヒューリスティックス、格納された個人の関節位置、及び、店舗棚上の商品の予め計算されたマップと組み合わせられた W h e n C N N の出力) を使用して個人当たりの在庫配列を生成する第 8 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 2 2 】

棚上に置き間違えられた商品の位置を識別し、追跡し、更新するための第 9 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 2 3 】

ピクセルベースの分析を使用して棚上の商品に対する変化 (取る / 置く) を追跡する第 1 0 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 2 4 】

店舗の在庫監査を実行する第 1 1 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 2 5 】

手の中の複数の商品を識別する第 1 2 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 2 6 】

ショッピングスクリプトを使用して店舗から商品画像データを収集する第 1 3 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 2 7 】

会員顧客からチェックアウトを実行し、集金を行う第 1 4 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 2 8 】

ゲスト顧客からのチェックアウトを実行し、集金を行う第 1 5 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 2 9 】

カート内の未払商品を特定し、損失防止を行う第 1 6 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 3 0 】

顧客のカート内で誤って識別された商品を店員が識別するのを支援するために、例えばカラーコードを使用して顧客を追跡する第 1 7 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 3 1 】

位置ベースのインプレッション、方向性インプレッション、A / B 分析、顧客認識、グループダイナミクス等を含む顧客ショッピング分析を生成する第 1 8 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 3 2 】

ショッピング分析を使用して目標顧客タッチバックを生成する第 1 9 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 3 3 】

様々な活動を視覚化するために店舗のヒートマップオーバーレイを生成する第 2 0 のシステム、方法、及びコンピュータ・プログラム製品。

【 0 3 3 4 】

手検出のための第 2 1 のシステム、方法、及びコンピュータ・プログラム。

【 0 3 3 5 】

商品分類のための第 2 2 のシステム、方法、及びコンピュータ・プログラム。

【 0 3 3 6 】

商品所有権決定のための第 2 3 のシステム、方法、及びコンピュータ・プログラム。

10

20

30

40

50

## 【 0 3 3 7 】

商品人物検出のための第 2 4 のシステム、方法、及びコンピュータ・プログラム。

## 【 0 3 3 8 】

商品存在物追跡のための第 2 5 のシステム、方法、及びコンピュータ・プログラム。

## 【 0 3 3 9 】

商品マルチカメラ人物一致のための第 2 6 の方法及びコンピュータ・プログラム。

## 【 0 3 4 0 】

実質的に本明細書に記載されているレジレス・チェックアウトのための第 2 7 のシステム、方法、及びコンピュータ・プログラム製品。

## 【 0 3 4 1 】

第 1 ~ 第 2 6 のシステムの何れかと、上記列挙した第 1 ~ 第 2 6 のシステムの何れかの他の 1 つまたは複数のシステムとの組合せ。

## 【 0 3 4 2 】

ここに記載されるのは、実空間のエリア内の被写体による在庫商品を置くこと及び取ることを追跡する方法であって：

## 【 0 3 4 3 】

各カメラの視野が少なくとも 1 つの他のカメラの視野と重なり合う複数のカメラを使用して、実空間内の対応する視野のそれぞれの画像シーケンスを生成すること；

## 【 0 3 4 4 】

複数のカメラから画像シーケンスを受け取り、第 1 の画像認識エンジンを使用して画像を処理し、被写体及び識別された被写体の実空間内の位置を識別する第 1 のデータセットを生成すること；

## 【 0 3 4 5 】

画像シーケンス内の画像内の識別された被写体の手の画像を含む有界ボックスを指定するために第 1 のデータセットを処理すること；

## 【 0 3 4 6 】

複数のカメラからの画像シーケンスを受信し、第 2 の画像認識エンジンを使用して識別された被写体の手の分類を生成するために、画像内の有界ボックスを処理すること、但し、前記分類は、識別された被写体が在庫商品を保持しているかどうか、棚に対する識別された被写体の手の位置を示す第 1 の近似度分類と、識別された被写体の身体に対する識別された被写体の手の位置を示す第 2 の近似度分類と、識別された被写体に関連するバスケットに対する識別された被写体の手の位置を示す第 3 の近似度分類と、可能性のある在庫商品の識別子とを含み； 及び、

## 【 0 3 4 7 】

識別された被写体の画像シーケンス内の画像セットの手の分類を処理し、識別された被写体による在庫商品を取ることを、及び、識別された被写体による在庫陳列構造上に在庫商品を置くことを検出すること、を有する。

## 【 0 3 4 8 】

この説明された方法では、第 1 のデータセットが、識別された各被写体について、実空間内の座標を有する候補関節のセットを含むことができる。

## 【 0 3 4 9 】

この説明された方法は、有界ボックスを指定するために第 1 のデータセットを処理することを含むことができ、各被写体の候補関節のセット内の関節の位置に基づいて有界ボックスを指定することを含む。

## 【 0 3 5 0 】

この説明された方法では、第 1 及び第 2 の画像認識エンジンの一方または両方が畳み込みニューラル・ネットワークを備えることができる。

## 【 0 3 5 1 】

この説明された方法は、畳み込みニューラル・ネットワークを使用して有界ボックスの分類を処理することを含むことができる。

10

20

30

40

50

## 【 0 3 5 2 】

非一時的データ記憶媒体を備えるコンピュータ可読メモリと、本明細書に記載の処理の何れかによって実空間のエリア内の被写体による在庫商品を置くこと及び取ることを追跡するために、コンピュータによって実行可能なメモリに記憶されたコンピュータ命令とを含むコンピュータ・プログラム製品が記載されている。

## 【 0 3 5 3 】

被写体の手を含む画像シーケンスを生成する複数のカメラと、複数のカメラに結合された処理システムであって、画像シーケンスを受信して時系列に手の分類を生成する手画像認識エンジンと、画像シーケンスから手の分類を処理して被写体による、在庫商品を置くこと及び取ることのうちの1つである行為を識別するロジックとを含む処理システムと、を含むシステムが記載されている。

10

## 【 0 3 5 4 】

このシステムは画像シーケンス内の画像内の被写体の関節の位置を識別し、識別された関節に基づいて被写体の手を含む対応する画像内の有界ボックスを識別するロジックを含むことができる。

## 【 0 3 5 5 】

追記に列挙するコンピュータ・プログラムは、本明細書に添付され、本願において提供されるシステムの特定の部分を実装するためのコンピュータ・プログラムの一例の一部を含む。追記には、被写体の関節及び在庫商品を識別するためのヒューリスティックスの例が含まれる。追記は、被写体のショッピングカート・データ構造を更新するためのコンピュータ・プログラム・コードを提示する。追記はまた、畳み込みニューラル・ネットワークのトレーニング中に学習率を計算するためのコンピュータ・プログラム・ルーチンを含む。追記には、各カメラからの画像フレーム毎、被写体毎、手毎のデータ構造における畳み込みニューラル・ネットワークから、被写体の手の分類結果を保存するためのコンピュータ・プログラム・ルーチンが含まれている。

20

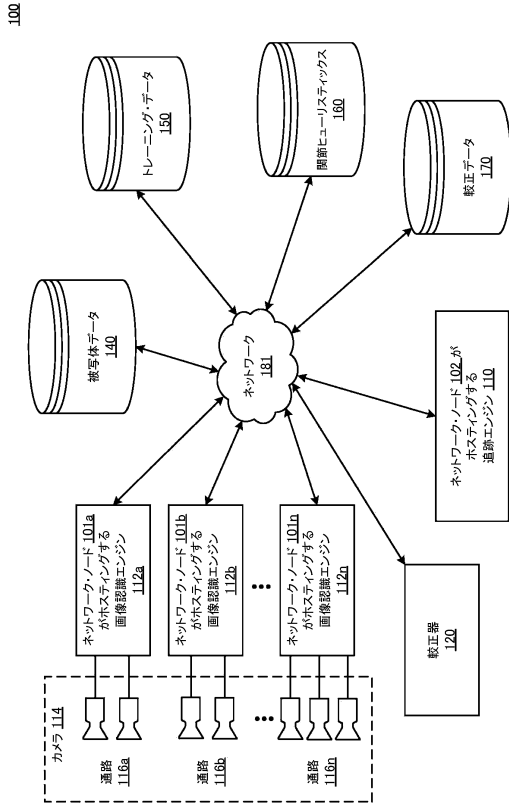
30

40

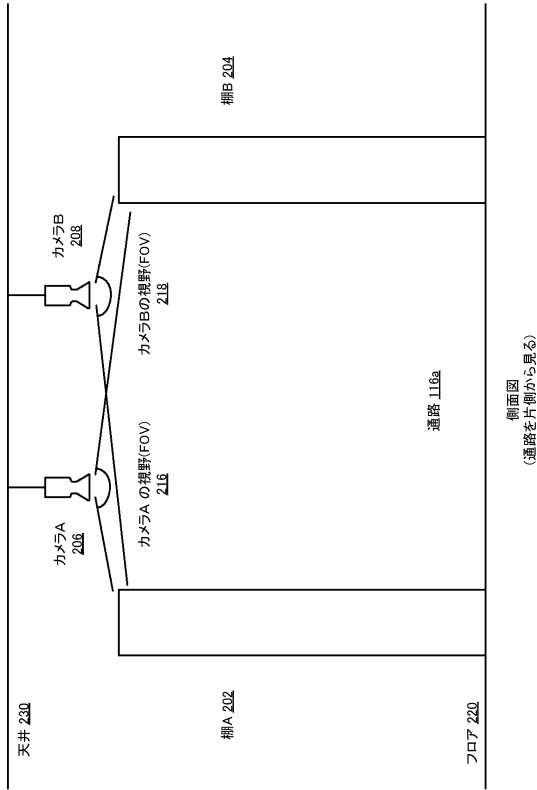
50

【図面】

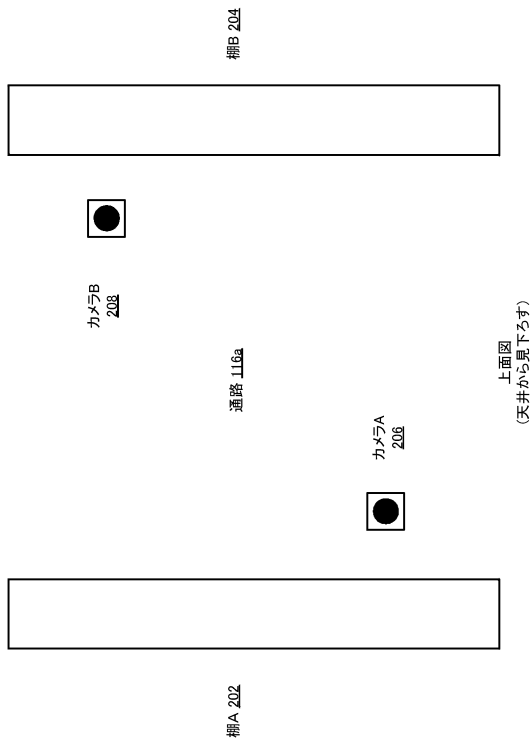
【図 1】



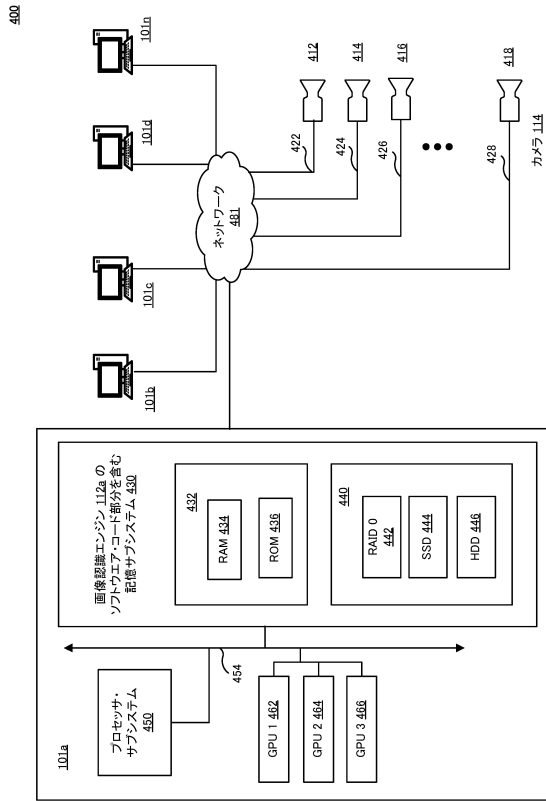
【図 2】



【図 3】



【図 4】



10

20

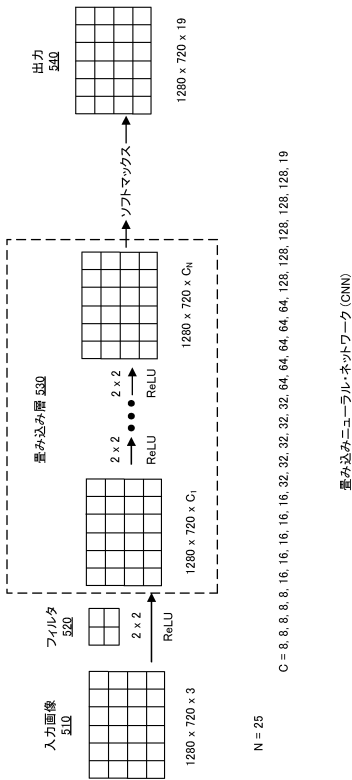
30

40

50

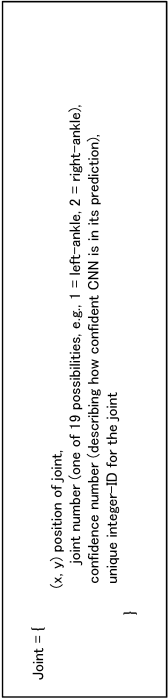
【図 5】

500

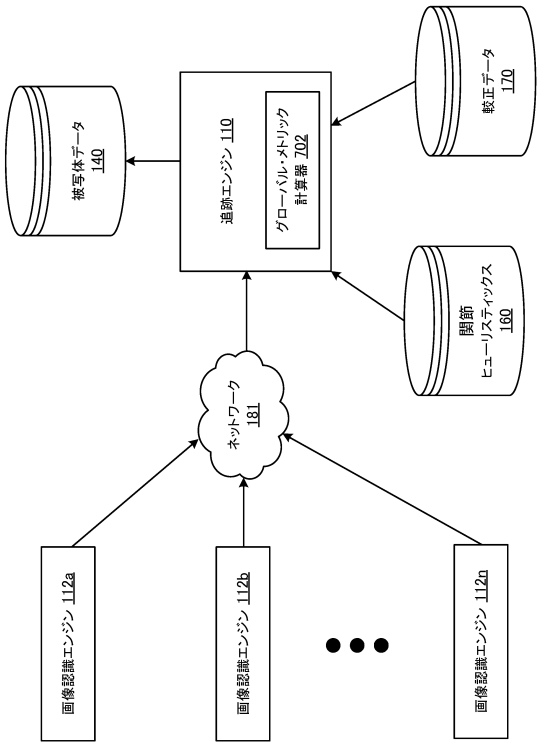


【図 6】

600

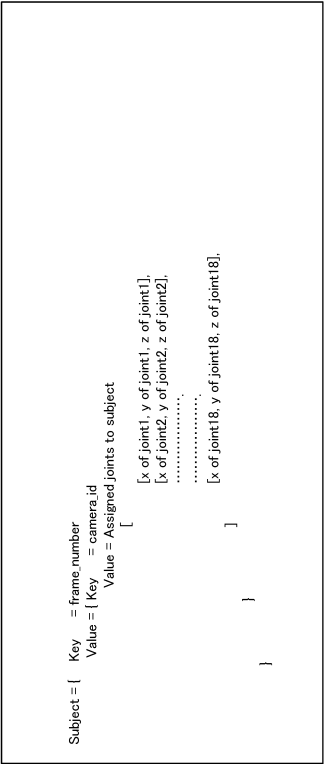


【図 7】



【図 8】

800



10

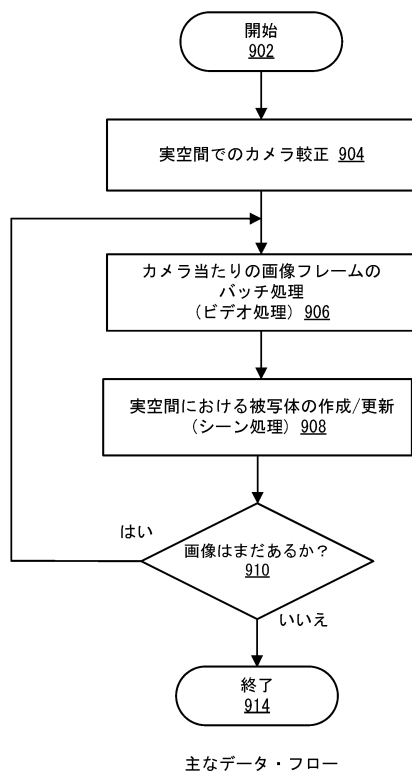
20

30

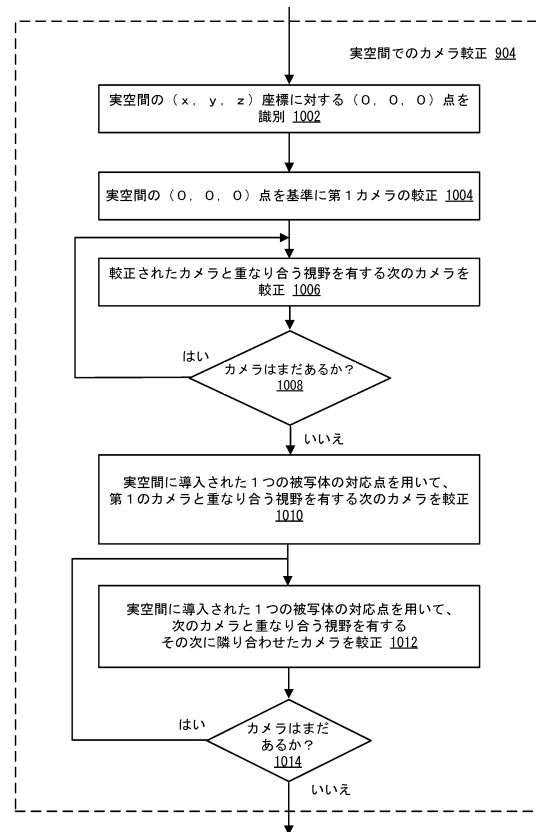
40

50

【図 9】



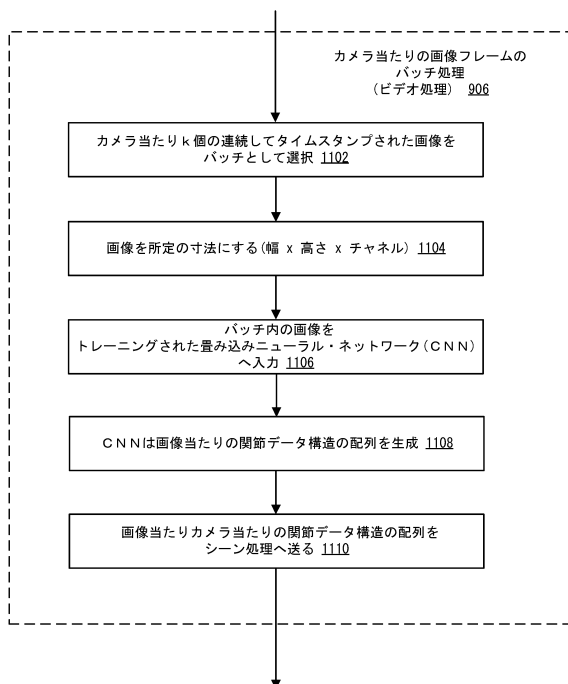
【図 10】



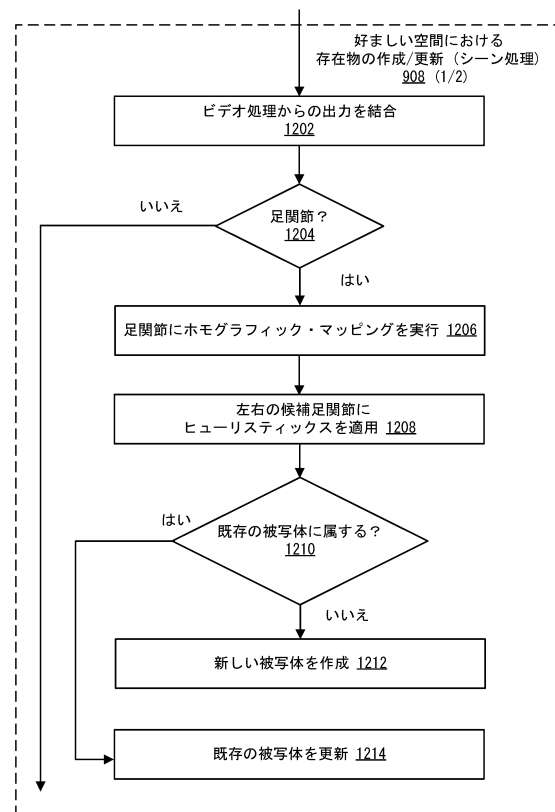
10

20

【図 11】



【図 12 A】

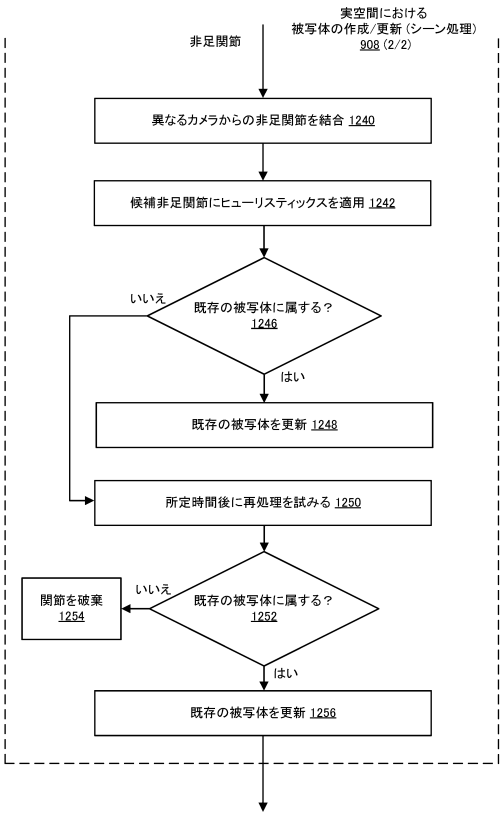


30

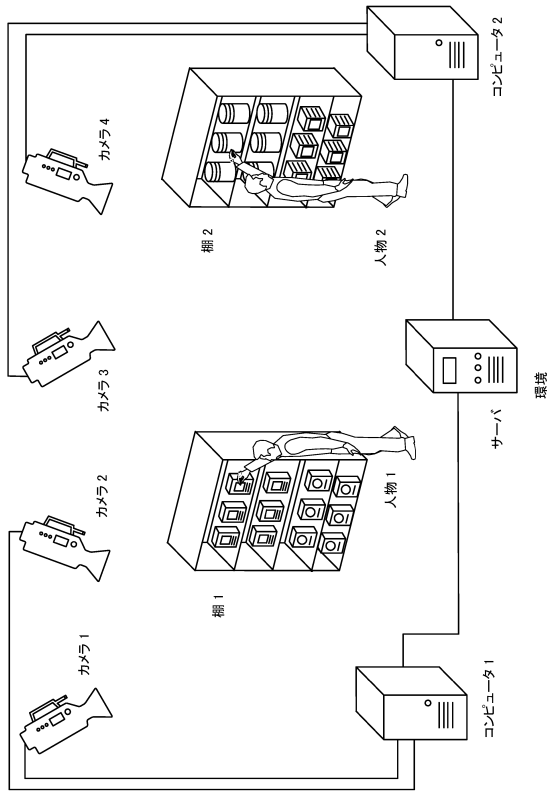
40

50

【図 1 2 B】



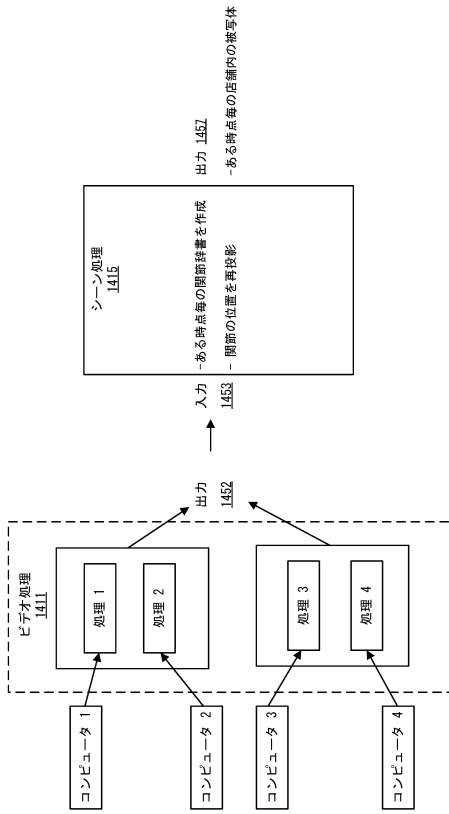
【図 1 3】



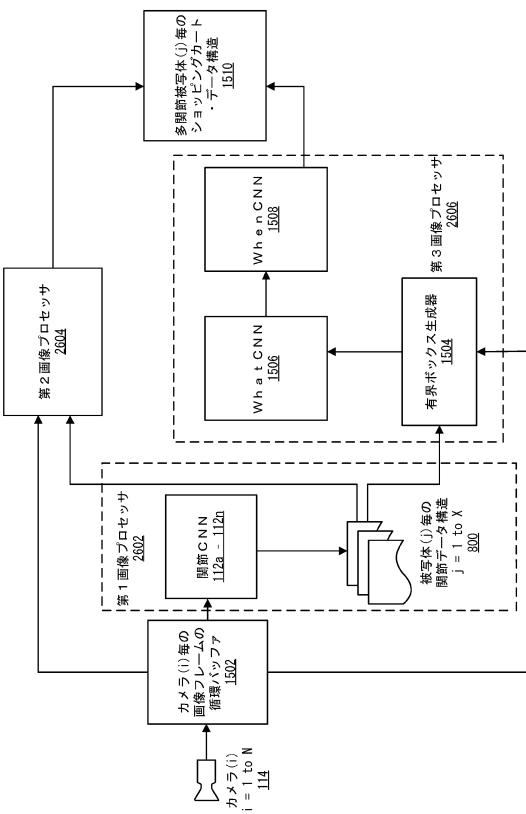
10

20

【図 1 4】



【図 1 5 A】



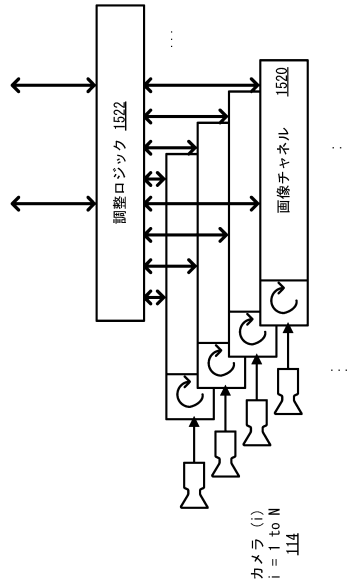
30

40

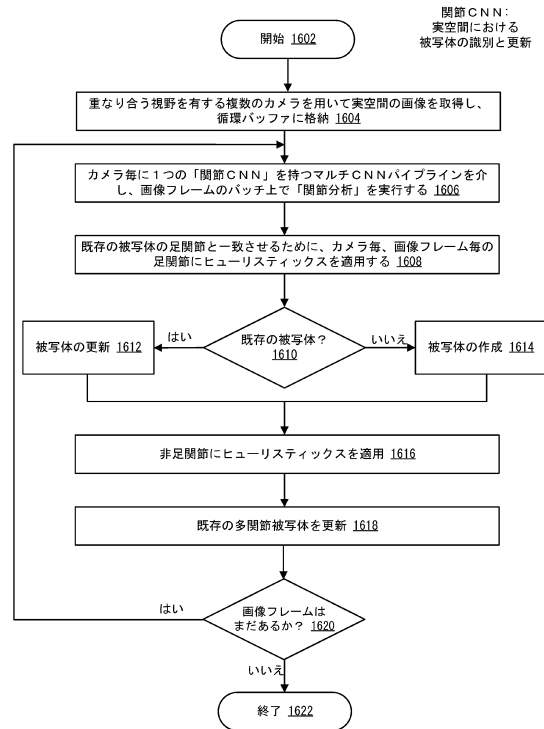
50



【図 15B】



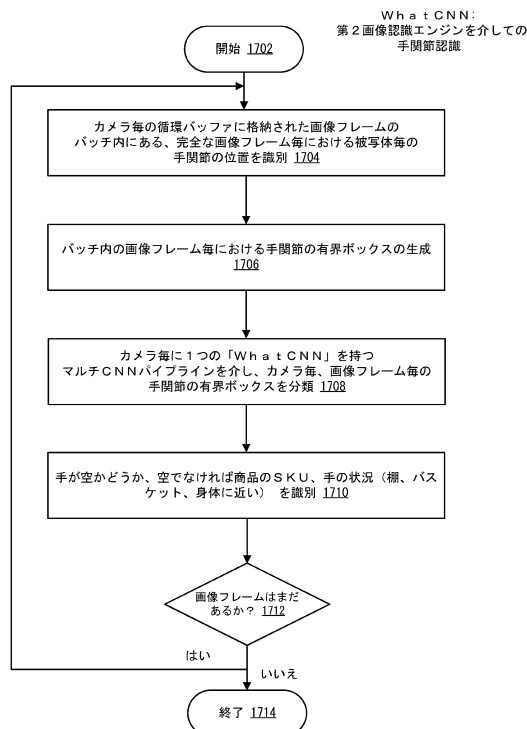
【図 16】



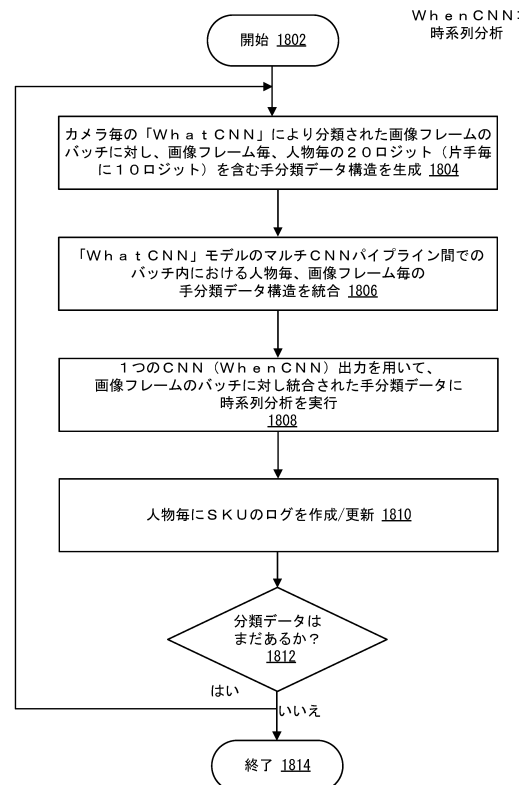
10

20

【図 17】



【図 18】

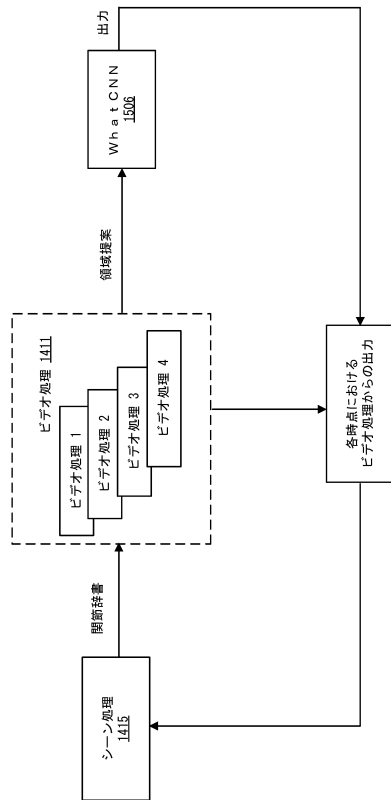


30

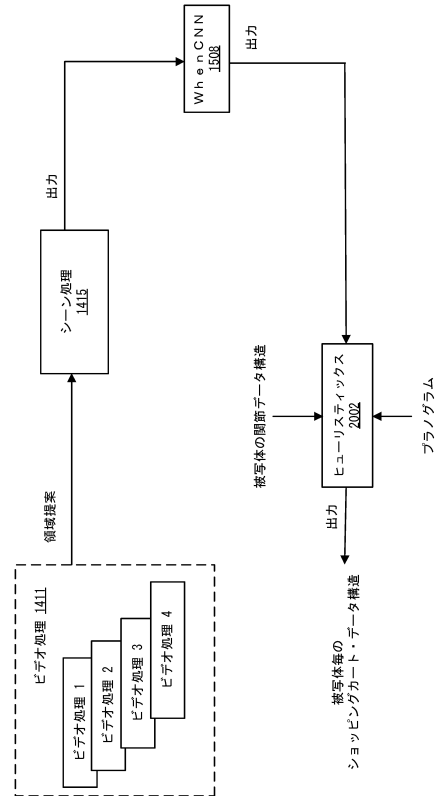
40

50

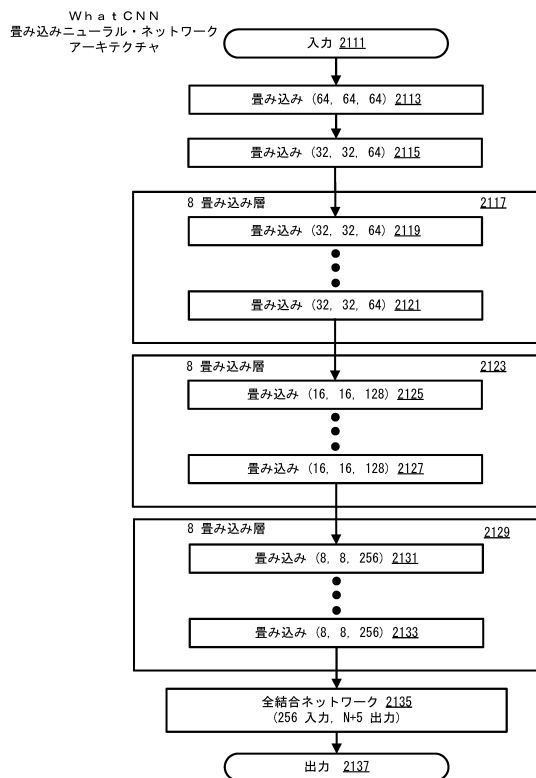
【図 19】



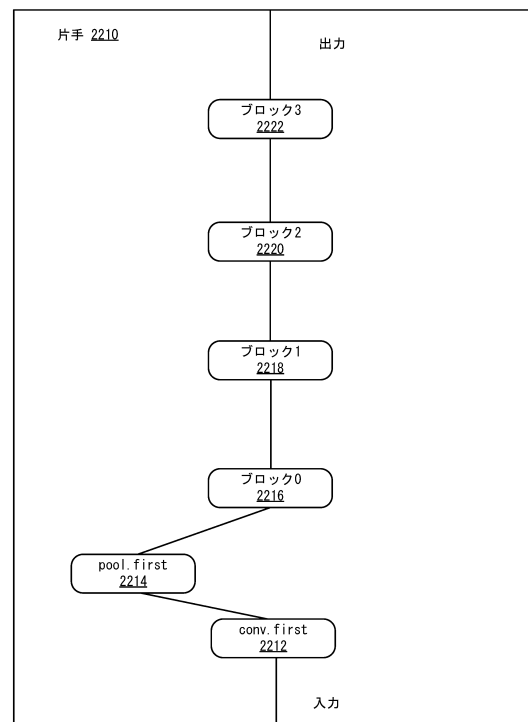
【図 20】



【図 21】



【図 22】



10

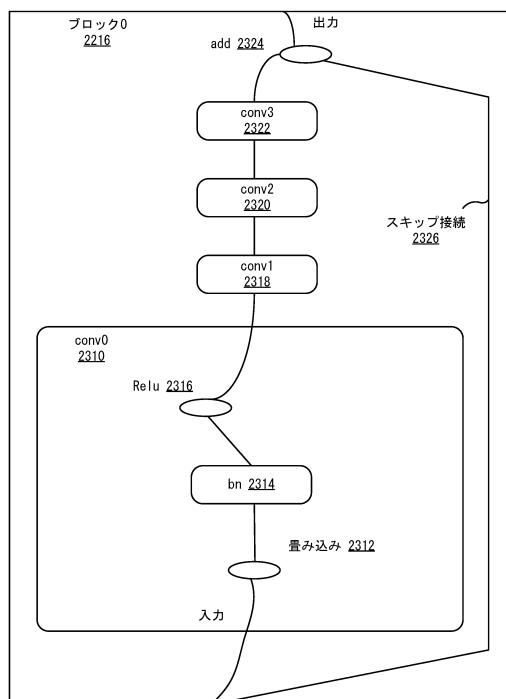
20

30

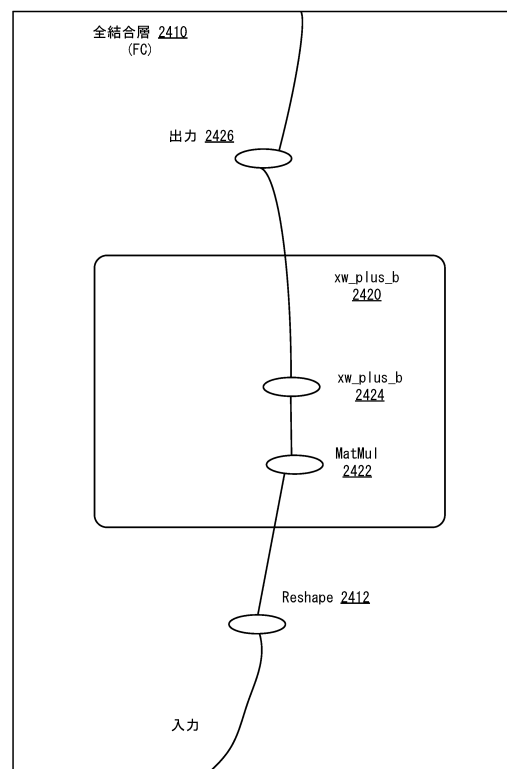
40

50

【 図 2 3 】



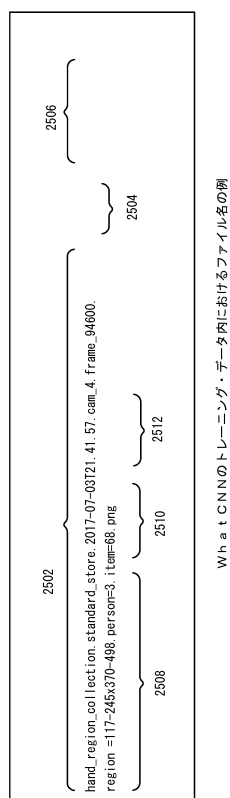
【圖 24】



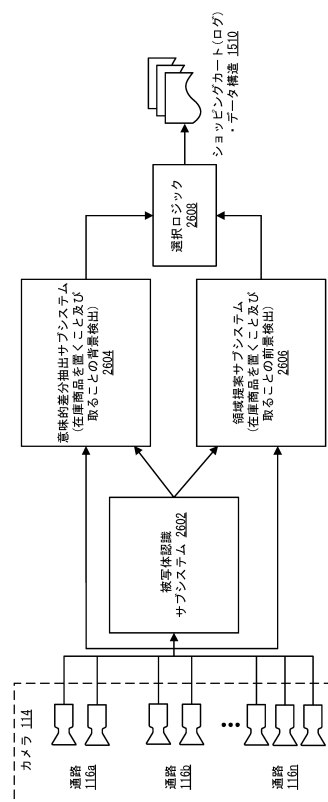
10

20

【 図 2 5 】



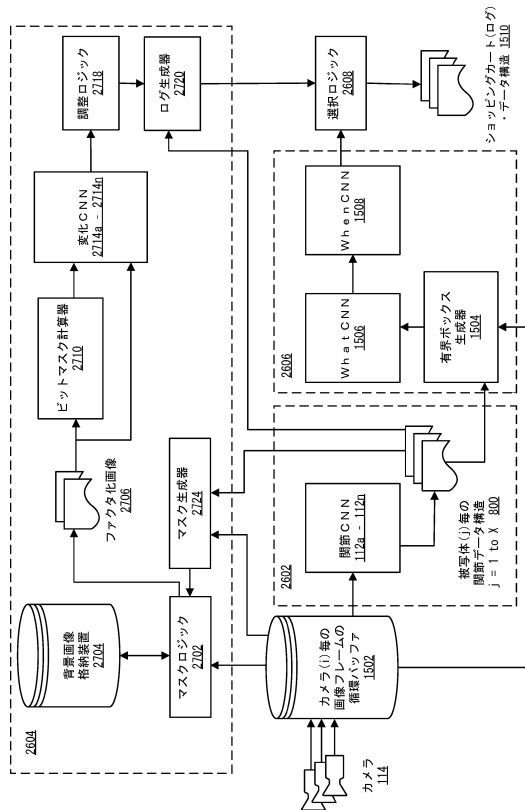
【 図 2 6 】



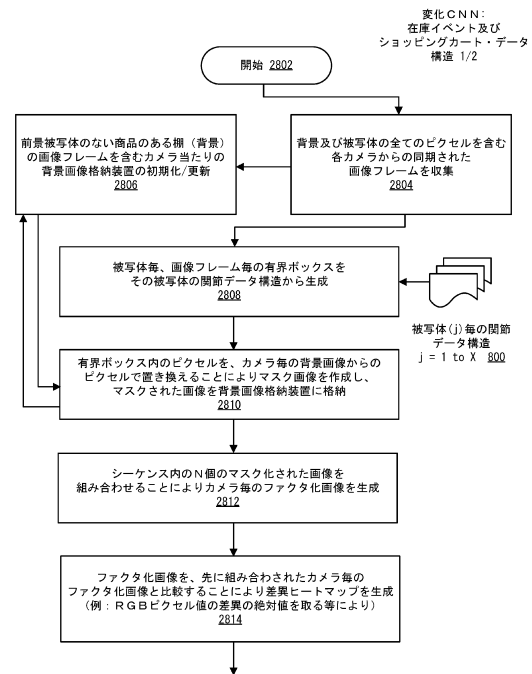
30

40

【 図 2 7 】



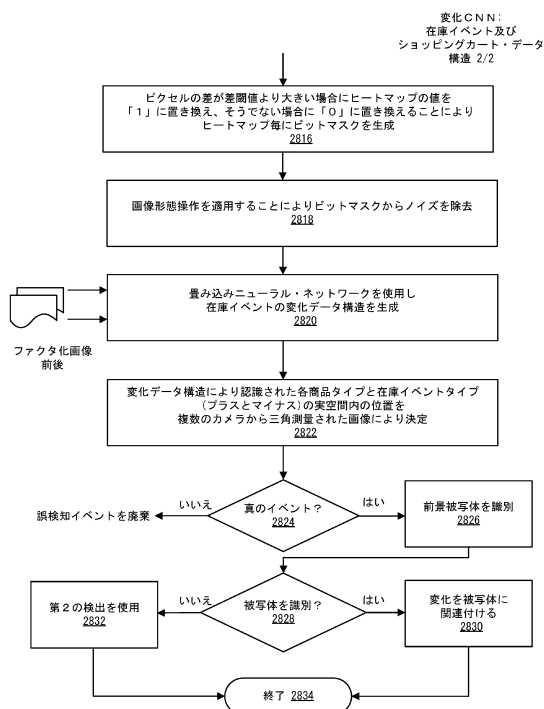
【 図 2 8 A 】



10

20

【 図 2 8 B 】



30

40

## フロントページの続き

(33)優先権主張国・地域又は機関

米国(US)

(31)優先権主張番号 15/907,112

(32)優先日 平成30年2月27日(2018.2.27)

(33)優先権主張国・地域又は機関

米国(US)

(31)優先権主張番号 15/945,466

(32)優先日 平成30年4月4日(2018.4.4)

(33)優先権主張国・地域又は機関

米国(US)

(31)優先権主張番号 15/945,473

(32)優先日 平成30年4月4日(2018.4.4)

(33)優先権主張国・地域又は機関

米国(US)

8 2 8 8 番ストリート 3 3

(72)発明者 ノヴァク, ジョン エフ.

アメリカ合衆国 カリフォルニア州 9 4 1 3 1 サンフランシスコ, バーンサイド アヴェニュー  
1 5 5

(72)発明者 オグル, ブランドン エル.

アメリカ合衆国 カリフォルニア州 9 4 1 5 8 サンフランシスコ, ナンバー 3 2 6 4 番ストリ  
ート 1 1 5 5

審査官 片岡 利延

(56)参考文献 米国特許出願公開第 2 0 1 4 / 0 2 1 9 5 5 0 ( U S , A 1 )

特開 2 0 1 2 - 1 2 3 6 6 7 ( J P , A )

特開 2 0 0 9 - 0 5 5 1 3 9 ( J P , A )

特開 2 0 0 4 - 3 4 8 6 1 8 ( J P , A )

特開 2 0 1 0 - 0 0 2 9 9 7 ( J P , A )

米国特許出願公開第 2 0 1 3 / 0 1 8 2 1 1 4 ( U S , A 1 )

米国特許出願公開第 2 0 1 7 / 0 2 0 6 6 6 9 ( U S , A 1 )

米国特許第 0 8 0 9 8 8 8 8 ( U S , B 1 )

Marcel Germann et al. , Space-time Body Pose Estimation in Uncontrolled Environments , [online] , 2011年 , <https://ieeexplore.ieee.org/document/5955367>Umar Iqbal et al. , PoseTrack: Joint Multi-Person Pose Estimation and Tracking , [online] , 2017年07月26日 , <https://ieeexplore.ieee.org/document/8099978>

(58)調査した分野 (Int.Cl. , D B 名)

G 0 6 T 7 / 2 0

G 0 6 T 7 / 0 0