

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)公開番号

特開2022-19339

(P2022-19339A)

(43)公開日 令和4年1月27日(2022.1.27)

(51)国際特許分類		F I		テーマコード(参考)	
G 0 6 T	7/00 (2017.01)	G 0 6 T	7/00	3 5 0 C	5 C 0 5 4
G 0 6 T	7/246(2017.01)	G 0 6 T	7/246		5 L 0 9 6
G 0 6 N	3/08 (2006.01)	G 0 6 N	3/08		
H 0 4 N	7/18 (2006.01)	H 0 4 N	7/18	D	
		H 0 4 N	7/18	G	
		審査請求	未請求	請求項の数	16
				OL	(全21頁) 最終頁に続く

(21)出願番号	特願2020-123119(P2020-123119)	(71)出願人	000001007 キヤノン株式会社 東京都大田区下丸子3丁目30番2号
(22)出願日	令和2年7月17日(2020.7.17)	(74)代理人	100126240 弁理士 阿部 琢磨
		(74)代理人	100124442 弁理士 黒岩 創吾
		(72)発明者	館 俊太 東京都大田区下丸子3丁目30番2号 キヤノン株式会社内
		(72)発明者	小川 修平 東京都大田区下丸子3丁目30番2号 キヤノン株式会社内
		(72)発明者	御手洗 裕輔 東京都大田区下丸子3丁目30番2号 最終頁に続く

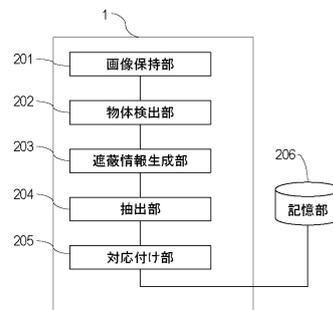
(54)【発明の名称】 情報処理装置、情報処理方法及びプログラム

(57)【要約】

【課題】 外見的特徴や姿勢が類似した物体が近接する場合においても安定して追尾を継続できる。

【解決手段】 上記課題を解決する本発明にかかる情報処理装置は、画像から少なくとも1つ以上の物体を検出する情報処理装置であって、遮蔽する物体と遮蔽された物体との遮蔽関係を示す画像特徴を学習した学習済みモデルに基づいて、前記画像から検出された各物体について、前記画像から検出された他の物体との遮蔽関係を推定する推定手段と、前記推定手段によって推定された遮蔽関係に基づいて、前記画像から検出された各物体について、前記画像と異なる時刻に撮像された画像において検出された物体との対応関係を特定する特定手段と、を有する。

【選択図】 図3



**【特許請求の範囲】****【請求項 1】**

画像から少なくとも 1 つ以上の物体を検出する情報処理装置であって、  
遮蔽する物体と遮蔽された物体との遮蔽関係を示す画像特徴を学習した学習済みモデルに基づいて、前記画像から検出された各物体について、前記画像から検出された他の物体との前記遮蔽関係を示す遮蔽情報を推定する推定手段と、  
少なくとも前記遮蔽情報に基づいて、前記画像から検出された各物体について、前記画像と異なる時刻に撮像された画像において検出された物体との対応関係を特定する特定手段と、を有することを特徴とする情報処理装置。

**【請求項 2】**

前記推定手段は、入力画像における物体の遮蔽された領域を推定する前記学習済みモデルに基づいて、前記画像から検出された各物体について、前記物体の遮蔽されている部分領域を示す前記遮蔽情報を推定することを特徴とする請求項 1 に記載の情報処理装置。

**【請求項 3】**

前記特定手段は、前記物体の画像特徴と、前記推定手段によって推定された前記遮蔽情報と、に基づいて、前記画像と異なる時刻に撮像された画像において検出された物体との対応関係を特定することを特徴とする請求項 1 または 2 に記載の情報処理装置

**【請求項 4】**

前記推定手段は、前記画像から検出された各物体について、他の物体によって遮蔽されていることを示す尤度である前記遮蔽情報を推定し、  
前記尤度を該物体の画像特徴に対応付けて保持する保持手段を更に有することを特徴とする請求項 1 乃至 3 のいずれか 1 項に記載の情報処理装置。

**【請求項 5】**

前記遮蔽情報は、前記画像の各領域について、遮蔽されている前記物体の領域にはより大きい前記尤度を示し、それ以外の領域にはより小さい前記尤度を示す情報であることを特徴とする請求項 4 に記載の情報処理装置。

**【請求項 6】**

前記特定手段は、前記画像から検出された各物体について、前記画像における位置と、前記画像から検出された画像特徴と、前記画像における前記遮蔽関係と、に基づいて、前記画像と異なる時刻に撮像された画像から検出された物体との対応関係を特定することを特徴とする請求項 1 乃至 5 のいずれか 1 項に記載の情報処理装置。

**【請求項 7】**

前記画像から前記物体毎の領域を取得する取得手段をさらに有し、  
前記推定手段は、前記取得手段によって取得された前記物体毎の領域について、遮蔽されている物体の有無を示す前記遮蔽情報を推定することを特徴とする請求項 1 乃至 6 のいずれか 1 項に記載の情報処理装置。

**【請求項 8】**

前記特定手段によって特定された、前記画像から検出された各物体と、前記画像と異なる時刻に撮像された画像において検出された物体と、の対応付けに基づいて、遮蔽されている物体があるか否かを判定する判定手段と、  
前記遮蔽されている物体を記憶する記憶手段と、をさらに有することを特徴とする請求項 1 乃至 7 のいずれか 1 項に記載の情報処理装置。

**【請求項 9】**

前記判定手段は、前記画像より前に撮像された画像において検出された物体のうち、前記画像から検出された各物体と対応しない物体を第 1 の物体として判定し、  
前記記憶手段は、前記画像が撮像された時点において、前記第 1 の物体が遮蔽されたことを記憶することを特徴とする請求項 8 に記載の情報処理装置。

**【請求項 10】**

前記判定手段は、前記画像から検出された各物体のうち、前記画像より前に撮像された画像において検出された物体と対応しない物体を第 2 の物体として判定し、

10

20

30

40

50

前記記憶手段は、前記画像から検出された前記第 2 の物体について、前記記憶手段によって前記画像より前に撮像された画像において遮蔽されていると判定された前記第 1 の物体との類似度が所定の閾値より大きい場合に、前記画像が撮像された時点において前記第 1 の物体は遮蔽されていないことを記憶することを特徴とする請求項 9 に記載の情報処理装置。

【請求項 1 1】

第 1 の画像において 2 つの物体を検出し、前記第 1 の画像の後で撮像された第 2 の画像において 1 つの物体を検出した場合、

前記推定手段は、前記第 2 の画像において検出された物体について他の物体を遮蔽していることを示す前記遮蔽情報を推定し、

前記特定手段は、前記推定手段によって推定された前記遮蔽情報に基づいて、前記第 1 の画像において検出された物体のうち他の物体を遮蔽している物体について、前記第 2 の画像において検出された物体と同一の物体であることを特定することを特徴とする請求項 1 乃至 1 0 のいずれか 1 項に記載の情報処理装置。

【請求項 1 2】

前記第 2 の画像より後に撮像された第 3 の画像から 2 つの物体を検出した場合、

前記推定手段は、前記第 3 の画像から検出された 2 つの物体のうち、前記第 2 の画像から検出された物体の画像特徴と対応付けられた物体について、他の物体を遮蔽していることを示す前記遮蔽情報を推定し、

前記特定手段は、前記第 1 の画像から検出された物体のうち他の物体によって遮蔽された物体について、前記第 3 の画像から検出された物体のうち前記第 2 の画像から検出された物体の画像特徴と対応付けられた物体とは異なる物体を、前記第 2 の画像において遮蔽された物体と同一の物体であることを特定することを特徴とする請求項 1 1 に記載の情報処理装置。

【請求項 1 3】

前記学習済みモデルは、ニューラルネットワークであることを特徴とする請求項 1 乃至 1 2 のいずれか 1 項に記載の情報処理装置。

【請求項 1 4】

前記学習済みモデルは、より短い時間間隔で撮像された複数の画像に基づいて、前記物体の遮蔽関係を示す画像特徴を学習させたモデルであることを特徴とする請求項 1 乃至 1 3 のいずれか 1 項に記載の情報処理装置。

【請求項 1 5】

コンピュータを、請求項 1 乃至 1 4 のいずれか 1 項に記載の情報処理装置が有する各手段として機能させるためのプログラム。

【請求項 1 6】

画像から少なくとも 1 つ以上の物体を検出する情報処理方法であって、

遮蔽する物体と遮蔽された物体との遮蔽関係を示す画像特徴を学習した学習済みモデルに基づいて、前記画像から検出された各物体について、前記画像から検出された他の物体との遮蔽関係を示す遮蔽情報を推定する推定工程と、

少なくとも前記遮蔽情報に基づいて、前記画像から検出された各物体について、前記画像と異なる時刻に撮像された画像において検出された物体との対応関係を特定する特定工程と、を有することを特徴とする情報処理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、被写体を追尾する技術に関する。

【背景技術】

【0002】

画像内の特定の被写体を追尾するための技術としては、輝度や色情報を利用するものやテンプレートマッチングなどが存在する。近年、Deep Neural Network

10

20

30

40

50

(以下DNNと省略)を利用した技術が、高精度な追尾技術として注目を集めている。

【0003】

非特許文献1は、画像内の特定の被写体を追尾するための方法の1つである。追尾対象が映った画像と、探索範囲となる画像を、重みが同一のConvolutional Neural Network(以下CNNと省略)にそれぞれ入力する。CNNから得られたそれぞれの特徴量同士の相互相関を算出することによって、探索範囲の画像中で追尾対象が存在する位置を特定するものである。このような追尾手法は追尾対象の位置を正確に同定できる一方、追尾対象に類似した物体が画面の上で重なるような場合に、誤った対象を追尾する失敗が発生し易い。

【0004】

これを回避するために特許文献1の手法に代表されるように、検出物体の領域の色特徴や奥行き情報からヒストグラムを作成し、その変化等を調べて物体が遮蔽されているか否かを判定する手法がある。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】米国特許出願第10185877(B2)号広報

【非特許文献】

【0006】

【非特許文献1】Bertinetto et al., Fully-Convolutional Siamese Networks for Object Tracking, arXiv 2016

【発明の概要】

【発明が解決しようとする課題】

【0007】

しかしながら、特許文献1に示される方法では、同じような姿勢の物体や外見的特徴の類似した物体が画面上で重なると、色やテクスチャといった特徴量のヒストグラムに差異が出にくいいため判定できないという課題がある。例えば、スポーツの集団競技等においては、狭い範囲に存在する複数の人物の服装や姿勢が同一になることも多く、異なる人物を同じ人物と見なして追尾する失敗が起こりうる。本発明は、このような課題に鑑みなされたものであり、外見的特徴や姿勢が類似した物体が近接する場合においても安定して追尾を継続することを目的とする。

【課題を解決するための手段】

【0008】

上記課題を解決する本発明にかかる情報処理装置は、画像から少なくとも1つ以上の物体を検出する情報処理装置であって、遮蔽する物体と遮蔽された物体との遮蔽関係を示す画像特徴を学習した学習済みモデルに基づいて、前記画像から検出された各物体について、前記画像から検出された他の物体との遮蔽関係を推定する推定手段と、前記推定手段によって推定された遮蔽関係に基づいて、前記画像から検出された各物体について、前記画像と異なる時刻に撮像された画像において検出された物体との対応関係を特定する特定手段と、を有する。

【発明の効果】

【0009】

本発明によれば、外見的特徴や姿勢が類似した物体が近接する場合においても安定して追尾を継続できる。

【図面の簡単な説明】

【0010】

【図1】物体検出の一例を説明する模式図

【図2】情報処理装置のハードウェア構成例を示す図

【図3】情報処理装置の機能構成例を示すブロック図

10

20

30

40

50

【図 4】情報処理装置が実行する処理手順を示すフローチャート

【図 5】情報処理装置の処理の結果例を示す図

【図 6】情報処理装置の処理の結果例を示す図

【図 7】情報処理装置が実行する処理手順を示すフローチャート

【図 8】遮蔽に関する情報の派生の例

【図 9】情報処理装置の機能構成例を示すブロック図

【図 10】情報処理装置の処理の結果例を示す図

【図 11】情報処理装置が実行する処理手順を示すフローチャート

【図 12】情報処理装置の学習処理の例を示す図

【図 13】情報処理装置の処理の結果例を示す図

10

【図 14】情報処理装置の処理の結果例を示す図

【図 15】情報処理装置の処理の結果例を示す図

【図 16】情報処理装置の学習処理の例を示す図

【発明を実施するための形態】

【0011】

<実施形態 1>

実施形態に係る情報処理装置を、図面を参照しながら説明する。なお、図面間で符号の同じものは同じ動作をすることで重ねての説明を省く。また、この実施の形態に掲載されている構成要素はあくまで例示であり、この発明の範囲をそれらだけに限定する趣旨のものではない。

20

【0012】

本実施形態では、動画もしくは連続撮影した静止画フレームから人物を検出し、追尾する機能について説明する。適用範囲は検出・追尾対象の物体のカテゴリを限定しないが、本実施形態 1 は対象を人物に限定する。本実施形態では、時間的に連続する画像毎に人物を検出し、連続する画像間でそれぞれの人物がどの人物と同一人物であるかを対応付けることで、人物の追尾を実現する。本実施形態では特に、スポーツイベントなどの撮影を想定し、人物の服装や移動方向等が類似しており、高頻度で近接・交差するとする。このような場合、各画像における人物の位置または服装の色といった外見的特徴に近い人物同士を対応付けるだけでは、誤った対応付けが発生しやすい。このような失敗をここでは誤マッチングと呼ぶ。

30

【0013】

本実施形態では撮影者から見て物体が重なっている時の、遮蔽関係のパターンを学習した学習済みモデルが出力する遮蔽に関する情報に着目する。学習済みモデルによって出力された遮蔽関係を、物体の対応関係の特定に併せて用いることで、手前にいる人物と奥にいる人物同士を対応付ける失敗を抑制し、追尾の精度を向上する。

【0014】

これを模式的に示した図が図 1 である。図 1 (A) の画像 2100, 22120, 2140 は同一の絵柄の 2 枚のトランプカードがテーブル上で交差していく様子を上から写した動画の 3 フレーム分の静止画を示している (時系列順に左から右に並んでいる)。画像 2100, 2120, 2140 を観察しただけでは各画像におけるカードがそれぞれどのように移動したかを確定することができない。一方で図 1 (B) は (A) よりも高フレームレートで同じ様子を撮影した例である。つまり、より短い時間間隔で撮像された画像群である。図 1 (B) の画像を時系列順に観察していけば、どちらのカードが次の画像でどこに移動したかを対応付けることができる。全体としては左側のカード 2201 が右側のカード 2202 の上を通過し、右側に移動したということを比較的容易に推定することができる。画像 2210 や画像 2230 に示すように、物体同士の交差の瞬間に過渡的に生じる見えを観察することで、画像 2220 においてどちらのカードが手前側を通過し、どちらが奥側にあるのかが判定可能となる。この判定に際しては、2.5 次元の奥行画像といった特別なセンサーやオプティカルフロー等の生成のコストの高い情報は必ずしも必要でない。物体同士が手前と奥で重なったときに、どのような見えが生じ易いかという、遮蔽関

40

50

係と見えの特徴 ( appearance feature ) とのパターン認識の問題として解くことができる。これは図 1 ( C ) および図 1 ( D ) に示す人物の交差のようなシーンでも同様である。本図 1 ( C ) ( D ) では人物の服装や姿勢等の見え、移動方向は同一であるとするとする。このような場合も、物体が交差する前後の見えの状態に着目して観察すれば、図 1 ( D ) の画像 2 4 2 0 では人物 2 4 0 1 が手前側に、人物 2 4 0 2 が奥側にいると判定する。以降の画像において、この遮蔽関係を維持したままであれば、人物 2 4 0 1 が人物 2 4 0 2 を一度遮蔽した場合に、手前側の人物 2 4 0 1 を追尾し、奥側の人物 2 4 0 1 の遮蔽関係と画像特徴を保持する。そして、遮蔽が解消したときには、人物 2 4 0 1 の追尾を継続しつつ、奥側にいた人物 2 4 0 2 を再び検出することが可能である。以上が本実施形態の原理の概要を示す説明である。詳細な処理については後述する。

10

## 【 0 0 1 5 】

図 2 は、本実施形態における、画像認識によって追尾対象を追尾する情報処理装置 1 のハードウェア構成図である。CPU H 1 0 1 は、ROM H 1 0 2 に格納されている制御プログラムを実行することにより、本装置全体の制御を行う。RAM H 1 0 3 は、各構成要素からの各種データを一時記憶する。また、プログラムを展開し、CPU H 1 0 1 が実行可能な状態にする。記憶部 H 1 0 4 は、本実施形態の処理対象となるデータを格納するものであり、追尾対象となるデータを記憶する。記憶部 H 1 0 4 の媒体としては、HDD、フラッシュメモリ、各種光学メディアなどを用いることができる。入力部 H 1 0 5 は、キーボード・タッチパネル、ダイヤル等で構成され、ユーザからの入力を受け付けるものであり、追尾対象を設定する際になどに用いられる。表示部 H 1 0 6 は、液晶ディスプレイ等で構成され、被写体や追尾結果をユーザに対して表示する。また、本装置は通信部 H 1 0 7 を介して、撮影装置等の他の装置と通信することができる。

20

## 【 0 0 1 6 】

図 3 は、情報処理装置の機能構成例を示すブロック図である。図 3 では CPU H 1 0 1 において実行される処理を、それぞれ機能ブロックとして示している。情報処理装置 1 は、画像取得部 2 0 1、物体検出部 2 0 2、遮蔽情報生成部 2 0 3、抽出部 2 0 4、対応付け部 2 0 5 を有し、外部の記憶部 2 0 6 に接続されている。記憶部 2 0 6 は情報処理装置 1 の内部にあってもよい。それぞれの機能を簡単に説明する。画像取得部 2 0 1 は、撮像装置によって特定の物体 ( 本実施形態では人物 ) を撮像した動画や連続静止画の画像を取得する。物体検出部 2 0 2 は、画像取得部 2 0 1 によって取得された画像から予め設定された所定の物体を示す画像特徴を検出する。例えば、さまざまな姿勢の人物の画像を用いて人体 ( 頭や動体 ) を示す画像特徴を予め学習した学習済みモデルに基づいて、画像における人物の領域を検出する。遮蔽情報生成部 2 0 3 は、遮蔽する物体と遮蔽された物体との遮蔽関係を示す画像特徴を学習した学習済みモデルに基づいて、画像から検出された各物体について、画像から検出された他の物体との遮蔽関係を示す遮蔽情報を推定する。遮蔽情報とは、注目物体が他の物体によって遮蔽されている可能性を表す尤度 ( 被遮蔽 / 遮蔽スコア ) である。例えば、ある物体について、他の物体によって遮蔽されている可能性が高ければ、被遮蔽 / 遮蔽スコアを 1 に近づける。ある物体について、他の物体を遮蔽している可能性が高ければ、被遮蔽 / 遮蔽スコアを 0 に近づける。このような遮蔽関係を示す被遮蔽 / 遮蔽スコアを、学習済みモデルを用いて推定する。抽出部 2 0 4 は、ある画像

30

40

## 【 0 0 1 7 】

図 4 は本実施形態の処理の流れを示したフローチャートである。以下の説明では、各工程

50

(ステップ)について先頭にSを付けて表記することで、工程(ステップ)の表記を省略する。ただし、情報処理装置はこのフローチャートで説明するすべての工程を必ずしも行わなくても良い。図4のフローチャートに示した処理は、コンピュータである図2のCPU H101により記憶部H104に格納されているコンピュータプログラムに従って実行される。

#### 【0018】

S301では、情報処理装置1が、各動画フレームについて繰り返すループ処理を開始する。S302では、画像取得部201が人物を撮像した動画や連続静止画の画像フレームを順次取得する。以降の処理はS301~S311まで各画像について順次処理がなされる。なお、画像取得部201は、情報処理装置に接続された撮像装置によって撮像された画像を取得してもよいし、記憶部H104に記憶された画像を取得してもよい。図5(A)中の動画フレーム3100, 3110, 3120, 3130, 3140が取得した画像フレームの例である。

10

#### 【0019】

次にS303では、物体検出部20が、所定の物体(ここでは人物)の画像特徴に基づいて、前記取得された画像から少なくとも1つ以上の所定の物体を検出する。画像内から物体を検出する公知技術としては、Liuによる手法等が挙げられる(Liu, SSD: Single Shot Multibox Detector, In: ECCV 2016)。画像内から候補物体を検出した結果を図5(A)に示す。図5(A)中の矩形枠3101, 3102, 3103, 3111, 3112, 3113, 3121, 3122, 3131, 3132, 3141, 3142が検出された物体領域を示すBounding Box(以下BB)である。

20

#### 【0020】

S304では、遮蔽マップ生成部203が、各画像について、領域毎に遮蔽されているか否か(遮蔽関係)についての遮蔽情報を示した遮蔽マップを生成する。遮蔽マップ生成部203が、各画像について、遮蔽されている物体のうちの見えている領域(被遮蔽物体領域)を推定する。ここでは各人物が他の人物と重なっているか、重なっている場合に奥側にいるか、手前側にいるかを判定し、その結果を遮蔽状態のスコア(尤度)として領域ごとに出力する。これは意味的領域分割の認識タスクの一種であり、Chenらの手法等の公知の手法を使って実現することができる。(Chen, DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs, 2016)。

30

#### 【0021】

図6(A)に遮蔽マップの生成処理を説明する模式図と結果の一例を示す。ニューラルネットワーク402は入力された画像から、入力画像の各画素について遮蔽状態を判定するニューラルネットワークである。RGB画像401が入力されると、ニューラルネットワーク402は画像中に人物がいるか否か、さらにその人物が遮蔽されているか否かを推定した結果を遮蔽マップ404として出力する。同マップは遮蔽されていない人物および人物以外の領域と推定された場合は0、遮蔽されている人物の領域には1、の被遮蔽スコアが出力される。遮蔽マップ404中の黒い領域ほど高い被遮蔽スコアであることを示す。すなわち、黒い領域は遮蔽された人物の領域であると推定されたことを示している。ニューラルネットワーク402は入力画像に対してこのような出力ができるように事前に学習を行っている(学習については後述する)。なお、図に示した遮蔽マップ404は推定結果として理想的な出力状態の一例を示したものである。

40

#### 【0022】

なお、RGB画像401のほかに、専用センサー等を使って2.5次元奥行画像405を別途取得するような派生的な形態も考えられる。前記奥行画像405を3チャンネルのRGB画像401と連結した4チャンネルの情報をRGB画像の代わりに画像入力として学習・認識する。これにより遮蔽領域の情報をより高精度にすることも可能である。

50

## 【 0 0 2 3 】

次に、S 3 0 5では、情報処理装置1が、S 3 0 3で検出された各物体について、S 3 0 6からS 3 0 7のループ処理を実行する。S 3 0 5 ~ S 3 0 8では、抽出部2 0 4が、生成された遮蔽マップから検出物体ごとに遮蔽関係を示す情報を抽出し、記憶部2 0 6に記憶する。S 3 0 6では、抽出部2 0 4が、遮蔽マップから検出物体毎に遮蔽関係を示す情報を抽出する。具体的には、図6 ( A )の人物検出枠4 0 7の中の遮蔽マップ4 0 4の被遮蔽スコアを平均する。この被遮蔽スコアが1に近いほどその物体は遮蔽されている可能性が高く、被遮蔽スコアが0に近いほどその物体は遮蔽されていない可能性が高いことを示す。なお、検出枠の位置のずれや、遮蔽マップ4 0 4にノイズが含まれることを想定して、図6 ( B )に示すように枠の中央付近を重視した重み付き平均で取得する。図中の演算4 1 2は各画像の部分領域毎ごとの要素積 ( アダマール積 ) を意味する。マップ4 1 3は中央にピークがあり、画像ブロックの総和が1となる2次元ガウス関数である ( 縦横サイズを人物検出枠に合わせて変形してある )。取得結果の例を図6 ( A )に記号o c cを付して被遮蔽スコア値4 0 9と4 1 0として示す。左側の検出枠は奥側にいる人物のため被遮蔽スコアが高く、右側の検出枠は手前側のため被遮蔽スコアが低いと判定されている。以上のような処理を先ほどの図5の入力画像に対して処理した結果例を図5 ( B ) ( 遮蔽マップ ) および図5 ( C ) ( 各枠の被遮蔽スコア推定結果 ) として示す。交差開始 ~ 終了の間、奥側に位置する人物3 1 0 1に対応する被遮蔽スコアは人物3 1 0 2のそれよりも相対的に高いことを示している。

10

## 【 0 0 2 4 】

S 3 0 7では、記憶部2 0 6が、抽出部2 0 4によって取得された各検出物体の被遮蔽スコアを記憶する。同時に、各検出物体の位置・サイズの情報、および色やテクスチャのヒストグラムといった物体の見えに関する特徴量、も記憶する。ここではこれら複数種類の数量を一括して検出物体の特徴量と呼ぶ。なお、見えに関する特徴量としてはこの他にニューラルネットワークの中間層情報等を利用してよい。(例えば " Hariharan , et . al , Hypercolumns for Object Segmentation and Fine-grained Localization , in CVPR 2015 ")。

20

## 【 0 0 2 5 】

S 3 0 8では、情報処理装置1が、各画像について繰り返すループと、各画像において検出された人物について繰り返すループを終了する。このループは画像毎に、その画像から検出された人物すべてについて遮蔽情報を取得したときに終了する。次に、S 3 0 9 ~ S 3 1 0では、対応付け部2 0 5が、前後の画像間の物体の対応付けを行う ( ただし一つ目の動画フレームの場合は過去のフレームがないためこれを行わない )。まず、S 3 0 9で、対応付け部2 0 5が記憶部2 0 6に記憶された過去の物体の特徴量である、被遮蔽スコア、位置サイズおよび見えの特徴量を取得する。次にS 3 1 0で、対応付け部2 0 5が過去の動画フレーム中に検出された物体と、現在処理しているフレーム中に検出された物体の対応付けを行う。

30

## 【 0 0 2 6 】

S 3 1 0における対応付け部2 0 5の詳細な処理フローを図7に示す。S 5 0 1で、対応付け部2 0 5は、まず現フレームの検出物体と一つ前のフレームで検出された物体の間で全組み合わせのペアを作る。前後のフレームでそれぞれn人とm人の人物が検出されていれば、全部でn x m個のペアが生成される。次に、S 5 0 2で、対応付け部2 0 5は全ての物体ペアについて類似度を算出する。類似度としては検出物体同士の特徴量の差分に基づいた指標を用いることができる。一例として過去の検出物体c 1と現在の検出物体c 2の類似度を下式のように算出する。

40

( 数式 1 )

$$L ( c_1 , c_2 ) = \begin{matrix} - W_1 & | & | & B B_1 - B B_2 & | & | \\ - W_2 & | & | & f_1 - f_2 & | & | \\ - W_3 & | & | & o c c_1 - o c c_2 & | & | \end{matrix}$$

ここで、B Bとは各物体の ( 中心座標値x、中心座標値y、幅、高さ ) の4変数をまとめ

50

たベクトルであり、 $f$  は各物体の特徴を示したものである。 $\|x\|$  は  $x$  の  $L^p$  ノルムである。 $occ$  は各物体の被遮蔽スコアである。 $W_1, W_2, W_3$  はそれぞれ経験的あるいは機械学習的に調整して設定される 0 以上のバランス係数である。ここで各特徴量のばらつきを事前に統計的に求めておいて各特徴量を正規化する等してもよい。物同士が交差する場合であっても、他の物体を遮蔽する側の人物を追尾することによって、被遮蔽側の人物が再び画像で確認されたときに、直前で他の物体を遮蔽する側の人物と対応付けようとすると数式 1 の 3 つめの項の値が小さくなり、類似度が低く算出される。つまり、この処理によって、遮蔽関係が異なる物同士はマッチングされる可能性が低くなり、追尾の誤マッチングが抑制できる。

【0027】

10

次に、S503において、対応付け部205が、過去の物体と現在の物体との類似度に基づいて物体間の対応関係を特定するための対応付け（マッチング）を行う。マッチングの方法にはいくつか存在する。例えば、類似度が高い候補同士から優先的にマッチングする方法や、ハンガリアンアルゴリズムを用いる方法等がある。ここでは前者を用いる。

【0028】

S503では、対応付け部205が、現フレームの全物体について対応付けが終了していなければS506で類似度最大のペアから同一人物として対応付けていく。対応付けの終わったペアの物体は対応付けの候補から省いていく。上記の処理の際に、その時点で残っているペアの中の最大の類似度の大きさが所定の閾値を下回った場合は、もはや類似した物体ペアが残っていないことを意味する。その場合はそれ以上無理に対応付けることなく（S505）、対応付けを終了する。

20

【0029】

以上の処理S301～S311を動画フレームごとに行う。その結果、図5(D)に結果例を示すように、動画中から人物を検出し、それぞれの物体がどこに移動したかの一連の追尾結果が得られる（フレーム間の同一の人物に記号A, B, Cで付して追尾の結果を示している）。

【0030】

<変形例>

本実施形態では物体ペア同士のマッチングの類似度として差分に基づき、被遮蔽スコアや見えといった各指標の距離を重み付け和した。ここで例えばKLダイバージェンスを使うことも考えられる。またメトリック学習を行ってより精度の高い距離指標を求めることも考えられる。また単一の類似度を一度だけ用いるのではなく、まず見えの特徴で類似度を判定し、条件を満たしたものは次に遮蔽状態のスコアの類似度に基づいて判定する、等のルールベースによる方法や段階的な判定方法も考えられる。またさらにニューラルネットやサポートベクトルマシンといった公知の識別器の手法を用い、説明変数を特徴量、目的変数を同一物体か否かの結果、として学習・識別し、この値によってマッチングを判定することも可能である。以上のようにフレーム間の物体間の対応付けは特定の形態に限定されない。

30

【0031】

またさらに別の派生形態として、被遮蔽スコアの推定値を安定させるために、下式のように過去のスコアを移動平均した値を用いる工夫も考えられる。

40

(数式2)

$$occEMA(t) = (1 - \alpha) \times occEMA(t-1) + \alpha \times occ(t)$$
 上式は指数移動平均値と呼ばれる値であり、 $occEMA(t)$  は時刻  $t$  の被遮蔽スコアの指数移動平均値、 $occ(t)$  は時刻  $t$  の被遮蔽スコア、 $\alpha$  は  $0 < \alpha < 1$  の係数である。過去の複数フレームで追尾ができていない物体については上式で指数移動平均値を算出しておき、類似度を比較する際には元の被遮蔽スコアではなく、指数移動平均被遮蔽スコアを用いる。これにより、交差時に複数のフレームにまたがって徐々に重畳状態が起こるような場合に、複数フレームの被遮蔽スコアの平均値に基づいてマッチングできるので、より物体間の対応付けが安定する。

50

## 【 0 0 3 2 】

またさらに別の派生形態として、マッチングの際に前後フレーム間の類似度だけでなく、 $n$ ステップ前の過去の複数のフレームの特徴量・位置を用いてマッチングを行うような形態も考えられる。この方法を用いることで、一度物体が遮蔽されて追尾できないフレームが発生しても、その後のフレームで遮蔽が解消されれば再び追尾が可能になる。この形態では例えば、 $n$ フレームまでさかのぼって物体の特徴量の平均値を求め、これに基づいて現フレームから検出された物体との類似度の算出を行う。もしくは、過去の $n$ フレームの物体と現フレームの物体間でそれぞれ類似度を求め、得られた $n$ 個の類似度の平均値が最も高い物体に対応付ける。また、過去だけでなく、 $n$ ステップの未来のフレームの結果も使って双方向で判定を行うことも考えられる。この形態は未来のフレームを処理するまで結果が判明しないため処理のリアルタイム性には劣るが、過去のみを見る方法よりも高精度である。

10

## 【 0 0 3 3 】

またさらに別の派生形態として、検出の失敗に対応するための形態が考えられる。物体検出・追尾においては物体の姿勢が特殊な形状に変化した、等の理由で物体検出の段階で一時的に失敗するようなことも起こり得る。このような未検出が起こると、フレーム間の対応付けの際に、前のフレームに存在した物体が、現フレームでは対応なしと判定される。すると追尾はそこで途切れることになる。このような失敗を防ぐために、以下のような工夫もありえる。すなわち、マッチングで未対応の人物が発生したら、その情報をリストに記憶しておき、次のフレームのマッチングのときに対応付けの候補に加える（一定時間が経過してもまだ未対応であれば物体自体がもう存在しないと判断し、リストから除去する。ここではこれをタイムアウト処理と呼ぶ）。

20

## 【 0 0 3 4 】

このように動画フレームをまたがる物体の対応付けについては種々のやり方が考えられ、特定の形態に限定されない。

## 【 0 0 3 5 】

< 遮蔽情報の形態のバリエーションおよび学習方法 >

本実施形態では、遮蔽マップとして、遮蔽されている物体のうちの見えている領域（被遮蔽物体領域）を推定した。この形態についても様々な派生形態が考えられる。一例を図8に示す。ここでは図8（B）に示すように、画像801のように奥側の物体の見えている領域を推定する以外でもよい。例えば、画像802のように奥側の物体の全領域を推定する。また、画像803のように、手前側の遮蔽物体の領域を推定する（図の領域440のように他物体と重なっていない物体も手前側領域として含めて推定している。ただし別の形態としてこのような単独の物体は手前側の領域に含めないことも考えられる）。また、画像801～803のように前景領域を推定するのではなく、画像804のように物体の中心や重心の位置を推定することも考えられる。画像804の場合においては被遮蔽物体の中心付近の領域に大きな正の値を、遮蔽物体の中心付近に小さな負の値を推定するようにする（ここでいう物体の中心領域は図示するようにガウス関数状の領域を推定させるような形態が考えられる）。

30

## 【 0 0 3 6 】

ここで遮蔽状態の情報の学習方法について図6（D）を用いて説明する。前述のChenらの手法等で示されるニューラルネット402は、入力画像であるRGB画像401に対して遮蔽物体の被遮蔽スコアマップ403を出力する。403の結果例を430に示す。CHENらの手法等は特定カテゴリ物体の前景領域を推定する手法であるが、ここでは遮蔽情報の教師値431を与えて、教師値431と同じようなマップが推定によって得られるようニューラルネット402の学習を行う。具体的には出力結果のマップ403と教師値431を比較し、交差エントロピーや二乗誤差などの公知の方法で損失値算出432を行う。損失値が漸減するように誤差逆伝搬法等でニューラルネット402の重みパラメータを調整する（この処理についてはChenらの手法と同一のため詳細は略す）。入力画像と教師値は十分な量を与える必要がある。重なった物体の領域の教師値を作成するの

40

50

はコストがかかるため、CGを用いることや、物体画像を切り出して重畳する画像合成の方法を用いて学習データを作成するようなことも考えられる。以上が学習方法になる。

【0037】

またさらに、本実施形態では上記で求めた物体の枠の中で取得して被遮蔽スコアと呼ぶ指標を求めた。遮蔽情報の取得の形態の様々な例を図8(C)に示す。図8(C1)は本実施形態の形態である。この他に、(C2)奥側の被遮蔽スコアと手前側の被遮蔽スコアの差分値を物体枠内で取得する、(C3)物体の中心のスコアを1点だけ参照する、等様々な考えられる。また、枠内で取得する際に、物体の枠内で取得する際に、他の物体枠と重なっている領域についてはどちらの物体の領域か判然としないために取得から省くような方法も考えられる。

10

【0038】

またさらに、上述の<遮蔽状態の推定>と<各物体の被遮蔽スコアの取得>を同時に行う方法も考えられる。例として、Liuの手法等で使われている公知な方法であるアンカーと呼ばれる手法があげられる。この手法では物体の候補枠の集合が求められるので、これを利用して各候補枠が遮蔽物体か被遮蔽物体かの被遮蔽スコアを推定し対応付けることが考えられる(この形態の詳細については実施形態3で示すのでここでは説明を略す)。

【0039】

またさらに、上で示したような複数の形態の遮蔽情報をそれぞれ取得し、これを遮蔽に関する多次元の特徴として後段の物体の対応付けに用いてもよい。もしくは前記の遮蔽に関する多次元の特徴から機械学習によって物体の遮蔽されている面積の割合を推定して用いてもよい。この場合は前記の遮蔽に関する多次元の特徴を説明変数とし、物体が遮蔽されている面積の割合を目標変数とし、ロジスティック回帰等の公知技術で回帰推定を行う等すればよい。

20

【0040】

<実施形態2>

本実施形態では実施形態1と同様に人物の検出と追尾を行う。ハードウェア構成は実施形態1の図2と同様である。本実施形態における機能構成例を示すブロック図は図9(A)になる。実施形態1の構成に新たに遮蔽状態判定部301が追加されている。実施形態1では追尾中に人物の枠は人物同士の重なりによって、人物の検出ができないことがある。例えば図10(A)中の動画フレーム4120に示すように、人物間で重なった面積が大きいときには、奥側の人物が検出できないことは多い。このような時に遮蔽状態判定部301が、人物は存在しているが被遮蔽状態にある、と判定する。

30

【0041】

実施形態1で説明したような物体検出部の一時的な検出の失敗による未検出と異なり、人物の集団が同じ方向に同じ速度で移動しているような場合、長時間未検出の状態が続く。さらに被遮蔽状態から脱した画面上の位置が、被遮蔽状態が開始した位置から離れることがある。このため被遮蔽状態であると積極的に判定し、推定した前記状態に応じた処理を行うことで追尾の成功率を高めることが望ましい。

【0042】

本実施形態も全体の処理フローは実施形態1の図4と同じであるが、S310の処理の詳細が下記のように異なる。ここでは、実施形態1と異なるS310の処理についてのみ説明する。図11を用いて遮蔽状態判定部301が行うS310処理の詳細なフローについて説明する。まずこれまでと同じようにS601で現フレームと前フレームで物体の対応付けを行う。S602で対応付けられなかった前フレームの物体がある場合、被遮蔽状態に入った可能性がある。そこでS603で当該物体のそれまでの被遮蔽スコアの高さが閾値以上かを調べる。これは動画フレームのフレームレートが十分に高ければ、遮蔽により未検出になる前後で被遮蔽スコアが高くなることが多いためである。さらにS604で当該物体の周辺領域で現フレームの物体の検出数の数が減っていないかを調べ、上記の二つの結果が真であれば当該物体は被遮蔽状態に入ったと推定し被遮蔽状態のリストに記憶する(S605)。被遮蔽状態のリストに記憶された物体については前回検出されたときの

40

50

特徴量と位置も合わせて記憶する。これによって、遮蔽が解消されて再び検出されたときに追尾できる可能性が向上する。

【 0 0 4 3 】

S 6 0 6 ~ S 6 1 0 は被遮蔽状態の物体が再出現したかどうかを判定する処理である。S 6 0 3 で対応付けられなかった現フレームの物体がある場合、被遮蔽状態を脱して再度検出できるようになった可能性がある。そこで S 6 0 7 で当該物体の被遮蔽スコアの高さが閾値以上かを調べる。さらに S 6 0 8 で当該物体の周辺領域で現フレームの物体の検出数の数が増えていないかを調べる。両方の結果が真で、且つ被遮蔽状態のリストに記憶されている物体のいずれかと当該物体が所定閾値以上に類似度が高い場合 ( S 6 0 8 )、当該物体は被遮蔽状態から脱して再度出現したと推定する。そのとき、対応付けた物体を被遮蔽状態のリストから除去する ( S 6 0 9 ) 被遮蔽状態のリストから除去された物体については、現在の入力画像から検出された特徴量と位置を取得する。

10

【 0 0 4 4 】

ここで、対応付けの処理の工夫として、例えば、フレーム間の物体のマッチングの際に、被遮蔽状態にある人物とのマッチングは距離による類似度のペナルティを減ずる。再出現を待つタイムアウトの時間を長く取る。遮蔽状態の物体との対応付けの閾値は、通常の物体間のマッチングよりも閾値を低く設定する、等が考えられる。

【 0 0 4 5 】

またさらに、ここでは二人の人物の重なりを想定して説明を行ったが、3人以上の人物の間で重なりが生じることもある。この場合は、遮蔽状態に入ったと判定されれば被遮蔽状態のリストに加えておき、再出現したら前フレームとの対応付けを行い、被遮蔽状態のリストから都度除去する。これにより3人以上についてもある程度の対応が可能である。

20

【 0 0 4 6 】

< 実施形態 3 >

本実施形態では、ユーザが指定した単一の物体を追尾する形態について説明する。ここでは追尾対象は人体等の特定カテゴリに限らず、ユーザが指定した不特定の物体を追尾する形態を扱う。例えば、犬などの動物や、車などの乗り物であってもよい。

【 0 0 4 7 】

機能ブロックの図は図 9 ( B ) になる。これまでの構成に新たに追尾物体指定部 3 0 2 が追加されている。ここで追尾物体指定部 3 0 2 と物体検出部 2 0 2 の機能は非特許文献 1 の方法を用いることで容易に実現することができる。追尾物体指定部 3 0 2 はユーザが動画フレーム中で追尾対象物体の枠位置を指定する機能部である。これにより追尾すべき物体の特徴が初期化される。物体検出部 2 0 2 は各動画中で最も対象物体と一致度の高い画像領域を同定する。同定した結果例を図 1 2 ( A ) に示す。図 1 2 ( A ) の動画フレーム 5 1 1 0 上の枠 5 1 1 1 がユーザによって指示された追尾物体の枠である。動画フレーム 5 1 2 0 ではこの物体が画面中で右側に移動しており、物体検出部 2 0 2 によって枠 5 1 2 1 として検出されている。非特許文献 1 の方法は物体の追尾手法として優れるが、類似物体間で容易に誤スイッチが生じる。そこで本実施形態ではこれまでの実施形態と同様に、追尾物体に対して遮蔽状態に関する情報を推定し、誤スイッチが生じていないかを判定する。

30

40

【 0 0 4 8 】

このために遮蔽情報生成部 2 0 3 として図 1 3 ( B ) に示すようなニューラルネット 6 3 0 0 を用いる。これは検出された追尾物体の画像 6 3 0 1 (ここでは処理の簡単のために正方形の画像に縦横比率を正規化している) を入力すると、画像パターンを見て、遮蔽されている ( Y e s ) がされていない ( N o ) かの分類結果 6 3 0 2 を出力する分類器である。遮蔽の有無の定義としては、物体の面積が何%以上遮蔽されているか否かとして定義する。この 2 クラスの値を教師値として与えてニューラルネット 6 3 0 0 を学習させる。この技術は通常の画像分類タスクと同様の広く公知な方法のため詳細を略す。また、教師値 ( 目標変数 ) を遮蔽の有無の 2 値ではなく遮蔽面積の割合として与えて回帰学習を行えば、推定結果 6 3 0 3 のように遮蔽の割合を推定することができる。この回帰学習には学

50

習時に与える損失値として二乗誤差等を用いる。

【 0 0 4 9 】

遮蔽情報生成部 2 0 3 で追尾物体候補の遮蔽度を推定した結果が図 1 4 ( A ) ( B ) である。図 1 4 ( A ) に示す物体の検出結果に対して、物体検出部 2 0 2 が図 1 4 ( B ) に符号 o c c を付して示したのが被遮蔽面積の推定値である。同図では被遮蔽スコアの変動幅は所定値（例えば 0 . 3 等の値）より小さく、追尾に失敗していないと判定できる（ここで、被遮蔽スコアだけでなく実施形態 1 で用いたような位置や見えの特徴量の類似度も併用して追尾の成功・失敗を判定してもよい）。

【 0 0 5 0 】

一方で図 1 4 ( C ) では、動画フレーム 7 2 2 0 から 7 2 3 0 にかけて物体 7 2 0 1 が物体 7 2 0 2 の向こう側を通過しており、その結果、物体検出部 2 0 2 が動画フレーム 7 2 3 0 における物体の位置を枠 7 2 3 1 として誤って推定している。この場合の遮蔽スコアは図 1 4 ( D ) に示すように 0 . 4 から 0 . 0 へと大きく変動しているため、交差によって誤追尾が発生したと判定することができる。誤追尾が発生したことが分かれば、そこで検出を止めたり、後段で修正する等の工夫を行うことができる。

10

【 0 0 5 1 】

以上が本実施形態の説明となる。

【 0 0 5 2 】

なお、遮蔽情報生成部 2 0 3 の学習は図 1 2 ( A ) 5 1 1 0 ~ 5 1 5 0 に示すように、不特定の物体について遮蔽状態が判定できるように様々な物体の遮蔽状態を推定できるように学習しておくことが望ましい。

20

【 0 0 5 3 】

なお他の派生の形態としては、図 1 3 ( B ) では、物体枠で切られた画像 6 3 0 1 を入力画像として示している。しかし、被遮蔽状態にあるか否かの判定には当該物体だけでなくその周辺を観察することが重要なため、入力画像としてはより広い範囲を入力することも考えられる（その場合、推定時にも同様の範囲を切り取って入力する）。

【 0 0 5 4 】

なお他の派生の形態としては、図 1 3 ( C ) に示すように、上述の L i u の手法のようなアンカーと言われる候補枠を使って物体の検出と遮蔽度の推定を同時に行う形態も考えられる。アンカー枠は図 1 3 ( D ) に示すような複数のサイズ・縦横比率の候補枠の集合である（ここでは 3 種類のアンカー枠を図示している）。アンカー枠は図 1 3 ( C ) の結果画像 6 4 5 0 に示すように、画像中の各ブロック領域に配置されている。ニューラルネット 6 4 0 0 は画像が入力されたら、各ブロック領域の各アンカーに当該物体があるか否かの被遮蔽スコアマップ 6 4 3 0 を生成する。被遮蔽スコアマップ 6 4 3 0 はアンカー枠の種類 3 個に対応した 3 枚のマップである。推定結果の例を図 1 3 ( C ) 6 4 5 0 に示す（以上の手法は広く公知のため詳細は上述の L i u の方法を参照されたい）。

30

【 0 0 5 5 】

ここで本実施形態の派生の形態として、物体が存在するか否かの推定と同時に、物体の被遮蔽スコアマップ 6 4 4 0 を生成する。これは各アンカー枠に、もしそこに物体がある場合、その被遮蔽割合がいくつになるかを推定したマップである。同マップもアンカーの種類の数に対応した 3 枚からなる（学習時には画像の各ブロックにおいて、各アンカー枠に被遮蔽スコアの教師値を与えてニューラルネット 6 4 0 0 を学習すればよい）。結果例を図 1 3 ( C ) 6 4 6 0 に示す。二つの推定マップを最終的に統合した例を統合結果例 6 4 7 0 として図示する。

40

【 0 0 5 6 】

上記の説明は物体検出の例になるが、非特許文献 1 の方法もアンカー候補枠ベースの手法であるため、物体を追尾しながら同時にその被遮蔽スコアを推定する派生形態を構成することが可能である。

【 0 0 5 7 】

< 実施形態 4 >

50

本実施形態では、ユーザが指定した単一の物体を追尾する形態について説明する。機能ブロックの図は実施形態3と同じで図9(B)である。これまでの実施形態では類似度を比較する際に、直前と直後のフレームで特徴量を比較することや、前後のnフレームを用いて比較すること等、ルールベースでフレーム間の物体の対応付けを行った。本実施形態では、この部分を機械学習に置き換えることでより精度の高い対応付けを行う。

#### 【0058】

リカレントニューラルネットは時系列データを処理して識別・分類等を行うことができる技術であり、Byeonらの方法などで公知なLong short term memoryネットワーク(以下LSTM)が代表的手法である。(Byeon et al., Scene labeling with LSTM recurrent neural networks, CVPR 2015)。当該手法で物体の特徴の経時的な変化を判別して物体間の対応付けを行うことができる。本実施形態の構成と結果例の模式図を図15に示す。ここでは1つの物体9102が追尾対象として指定され、Bertinettoら等の手法で追尾されている(t=2の動画フレームで誤スイッチが起こっている)。図15(C)のLSTMユニット9501~9504は、各時刻で追尾している物体の特徴9401~9404を受け取って、追尾が成功しているか、失敗しているかを判定して出力9701~9704として出力する。ここでは図示上LSTMユニットを複数書いているが、ここでは複数のユニットが存在するのではなく同一のユニットの各時刻の状態を示している。各時刻のLSTMユニットは次の時刻のLSTMユニットにリカレント入力9802を送る。LSTMユニットはその時点の物体の特徴とそれまでの過去の情報を含むリカレント入力9802を元に内部状態を必要に応じて変更する。これにより物体のパターンが経時的にどのように変化しているかを踏まえた上で現時点の追尾が成功しているか否かを判断することができる。

10

20

#### 【0059】

LSTMユニットへの入力の特徴量は実施形態3で説明したニューラルネットの特徴量などを用いることができる。例えば図13(B)の物体の被遮蔽スコアを判定するニューラルネット6310の最終層6320への入力値を用いる。ここでは前記層の出力値(1値のスカラー)でなく入力値(多次元ベクトル)を用いている。これは遮蔽状態を判断するのに用いたのと同じ多次元特徴を用いることで、遮蔽に関する多種の情報をLSTMに取り込むためである。これにより様々な遮蔽のパターンを判定できることが期待できる。

30

#### 【0060】

LSTMユニットの学習時には、教師値として各瞬間の追尾が成功しているか失敗しているかを与え、LSTMの各重みパラメータを調整する。また別の形態として図16(D)に示すように、追尾の成功・失敗ではなく、教師値として遮蔽状態にあるか否かを与えて学習すれば、被遮蔽状態にあるか否かを判定させることも可能である。

#### 【0061】

また別の形態として、実施形態3で説明した派生の形態と同様に、追尾物体をアンカー枠ベースで検出し、9401として図13(C)の物体の位置および被遮蔽スコアを同時に判定するニューラルネット6400の特徴量6420を使ってもよい。この形態であれば、物体の追尾や検出と当該物体の被遮蔽スコアを同時・高速に判定することができる。

40

#### 【0062】

<実施形態5>

本実施形態では、ユーザが指定した単一の物体を追尾する形態について説明する。基本機能構成は実施形態1と同様である。本実施形態では物体の遮蔽情報として、相対的な物体間の遠近情報を用いる。

#### 【0063】

図16(A1)にその例を示す。ここでは学習画像としてRGB画像801を用意する。さらにレーザーレンジファインダー装置やステレオ計測等によりRGB画像801に対応した距離画像833が得られている。距離画像833はカメラからの距離の絶対値をグレースケールで表したものであり、白い色ほど近い距離の物体を意味する。本実施形態では

50

R G B 画像 8 0 1 を入力画像とし、距離画像を教師値 8 3 1 として、ニューラルネット 4 0 2 の重みを学習する。ただし絶対値としての距離画像 8 3 1 と全く同じ出力結果 8 3 0 を得ることはパターン認識としては比較的難しい問題であり、本実施形態に用いる遮蔽情報としてはそこまで高精度であることを必要としない。そこで本実施形態では近傍の物体間の相対的な遠近関係を推定するような学習を行う。

【 0 0 6 4 】

例えば同図の出力結果 8 3 0 に示すように、人物 8 0 1 1 と 8 0 1 2 の距離の推定結果 8 3 0 1 と 8 3 0 2 は絶対値としては正しくない。人物 8 0 1 1 と離れた人物 8 0 1 3 に対応する推定結果 8 3 0 1 と 8 3 0 3 も正しくない遠近関係になっている。しかし近傍の二人の人物 8 0 1 1 と 8 0 1 2 の、遠近の順序関係だけに限定すれば、正しい結果である。このように < 局所の物体間 > の < 遠近順序の関係 > は正しく推定できるように学習し、これらを物体の遮蔽情報として集計して用いる。

10

【 0 0 6 5 】

以上は、学習時の損失値計算に以下の工夫を施すことで実現される。図 1 6 ( A 2 ) に図 1 6 ( A 1 ) の教師値 8 3 1 上の記号 \* の付近の領域を拡大した教師値領域 8 3 1 a を示す。対応する出力結果の領域 8 3 0 a も示す。ここで領域 8 3 1 a 上の各画素  $i$  と画素  $j$  に注目し、その遠近関係が正しいか否かで当該画素ペアの損失を求める。ここでは領域 8 3 0 a 上の画素  $i$  と画素  $j$  の遠近関係は教師値と一致するので損失は発生しない。対してもし領域 8 3 0 b のような推定結果であった場合は、遠近関係が正しくないので損失を計上する。このような判断を、所定距離内にある全画素ペアで行う。最終的に遠近関係を誤ったペア数を全ペア数で割った値を損失値の総計とする。このようにして学習したニューラルネット 8 0 2 が学習終了し、推定した距離の出力結果 8 3 4 を図 1 6 ( B ) に示す。

20

【 0 0 6 6 】

次に、相対的な距離の出力結果 8 3 4 を集計して物体の被遮蔽尤度を求める。ここでは別途検出しておいた人物検出枠 8 3 5 1 と 8 3 5 2 を用いて検出枠ごとに集計する。各枠内でそれぞれの距離の値を平均し、 $d_{ave1}$  と  $d_{ave2}$  とする。次にこの距離の値を隣接した物体枠間で比較して正規化して被遮蔽尤度のスコア値  $occ$  へと変換する。例えば下式で変換する。

( 数式 3 )

$$\begin{aligned} occ_i &= \text{Sigmoid}(\text{Log}(d_{ave_i} / d_{ave_j})) \\ &= 1 / (1 + d_{ave_j} / d_{ave_i}), \\ occ_j &= 1 / (1 + d_{ave_i} / d_{ave_j}), \end{aligned}$$

30

ただし

$$\text{Sigmoid}(x) = 1 / (1 + \exp(-x)).$$

ここで  $i$  と  $j$  は重なり部分のある二つの隣接した検出物体枠である。3 つ以上の物体が重なっている場合は、それぞれ上記の式で被遮蔽スコア  $occ_i$  を求め、そのうちの最大値をその物体の被遮蔽スコアとしてもよい。

【 0 0 6 7 】

以上が相対的な距離推定を行い、被遮蔽スコアを集計するまでの処理内容となる。被遮蔽スコアを用いた追尾処理は実施形態 1 と同様になるためここでは割愛する。

40

【 0 0 6 8 】

なお派生的な学習の工夫として下記のようなものが考えられる。( 1 ) 距離の教師値の差分が所定閾値 以上のペアのみに限定して損失を集計する。これにより距離画像の観測時のノイズに対しロバストに学習できる。( 2 ) ( 1 ) を行い、且つマージン領域を設定する。例えばペアの遠近関係が正しいか正しくないかのみならず、遠近関係が正しく、且つ所定閾値 以上値が相対的に離れていない場合に損失を発生させる。( 3 ) 距離の教師値の差分が閾値 未満の画素ペアに対する出力値が、閾値 以上に大きなケースも誤りとして損失を与える。これによりノイズ的な出力を抑制する。

【 0 0 6 9 】

以上、さまざまな形態があり得るが、相対的・局所的に距離を学習できるような形態であ

50

ればいずれでもよく、一つの形態に限定されない。本発明は、以下の処理を実行することによっても実現される。即ち、上述した実施形態の機能を実現するソフトウェア（プログラム）を、データ通信のネットワーク又は各種記憶媒体を介してシステム或いは装置に供給する。そして、そのシステム或いは装置のコンピュータ（またはCPUやMPU等）がプログラムを読み出して実行する処理である。また、そのプログラムをコンピュータが読み取り可能な記録媒体に記録して提供してもよい

【符号の説明】

【0070】

1 情報処理装置

201 画像取得部

202 物体検出部

203 遮蔽情報生成部

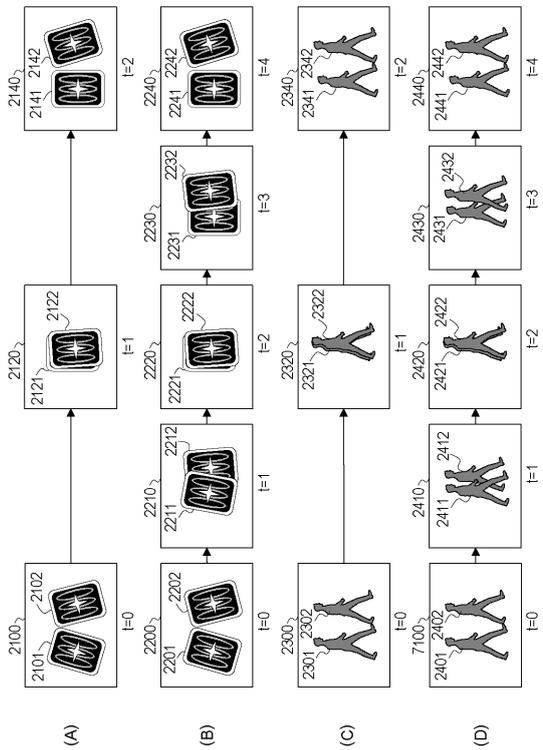
204 特徴量取得部

205 対応付け部

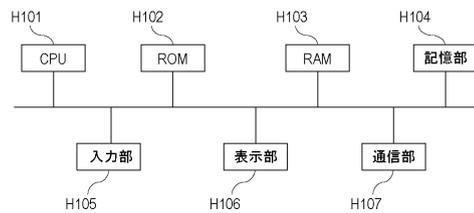
206 記憶部

【図面】

【図1】



【図2】



10

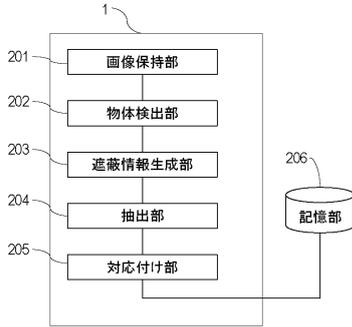
20

30

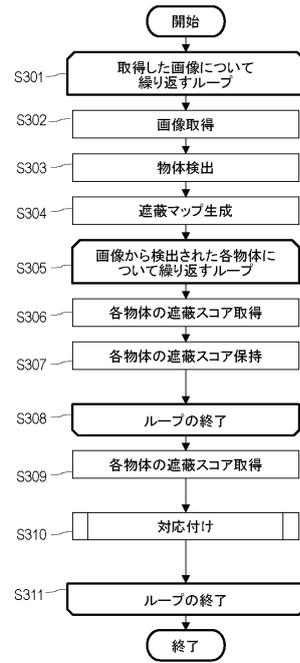
40

50

【 図 3 】



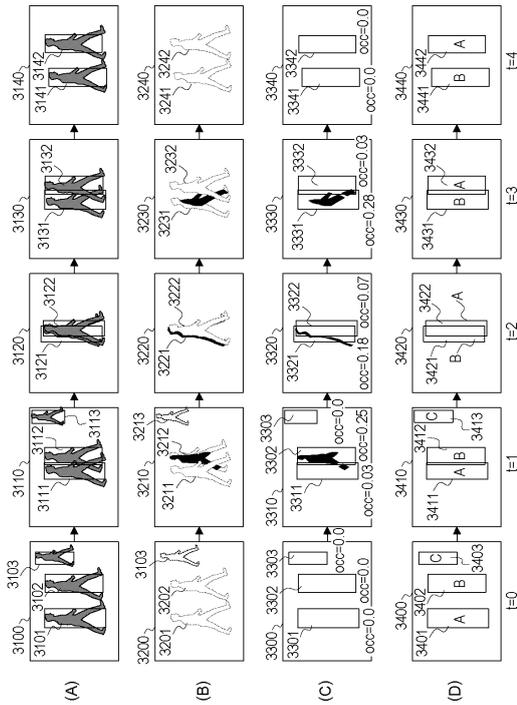
【 図 4 】



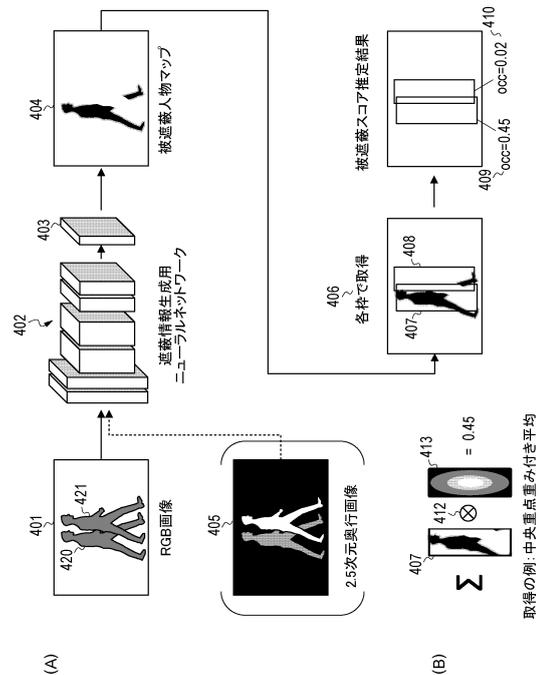
10

20

【 図 5 】



【 図 6 】

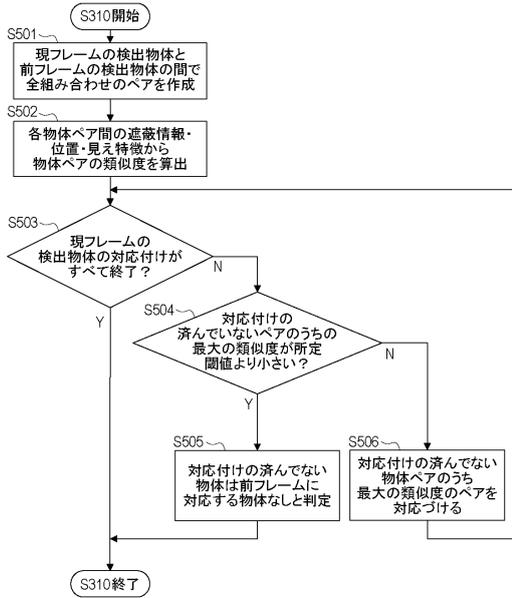


30

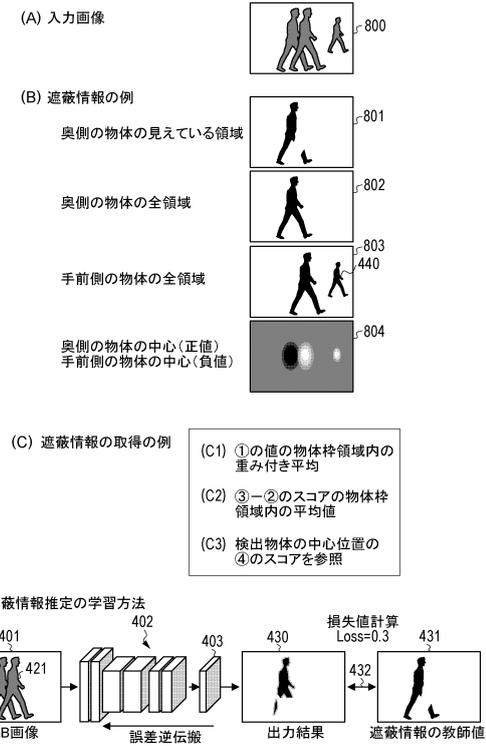
40

50

【 図 7 】



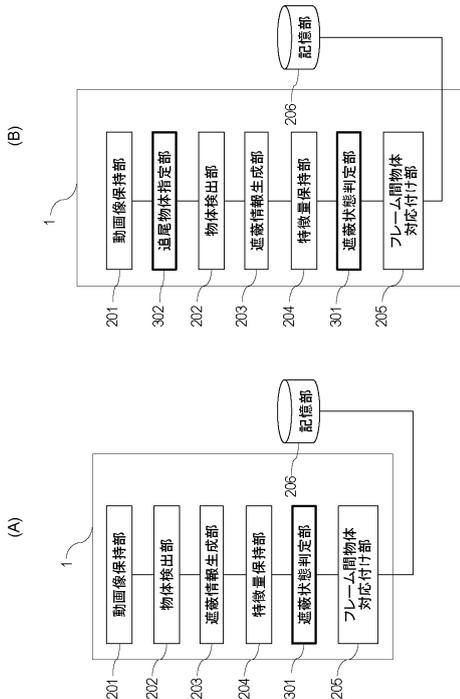
【 図 8 】



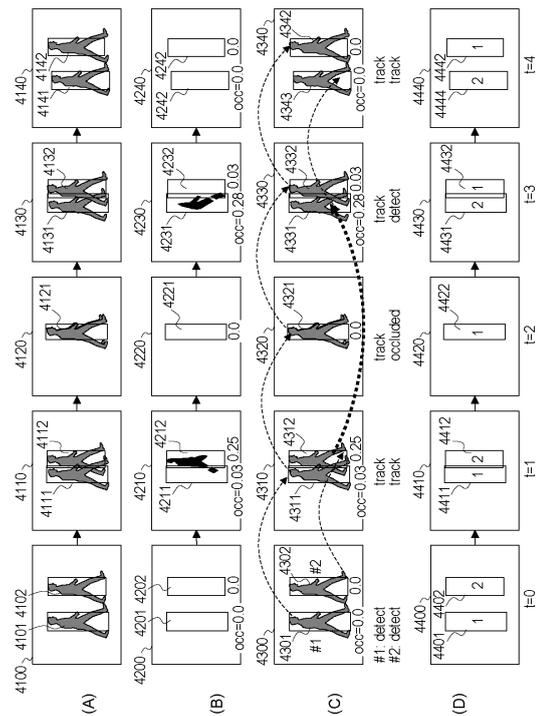
10

20

【 図 9 】



【 図 10 】

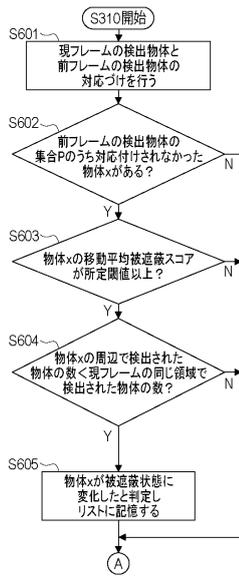


30

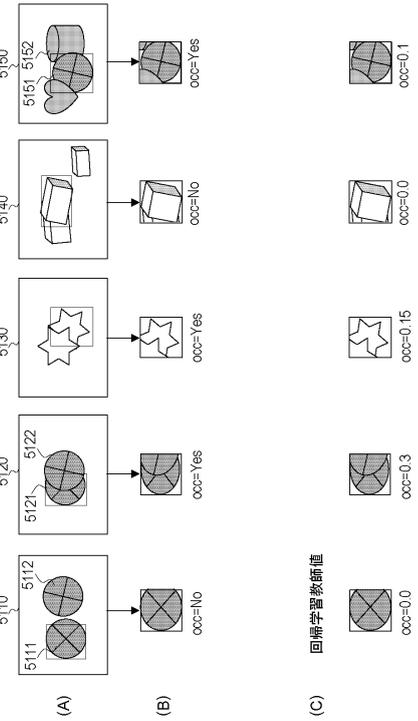
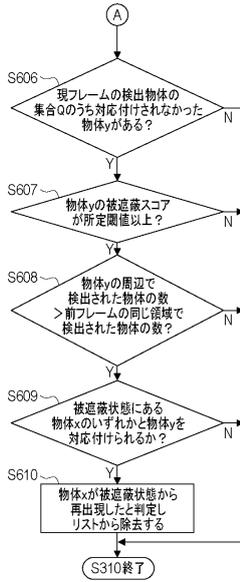
40

50

【 図 1 1 】



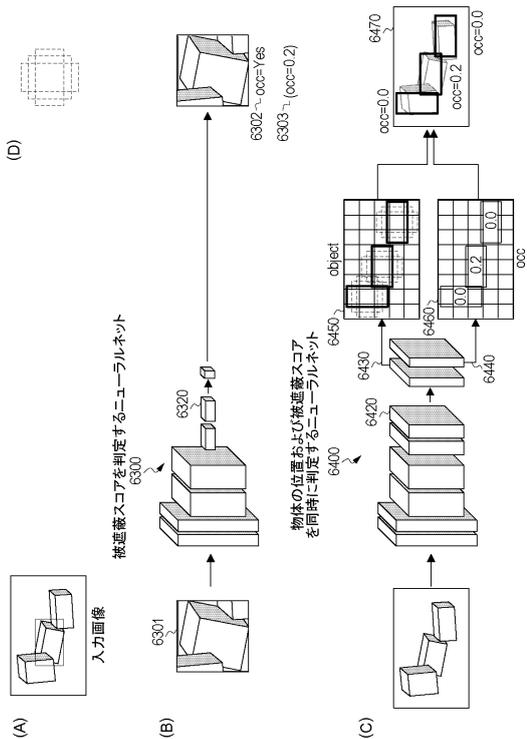
【 図 1 2 】



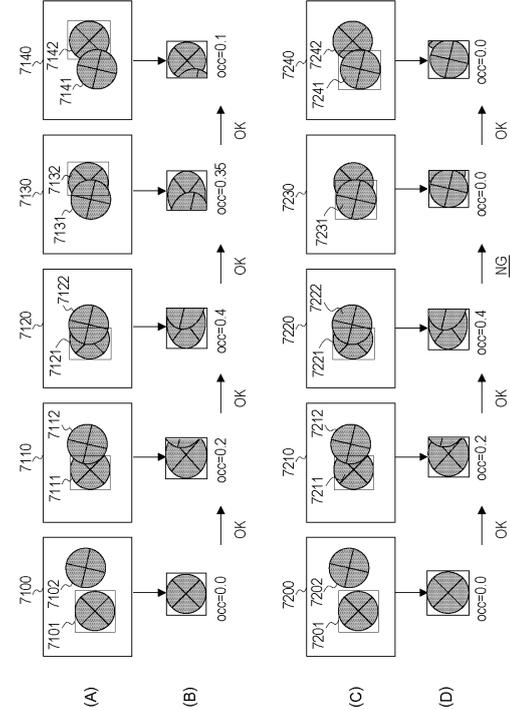
10

20

【 図 1 3 】



【 図 1 4 】

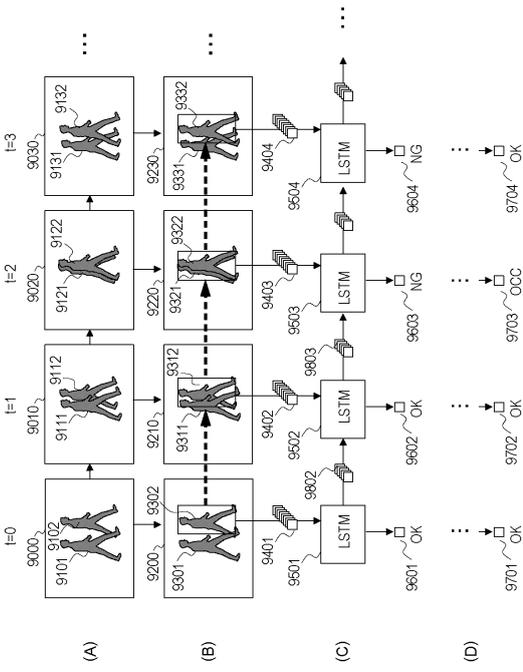


30

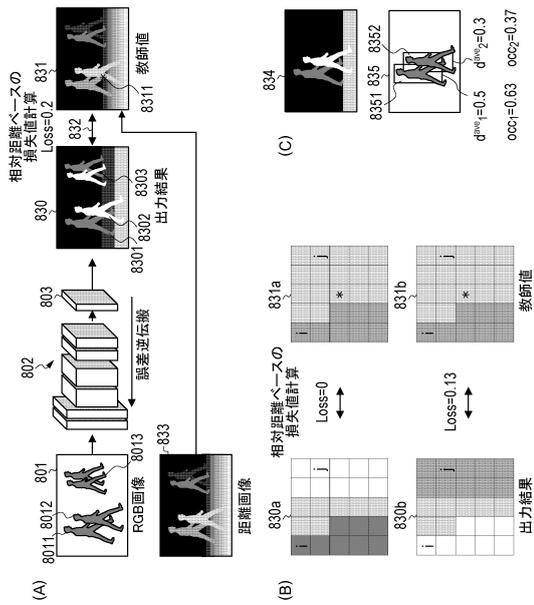
40

50

【 図 15 】



【 図 16 】



10

20

30

40

50

---

フロントページの続き

(51)国際特許分類

F I

H 0 4 N

7/18

K

テーマコード (参考)

キヤノン株式会社内

F ターム (参考)

5C054 CA04 CC02 FC12 FC13 HA19

5L096 AA02 AA06 AA09 BA18 CA04 DA02 FA15 FA35 FA59 FA69

FA77 GA40 HA05 HA11 JA11