

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7490359号
(P7490359)

(45)発行日 令和6年5月27日(2024.5.27)

(24)登録日 令和6年5月17日(2024.5.17)

(51)国際特許分類 F I
G 0 6 T 7/00 (2017.01) G 0 6 T 7/00 3 5 0 B

請求項の数 23 (全32頁)

(21)出願番号	特願2019-233229(P2019-233229)	(73)特許権者	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1年12月24日(2019.12.24)	(74)代理人	100126240 弁理士 阿部 琢磨
(65)公開番号	特開2021-103347(P2021-103347 A)	(74)代理人	100223941 弁理士 高橋 佳子
(43)公開日	令和3年7月15日(2021.7.15)	(74)代理人	100159695 弁理士 中辻 七朗
審査請求日	令和4年12月14日(2022.12.14)	(74)代理人	100172476 弁理士 富田 一史
		(74)代理人	100126974 弁理士 大朋 靖尚
		(72)発明者	館 俊太 東京都大田区下丸子3丁目30番2号キ 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)【特許請求の範囲】

【請求項1】

入力画像に含まれる複数の物体を検出する情報処理装置であって、
前記入力画像から抽出された画像特徴を入力すると前記物体が存在する可能性を示す尤度を領域毎に得るために用いる尤度マップを出力する学習モデルに、複数の異なる結合重みパラメータを設定することにより複数の前記尤度マップを生成する生成手段と、
前記入力画像に含まれる前記物体の正解としての位置の情報を取得する取得手段と、
前記複数の尤度マップのうち同一の尤度マップ内の、第1の注目領域の尤度と、該第1の注目領域の近傍の領域に対応付けられた前記尤度との差に基づいて、該第1の注目領域の近傍の領域に対応付けられた前記尤度を下げるための第1の損失関数と、前記複数の尤度マップを用いて得られた尤度を領域毎に統合して得た尤度のうち、前記取得された位置を含む領域で得られた当該尤度を上げるための第2の損失関数と、に基づいて前記尤度マップごとに前記結合重みパラメータを更新する学習手段と、を有することを特徴とする情報処理装置。

【請求項2】

前記第1の損失関数は、前記第1の注目領域の尤度と、該第1の注目領域の近傍の領域に対応付けられた前記尤度との差が小さいほど、該第1の注目領域の近傍の領域に対応付けられた前記尤度をより下げ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】

前記学習手段は、前記複数の尤度マップの前記取得された前記物体の位置を含む第2の

10

20

注目領域に対応付けられた前記尤度に基づいて、前記第2の注目領域に対応付けられた尤度を調整する第3の損失関数に更に基づいて、前記尤度マップごとに前記結合重みパラメータを更新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】

前記第3の損失関数は、前記取得された物体の位置についての情報に基づいて前記入力画像に含まれる前記物体の位置が1つである場合は、前記複数の尤度マップのうちいずれか1つの尤度マップの前記第2の注目領域に対応付けられた尤度がより大きくなるように調整し、前記取得された物体の位置についての情報に基づいて前記入力画像に含まれる前記物体の位置が複数ある場合は、前記物体の数に応じて前記尤度マップの前記第2の注目領域に推定された尤度がより大きくなるように調整する損失値を決定することを特徴とする請求項3に記載の情報処理装置。

10

【請求項5】

前記第3の損失関数は、前記複数の尤度マップのうち2つの尤度マップ間において前記入力画像における前記物体の位置を含む前記第2の注目領域に対応付けられた尤度に基づいて、前記尤度マップの少なくとも一方における前記第2の注目領域に対応付けられた尤度を下げるように調整する損失値を決定することを特徴とする請求項3または4に記載の情報処理装置。

【請求項6】

前記学習手段により更新された前記結合重みづけパラメータを設定した前記学習モデルによって生成された前記尤度マップに基づいて、前記尤度が閾値より大きい領域に存在する前記物体を検出する検出手段と、を更に有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。

20

【請求項7】

前記生成手段は、前記抽出された画像特徴のうち、前記結合重みパラメータが更新された学習済みモデル毎に異なる画像特徴の組み合わせを入力することによって前記尤度マップを生成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】

前記入力画像から画像特徴を複数の異なる組み合わせで抽出する抽出手段を更に有し、前記生成手段は、前記抽出された異なる組み合わせの画像特徴を前記学習モデルに入力することで前記尤度マップを生成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。

30

【請求項9】

前記物体とセンサとの距離を示す奥行き情報を取得する取得手段を更に有し、前記生成手段は、前記奥行き情報をさらに前記学習モデルに入力することによって、複数の前記尤度マップを生成することを特徴とする請求項1乃至8のいずれか1項に記載の情報処理装置。

【請求項10】

入力画像に含まれる複数の物体を検出する情報処理装置であって、前記入力画像から抽出された画像特徴を入力すると前記物体が存在する可能性を示す尤度を領域毎に得るための尤度マップを出力する学習済みモデルに、複数の異なる結合重みパラメータを設定することにより生成された複数の前記尤度マップを取得する第1の取得手段と、

40

前記取得された複数の尤度マップに基づいて、前記尤度が閾値より大きい注目領域を前記入力画像に含まれる前記物体として検出する検出手段と、を有し、

前記学習済みモデルは、前記複数の尤度マップのうち同一の尤度マップ内の、第1の注目領域の尤度と、該第1の注目領域の近傍の領域に対応付けられた前記尤度との差に基づいて、該第1の注目領域の近傍の領域に対応付けられた前記尤度を下げるための第1の損失関数と、前記複数の尤度マップのうち所定の2つの尤度マップ間で共通である第2の注目領域に対応付けられた尤度に基づいて、前記尤度マップの少なくとも一方における前記第2の注目領域に対応付けられた尤度を調整するための第3の損失関数と、に基づいて、

50

前記第 1 と第 3 の損失関数によって出力された損失値を小さくするように前記尤度マップごとに前記結合重みパラメータを学習させた学習モデルであることを特徴とする情報処理装置。

【請求項 1 1】

前記第 3 の損失関数は、前記複数の尤度マップのうち 2 つの尤度マップ間において前記入力画像における前記物体の位置を含む前記第 2 の注目領域に対応付けられた尤度に基づいて、前記尤度マップの少なくとも一方における前記第 2 の注目領域に対応付けられた尤度を下げるための損失値を決定する関数であることを特徴とする請求項 1 0 に記載の情報処理装置。

【請求項 1 2】

前記入力画像に含まれる前記物体の位置を取得する第 2 の取得手段を更に有し、
前記学習済みモデルは、前記複数の尤度マップを用いて得られた尤度を領域毎に統合して得た尤度のうち、前記取得された位置を含む領域で得られた当該尤度を上げるための第 2 の損失関数に更に基づいて、前記第 1 と第 2 と第 3 の損失関数によって出力された損失値を小さくするように前記尤度マップごとに前記結合重みパラメータを学習させた学習モデルであることを特徴とする請求項 1 0 または 1 1 に記載の情報処理装置。

【請求項 1 3】

前記第 2 の損失関数は、前記複数の尤度マップを統合した結果に対して、前記取得された位置に対応付けられた前記尤度が閾値より小さい場合は大きい損失値を出力することを特徴とする請求項 1 2 に記載の情報処理装置。

【請求項 1 4】

前記第 2 の取得手段は、前記物体の位置から前記物体の数を取得し、
前記第 2 の損失関数は、前記取得された前記物体の数に基づいて、前記複数の尤度マップを統合した結果に対して、前記尤度が閾値より大きい領域が前記取得された数と一致しない場合は大きい損失値を出力することを特徴とする請求項 1 2 または 1 3 に記載の情報処理装置。

【請求項 1 5】

前記損失関数の出力する損失値が収束するように前記学習済みモデルのパラメータを更新することによって学習する学習手段を更に有することを特徴とする請求項 1 0 乃至 1 4 のいずれか 1 項に記載の情報処理装置。

【請求項 1 6】

前記入力画像から抽出された画像特徴に基づいて、前記物体の数を特定する特定手段を更に有し、
前記第 1 の取得手段は、前記特定された前記物体の数に応じて取得する前記尤度マップの数を調整することを特徴とする請求項 1 0 乃至 1 5 のいずれか 1 項に記載の情報処理装置。

【請求項 1 7】

入力画像に含まれる複数の物体を検出する情報処理装置であって、
前記物体を撮像した画像における該物体それぞれの領域を取得する取得手段と、
前記入力画像から抽出された画像特徴を入力すると前記物体が存在する可能性を示す尤度を領域毎に得るための尤度マップを出力する学習モデルに、複数の異なる結合重みパラメータを設定することにより複数の前記尤度マップを生成する生成手段と、
前記複数の尤度マップのうち同一の尤度マップ内の、第 1 の注目領域の尤度と、該第 1 の注目領域の近傍の領域に対応付けられた前記尤度との差に基づいて、該第 1 の注目領域の近傍の領域に対応付けられた前記尤度を下げるための第 1 の損失関数と、前記複数の尤度マップのうち所定の 2 つの尤度マップ間で共通である第 2 の注目領域に対応付けられた尤度に基づいて、前記尤度マップの少なくとも一方における前記第 2 の注目領域に対応付けられた尤度を調整するための第 3 の損失関数と、に基づいて前記尤度マップごとに前記結合重みパラメータを更新する学習手段と、

前記学習手段により更新された前記結合重みづけパラメータを設定した前記学習モデル

10

20

30

40

50

によって生成された前記尤度マップに基づいて、前記入力画像において物体毎に対応する領域を推定する推定手段と、を有することを特徴とする情報処理装置。

【請求項 18】

入力画像に含まれる複数の物体を検出する情報処理装置であって、

前記入力画像から抽出された画像特徴を入力すると前記物体が存在する可能性を示す尤度を領域毎に対応付けた複数の尤度マップと、前記入力画像に含まれる前記物体の位置についての情報と、を取得する取得手段と、

前記複数の尤度マップを用いて得られた尤度を領域毎に統合して得た尤度のうち、前記取得された位置を含む領域で得られた当該尤度を上げるための第2の損失関数に基づいて、前記尤度マップごとに該尤度マップの領域毎の尤度を決定するためのパラメータを更新する学習手段と、有することを特徴とする情報処理装置。

10

【請求項 19】

入力画像に含まれる複数の物体を検出する情報処理装置であって、

前記入力画像から抽出された画像特徴を入力すると前記物体が存在する可能性を示す尤度を領域毎に得るための複数の尤度マップを生成する生成手段と、

前記複数の尤度マップのうち同一の尤度マップ内の、第1の注目領域の尤度と、該第1の注目領域の近傍の領域に対応付けられた前記尤度との差に基づいて、該第1の注目領域の近傍の領域に対応付けられた前記尤度を下げる第1の調整と、前記複数の尤度マップのうち所定の2つの尤度マップ間で共通である第2の注目領域に対応付けられた尤度に基づいて、前記尤度マップの少なくとも一方における前記第2の注目領域に対応付けられた尤度を調整する第2の調整と、を行うことによって前記尤度マップごとに尤度を調整する調整手段と、

20

前記調整された複数の尤度マップに基づいて、前記尤度が閾値より大きい領域を前記入力画像に含まれる前記物体として検出する検出手段と、を有することを特徴とする情報処理装置。

【請求項 20】

入力画像に含まれる複数の物体を検出する情報処理装置であって、

前記入力画像から抽出された画像特徴を入力すると前記物体が存在する可能性を示す尤度を領域毎に得るために用いる尤度マップを出力する学習モデルに、複数の異なる結合重みパラメータを設定することにより複数の前記尤度マップを生成する生成手段と、

30

前記入力画像に含まれる前記物体の正解としての位置の情報を取得する取得手段と、

前記複数の尤度マップのうち同一の尤度マップ内の、第1の注目領域に対応付けられた尤度が閾値より大きい場合、該第1の注目領域の近傍の領域に対応付けられた前記尤度をより下げるための第1の損失関数と、前記複数の尤度マップを用いて得られた尤度を領域毎に統合して得た尤度のうち、前記取得された位置を含む領域で得られた当該尤度を上げるための第2の損失関数と、に基づいて前記尤度マップごとに前記結合重みパラメータを更新する学習手段と、を有することを特徴とする情報処理装置。

【請求項 21】

前記第1の損失関数は、前記複数の尤度マップのうち同一の尤度マップ内の、第1の注目領域に対応付けられた尤度が閾値より大きい場合は該第1の注目領域の近傍の領域に対応付けられた前記尤度をより下げるための損失値を大きくするように、前記第1の注目領域に対応付けられた尤度が閾値より小さい場合は該第1の注目領域の近傍の領域に対応付けられた前記尤度を下げるための損失値をより小さくなるように損失値を出力する損失関数であることを特徴とする請求項 20 に記載の情報処理装置。

40

【請求項 22】

コンピュータを、請求項 1 乃至 21 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 23】

入力画像に含まれる複数の物体を検出する情報処理方法であって、

前記入力画像から抽出された画像特徴を入力すると前記物体が存在する可能性を示す尤

50

度を領域毎に得るための尤度マップを出力する学習モデルに、複数の異なる結合重みパラメータを設定することにより複数の前記尤度マップを生成する生成工程と、

前記入力画像に含まれる前記物体の位置についての情報を取得する取得工程と、

前記複数の尤度マップのうち同一の尤度マップ内の、第1の注目領域の尤度と、該第1の注目領域の近傍の領域に対応付けられた前記尤度との差に基づいて、該第1の注目領域の近傍の領域に対応付けられた前記尤度を下げるための第1の損失関数と、前記複数の尤度マップを用いて得られた尤度を領域毎に統合して得た尤度のうち、前記取得された位置を含む領域で得られた当該尤度を上げるための第2の損失関数と、に基づいて前記尤度マップごとに前記結合重みパラメータを更新する学習工程と、を有することを特徴とする情報処理方法。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、画像中の複数の対象物体を認識する技術に関する。

【背景技術】

【0002】

従来、画像中の特定の被写体を検出する方法が知られている。例えば非特許文献1ではニューラルネットワークを用いて被写体の位置・サイズ・カテゴリ等の認識を行う。非特許文献1に類する手法は入力画像を多層ニューラルネットワーク（深層ニューラルネットワークとも呼ばれ、以下DNNと略する）で処理する。これらの手法の一般的な形態は下記のようなものである。まず入力画像を $W \times H$ のサイズに変換しDNNに入力する。次にこの入力画像に対して畳み込み等の演算を複数回繰り返し、画像を中間的に $w \times h$ ($w < W, h < H$)の解像度の画像特徴へと変換する。DNNの後段の層ではこの特徴の $w \times h$ の各位置に対象物体が存在するか、しないかを判定する。存在すると判定した場合はさらに別途DNNで推定したサイズや精確な位置推定結果等、と合わせて検出結果を出力する。

20

【先行技術文献】

【非特許文献】

【0003】

【文献】J. Redmon, A. Farhadi, YOLO9000: Better, Faster, Stronger, CVPR, 2017

30

【発明の概要】

【発明が解決しようとする課題】

【0004】

従来、物体の有無を判定する単位である1つのブロックの中に同種の複数の物体が隣接して存在する場合、「物体が1つ存在する」という誤検出を起こすことがあった。つまり、近接した同一種類の物体同士を分離して認識することは難しかった。

【0005】

この課題は入力画像あるいは前述の中間的な画像特徴の解像度を上げることで回避できるが、その場合DNNの処理量および後処理（多重検出抑制等）の演算量が大きく増大する。

40

【0006】

本発明は、上記の課題に鑑みてなされたものであり、複数の物体が隣接して存在する場合でも、そのそれぞれを精度よく検出することを目的とする。

【課題を解決するための手段】

【0007】

上記の目的を達成する本発明に係る情報処理装置は、入力画像に含まれる複数の物体を検出する情報処理装置であって前記入力画像から抽出された画像特徴を入力すると前記物体が存在する可能性を示す尤度を領域毎に得るために用いる尤度マップを出力する学習モデルに、複数の異なる結合重みパラメータを設定することにより複数の前記尤度マップを生成する生成手段と、前記入力画像に含まれる前記物体の正解としての位置の情報を取得

50

する取得手段と、前記複数の尤度マップのうち同一の尤度マップ内の、前記複数の尤度マップのうち同一の尤度マップ内の、第1の注目領域の尤度と、該第1の注目領域の近傍の領域に対応付けられた前記尤度との差に基づいて、該第1の注目領域の近傍の領域に対応付けられた前記尤度を下げるための第1の損失関数と、前記複数の尤度マップを用いて得られた尤度を領域毎に統合して得た尤度のうち、前記取得された位置を含む領域で得られた当該尤度を上げるための第2の損失関数と、に基づいて前記尤度マップごとに前記結合重みパラメータを更新する学習手段と、を有することを特徴とする。

【発明の効果】

【0008】

本発明によれば、複数の物体が隣接して存在する場合でも、そのそれぞれを精度よく検出できる。

10

【図面の簡単な説明】

【0009】

【図1】情報処理装置の機能構成例を示すブロック図

【図2】情報処理装置が実行する処理を説明するフローチャート

【図3】特徴抽出部が実行する処理を説明するフローチャート

【図4】画像特徴の模式図

【図5】尤度マップの更新と統合の模式図

【図6】尤度マップ生成部の概念図

【図7】尤度マップの結合重みの概念図

20

【図8】情報処理装置の機能構成例を示すブロック図

【図9】情報処理装置が実行する処理を説明するフローチャート

【図10】損失関数の模式図

【図11】情報処理装置の機能構成例を示すブロック図

【図12】特徴の集計の模式図

【図13】特徴の集計の詳細

【図14】情報処理装置が実行する処理を説明するフローチャート

【図15】情報処理装置の機能構成例を示すブロック図

【図16】情報処理装置が実行する処理を説明するフローチャート

【図17】GUIの一例を示す図

30

【図18】情報処理装置の機能構成例を示すブロック図

【図19】情報処理装置が実行する処理を説明するフローチャート

【図20】尤度マップの一例を示す図

【図21】情報処理装置のハードウェア構成例を示すブロック図

【発明を実施するための形態】

【0010】

<実施形態1>

本実施形態は物体検出において複数の物体が近接・重畳している場合であっても、頑健に検出する一手法について説明する。ここでは物体の顔の検出を行う情報処理装置の例を説明する。ただし本発明は物体の顔に限定することなく各種物体検出に対して適用可能である。

40

【0011】

複数の物体が含まれる画像から物体それぞれの検出する場合、物体同士が隣接していると1つの物体として検出することや一方の物体を検出しないことがある。そのようなケースに対し、本実施形態では、物体が存在する可能性を示す尤度を入力画像の各領域に対応づけられた尤度マップを2枚以上用意し、複数の尤度マップから物体の位置を検出する。複数の尤度マップは、それぞれの異なる位置にある物体を検出するように、損失関数あるいはニューロン間の結合を用いて、尤度マップ内または尤度マップ間の重みを学習する、あるいは尤度マップを用いて得られる尤度を更新する。

【0012】

50

本発明の基本的な機能構成図である図1を用いて説明する。なお、これ以降図面を参照しながら説明する際に、図面間で同一の符号は同一の構成モジュールを意味する。特筆すべき差異がない場合は重ねての説明を省く。

【0013】

図21は、情報処理装置のハードウェア構成例を示すブロック図である。中央処理ユニット(CPU)211は、RAM213をワークメモリとして、ROM212や記憶装置214に格納されたOSやその他プログラムを読みだして実行し、システムバス219に接続された各構成を制御して、各種処理の演算や論理判断などを行う。CPU211が実行する処理には、実施形態の情報処理が含まれる。記憶装置214は、ハードディスクドライブや外部記憶装置などであり、実施形態の情報処理にかかるプログラムや各種データを記憶する。入力部215は、カメラなどの撮像装置、ユーザー指示を入力するためのボタン、キーボード、タッチパネルなどの入力デバイスである。なお、記憶装置214は例えばSATAなどのインターフェースを介して、入力部215は例えばUSBなどのシリアルバスを介して、それぞれシステムバス219に接続されるが、それらの詳細は省略する。通信I/F216は無線通信で外部の機器と通信を行う。表示部217はディスプレイである。センサ218は画像センサや距離センサである。センサで計測した結果を画像として記憶装置214に記憶する。尚、CPUはプログラムを実行することで各種の手段として機能することが可能である。なお、CPUと協調して動作するASICなどの制御回路がこれらの手段として機能しても良い。また、CPUと画像処理装置の動作を制御する制御回路との協調によってこれらの手段が実現されても良い。また、CPUは単一のもの

10

20

【0014】

図1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1を用いて各機能構成について説明する。情報処理装置1は、画像入力部101、特徴抽出部102、尤度マップ生成部104、統合部107、出力部108、記憶部109から構成される。特徴抽出部102は、入力画像を処理して画像に含まれる様々な物体の特徴を示す画像特徴103を生成する。尤度マップ生成部104は、画像特徴103を入力すると、特定の物体が存在する可能性を示す尤度を領域毎に示す尤度マップを出力する学習モデルに基づいて、複数の尤度マップを生成する。ここでは、学習モデルの層間の結合重みづけパラメータのセットを異なる組み合わせで用意する。このパラメータセットは、後述する損失関数を用いて、隣接する同種類または同じくらいの大きさの物体が隣接していてもそれぞれを異なるマップで検出できるように学習させたものである。統合部107は、生成された複数の前記マップを統合して、特定の物体が存在する位置を示すマップを出力する。出力部108は、統合結果から、入力画像において認識対象となる物体が存在する位置を出力する。記憶部109は、情報処理装置がパラメータを学習する際に出力の目標値である教師値を記憶する。教師値は、予め用意され記憶されたデータを取得する構成以外にも、ユーザーや外部装置によって画像において特定の物体が存在する位置を示した教師値を入力できる構成

30

40

【0015】

<認識処理の動作>

次に情報処理装置が実行する処理を説明するフローチャートである図2を用いて、処理手順を説明する。以下の説明では、各工程(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1はこのフローチャートで説明するすべてのステップを必ずしも行わなくても良い。

【0016】

本実施形態における、情報処理装置は、入力画像に含まれる複数の物体を検出する。例えば、物体が複数人映っている画像から、物体の存在する位置を検出する。そのために、

50

入力画像から物体を示す画像特徴を抽出する。次に、抽出された画像特徴に基づいて、物体が存在する可能性を示す尤度を出力するマップを少なくとも2つ以上生成する。ここで生成されるマップは、それぞれ異なる位置に存在する物体を検出できるように、それぞれ異なる検出結果（尤度）を出力するマップになるよう更新する。更新の方法は後述する。生成された複数のマップにおいて所定の値より大きい尤度が出力された注目領域について、同一のマップまたは異なるマップにおける注目領域の近傍の領域における尤度の値に基づいて損失値を出力する損失関数に基づいてマップを更新する。

【0017】

まずS1で、画像入力部101が認識対象の物体が映った入力画像を入力する。入力画像は、RGBカラー画像以外でも、白黒画像やグレースケールの濃淡画像でも良い。また、カメラで撮像する画像以外でも、赤外線カメラによる赤外線写真や、LidarやToFを代表とするアクティブ距離センサで得た距離画像でも良い。次にステップS2で特徴抽出部102が、入力画像から画像特徴を抽出する。画像特徴としては画素の色やテクスチャなどを集計した特徴ベクトルなど、公知の様々な方法が考えられる。マップ状の特徴であれば何れの方法でも適応でき、特定方法に限定されない。本実施形態の特徴抽出部102は多層ニューラルネット102aを備えることとする。具体的には以下のような手順を行って画像特徴を抽出する。

【0018】

<画像特徴の抽出手順>

S2において、多層ニューラルネット102aによってマップ状の高次元画像特徴を抽出する方法について説明する。特徴抽出部が実行する処理を説明するフローチャートを図3に示す。まずステップS201においてニューラルネットが画像特徴を格納するための3次元配列 $F(x, y)$ を初期化する(x, y は特徴の画素に関する添え字である)。次にステップS202からS207でニューラルネットの各層が入力画像に対して演算処理を行う。

【0019】

S203において、本実施形態のニューラルネット102aは図4に示すような各層の演算処理を行う。ここでいう演算処理とは、ニューラルネットワークの各層で、後段の検出処理において物体を検出するための画像特徴を入力画像から抽出する処理である。ニューラルネット102aは入力画像401に対して複数回の畳み込みフィルタ処理を行う畳み込み処理402a、402b、402cを備える。さらに各畳み込みの後に行う活性化関数処理を備える(図では略している)。さらにプール処理404a、404bを備える。L番目の層の畳み込みフィルタ処理および活性化関数処理を数式で表すと下記のようになる。

数式1

$$f_L(x, y, CHOUT) = \begin{cases} CHIN(x, y) \times w^L(x, y, CHIN, CHOUT) + B^L_{CHOUT} \\ \text{if } x < 0 \\ x \text{ Otherwise} \end{cases}$$

【0020】

ここで $f_L(x, y, z)$ はL番目の層が出力する特徴マップの結果で、z枚のマップからなる。(図4中に単位chとして付した数字は特徴マップの枚数である。)(\cdot)は半波整流よりなる活性化関数、 $w^L(x, y, CHIN, CHOUT)$ (ただし $x, y \in \{-K, \dots, 0, \dots, K\}$)はL番目の層の畳み込みの重みパラメータ、 B^L はL番目の層のバイアス項である。CHINはL-1番目の層が出力する特徴マップの番号、CHOUTはL番目の層が出力する特徴マップの番号を表す。なお上式ではRGB3チャンネルからなる入力画像 $I(x, y)$ は特徴マップ $f_0(x, y, z)$ として扱うものとする。

【0021】

10

20

30

40

50

なおここでは畳み込みの前後で特徴マップの x, y 方向のサイズが変化しないように、畳み込み処理の前に特徴マップ f_{L-1} の周囲の画素に 0 値を充填してから畳み込むものとする（パディング処理と呼ばれる）。

【0022】

プール処理 404a、404b は特徴マップを所定の局所領域ごとに代表値で代表させることでマップのサイズを縮小する処理である。プール処理は CNN の認識性能をロバストにする効果がある反面、結果の解像度が落ちるといった性質がある。図 6 に示す例ではプール処理 404a、404b はそれぞれ特徴マップを 2×2 画素ごとに統合して 2 分の 1 の解像度の特徴マップに縮小する処理である。

【0023】

ここまでの演算処理により、CNN の各層の段階でそれぞれ特徴マップ 403a、403b、403c が生成される。以上はすべて CNN に関する一般的な技術であり、上記の非特許文献 1、また下記の非特許文献 2、非特許文献 3 等で広く公知であるため、これ以上の詳細な説明は略す。必要に応じて先行文献を参照されたい。（非特許文献 2：A. Krizhevsky et al., ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012）。（非特許文献 3：M. D. Zeiler, R. Fergus, Visualizing and Understanding Convolutional Networks, ECCV 2014）。

【0024】

なお CNN は非特許文献 2 で行っているような大規模データベースを用いた分類タスクで重みパラメータをあらかじめ学習しておくものとする。この結果 CNN の入力画像に近い低層の特徴マップ 403a は、画像中の線分の傾きのような単純な模様によく反応するマップが生成される。そして後段の高層の特徴マップ 403c ではより広い範囲の画像パターンを集約し、複雑なパターン形状に反応するようなマップが生成される。CNN の上記の性質は非特許文献 3 等で広く公知である。

【0025】

次に S204 では、ニューラルネット 102a が現在処理中の層が所定の層か否かを判定し、所定の層であれば S205 に進み、解像度を揃えてから（S205）、特徴マップ 403 を特徴の配列に連結して追加する（S206）。所定の層とは、設計時点であらかじめ取ってくと決めている層であって、一部でも良いし、全部の層を連結しても良い。所定の層ではない層に対しては、連結処理を行わない。これを繰り返して最終的に特徴マップが複数個連結された画像特徴 103 が得られる。上記の処理は数式では下記のように表される。すなわち、画素ごとに画像特徴を抽出し、特徴と画素とを対応づけた特徴マップを生成する。

数式 2

$$F(x, y) = [f_1(x, y)^T, f_2(x, y)^T, \dots, f_M(x, y)^T]^T$$

ここで f_1, f_2, \dots, f_M は抽出された所定の特徴マップであり、特徴 F は上記特徴マップを Z 次元方向に連結したものである。 x, y は特徴マップの画素の座標である。

【0026】

なお一般的に CNN はプール処理を行うため、特徴マップの解像度は層によって異なっている。そのため前記の連結の前に連結特徴生成部 204 がアップサンプル処理 405a、405b（あるいはダウンサンプル処理）を行って解像度を変更し、各特徴マップを所定の解像度に揃える（S205）。ここでの解像度変更は値のコピーや線形補間といった一般的な方法を行えばよい。図 6 では各特徴マップ 403b と 403c をそれぞれ 2 倍と 4 倍にアップサンプルすることで、特徴マップ f_1, f_2, f_3 を入力画像 $I(x, y)$ と同じ解像度に揃えている。以上の結果、画像特徴 103 として特徴 $F(x, y)$ が得られた。

【0027】

なお本形態では、入力画像 $I(x, y)$ と特徴 $F(x, y)$ の解像度は同一である。し

10

20

30

40

50

かし仮にアップサンプル（ダウンサンプル）の処理の倍率を変更すれば、入力画像 I (x , y) よりも細かい（粗い）解像度の特徴を抽出することも可能である。なお本発明の説明においては特徴 F (x , y) の個々の画素を一般化した名称として以降は「領域ブロック」と呼ぶこととする。以上が画像特徴を生成する S 2 の説明になる。

【 0 0 2 8 】

< 尤度マップの生成 >

S 3 では、尤度マップ生成部 1 0 4 が、入力画像から抽出された画像特徴を入力すると物体が存在する可能性を示す尤度を領域毎に対応付けた尤度マップを出力する学習モデルに、異なる結合重みパラメータを設定することにより複数の尤度マップを生成する。ここでは、尤度マップ生成部 1 ~ N が、前段で得られた画像特徴に基づき、N 枚の顔の尤度マップを生成する。尤度マップを複数生成するのは近接・重畳した複数の被写体でもなるべく漏れなく検出するためであり、基本的に各尤度マップの機能は同質である。なお、尤度マップ生成部 1 ~ N は、それぞれが出力する尤度マップが、複数のマップで同じ物体を重複して検出しないように、かつ同じマップで隣接する物体が検出されないように、それぞれの異なるパラメータセットを学習済みであるものとする。学習方法については、学習処理の部分で後述する。

10

【 0 0 2 9 】

具体的には各尤度マップ生成部が、領域ごとに 1 x 1 サイズの畳み込み演算を行い、特定の物体（ここでは人物を検出したいので顔の特徴）があるか否かを示すスコア値（以降これを尤度スコアと呼ぶ）を算出する（下式）。ここではあらかじめ学習によって決定された重みパラメータ w , b を用いて各領域の特徴の畳み込み演算処理を行う。これによりその領域に物体が存在する可能性を示す尤度を取得する。

20

数式 3

$$v_k(x, y) = g(\sum_j w_{kj} \cdot f_{xyj} + b_k),$$

【 0 0 3 0 】

ここで v k は k 番目の尤度マップの尤度スコア (k = 1 , 2 , ... , N)、f x y j は高次元の特徴 F (x , y) の j 番目の次元の要素、w と b は縦横 1 x 1 サイズの畳み込みのパラメータである。w と b は画像特徴を元に、それぞれが独自に顔の尤度スコアのマップを生成するよう後述の学習処理によってあらかじめ調整されている。

【 0 0 3 1 】

関数 g (·) はスコア値を規格化する関数であり、

数式 4

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

等によって定義される。ここでは関数 g は尤度スコアを 0 ~ 1 の範囲に規格化するものである。関数 g としては上記の形態に限らず規格化のためのその他の適当な非線形関数を用いてもよい。ここでは、尤度が高い（1に近い）ほど顔が存在する可能性が高い。

30

【 0 0 3 2 】

< 尤度マップの更新処理 >

S 4 ~ S 7 は尤度マップを更新する処理である。なお、尤度マップは、入力画像から各領域の尤度（対象物体が存在する可能性を示すスコア）を出力する過程で用いるマップであり、各領域には尤度を出力するための内部パラメータが設定されている。前述の学習の結果、尤度マップ内の内部パラメータは、尤度マップ毎に異なっており、同一の入力画像における共通する位置の領域に異なる尤度が出力される仕組みになっているものとする。

40

【 0 0 3 3 】

前段で得られた複数の尤度マップは、尤度マップ間で重複して同一の被写体を検出していたり、一方でどのマップでも尤度スコアが低い被写体があったり、等があり得る。そのため本処理ステップでは尤度マップの出力の調整を行う。

【 0 0 3 4 】

ただし尤度マップの結果が良好な場合必ずしも本ステップは必要ではない。また後述するように、本ステップと同様の機能が前段の D N N の中で一体的に実現されるケースもあ

50

り得る。ただしここでは説明の便宜上、本処理を1つの独立した処理モジュールとして扱い、尤度マップ生成部104が各尤度マップの出力を調整する形態として説明する。

【0035】

まず図1に示すように各尤度マップの間および各尤度マップの内部にはマップ間の結合105およびマップ内の結合106が備わっている。この結合を通じて、各マップの各領域ブロックは周囲のマップや領域の状態に応じて自身の状態を更新する。この結合の具体的な形態例の1つを図5に示す。図に例示するように、マップ間および同一マップ内のブロックの間にマップ間抑制結合23およびマップ内抑制結合24が備わっている。なお、図5の入力画像は、画像の左側にいる人物2人が重なっている。この2人の人物が同じ領域(ブロック)内で検出される例を示したのが尤度マップ群22の尤度マップ2と3の出力である。この場合、検出結果の見え方としては、画像座標(5, 3)に2つ分のスコア(閾値より大きい尤度)が検出されていることがわかる。異なる尤度マップでは、異なる人物に対して尤度が高くなるように学習済みの状態で、同じ画像座標に閾値より大きい尤度が検出されているのは、同じ人物のダブルカウントではなく、同じまたは隣接領域に存在する物体2つ分を検出したことになる。従来の方法では、入力画像に対して1枚の尤度マップを出力するため、隣接する物体がある領域には1人分の検出結果しか得られなかった。しかし、本実施形態では、同じ領域(ブロック)に隣接する物体がある場合でも、別々の尤度マップがそれぞれの尤度を算出するため、隣接する物体でも検出が可能になる。

10

【0036】

この結合の入力信号の総和を下記式に基づいて算出し(S5)、各尤度マップ v_k をそれぞれ更新する(S6)。なお、 α を含む重みは後述の学習処理の際に、一緒に学習される。色々な箇所で損失値の計算がなされ、すべての損失の総和を最小化するように学習が行われる。 β を含む全重みは、上記のような損失総和を最小にするように適宜調整される。 z_{INTRA} と z_{INTER} はそれぞれ尤度マップ内、尤度マップ間、の反応の出方を調べ、それらの影響を加味して反応を増強するか減少するかを決める調整値になる。

20

数式5

$$z_{INTRA_k}(x, y) = \sum_{(x', y') \in R} z_{INTRA_{k'}}(x', y') \cdot v_k(x + x', y + y')$$

$$z_{INTER_k}(x, y) = \sum_{k'} z_{INTER_{k'k}} \cdot v_{k'}(x, y)$$

数式6

$$v_k^{NEW}(x, y) = g(\alpha v_k + z_k^{INTRA} + z_k^{INTER} + \beta)$$

30

【0037】

ここで z_k^{INTRA} と z_k^{INTER} はマップ内・マップ間の入力信号の総和、 z_k^{INTRA} と z_k^{INTER} は結合重みの強度である。 z_k^{INTRA} と z_k^{INTER} は、それぞれ尤度マップ内、尤度マップ間、の反応の出方を調べ、それらの影響を加味して反応を増強するか減少するかを決める調整値になる。調整前の値 $v_k^{NEW}(x, y)$ にZの二つの値を足し、発散しないよう $g(\cdot)$ でゲイン調整する。Rは近傍のブロックの集合であり、同一マップ内で相互に結合する範囲を示している。および β は尤度マップの出力の調整値である。パラメータ α, β は後述の学習処理によってあらかじめ適切に調整されているとする。

40

【0038】

上記の更新処理は複数回繰り返してもよいし、一度のみでもよい。図6にこの更新処理についてのブロック図を2つの形態例で示す。図6(A)は再帰的な結合処理を示したものである。図6(B)は図6(A)の再帰結合の繰り返しを3回に限定し、結合を全て等価なフィードフォワード結合に置き換えたものである。図6(B)のような形態は特に学習時やハードウェア実装時に好適な場合がある。

【0039】

ここまでのS4~S7の処理は、複数の尤度マップを参照して、1つの推定結果を得るための統合処理の一例を示すものである。複数のマップで示された尤度を、統合前に互いに参照することによって、学習モデルの学習が十分に進んでいない状態でも、ルールペー

50

スで物体の位置が検出可能である。＜尤度マップの注目領域の尤度＞を数式 5 に基づいて、＜相互に入力＞して出力を調整する。1つの物体が2か所以上に検出されることや、2つ以上の物体が融合してしまうことといった状態を防ぐことを目的とした処理である。

【0040】

上記に相当する機能が実現されるのであれば、別の形態例として図7のような形態も本発明の適用例の1つである。ここではマップ内・マップ間の結合重みに相当する処理としてニューラルネットの畳み込み処理24を用いており、再帰的な結合は用いない。畳み込み処理24は3チャンネルの尤度マップを入力とし、3チャンネルの尤度マップを出力としている。これにより図5のマップ内・マップ間の結合による出力調整と等価な出力調整処理を実現することが可能である。なお図中では例として黒丸が畳み込みの負の係数の重みを、白丸が正の係数の重みを示している。

10

【0041】

＜統合処理＞

S8は、統合部107が、各尤度マップに分散している検出結果を統合し、統一した結果を生成するステップである。本ステップではまず統合部107が各尤度マップkの各領域ブロックの尤度スコアを調べ、スコアが所定の閾値 τ_k よりも大きい場合に顔が検出されたと判断して変数 d_k に1の値をセットする(下式)。この処理によって、この処理によって、尤度が所定の値より大きい領域を物体の存在する候補領域とする。

数式7

$$d_k(x, y) = 1 \quad \text{if } v_k(x, y) > \tau_k$$
$$d_k(x, y) = 0 \quad \text{Otherwise}$$

20

【0042】

上記のdを要素とする3次元の行列(x、y、およびkの3次元)をここでは尤度マップDとする。さらにこの時、非最大値抑制処理も併せて行う。具体的にはk番目の検出結果 d_k において、所定の距離以内に複数の物体が近接して検出された場合は多重検出であると判断する。そして、尤度スコアが最も高い物体のみを真の検出として残し、低い方を偽として削除する(非最大値抑制処理は非特許文献1等で一般的な公知の方法であるのでここでは詳細を省略する)。

【0043】

なお各尤度マップの中では非最大値抑制処理を行う一方で、各尤度マップ1~Nの間については非最大値抑制処理を行わない。複数の尤度マップの同一位置に複数の反応が生じた場合は、近接した位置に複数の物体が存在すると判断していずれの検出結果も残す。

30

【0044】

なおさらにこのとき、被写体のより詳細な位置を推定してもよい。具体的な例の1つとしては例えば尤度マップの尤度スコア値からサブピクセル推定の方法で行う(サブピクセル推定は各尤度マップでそれぞれ独立に行う)。サブピクセル推定の詳細については非特許文献5等で広く公知であるためそちらを参照されたい。(非特許文献5; Psarakis & Evangelidis, An Enhanced Correlation-Based Method for Stereo Correspondence with Sub-Pixel Accuracy, 2005)。

40

【0045】

また他の詳細な推定の方法としては位置やサイズを回帰推定するマップを別途用意し、マップの値に基づいて物体の位置・サイズを微調整してもよい(位置およびサイズ推定のための推定マップは各尤度マップkで用意する必要がある)。本発明においてこれらの工夫の有無は発明の本質に関わらないため詳細を省略する。非特許文献1等で公知なため必要に応じて参照されたい。

【0046】

以上のようにして尤度マップを統合した結果をまとめ、例えば図5の統合結果23のような検出結果のリストとして出力する。図には検出した物体の位置と尤度スコアからなるリストの例を示している。

50

【 0 0 4 7 】

最後にステップ S 9 で結出力部 1 0 8 が上記の統合結果に基づき顔の枠等を表示デバイス等へ出力する。以上で情報処理装置の認識動作が終了する。

【 0 0 4 8 】

なお統合処理の他の派生の形態としては、1枚ずつ尤度マップを調べるのではなく、一旦全マップを重み付き和等して1枚のマップにしてから尤度スコアを調べる、等の形態も考えられる。また非最大値抑制の有無やその方法についてもさまざまな選択肢がある。また尤度スコア値の閾値 k についても1段階のみでなく2段階の閾値を用いる等も考えられる。このように尤度マップの統合の形態については複数考えられ、特定の形態に限定されない。

10

【 0 0 4 9 】

< 学習処理の動作 >

次に本認識装置の学習動作について説明する。学習動作の際の情報処理装置の機能構成例を図8に示す。ここでは学習に必要な損失値算出部 2 1 0 が追加されている。

【 0 0 5 0 】

学習の処理を説明するフローチャートは図9(A)である。まず、図9(A)のステップ S 2 1 で、画像入力部 1 0 1 が、学習画像のセット(バッチデータ)を選択して画像を入力する。また同時に、記憶部 1 0 9 が各学習画像に対応する教師値を統合部 1 0 7 および損失算出部 2 1 0 へ入力する(ステップ S 2 2)。教師値は、各画像中に存在する物体の中心位置を示したものである。なお、教師値は検出する物体やタスクによって異なる。例えば、人物を検出する場合は、人の顔の中心位置に G T を与える。物体を検出する場合は重心位置等に G T を与える。具体的には、図 1 0 (A) の入力画像 1 0 0 1 に対する教師値は 1 0 0 2 であって、人物の顔の中心位置がある領域にラベルを付けた教師値になっている。教師データは、入力画像に対して正解の位置の座標のみを対応付けたデータでもよいし、人物の顔の中心位置には 1、それ以外の位置には 0 を入れたマップ形式のデータでもよい。

20

【 0 0 5 1 】

次に S 2 3 で特徴生成部 1 0 2 が、入力画像についての特徴を生成し、ついで尤度マップ生成部 2 0 4 a ~ 2 0 4 c が尤度マップを生成し、統合部 1 0 7 がそれらを統合した結果を生成する。尤度マップの統合は、予め決められた重みで統合してもよいし、統合方法を学習してもよい。

30

【 0 0 5 2 】

次に学習に必要なとなる認識結果の損失値の計算を行う。本実施形態の学習においては最終の統合結果についての損失値と、中間の生成物である尤度マップについての損失値の両方を学習計算に用いることとする。

【 0 0 5 3 】

< 統合した尤度マップに対する損失関数 >

まず統合した尤度マップに対する損失関数の方法は以下である。

【 0 0 5 4 】

S 2 4 では、統合部 1 0 7 が、統合した尤度マップに対する損失関数と予め与えた教師値とを比較して物体の検出位置についての損失関数(第2の損失関数)を用いた損失値を出力する。前記複数の尤度マップを統合した結果に対して、前記取得された位置に対応付けられた前記尤度が閾値より小さい場合は前記尤度を上げるための損失関数である。なお、損失値はさまざまな形態の利用が可能である。ここでは例として下式のように二乗誤差を用いて、領域ブロック (x, y) ごとに損失値を計算し、総和する。まず、教師データが示す物体の正解位置と、学習モデルに入力画像を入力することで推定された物体の位置と検出された物体の数との差を求める。

40

数式 8

$$L o s s ^{-}(D, T^{NUM}) = \sum_{x, y} (k d_k(x, y) - T^{NUM}(x, y))^2$$

【 0 0 5 5 】

50

ただしDは統合した尤度マップ、 $T^{NUM}(x, y)$ は教師値であり、領域ブロック (x, y) に顔の中心が位置する物体の総数を与える。上式の損失値を使って教師あり学習を行うことで、各尤度マップの反応結果が真の物体の数となるべく一致して反応するよう、ネットワークの重みパラメータが調整される（学習における重みパラメータの更新の方法については後述する）。つまり、数式8の損失関数は、正解位置に高い尤度が検出されなかった場合、または間違った位置に高い尤度が検出された場合は、すべての尤度マップに対して同程度の損失値を出力する。

【0056】

また損失関数の別の形態として下式のように交差エントロピーを用いることも可能である。数式9によれば、GTで示された正解位置に対応付けた各マップの尤度を比較する際に、いずれかの尤度マップを用いて物体の位置が検出できた場合がある。その場合は、その正解位置に対応づけられた尤度が低い他の尤度マップについては、その正解位置に物体が存在する可能性を示す尤度を小さくするように、各尤度マップ生成部（学習モデル）の層間の重みパラメータを学習する。ある領域に複数の物体が含まれる場合は、その物体の数に応じて大きな尤度を出力する尤度マップを準備する必要がある。そのため、GTから物体の数だけ大きい尤度（1.0等）を示すGTマップを生成し、もしGTの数よりも大きな尤度を示す尤度マップの数が少ない場合は、より物体を積極的に周囲の画像特徴をかくしゅうする。数式9の損失関数を用いることで、数式8よりもより細かい学習ができる。

数式9

$$Loss(V, T) = \sum_{x, y} \{ t_k(x, y) \log(v'_k(x, y)) - (1 - t_k(x, y)) \log(1 - v'_k(x, y)) \}$$

【0057】

ただしTは教師値であり、0か1かを要素tの値として持つ3次元の行列である。Tの各要素は各領域・各マップに物体（の中心）が存在するか否かを示している。

【0058】

なお、ここでは物体が同一ブロック (x, y) 内に複数存在する場合にも適切に尤度マップの損失値を計算するため下記のような工夫を設ける。まず、あるブロック (x, y) にn個の物体が存在する場合、教師値 $T(x, y)$ の値として、先頭にn個の1の値、残りを $N - n$ 個の0の値からなるN要素のベクトルtを与える。次に、尤度マップ $v_k(x, y)$ の尤度スコア値を降順にソートし、この値を $v'_k(x, y)$ とする（以降tおよび v' を<ソート教師値>および<ソートスコア値>と呼ぶ）。このようにしてから、数式9を用いて損失値を計算する。

【0059】

具体例を1つ示す。仮にいま $N = 4$ 枚の尤度マップがあり、真値としてブロック (x, y) に2つの物体が存在するとする。さらに認識結果として当該ブロック (x, y) の尤度マップの尤度スコア値が

数式10

$$V(x, y) = [v_1(x, y), v_2(x, y), v_3(x, y), v_4(x, y)]^T \\ = [0.1, 0.7, 0.5, 0.9]^T$$

と得られているとする。このとき、尤度スコア値をソートしたソートスコア値、およびソート教師値はそれぞれ

数式11

$$V'(x, y) = [0.9, 0.7, 0.5, 0.1]^T, \\ T(x, y) = [1, 1, 0, 0]^T$$

となる。当該ブロックの損失値は

数式12

$$Loss = \sum_k t_k \log(v'_k) - (1 - t_k) \log(1 - v'_k) \\ = 0.105 + 0.358 + 0.693 + 0.105 \\ = 1.261$$

10

20

30

40

50

と算出される（添え字を一部省略している）。各尤度マップが全体として全認識対象を過不足なく検出していれば、上記の損失値はゼロとなる。その際どの尤度マップがどの物体を検出したかは問わないことに注意されたい。この損失値は特定の対応関係（例えば尤度マップ1が前側の物体、尤度マップ2が後側の物体に反応するといった関係）を特に設けず、全体として検出精度が上がるよう各尤度マップ生成部のパラメータセットを学習させることを意味する。以上が尤度マップの〈統合結果〉についての損失値となる。

【0060】

なおここで示したようにDNNにおける損失値は、さまざまな形態の損失関数を採用することが可能である。本発明の適用対象は特定の形態の損失値に限定されない（なお学習の計算の都合上、損失関数は解析的に微分計算できる形が好適である）。

10

【0061】

＜尤度マップの損失値算出＞

次に、それぞれの尤度マップについて、尤度マップ内の各領域に対応づけられたスコアに基づいた損失値の算出の方法について説明する。まずS25では尤度マップ内の各領域に対応づけられたスコアに基づいた損失値を計算し、S26では複数の尤度マップ間の対応する注目領域に対応づけられたスコアに基づいて損失値を計算する。模式図を図10に示す。この二種類の損失値を適切に設計すれば、

- (1) 近接した複数の物体に対して複数の尤度マップが分担して反応する
- (2) 1つの物体に対して1つ以上の尤度マップが反応しない

の二つの性質を持った尤度マップ生成部106のパラメータセットを複数パターン学習で得ることが可能である。

20

【0062】

まず単一の尤度マップについての損失関数（第1の損失関数）を下式のように定義する。複数の尤度マップのうち同一の尤度マップ内の第1の注目領域について、該注目領域の近傍の領域に対応付けられた尤度と該注目領域に対応づけられた尤度との差に基づいて、近傍領域の尤度を下げるための損失関数を用いる。第1の損失関数は、注目領域の尤度と近傍領域の尤度との差が小さいほどより大きい損失値を出力し、近傍領域の尤度を小さくする方向の調整をする。また、第1の損失関数は、注目領域の尤度と、近傍の領域の尤度との差が大きいほど、より小さい損失値を出力するか、または損失値を与えない。このような損失関数を用いることで、同一マップ内の隣接する領域で得られる尤度のコントラストがよりはっきりし、1つの領域において1つの物体を検出するための尤度マップを得られる。言い換えれば、1つのマップ内の特定領域において複数の物体を検出することを抑制できる。

30

数式13

$$Loss_{INTRA} = - \sum_{k \in R} \log \left(\frac{1}{\sum_{k \in R} \exp(-\frac{(x - x_k)^2 + (y - y_k)^2}{\sigma_k^2})} \right)$$

ただし関数は

数式14

$$f(x, y) = \frac{1}{\sigma_1} \exp\left(-\frac{(x - x_1)^2 + (y - y_1)^2}{\sigma_1^2}\right) - \frac{1}{\sigma_2} \exp\left(-\frac{(x - x_2)^2 + (y - y_2)^2}{\sigma_2^2}\right)$$

40

である。図10(A)に示すような、正のピークとピークの近傍に負の窪みを持つメキシカンハット型の関数1003である（ σ_1 , σ_2 , x_1 および y_1 は関数形状を決定する定数のパラメータである）。本損失値は尤度マップに対して $f(x, y)$ を畳み込みカーネルとして畳み込んだ結果の総和の値である。Rは畳み込みを行う領域の範囲である。

【0063】

複数の反応が近接して同時に一枚の尤度マップ上に生じると、損失値 $Loss_{INTRA}$ は大きな値を取る。例えば仮に図10(A)のように物体2人が近接して写っている入力画像1001が入力されたとする。尤度マップ生成部が生成した尤度マップ群1004では、尤度マップ1が両方の物体に対して反応している（反応の強さをグレースケールの濃淡で表す）。一方の尤度マップ群1005では尤度マップ1と尤度マップ2に分散して反

50

応が生じている。この場合メキシカンハット型の関数の性質のため、前者では大きな、後者では小さな損失値が算出される。

【0064】

上記は「近接した複数の物体が一枚の尤度マップで同時に検出される」ことを抑制する損失関数の一形態であるが、本発明が適用可能な形態としてはこの他の形態も考えられる。例えば下式のような損失値の定義もあり得る。

数式15

$$Loss_{INTRA} = \sum_k (x, y) V_k(x, y) - t_h$$

【0065】

ここで (\cdot) は半波整流の関数である。同損失値は、各尤度マップの尤度スコアの総和が所定閾値 t_h を越えるとペナルティを与える。すなわち、複数の尤度マップのうち同一の尤度マップ内の、第1の注目領域について、該注目領域に対応付けられた尤度が閾値より大きい場合は該注目領域の近傍の領域に対応付けられた尤度を下げるとの損失関数を用いる。また、該注目領域に対応付けられた尤度が閾値より小さい場合は、第1の注目領域の近傍の領域に対応付けられた尤度を下げるとの損失値をより小さくなるように損失値を出力する。このためこの損失値を用いて学習を行うと、どれか1つの尤度マップだけが突出して物体に反応するような動作が抑制される。以上が単一の尤度マップに対して定義される損失値の説明である。

10

【0066】

次に複数の尤度マップ間の対応する注目領域に対応づけられたスコアに基づいて損失値を出力する損失関数(第3の損失関数)の例を示す。第3の損失関数は、複数の尤度マップのうち所定の2つの尤度マップにおいて尤度を比較する。入力画像の物体の位置(第2の注目領域)に大きな尤度に対応づけられた場合に、2つの尤度マップで競合する(対応領域の)尤度を下げるように学習モデルの重みパラメータを学習するための損失関数である。なお、第2の注目領域は、取得された教師値によって示された物体が存在する位置を含む各尤度マップの対応する領域である。ただし、注目領域に2つ以上の物体が存在する場合は、物体の数と尤度マップの数に応じて損失値が異なる。注目領域に物体が1つだけ存在する場合は、第2の損失関数は、2つの尤度マップにおいて共通する位置の注目領域に対応づけられた尤度に基づいて、2つの尤度マップにおける注目領域の少なくとも一方に対して大きな損失値を出力する。注目領域に物体が2つ以上存在する場合は、物体の数 m と同じ数の尤度マップに対しては注目領域に推定された尤度が低いときはその尤度を大きくするような損失値が決定される。物体の数 m より尤度マップが多い場合は、注目領域を尤度の大きい順にソートしたときの $m + 1$ 番目以降の尤度マップの注目領域に対し、尤度が小さくなるような損失値を決定する。また、入力画像の物体が存在しない位置に対応する尤度マップの領域に対して大きい尤度が推定された場合は、その尤度を小さくするために大きな損失値が決定される。つまり、第3の損失関数は、取得された物体の位置に基づいて入力画像に含まれる物体の位置が1つである場合は、複数の尤度マップのうちいずれか1つの尤度マップの第2の注目領域に対応付けられた尤度がより大きくなるようにする。また、取得された物体の位置についての情報に基づいて入力画像に含まれる前記物体の位置が複数ある場合は、物体の数に応じて尤度マップの第2の注目領域に推定された尤度が所定の閾値より大きくなるように調整するような損失値を決定する。損失値が小さいほど学習が進んでいると判断できるため、後の学習処理では算出された損失値を小さくする(または収束させる)ためにパラメータを調整する。この第2の損失関数によって、異なる尤度マップで同一物体を検出しないようにし、異なる尤度マップで異なる物体を検出できるように検出する対象を役割分担させることを目指す。なお、所定の2つの尤度マップとは、生成された複数の尤度マップからすべての組み合わせを指す。

20

30

40

数式16

$$Loss_{INTER} = - \sum_{x, y} k_k(T(x, y)) \log \{ k(V'(x, y)) \}$$

と定義する。ただし、 k はソフトマックス関数

50

数式 17

$$T(x) := \exp(x_i) / \sum_{j=1 \text{ to } N} \exp(x_j), \quad X = [x_1, \dots, x_N]^T$$

である。T(x, y) および V'(x, y) は N 個の要素からなるベクトルであり、先掲の <ソート教師値> および <ソートスコア値> と同じものである。

【0067】

上記の損失関数の結果の例を図 10 (B) に示す。ここでは入力画像 1006 のように単一の物体が写った画像が入力されている。これに対して反応結果 1008 のようにマップ 1 とマップ 2 の両方が同時に反応した場合、損失値 Loss^{INTER} は大きな値をとる。対して反応結果 1009 のように正しくいずれか 1 つの尤度マップのみが反応している場合は、同損失値は小さな値をとる。

10

【0068】

上記ではソフトマックス関数およびソートされた尤度スコアを用いたが、これは実現例の 1 つを示すのみである。「1 つの物体に対してなるべく 1 つのマップの領域しか反応しない」ことを促進するような損失関数の設計であればさまざまな形態の採用が可能である。以上が尤度マップに関する損失値の定義である。

【0069】

<学習パラメータの更新>

次にこのようにして得られた各種の損失関数が出力した損失値を使って、各マップ生成部のパラメータを更新する。本形態で学習更新の対象となるパラメータの 1 つは図 8 の尤度マップ生成手段 204a, 204a, 204c, の、それぞれの重みパラメータ w_k, b_k である。さらに各尤度マップ生成手段のマップ内結合とマップ間結合の重みパラメータ $w_{k'}$, $b_{k'}$ 、および結合調整パラメータ α_k , β_k である (ここで k は k 番目の尤度マップ生成手段を表す添え字である)。図 8 では右上方向の矢印を付して学習対象のパラメータを示す。

20

【0070】

上記各パラメータはそれぞれ乱数で初期化してから学習を開始する。特に各尤度マップはそれぞれ同質のマップであり、マップ間に競合的な損失値を与えて学習させることで、対象に対する反応が各マップに分散するように誘導することを企図している。もし各マップの重みパラメータ w_k, b_k の初期値が同一であると、常に同じ反応となって競合し、適切に学習が進まない。そのため各尤度マップは必ず異なる値で初期化する必要がある。

30

【0071】

前述の方法で算出した損失値の総和の値を E とし、E を入力画像のバッチセットごとに算出し、これを減らすような勾配の方向に各パラメータを微小に更新すればよい。具体的には下式のようにパラメータの値を各々更新する (S27)。例えば、E が所定の値より小さい値に収束するまで、パラメータセットを更新する。

数式 18

$$E = Loss + \alpha_1 Loss^{INTRA} + \alpha_2 Loss^{INTER},$$

$$w^{t+1} = w^t + \eta \frac{\partial E}{\partial w^t},$$

$$b^{t+1} = b^t + \eta \frac{\partial E}{\partial b^t},$$

$$\alpha_k^{t+1} = \alpha_k^t + \eta \frac{\partial E}{\partial \alpha_k^t},$$

$$\beta_k^{t+1} = \beta_k^t + \eta \frac{\partial E}{\partial \beta_k^t}.$$

40

【0072】

ただし Loss は尤度マップの統合結果に対して算出された損失値、 α_1 , α_2 は各種の損失値のバランスを調整するための定数、 η , η' は適当に設定された 1 以下の微小な係数 (学習係数) である。なおここで各偏微分 $\frac{\partial E}{\partial x}$ の値はニューラルネットワークの一般的な方法である誤差逆伝搬法を用いて求める (誤差逆伝搬法については非特許文献 4 などに広く公知のためここでは省略する。また上式は見易さのため添え字を一部省略している) (非特許文献 4 : Y. LeCun et al. Handwritten

50

digit recognition with a back-propagation network. 1990.)。

【0073】

なお尤度マップ生成部の結合重みは再帰的な結合を含むが、の学習には再帰的ネットワークの学習において一般的な方法を併せて用いるものとする（たとえば図6(B)のように有限繰り返し数の処理ブロックに展開した上で、誤差逆伝播法で更新する）。

【0074】

なおここで学習対象としなかった特徴生成部102のニューラルネット102aの重みについても、同様に誤差逆伝播法で学習してもよい（これは入力から出力まで一貫して重みパラメータを学習する形態でありEnd-to-end学習と呼ばれる）。

10

【0075】

なおさらに、統合部107の検出閾値パラメータである k 等を学習対象パラメータに加えるような形態なども考えられる。

【0076】

<派生の形態>

ここまで各処理ブロックの機能モジュールやその学習形態について順を追って説明を行ってきた。ここでは考えられるその他の派生の形態についていくつかの例を加える。

【0077】

例えば損失値の計算において、統合結果の損失値と尤度マップの損失値の算出方法についてそれぞれ述べたが、他の形態として、どちらか片方のみを用いたり、部分的に用いたり、学習の進み具合に応じてこれらを切り替えたりといった形態も考えられる。

20

【0078】

また例えば、ここまでは物体を検出する際には、物体の中心位置を基準位置として学習し、検出したが、この基準位置を変えることもできる。例えば（尤度マップの数は増えるが）、物体の上下左右端を基準位置としてそれぞれを推定するようなマップを学習し、検出するような形態でもよい。

【0079】

また、本実施形態ではN個の尤度マップ生成部、およびその結果としてのN枚の尤度マップを用いて対象を認識したが、この数Nを認識時に動的に変更するような方法も考えられる。例えば画像内に対象が多数重畳しているときは尤度マップ数が多いほうが、検出精度が高くなると考えられるが、物体の数に対してマップ数が多すぎるとノイズ状の反応が却って増えることや、余計な演算量が増えることがある。そのため尤度マップ数を適切に増減するような形態も考え得る。

30

【0080】

これを説明するために図8に追加的にマップ数決定部211と重みパラメータ提供部212を示す。マップ数決定部211は画像特徴103に基づいて画像シーンを考慮して最も良好な結果が得られるようにマップの数 n を決定する（例として対象物体で混雑した画像に対しては大きな n を与える等）。次にその結果を受けた重みパラメータ提供部212は1~ n の尤度マップ生成部に対して n 個の重みパラメータの提供を行う。

【0081】

40

マップ数決定部211がマップ数を決定するやり方として例えば以下の形態が考えられる。まず1個、2個、...、N個の尤度マップ生成部からなる、N通りの異なる設定の情報処理装置を用意し、各個に学習を済ませておく。次に入力画像 x が与えられたときの検出結果の精度を各N通りの設定について調べ、精度の良し悪しの値を記憶しておく（検出精度の良さを測る基準として例えば先に挙げた統合部107の統合結果の損失値などを使えばよい）。

【0082】

マップ数決定部211は画像 x の画像特徴103を説明変量とし、検出結果の精度を目標変量として、各N通りの情報処理装置の検出精度を推定する回帰学習を行う（ニューラルネットやサポートベクトル回帰等の一般的な方法を用いる）。認識時にはこの回帰器の

50

推定結果を用いて、マップ数決定部 2 1 1 が各 N 通りの設定の検出精度の期待値を調べ、最も期待値の高かったマップ数 n を採用する。なおこのときに、推定検出精度に使用マップの少なさ (= 総計算量)、も考慮するような合成指標を使って精度と計算量から使用マップ数を決定してもよい。以上のようにすることで、動的に尤度マップの構成を変更することが可能である。

【 0 0 8 3 】

またさらに別の派生の形態は以下のものである。これまで、本実施形態では説明の便宜上、各処理モジュールを明確に区別できるものとして説明してきた。ここで考えられる他の形態としては、ニューラルネットの各部が本実施形態の機能モジュールと同等機能を持ち、それらが境目なく結合した形態である。

【 0 0 8 4 】

例えば、特徴の生成部 1 0 2、複数の尤度マップ生成部 1 0 4、尤度マップ間の結合 1 0 5 や、マップ内の結合 1 0 6、および統合部 1 0 7、の各機能が、DNN の各層の上に分散して実現するような形態が考えられる。このような機能を実現するためには、どの層にどの機能的役割を実現させるかをある程度決めた上で、各機能が十分実現できる程度の層数および入出力チャンネル数、および層間の結合関係、を持った DNN を用意する。次に実現させたい機能的な役割に応じて中間層や最終層に対する損失関数を適切に設計し、学習データと教師値を与えて学習させる。学習がうまくいかない場合は中間層の出力を調べ、機能実現のために特徴表現や識別力が足りていないケース、出力調整に必要な繰り返し数 (層数) が足りていないケース、等を同定して適宜必要な箇所のニューロンの層数やチャンネル数を増やす。1 つの形態例として以上のような方法が考えられる。

【 0 0 8 5 】

またさらに別の派生の形態は以下のものである。同一種類の複数の物体に対して複数の尤度マップが反応することが本発明の要諦の 1 つであるが、ここでの「同一種類」とは物体カテゴリーの部分集合であってもよい。例えば、物体を見えのサイズ、アスペクト比や姿勢で分けたものをそれぞれ異なるカテゴリーとしてもよい。またさらに、任意の複数のカテゴリーをまとめた上位集合を作り、これを 1 つの種類であると見なす形態であってもよい。例えば、犬と猫を 1 つのカテゴリーとする等がこれに当たる。

【 0 0 8 6 】

例を示すと、種類の異なる複数のカテゴリー A, B, C の物体があり、本情報処理装置に複数の尤度マップ X, Y があるとすると、もし尤度マップ X, Y どちらもカテゴリー A, B, C の物体全てを検出するように学習を行うのであれば、これは本発明の一形態である。さらに、尤度マップ X はカテゴリー A, B に対して、尤度マップ Y はカテゴリー B, C に対して、それぞれ反応するように学習するような形態も、カテゴリー B に対しては複数の尤度マップが反応するという点において本発明の形態の 1 つである。また尤度マップ X, Y とともにカテゴリー A, B, C の被写体すべてに反応するが、尤度マップ X は特にカテゴリー A に優先的に反応し、尤度マップ Y はカテゴリー B に優先的に反応する、というように尤度マップを部分的に特性付けるような形態も考えられる。

【 0 0 8 7 】

以上ここまで本発明を適用して得られる派生的な形態について説明を加えてきた。本発明をなす根幹を述べると、< 同一種類の複数の物体が存在している > ときに、それらが < 複数の同質の尤度マップ上に分散して検出される > よう設計された各機能モジュールあるいはそれらの重みパラメータ、とまとめられる。上記目的に適う機能モジュール、重みパラメータ、(あるいは適切にパラメータを学習するための損失関数) は本発明の実施形態に含まれ、本発明の実現形態は特定の形態のみに限定されない。以上で実施形態 1 の説明を終える。

【 0 0 8 8 】

< 実施形態 2 >

第二の実施形態では、第一の実施形態と同様に物体の検出を目的とする。情報処理装置の基本的な構成例のブロック図は図 1 1 である。実施形態 1 と異なる点は各マップ 3 0 4

10

20

30

40

50

a ~ 304c に新たに複数の特徴集計部 303a ~ 303c が加わっている点である。本実施形態は各尤度マップ生成部に与える特徴にバリエーションを加えた形態となっている。これにより各尤度マップの出方に変化が生じ、近接した物体の分離度が向上する（理由は後述する）。

【0089】

なお本実施形態では各尤度マップ生成部 304 の間および内部の結合は除いており、実施形態 1 の処理フロー中で行った尤度マップの更新は行わない。ただし派生的な形態として実施形態 1 と同様に上記結合を構成に含めて尤度マップの更新を行うことも考えられる。また、ハードウェア構成は実施形態 1 と同様に図 21 のような構成を用いる。

【0090】

画像特徴の抽出処理の説明のための模式図を図 12 に示す。図 12 (A) はこれまでの画像特徴の形態図である。ニューラルネットの各階層の出力結果を連結して一種類の階層特徴

$F(x, y) = [f_1(x, y)^T, f_2(x, y)^T, f_3(x, y)^T]^T$
を生成して用いている。

【0091】

図 12 (B) は本実施形態で開示する画像特徴の抽出処理の形態である。階層特徴の生成時の集計の方法を N 通りに変更することで、

$F_k(x, y) = [f_{k1}(x, y)^T, f_{k2}(x, y)^T, f_{k3}(x, y)^T]^T$
($k = 1, 2, \dots, N$) と、N 通りの特徴を生成している。

【0092】

集計方法の具体例を図 13 に示す。同図は特徴集計部 303 が CNN の第 j 層の出力の特徴 f_{raw_j} に 4 通りのサブサンプルを行って 4 通りの画像特徴 $f_{1j} \sim f_{4j}$ を生成している。すなわち、モデル毎に異なる画像特徴が入力されるようになっている。ここではサブサンプルにより特徴マップの縦横それぞれの解像度を 2 分の 1 にしているが、 2×2 の領域ブロック範囲 ($Range(k, j)$ と記号を付して示す) の位相を都度変更しながらサブサンプルを行う。これにより各特徴 f_{kj} がそれぞれ微妙に異なるバリエーションを持った画像特徴となっている。

【0093】

本実施形態 2 の処理のフローのうち、特に画像特徴抽出部分（実施形態 1 でのステップ S2 に相当する）について詳細化したフローを図 14 に示す。本フローでは特徴集計部 303 が、ステップ S42 ~ S49 のループで N 通りの画像特徴の集計および生成を行う。ステップ S43 ~ ステップ S46 では、尤度マップの番号 k に応じて異なる集計範囲 $Range(k, j)$ を設定する（ステップ S44）。同範囲でサブサンプルを行って特徴 f_{kj} を生成する（ステップ S45）。さらにステップ S47 では $f_{kj}(x, y)$ を連結して階層特徴 $F_k(x, y)$ を生成する。そして、尤度マップ生成部 k へと出力する（ステップ S48）。

【0094】

このようにサブサンプルの集計パターンを様々に変更することで、位相等の微妙に異なる特徴を複数の尤度マップ生成部に提供することができる。異なる特徴に基づいて物体の尤度スコアをそれぞれ判定するため、単一の特徴に基づいて判定するのに比較し、近接した物体パターンを分離・検出できる可能性が高い。なお同様の方法として、ニューラルネット 102a で高解像の入力画像を処理して解像度の高い画像特徴を得て用いることも考えられるが、ニューラルネット 102a の計算量の増大を伴う。本実施形態の形態は特徴の集計方法を変えるだけのため、計算量を特段に増やすことなく同種の効果が得られる。

【0095】

なお上記は形態の一例であり、集計の方法の変化の付け方は他にもあり得る。（ 2×2 ではなく 1×2 と 2×1 といった非正方形の範囲を用いる。サブサンプルの他に最大値プーリングや平均値サンプルを行う。一部分岐した DNN を用いて特徴ごとに一部だけ異なる層の特徴マップを連結させる、等）。またベースとなるニューラルネットの特徴はここ

10

20

30

40

50

では階層型の画像特徴を用いているが、適用可能な形態はこれに限定されない。

【0096】

なお学習時は実施形態1と同様に損失を計算し、誤差逆伝搬の方法で各重みを更新すればよい。

【0097】

<実施形態3>

本実施形態では、本情報処理装置への情報入力部の1つとして画像の奥行き情報（以下2.5次元情報）を加え、これを利用する形態について説明する。さらにユーザーの指示を受ける表示切替指示部を設け、ユーザーの介在に基づいて認識結果の提示の仕方を切り替える形態について説明する。図15に機能構成図を示す。また、ハードウェア構成は実施形態1と同様に図21のような構成を用いる。

10

【0098】

2.5次元情報の利用方法としては学習時に用いる場合と、認識時に用いる場合の2種類が考えられる。

【0099】

<学習時の2.5次元情報の利用>

学習時の2.5次元情報の利用の仕方の1つとして、以下に例を挙げる。まず処理フローの図9(B)に示すように、観測した2.5次元情報のマップを真値の一種として与える（ステップS33）。ここでは2.5次元情報のマップを変数 $Dist(x, y, d)$ として表す。2.5次元情報のマップ $Dist(x, y, d)$ は簡単のために尤度マップと同じ画像解像度を持った3次元の行列であるとする（なお奥行き方向 d は $d = 0, 1, \dots, d_N$ とあらかじめ離散化している。 d_N は最大の距離に対応する適当な値である）。行列 $Dist$ の要素のうち、物体が存在する位置・奥行きに当たる要素には1、それ以外の要素には0が入っているとす。次に下式のようにマップ内の損失値の計算において2.5次元の値を利用する（ステップS36）。

20

数式19

$$Loss_{INTRA}(x, y) = - \sum_k v'_k(x, y, d) \cdot \sum_d Dist(x, y, d)$$

30

ただし、ここで $'$ はメキシカンハット関数を奥行き方向に拡張した下記の関数である。

数式20

$$v'_k(x, y) := \frac{1}{2} \exp\left(-\frac{x^2 + y^2 + d^2}{\sigma^2}\right)$$

(σ はスケール調整の定数)

$v'_k(x, y, d)$ は尤度マップの尤度スコア $v_k(x, y)$ を奥行き方向に複製して便宜的に3次元に拡張した変数である。（ $v'_k(x, y, d) := v_k(x, y), d = 0, 1, \dots, d_N$ ）

上記式は、画像の像面上の距離、および奥行き上の距離、が両方共に近い物体を1つの尤度マップで検出することに対してペナルティを与えることを意味する。この損失値を使って学習することにより、奥行きおよび画像面上の距離の近い物体はなるべく各尤度マップに分散して反応が生じるように誘導される。

40

【0100】

なお同様に距離情報を利用した派生の形態として、手前側の物体を大きな番号の尤度マップで検出し、遠い側の物体を小さな番号の尤度マップで検出したときにペナルティを与えるような損失値、といった形態も考えられる。ただしマップの数を大きく超える多数の物体が1列棒状に並んでいるような場合（集合写真等では頻繁に起こりえる）、このように距離と尤度マップを密接に対応させて学習させると、物体の検出分離度は逆に悪くなるので注意が必要である。本実施形態の数式19の形態のように、奥行き情報を相対的にのみ用いる方法が本発明においてはより好適であると考えられる。

50

【0101】

<認識時の2.5次元情報の利用>

認識時に奥行き情報を利用する形態の1つとしては、2.5次元情報のマップを画像特徴の1つとして連結して認識に用いることである。

【0102】

また他の形態の1つは、認識時に、ユーザーの指示部である表示切替指示部4001を用いて、ユーザーの指示と奥行き情報に基づいて表示を切り替える形態である。

【0103】

この後者の形態について詳細を述べる。入力画像を示す図17(A)および結果の表示例を示す図17(B)~(G)を利用して説明する。まず認識処理が開始されると図17(B)のような画像が入力される(ステップS61)。次に尤度マップが生成される(ステップS62)。次に尤度マップを統合することで、例として図17(C)のような検出枠の結果が得られる(ステップS63)。この結果をそのまま表示すると視認性が低い。そのため、考えられるユーザーインターフェースとしては2.5次元情報入力部4002から入力された奥行き情報を用いて、検出物体のうちもっとも距離の近い物体を判定してその検出枠を表示することが考えられる(ステップS64)。同時に拡大表示窓302aにその拡大結果も表示する(ステップS65)。結果例を図17(D)に示す。次に表示切替指示部4001の1つの形態例である左右矢印状のボタン301aによってユーザーが指示を与え、表示の切り替えを行う(ステップS66、ステップS67)。このとき、ユーザーの左右ボタンの押し下げに応じ、被写体の左右位置順に結果を切り替えて表示する形態が考えられる。このように表示を切り替えた例を図16(E)に示す。また物体の遠近順に切り替えて表示する、といった形態もあり得る。また、ユーザー指示部を用いずに、表示を一定時間ごとに自動的に切り替えるような形態も考えられる。

【0104】

なお派生の形態として、表示するのは上述のような検出枠等でなく、尤度マップの尤度スコア値を濃淡マップとして表示するようなことも考えられる。このようにマップの値をユーザーの指示に応じて表示することで、ニューラルネットの認識結果の内部状態の情報を提示するインターフェースとして利用することが考えられる。図17(F)にその結果の例を示す。ここでは尤度マップの尤度スコア値を値の大きさに応じた灰色~黒色の矩形で示している。さらにユーザーが左右矢印状のボタン301bを押し下げて指示を与え、表示する尤度マップの結果を切り替えて表示した例を図17(G)に示す。図17(F)と図17(G)は2つの異なる尤度マップの尤度スコアを切り替えて示している。そのため尤度スコア値の濃淡の一部には重なった領域が存在し、それぞれの値が異なっている点に注意されたい(例えば記号303を付した矩形)。ここでは尤度マップの尤度スコア値をそのまま全て表示しているが、検出結果と2.5次元情報を併せて用いることで、各検出された物体ごとに尤度スコア値を分けて表示する、等も考えられる。この機能によって、ユーザーは学習済みモデルの学習が十分に進み、物体を正確に検出できていることを確認できる。

【0105】

<実施形態4>

実施形態4では画像の意味的領域分割(セマンティック・セグメンテーション)を行う情報処理装置について説明する。本発明が物体毎の物体検出タスクのみならず、領域ラベリングのタスク等種々のタスクに対しても広く適用可能であることを示す。

【0106】

重畳・近接した物体の領域を正しく同定・分離するのは一般に困難な課題である。非特許文献6などはこのためにまず画像中の物体の検出を行い、同定された物体領域を入力として再度DNNの処理を行い、各画素が前景か否かを判別して物体の前景領域を生成している。この方法は領域分割の精度が高いが、物体ごとにDNNの処理を行うため演算量は多い(非特許文献6:K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask R-CNN, ICCV, 2017)。

10

20

30

40

50

【0107】

本実施形態の情報処理装置は物体ごとにDNNの処理を行わない。非特許文献6が物体ごとに前景領域を判別するマップを用意するのに対して、本発明の実施形態はN個の尤度マップのみを用いる(マップの数Nは画像中に同時に出現する対象物体の数より小さいことを想定している)。同マップを以降領域尤度マップと呼ぶ。本実施形態では、物体領域として物体の領域を同定することを目的とする。また複数の物体については、物体ごとに領域を分離・同定することを目的とする。

【0108】

<学習動作>

学習時には、教師値として図20(B)のように領域のラベルの真値、 $l(x, y)$ {0, 1, ..., L}を用意する。(図では物体領域の色の違いでラベルの値の違いを表現している)。ラベル $l(x, y)$ が0の領域ブロックは物体が存在しない領域である。ラベルが1, 2, ..., Lの領域は物体領域であり、異なる数値でそれぞれ画像中の異なる物体を意味している。物体毎の領域を示した領域情報が本実施形態の教師データである。すなわち、教師値とは、各画素にどの物体が存在するかを示すラベル(例えば、左の人物は1、中央の人物は2、左の人物は3、人物がいない領域は0)を持った画像情報である。

10

【0109】

各領域尤度マップは、物体領域($l(x, y) > 0$)に対しては大きな尤度スコア値が、それ以外の領域($l(x, y) = 0$)に対しては小さな尤度スコアが出るように重みパラメータを学習する。用いる損失関数としては実施形態1の物体検出タスクと同様の交差エントロピー等であればよい。具体的には例えば実施形態1の数式8や数式9を用いる。また実施形態1の物体検出タスクでは物体の中心の領域ブロックに正の教師値を与えて学習したのに対して、領域判別タスクでは物体の領域に対応するブロックすべてに正の教師値を与えて学習する。

20

【0110】

ここで本発明の特性である<複数の尤度マップが分散協調して認識を行う>ことを実現する形態として、さらに以下のような二つの特性を実現する損失関数の項を加える。

- (1) 1枚の領域尤度マップは、近接や重畳している複数の物体領域に同時に反応しない
- (2) 複数の領域尤度マップは、同一物体の領域に対して同時に反応しない

まず上記(1)を実現する損失関数を説明する。形態としては種々あり得るが、例えば下式のようなものである。

30

数式21

$$Loss_{INTRA} = \sum_{x, y} \sum_R \{ 1 - (l(x, y) - l(x + x', y + y')) \} \times v_k(x, y) \times v_k(x + x', y + y')$$

上記の損失関数は、1枚の領域尤度マップが、異なる複数の物体の領域に対して反応した場合に損失値のペナルティを与える。ただし δ はディラックのデルタ関数であり、2つの領域のラベルが同じ値の時に1、異なる時に0を返す。 v_k はk番目の領域尤度マップの物体領域の尤度スコア値である。またここでRは同時反応を抑制する所定の近傍のブロックの範囲であり、この範囲の外であれば異なる物体の領域に反応してもペナルティを与えない。

40

【0111】

次に先述の(2)の特性を実現する損失関数についてであるが、これは実施形態1の数式16等を使えば実現できる。

【0112】

以上に述べた損失関数を用いて、損失値の総和を下げるように学習対象の各パラメータを学習更新する。学習が進めば、領域尤度マップが正しく物体の領域に反応し、且つ近接・重畳した複数の被写体の領域はなるべく異なる複数の領域尤度マップに分散して検出されるようになる。

【0113】

50

< 認識動作 >

図 18 は領域尤度マップを生成する情報処理装置の機能構成例のブロック図である。基本的な構成は実施形態 1 のものとほぼ同一であり、同一の処理を行うモジュールには同じ番号を付している。実施形態 1 と異なる点の 1 つは物体位置推定部 500 が新たに加わっている点である。また他の異なる点の 1 つとしては尤度マップ生成部 104 および統合部 107 が検出する対象が物体の中心位置（実施形態 1）か、物体の前景領域か（本実施形態）の違いがある。処理を説明するフローチャートは図 19、処理の過程と結果の一例は図 20 に示す。また、ハードウェア構成は実施形態 1 と同様に図 21 のような構成を用いる。

【0114】

ここで、領域尤度マップを用いた認識処理の流れを簡単に説明する。これまでの実施形態と同様に、まず、情報処理装置の画像入力部 101 が認識対象となる入力画像を入力する（ステップ S71）。次に、特徴抽出部 102 が、入力画像から画像特徴 103 を抽出する（ステップ S72）。複数の領域尤度マップからなる尤度マップ生成部 104 が対象被写体の領域ブロックか否かを示す尤度スコアのマップを生成する（ステップ S73）。領域尤度マップ 1 の結果例を図 20（C）、同領域尤度マップ 2 の結果例を図 20（D）に示す。図では尤度スコア値の大きさをグレースケールの矩形で表示している（黒いほどスコアが高い）。なお各領域尤度マップは近接した物体が異なるマップ上に分散して検出されるように、尤度マップ生成部 104 が用いる学習済みモデルはあらかじめ学習が施されているものとする（その方法については後述する）。

【0115】

ステップ S74 ~ ステップ S76 は領域尤度マップの統合処理になる。まずは統合部 107 が所定閾値以上の尤度スコア値を含む領域ブロックを物体の領域とする（ステップ S74）。図 20（E）に物体が存在する領域と判定された領域の例を示す。ここでは領域尤度マップ 1 で物体が存在する領域と判定した領域は黒の矩形で、同じく領域尤度マップ 2 で物体を検出した領域は灰色で示している（なお領域尤度マップ 1 と 2 の両方が物体領域とした箇所は、尤度スコア値がより高い方のマップの色で示している）。つまり、図 20（E）の物体領域マップは前述した尤度の大きさに応じた色分けではなく、領域尤度マップ毎に閾値以上の尤度を検出した位置を示すマップである。

【0116】

次にステップ S75 では、物体位置推定部 500 が物体の位置検出を行い、物体の位置の情報を統合部 107 に提供する。物体の検出の方法はこれまで実施形態 1 や非特許文献 1 等に開示されるような方法を別途行うものとする。ここでは実施形態 1 の方法を用いて、画像特徴 103 に基づいて検出したとする。図 20（F）に記号 501f ~ 503f を付して検出された物体の検出枠の例を示す。

【0117】

次にステップ S76 では、統合部 107 が、物体が存在する領域を個々の物体の領域に分割する。方法としてはまず検出枠と領域尤度マップとを対応づける。ここでは各枠内の領域のうち物体と判定した領域の数が最も多かった領域尤度マップを各枠に対応させる。（例として図 20（F）の検出枠 503f の場合、領域尤度マップ 1 を対応させる。次に各検出枠に、対応する領域尤度マップの物体領域を各物体の領域として決定する（例として図 20（G）の領域 503g）。

【0118】

最後に、ステップ S77 で、出力部 108 が、各物体が存在する領域を示す結果を出力する。例えば図 20（G）のように一体ずつ物体領域 503g を表示してもよい。その際、複数の物体が重畳している領域（例えば図 20（H）に符号 504h を付して示した灰色の矩形の領域）は各領域尤度マップの尤度スコア値の大小や 2.5 次元情報を用いて前側か後側かを推定してもよい。この場合被遮蔽領域に矩形 504h のような色の変化をつけて表示してもよい。また図 20（I）のように全ての物体領域を重ね、物体ごとに領域の色を変える等して表示してもよい。実施形態 3 と同様にこれらの表示を切り替えるユー

10

20

30

40

50

ザー指示部等を備えてもよい。

【0119】

なお領域尤度マップの統合にはここで述べた以外に様々な細かな派生的工夫が考えられる。例えば物体領域の数ではなく、尤度スコア値の大小を対応づけの判定に用いる、2.5次元情報を用いて前側の物体の検出枠から領域の対応づけを行っていく。対応づけ済みの領域は取り除いて重畳領域を正確に対応づける。領域尤度マップの純度の高い（他のマップとの混在が少ない）枠の領域から対応づけしていく、等である。また物体の検出枠を使わずに、領域尤度マップのうち孤立した物体領域をそれぞれ別個の物体の領域とする簡便な方法等もあり得る（遮蔽された物体が複数の領域に過剰分割される可能性もあるため注意を要する）。このように様々な形態の方法が考え得るが、本発明の根幹に関わりのない表層的な相違のため、ここでは詳細を略す。

10

【0120】

以上、本発明が物体検出タスクのみならず領域判別タスクにも適用可能であることを示した。特に、複数の領域尤度マップを同時に使うことで、同種の物体が近接・重畳した場合に適することを示した。またこれは非特許文献6のような、物体の検出を行ってから物体ごとにDNNで領域判定を行う演算量の大きな方法と異なる形態であることを示した。

【0121】

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。

20

【符号の説明】

【0122】

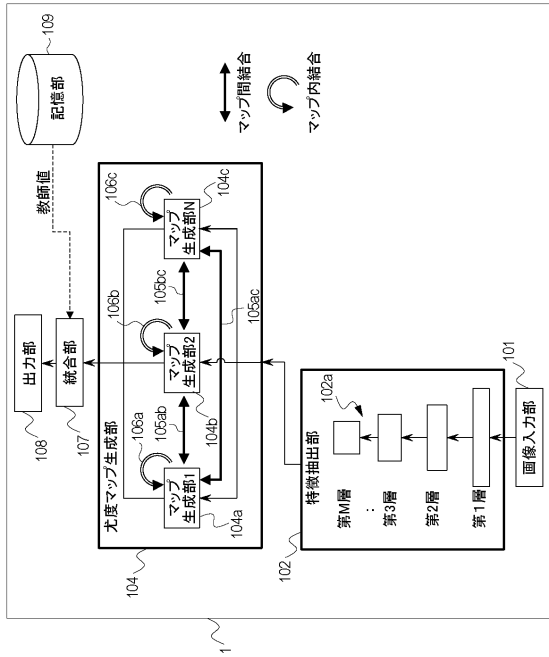
- 101 画像入力部
- 102 特徴抽出部
- 103 画像特徴
- 104 尤度マップ生成部
- 105 マップ間結合経路
- 106 マップ内結合経路
- 107 統合部
- 108 出力部
- 109 記憶部

30

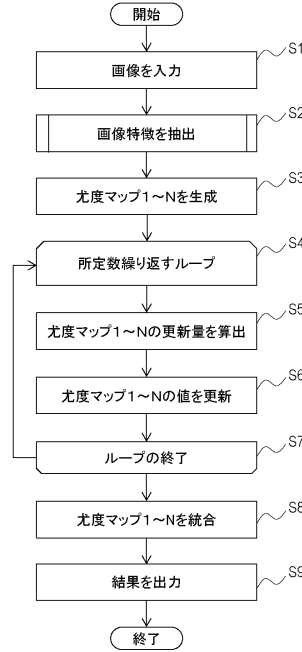
40

50

【図面】
【図 1】



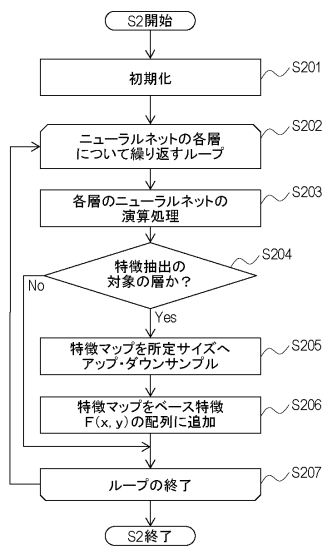
【図 2】



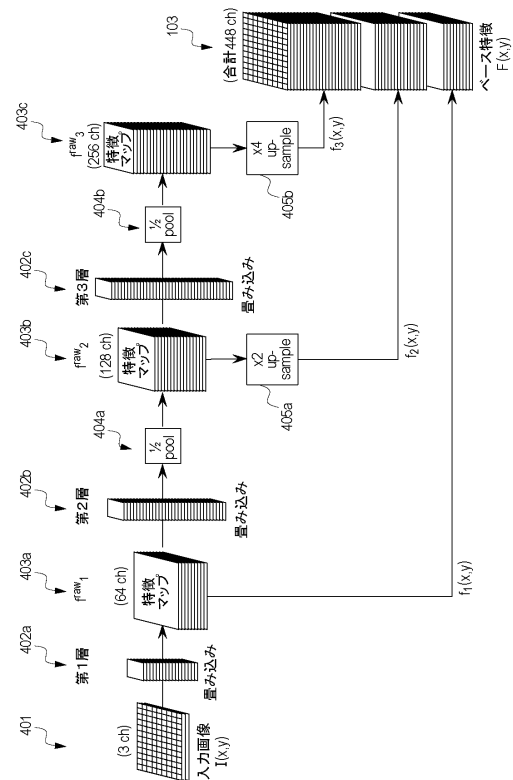
10

20

【図 3】



【図 4】

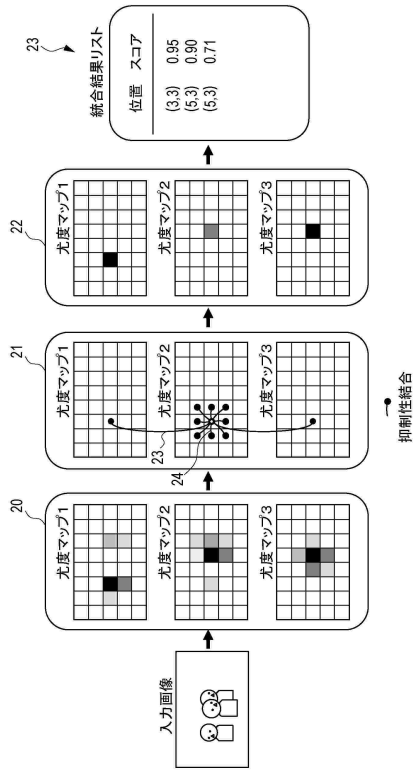


30

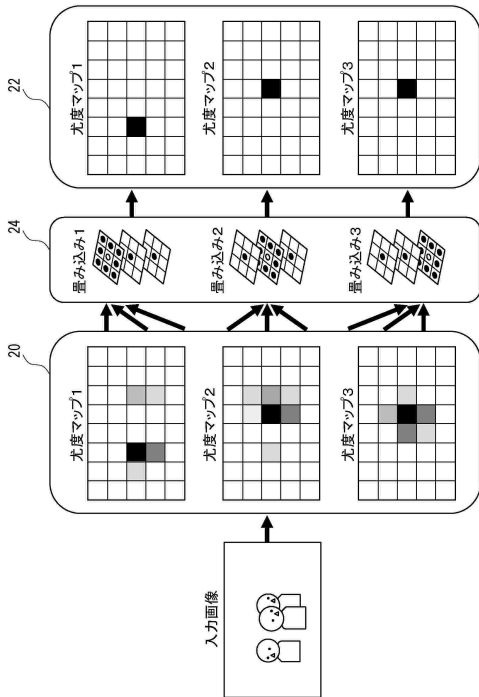
40

50

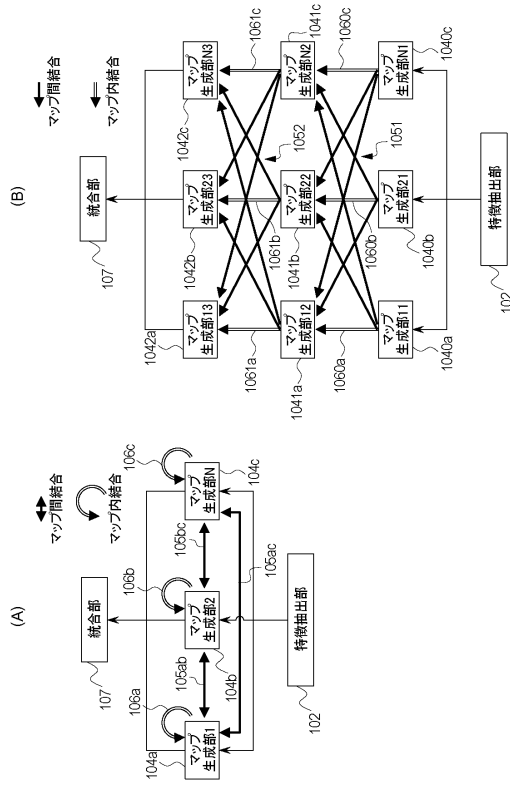
【図5】



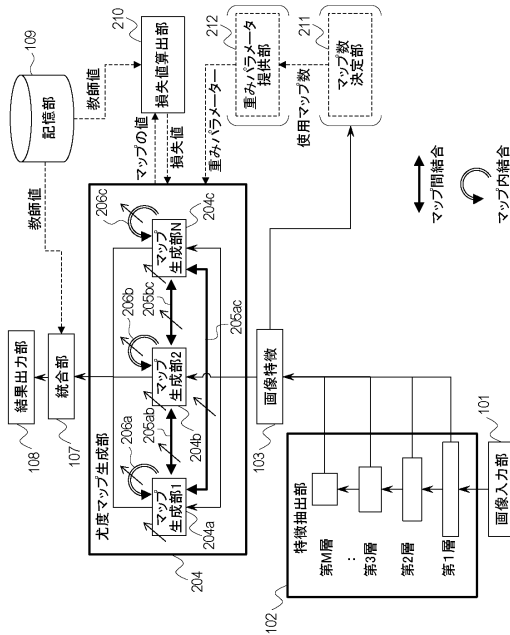
【図7】



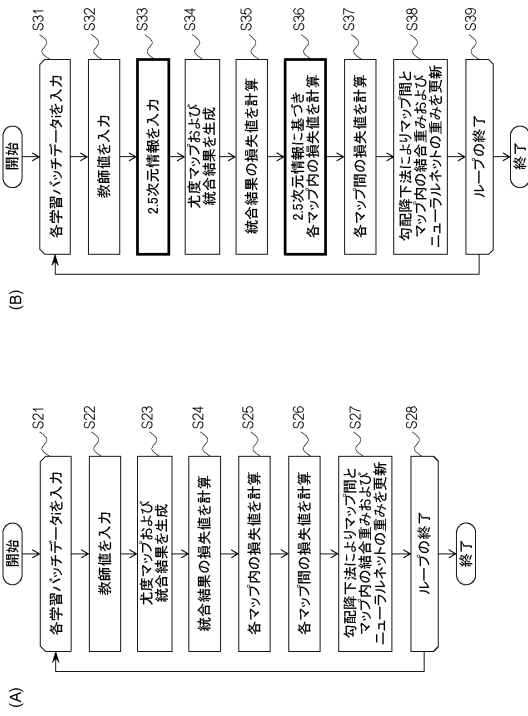
【図6】



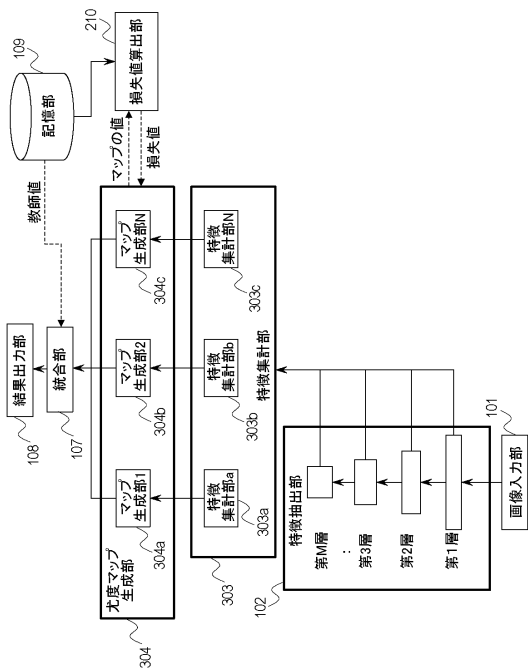
【図8】



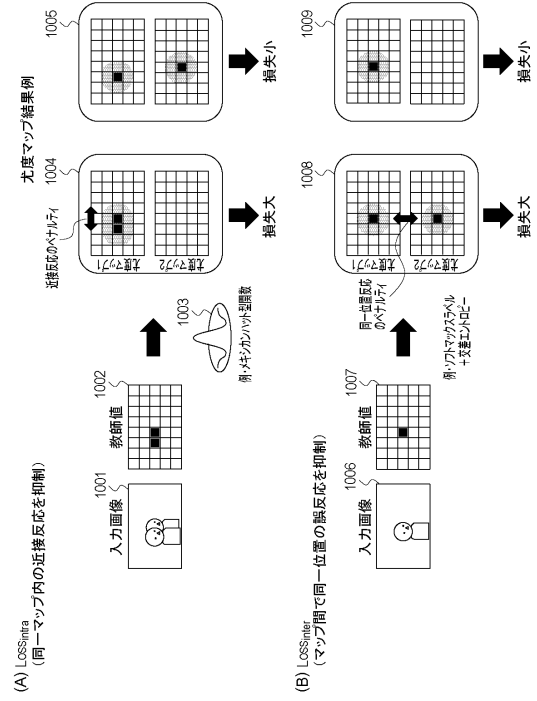
【図 9】



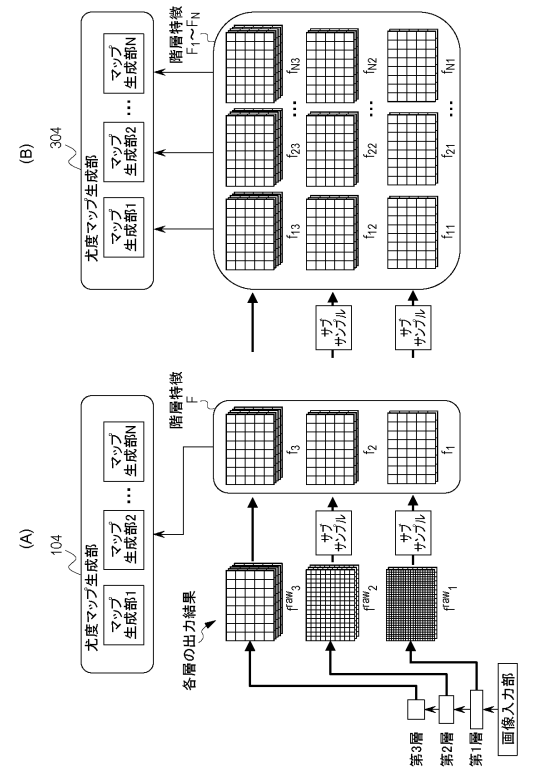
【図 11】



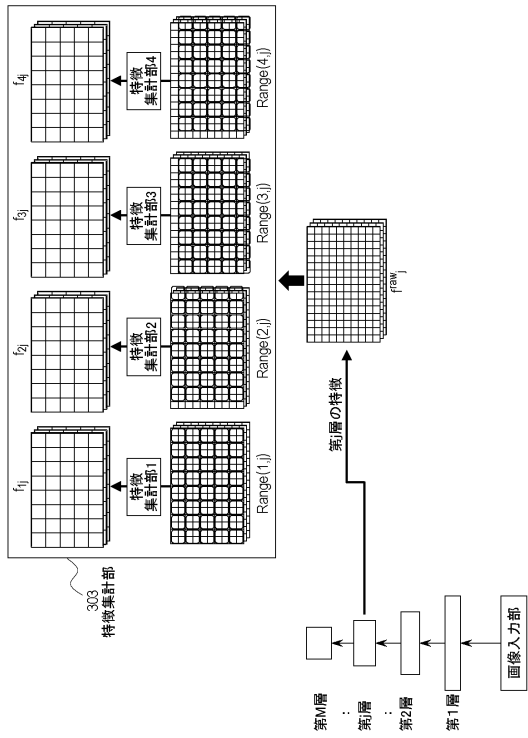
【図 10】



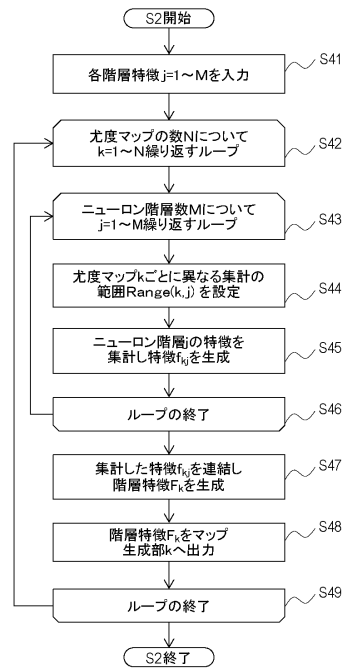
【図 12】



【図13】



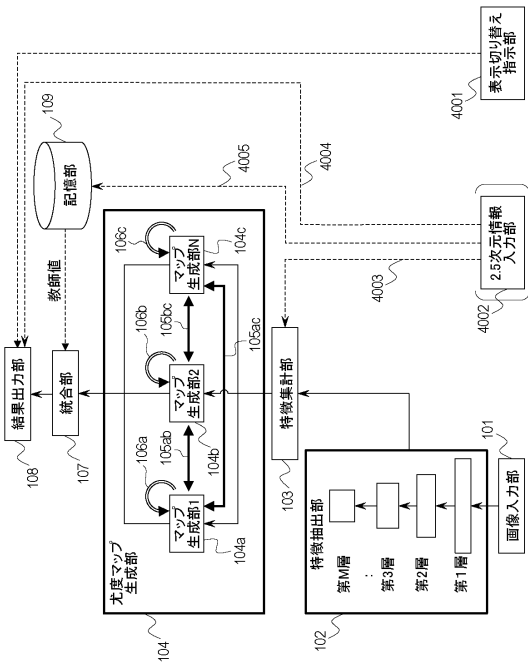
【図14】



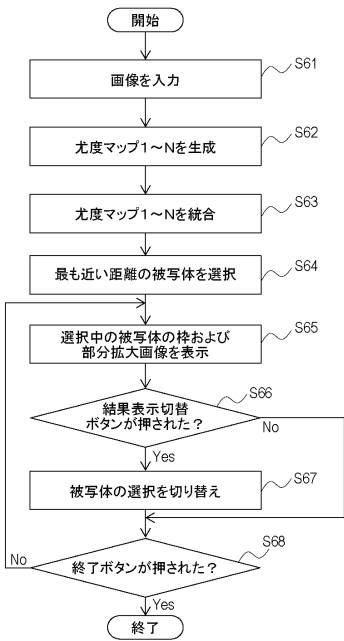
10

20

【図15】



【図16】

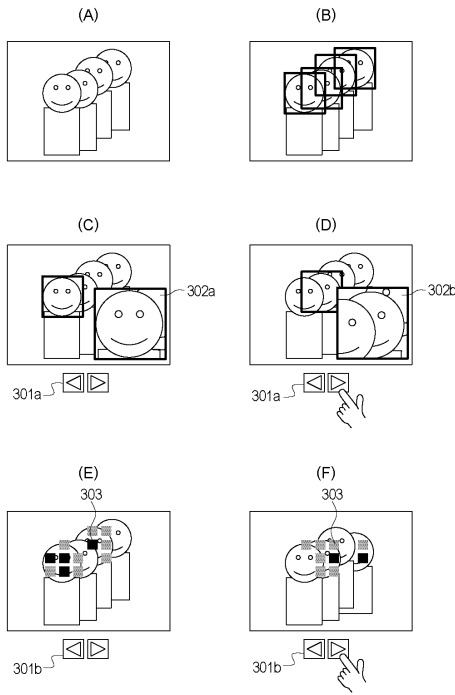


30

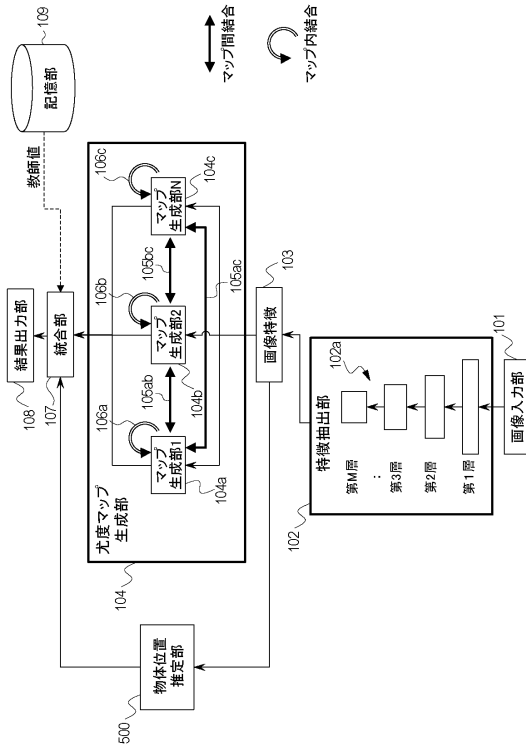
40

50

【図 17】



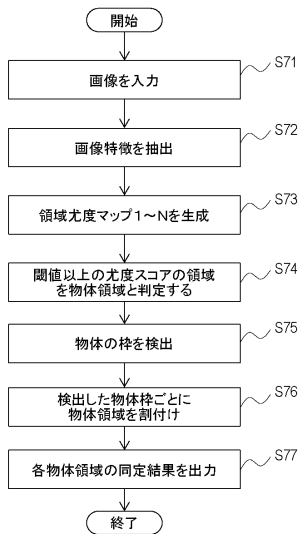
【図 18】



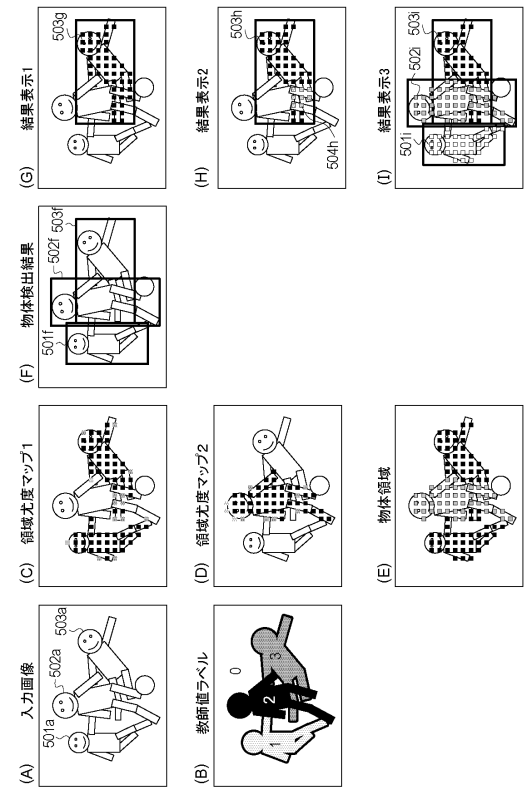
10

20

【図 19】



【図 20】

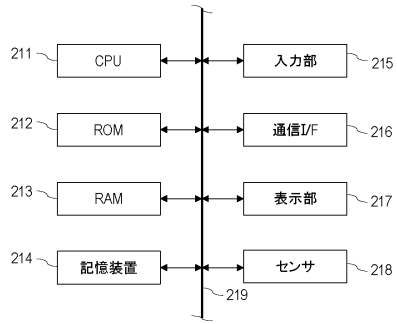


30

40

50

【図 21】



10

20

30

40

50

フロントページの続き

ヤノン株式会社内

審査官 吉川 康男

- (56)参考文献 特開 2 0 1 9 - 0 3 2 7 7 3 (J P , A)
A I による三次元モデル自動構成に基づく V R 地震シミュレーションシステム , 第 2 4 回
画像センシングシンポジウム S S I I 2 0 1 8 IS1-28 , 2018年06月13日
Joseph Redmon;Santosh Divvala;Ross Girshick;Ali Farhadi , You Only Look Once: Unified, R
eal-Time Object Detection , 2016 IEEE Conference on Computer Vision and Pattern Recogn
ition (CVPR) , IEEE , 2016年06月27日 , 779-788 , <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=7780460>
畳み込みニューラルネットワークを用いた歩行者検出の高速化 , 電子情報通信学会技術研
究報告 V o l . 1 1 7 N o . 5 0 5 N L P 2 0 1 7 - 1 0 3 , 2018年03月06日
Burak Uzkent;Aneesh Rangnekar;Matthew J. Hoffman , Aerial Vehicle Tracking by Adaptive
Fusion of Hyperspectral Likelihood Maps , 2017 IEEE Conference on Computer Vision and
Pattern Recognition Workshops (CVPRW) , IEEE , 2017年07月21日 , 233-242 , <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8014769>
- (58)調査した分野 (Int.Cl. , D B 名)
G 0 6 T 7 / 0 0