

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7683723号  
(P7683723)

(45)発行日 令和7年5月27日(2025.5.27)

(24)登録日 令和7年5月19日(2025.5.19)

(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
G 0 6 N	3/08	(2023.01)	G 0 6 N	3/08	
G 0 6 N	20/00	(2019.01)	G 0 6 N	20/00	

請求項の数 10 (全23頁)

(21)出願番号	特願2023-555935(P2023-555935)	(73)特許権者	000004237 日本電気株式会社 東京都港区芝五丁目7番1号
(86)(22)出願日	令和3年10月26日(2021.10.26)	(74)代理人	100124811 弁理士 馬場 資博
(86)国際出願番号	PCT/JP2021/039520	(74)代理人	100088959 弁理士 境 廣巳
(87)国際公開番号	WO2023/073813	(74)代理人	100097157 弁理士 桂木 雄二
(87)国際公開日	令和5年5月4日(2023.5.4)	(74)代理人	100187724 弁理士 唐鎌 睦
審査請求日	令和6年4月5日(2024.4.5)	(72)発明者	朴 君 東京都港区芝五丁目7番1号 日本電気株式会社内
		審査官	藤原 敬利

最終頁に続く

(54)【発明の名称】 画像処理システム

(57)【特許請求の範囲】

【請求項1】

画像から互いに相違する複数の推論タスクを行う学習済みモデルを生成する学習部を含み、

前記学習済みモデルは、

前記画像から前記複数の推論タスクに共通な第1の特徴量を抽出する第1のコンポーネントと、

前記推論タスクに対応して設けられ、前記第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出する第2のコンポーネントと、

前記推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成する第3のコンポーネントと、

前記推論タスクに対応して設けられ、前記第3の特徴量から対応する推論タスクの推論結果を出力する第4のコンポーネントと、

を含む画像処理システム。

【請求項2】

前記第3のコンポーネントは、前記推論タスク毎に抽出された第2の特徴量のうちの1つを基準特徴量とし、前記基準特徴量以外の前記第2の特徴量のサイズを前記基準特徴量のサイズに合わせて変更し、前記基準特徴量のサイズに合わせてサイズを変更した後の前記基準特徴量以外の前記第2の特徴量と前記基準特徴量とを結合して前記第3の特徴量を生成し、前記推論タスク毎に、前記第3の特徴量のサイズを前記第4のコンポーネントの

10

20

入力サイズに合わせて変更して前記第 4 のコンポーネントへ出力する、  
請求項 1 に記載の画像処理システム。

【請求項 3】

前記第 3 のコンポーネントは、前記推論タスクに対応するサブコンポーネントを含み、  
前記サブコンポーネントは、対応する前記推論タスクの前記第 2 の特徴量を基準特徴量とし、  
対応する前記推論タスク以外の前記第 2 の特徴量のサイズを前記基準特徴量のサイズ  
に合わせて変更し、前記基準特徴量のサイズに合わせてサイズを変更した後の前記対応する  
前記推論タスク以外の前記第 2 の特徴量と前記基準特徴量とを結合して前記第 3 の特徴  
量を生成し、前記第 4 のコンポーネントへ出力する、  
請求項 1 に記載の画像処理システム。

10

【請求項 4】

前記学習部は、複数の学習段階に分けて前記学習済みモデルの学習を行い、  
前記複数の学習段階は、少なくとも、  
前記複数の推論タスクのうち何れか 1 つを学習対象タスクとし、前記学習対象タスク  
以外の推論タスクに係る前記第 2 のコンポーネントおよび前記第 4 のコンポーネントと前  
記第 1 のコンポーネントのパラメータを固定して、前記学習対象タスクに係る前記第 2 の  
コンポーネントおよび前記第 4 のコンポーネントのパラメータを学習する第 1 の学習段階  
と、

前記第 1 のコンポーネントのパラメータを固定して、前記複数の推論タスクのそれぞれ  
に係る前記第 2 のコンポーネントおよび前記第 4 のコンポーネントのパラメータを学習する  
第 2 の学習段階と、を含む、  
請求項 1 乃至 3 の何れかに記載の画像処理システム。

20

【請求項 5】

前記推論タスクに対応して設けられた第 4 のコンポーネントは、前記第 3 の特徴量を構  
成する複数の前記第 2 の特徴量のうち、対応する前記推論タスクの第 2 の特徴量の優先度  
合いを定める重みをそれ以外の第 2 の特徴量の優先度合いを定める重みより大きくする、  
請求項 1 乃至 4 の何れかに記載の画像処理システム。

【請求項 6】

学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を出力  
する推論部を含み、

30

前記学習済みモデルは、  
前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出する第 1 のコンポーネ  
ントと、

前記推論タスクに対応して設けられ、前記第 1 の特徴量から対応する推論タスクに固有  
な第 2 の特徴量を抽出する第 2 のコンポーネントと、

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成する第 3 の  
コンポーネントと、

前記推論タスクに対応して設けられ、前記第 3 の特徴量から対応する推論タスクの推論  
結果を出力する第 4 のコンポーネントと、  
を含む画像処理システム。

40

【請求項 7】

コンピュータによる画像処理方法であって、

前記コンピュータが、画像から互いに相違する複数の推論タスクを行う学習済みモデル  
を生成し、

前記生成では、前記コンピュータが、前記学習済みモデルに、  
前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、  
前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量  
を抽出させ、

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、  
前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させ

50

る、

画像処理方法。

【請求項 8】

コンピュータによる画像処理方法であって、

前記コンピュータが、学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を推定して出力し、

前記推定では、前記コンピュータが、前記学習済みモデルに、

前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、

前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、

前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、

画像処理方法。

【請求項 9】

コンピュータに、画像から互いに相違する複数の推論タスクを行う学習済みモデルを生成する処理を行わせるためのプログラムであって、

前記生成では、前記学習済みモデルに、

前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、

前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、

前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、

プログラム。

【請求項 10】

コンピュータに、学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を推定して出力する処理を行わせるためのプログラムであって、

前記推定では、前記学習済みモデルに、

前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、

前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、

前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、

プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、画像処理システム、画像処理方法、および、記録媒体に関する。

【背景技術】

【0002】

単一の多層ニューラルネットワーク DNN (Deep Neural Network) で複数のタスクを同時に学習および推定する手法がある。この手法はマルチタスク学習と呼ばれる。マルチタスク学習は、タスク数に比例して増加する学習および推定時間を削減することができる。これにより、複数のタスクから得られる情報が必要な人物画像解析などのアプリケーションにおいてマルチタスク学習は有効な手法の一つになっている。

【0003】

マルチタスク学習の一例が特許文献 1 に記載されている。特許文献 1 に記載の技術（以下、本発明に関連する技術と記す）では、DNN は、人物の顔が写っている画像から複数

10

20

30

40

50

のタスクに共通な特徴量  $x^L$  を抽出する。次に、DNNは、特徴量  $x^L$  から顔の表情を識別するタスクに固有な特徴量を抽出して推定結果  $y^c$  を出力すると共に、それと並行して、特徴量  $x^L$  から顔領域の目や鼻の位置を推定するタスクに固有な特徴量を抽出して推定結果  $y^r$  を出力する。

【先行技術文献】

【特許文献】

【0004】

【文献】特開2018-55377号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、本発明に関連する技術では、画像から全てのタスクに共通な特徴量を抽出し、この共通な特徴量からタスクに固有な特徴量を抽出して各タスクの推定結果を推定するように構成されている。そのため、或るタスクに固有な特徴量を他のタスクの推定に利用できないという課題がある。

【0006】

本発明は、上述した課題、すなわち、複数のタスク間でタスク固有な特徴量を相互に利用できない、という課題を解決する画像処理システムを提供することにある。

【課題を解決するための手段】

【0007】

本発明の一形態に係る画像処理システムは、  
画像から互いに相違する複数の推論タスクを行う学習済みモデルを生成する学習部を含み、

前記学習済みモデルは、

前記画像から前記複数の推論タスクに共通な第1の特徴量を抽出する第1のコンポーネントと、

前記推論タスクに対応して設けられ、前記第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出する第2のコンポーネントと、

前記推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成する第3のコンポーネントと、

前記推論タスクに対応して設けられ、前記第3の特徴量から対応する推論タスクの推論結果を出力する第4のコンポーネントと、  
を含むように構成されている。

【0008】

本発明の他の形態に係る画像処理システムは、

学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を出力する推論部を含み、

前記学習済みモデルは、

前記画像から前記複数の推論タスクに共通な第1の特徴量を抽出する第1のコンポーネントと、

前記推論タスクに対応して設けられ、前記第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出する第2のコンポーネントと、

前記推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成する第3のコンポーネントと、

前記推論タスクに対応して設けられ、前記第3の特徴量から対応する推論タスクの推論結果を出力する第4のコンポーネントと、  
を含むように構成されている。

【0009】

本発明の他の形態に係る画像処理方法は、

画像から互いに相違する複数の推論タスクを行う学習済みモデルを生成し、

10

20

30

40

50

前記生成では、前記学習済みモデルに、  
 前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、  
 前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、  
 前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、  
 前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、ように構成されている。

## 【 0 0 1 0 】

本発明の他の形態に係る画像処理方法は、  
 学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を推定して出力し、

10

前記推定では、前記学習済みモデルに、  
 前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、  
 前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、  
 前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、  
 前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、ように構成されている。

## 【 0 0 1 1 】

本発明の他の形態に係るコンピュータ読み取り可能な記録媒体は、  
 コンピュータに、画像から互いに相違する複数の推論タスクを行う学習済みモデルを生成する処理を行わせるためのプログラムであって、

20

前記生成では、前記学習済みモデルに、  
 前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、  
 前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、  
 前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、  
 前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、

プログラムを記録するように構成されている。

30

## 【 0 0 1 2 】

本発明の他の形態に係るコンピュータ読み取り可能な記録媒体は、  
 コンピュータに、学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を推定して出力する処理を行わせるためのプログラムであって、

前記推定では、前記学習済みモデルに、  
 前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、  
 前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、  
 前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、  
 前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、

プログラムを記録するように構成されている。

40

## 【発明の効果】

## 【 0 0 1 3 】

本発明は、上述したような構成を有することにより、複数のタスク間でタスク固有な特徴量を相互に利用することができる。このため、複数のタスクのそれぞれにおいて、当該タスク固有の特徴量と他のタスク固有の特徴量とを考慮した学習および推定が可能となる。

## 【図面の簡単な説明】

## 【 0 0 1 4 】

【図 1】本発明の第 1 の実施形態に係る画像処理装置のブロック図である。

50

【図 2】本発明の第 1 の実施形態に係る画像処理装置における学習フェーズの動作の一例を示すフローチャートである。

【図 3】本発明の第 1 の実施形態に係る画像処理装置における推定フェーズの動作の一例を示すフローチャートである。

【図 4】本発明の第 1 の実施形態で使用するモデルの一例を示す構成図である。

【図 5】本発明の第 1 の実施形態で使用するモデルのコンポーネント C M 3 の一例を示す構成図である。

【図 6】本発明の第 1 の実施形態で使用するモデルのコンポーネント C M 3 の他の例を示す構成図である。

【図 7】本発明の第 1 の実施形態で使用するモデルの機械学習に用いられる訓練データのリストの一例を示す図である。

10

【図 8】本発明の第 1 の実施形態に係る画像処理装置における学習部の学習処理の一例を示すフローチャートである。

【図 9】本発明の第 2 の実施形態に係る画像処理装置のブロック図である。

【図 10】本発明の第 3 の実施形態に係る画像処理装置のブロック図である。

【発明を実施するための形態】

【0015】

次に、本発明の実施の形態について、図面を参照して詳細に説明する。

[第 1 の実施の形態]

図 1 は、本発明の第 1 の実施形態に係る画像処理装置 10 のブロック図である。この画像処理装置 10 は、画像から互いに相違する複数の推論タスクを行うように構成されている。図 1 を参照すると、画像処理装置 10 は、カメラ I / F (インターフェース) 部 11 と通信 I / F 部 12 と操作入力部 13 と画面表示部 14 と記憶部 15 と演算処理部 16 とを含んで構成されている。

20

【0016】

カメラ I / F 部 11 は、有線または無線により画像サーバ 17 に接続され、画像サーバ 17 と演算処理部 16 との間でデータの送受信を行うように構成されている。画像サーバ 17 は、有線または無線によりカメラ 18 に接続され、カメラ 18 で撮影されたそれぞれ撮影時刻が異なる複数の画像を過去一定期間分蓄積するように構成されている。カメラ 18 は、例えば、数百万画素程度の画素容量を有する CCD (Charge - Coupled Device) イメージセンサや CMOS (Complementary MOS) イメージセンサを備えたカラーカメラや白黒カメラであってよい。カメラ 18 は、防犯・監視の目的のために多くの人が行きかう街頭、屋内などに設置されたカメラであってよい。或いはカメラ 18 は、車などの移動体に搭載されて移動しながら同一或いは異なる撮影領域を撮影するカメラであってよい。カメラ 18 は 1 台に限定されず、異なる場所から異なる撮影領域を撮影する複数台のカメラであってよい。

30

【0017】

通信 I / F 部 12 は、データ通信回路から構成され、有線または無線によって図示しない外部装置との間でデータ通信を行うように構成されている。操作入力部 13 は、キーボードやマウスなどの操作入力装置から構成され、オペレータの操作を検出して演算処理部 16 に出力するように構成されている。画面表示部 14 は、LCD (Liquid Crystal Display) などの画面表示装置から構成され、演算処理部 16 からの指示に応じて、各種情報を画面表示するように構成されている。

40

【0018】

記憶部 15 は、ハードディスクやメモリなどの記憶装置から構成され、演算処理部 16 における各種処理に必要な処理情報およびプログラム 151 を記憶するように構成されている。プログラム 151 は、演算処理部 16 に読み込まれて実行されることにより各種処理部を実現するプログラムであり、通信 I / F 部 12 などのデータ入出力機能を介して図示しない外部装置や記録媒体から予め読み込まれて記憶部 15 に保存される。記憶部 15 に記憶される主な処理情報には、画像情報 152、モデル 153、および推定結果情報 1

50

54がある。

【0019】

画像情報152は、カメラI/F部11を通じて画像サーバ17から取得されたカメラ18のフレーム画像である。

【0020】

モデル153は、カメラ18のフレーム画像から互いに相違する複数の推論タスクを同時に学習および推定する機械学習モデルである。モデル153は、例えば、DCNN(Deep Convolutional Neural Network)を用いて構成してよい。本実施形態では、モデル153は、物体検出、ポーズ推定、セマンティックセグメンテーション推定の3つの推論タスクを行うようにパラメータが学習される。パラメータが学習されたモデルを学習済みのモデルと呼び、学習前のモデルと区別する。

10

【0021】

物体検出は、画像内のクラスと物体位置を検出する。物体検出の結果は、クラス名、クラスの推定信頼度、物体位置を表すバウンディングボックス(以下、矩形と記す)を含む。検出するクラスは、例えば人物であってよい。但し、検出するクラスは人物に限定されず、動物や物であってよい。

【0022】

ポーズ推定は、画像内の人物の骨格情報を推定する。人物の骨格情報は、人体を構成する関節の位置を表す情報を含んでいる。関節は、首や肩などの関節のみならず、目や鼻などの顔のパーツも関節に含めてよい。ポーズ推定の結果は、関節名(関節ID)、関節の位置、関節の信頼度を含む。

20

【0023】

セマンティックセグメンテーション推定は、画像内の各ピクセルのクラスを推定する。セマンティックセグメンテーション推定の結果は、各ピクセルのクラスを含む。推定するクラスは、物体検出で検出するクラスと同じである。

【0024】

推定結果情報154は、学習済みのモデル153を用いて画像から推定した結果を表す情報である。推定結果情報154は、物体検出結果、ポーズ推定結果、および、セマンティックセグメンテーション推定結果を含む。

【0025】

演算処理部16は、MPUなどの1以上のプロセッサとその周辺回路を有し、記憶部15からプログラム151を読み込んで実行することにより、上記ハードウェアとプログラム151とを協働させて各種処理部を実現するように構成されている。演算処理部16で実現される主な処理部には、取得部161、学習部162、および、推定部163がある。

30

【0026】

取得部161は、カメラI/F部11を通じて画像サーバ17から、カメラ18で撮影された動画を構成するフレーム画像あるいはそれをダウンサンプリングしたフレーム画像を取得し、記憶部15に画像情報152として保存するように構成されている。取得されたフレーム画像には、カメラIDおよび撮影時刻が付加されている。フレーム画像の撮影時刻は、フレーム毎に相違する。

40

【0027】

学習部162は、訓練データを用いて、モデル153に上記3つの推論タスクを同時に学習させるように構成されている。即ち、学習部162は、画像から上記3つの推論タスクを行う学習済みモデル153を生成する。学習部21は、上記生成では、モデル153に、画像から上記3つの推論タスクに共通な第1の特徴量を抽出させ、次に、推論タスク毎に、第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出させ、次に、推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成させ、次に、推論タスク毎に、第3の特徴量から対応する推論タスクの推論結果を出力させる。

【0028】

推定部163は、学習済みモデル153を用いて、画像から上記3つの推論タスクの推

50

論結果を推定して出力するように構成されている。推定部 31 は、上記推定では、学習済みモデル 153 に、先ず、画像から上記 3 つの推論タスクに共通な第 1 の特徴量を抽出させ、次に、推論タスク毎に、第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、次に、推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、次に、推論タスク毎に、第 3 の特徴量から対応する推論タスクの推論結果を出力させる。

#### 【0029】

次に、画像処理装置 10 の動作を説明する。画像処理装置 10 のフェーズは、学習フェーズと推定フェーズとに大別される。学習フェーズは、モデル 153 を機械学習するフェーズである。推定フェーズは、学習済みのモデル 153 を用いて、画像から上記 3 つの推論タスクの推論結果を推定し、出力するフェーズである。

10

#### 【0030】

図 2 は学習フェーズの動作の一例を示すフローチャートである。図 2 を参照すると、先ず、取得部 161 は、カメラ I / F 部 11 を通じて画像サーバ 17 からカメラ 18 で撮影されたフレーム画像を取得し、記憶部 15 に画像情報 152 として保存する（ステップ S1）。次に、学習部 162 は、モデル 153 の機械学習に用いる訓練データを作成する（ステップ S2）。次に、学習部 162 は、訓練データを用い、入力を画像とし、出力を上記 3 つの推論タスクの推定結果とするモデル 153 を機械学習させ、学習済みのモデル 153 を生成する（ステップ S3）。

#### 【0031】

図 3 は推定フェーズの動作の一例を示すフローチャートである。図 3 を参照すると、先ず、取得部 161 は、カメラ I / F 部 110 を通じて画像サーバ 17 からカメラ 18 で撮影されたフレーム画像を取得し、記憶部 15 に画像情報 152 として保存する（ステップ S11）。

20

#### 【0032】

次に、推定部 163 は、学習済みのモデル 153 を用いて、画像情報 152 に含まれるフレーム画像から上記 3 つの推論タスクの推定結果を同時に推定する（ステップ S12）。次に、推定部 163 は、推定した 3 つの推論タスクの推定結果を画面表示部 14 に表示し、または / および、通信 I / F 部 12 を通じて外部装置へ送信する（ステップ S13）。

#### 【0033】

続いて、モデル 153 と学習部 162 を詳細に説明する。

30

#### 【0034】

先ず、モデル 153 の詳細を説明する。

#### 【0035】

図 4 は、モデル 153 として利用できるマルチタスクモデルの一例を示す構成図である。この例のモデル 153 は、8 個のコンポーネント CM から構成され、全体が 1 つの多層ニューラルネットワークになっている。

#### 【0036】

コンポーネント CM1 は、多層ニューラルネットワークの低層側に設けられ、画像を入力し、全てのタスクに共通な低次元特徴量 FM1 を抽出するように構成されている。コンポーネント CM1 は、バックボーンとも呼ばれる。コンポーネント CM1 によって抽出された特徴量 FM1 は、低次元特徴マップとも称する。コンポーネント CM1 は、1 以上の畳み込み層を含んで構成してよい。例えば、コンポーネント CM1 は、SSD (Single Shot MultiBox Detector) の構成要素である VGG-16 を使用してよい。或いは、コンポーネント CM1 は、例えば、OpenPose の構成要素である VGG-19 を使用してよい。或いは、コンポーネント CM1 は、例えば、SegNet の構成要素であるエンコーダを使用してよい。或いは、コンポーネント CM1 は、例えば、SSD や OpenPose や SegNet 以外のモデルのバックボーンを使用し

40

#### 【0037】

50

コンポーネントCM2-1は、コンポーネントCM1から特徴量FM1を入力し、物体検出タスクに固有な高次の特徴量FM2-1を抽出するように構成されている。コンポーネントCM2-1は、1以上の畳み込み層を含んで構成してよい。例えば、コンポーネントCM2-1は、SSDの構成要素である特別な畳み込み層(Extra Feature Layers)を使用してよい。但し、コンポーネントCM2-1は上記に限定されず、SSD以外の物体検出モデルにおいて物体検出タスクに固有な高次の特徴量を抽出する畳み込み層を使用してよい。

【0038】

コンポーネントCM2-2は、コンポーネントCM1から特徴量FM1を入力し、ポーズ推定タスクに固有な高次の特徴量FM2-2を抽出するように構成されている。コンポーネントCM2-2は、1以上の畳み込み層を含んで構成してよい。例えば、コンポーネントCM2-2は、OpenPoseの構成要素である、キーポイントの位置を表すPart Confidence Mapを生成する畳み込み層と、キーポイント間の関連度を表すPart Affinity Fieldsを生成する畳み込み層と、生成されたPart Confidence MapとPart Affinity Fieldsと抽出元の特徴量FM1とを結合(concatenate)する層(これによって結合して得られた特徴マップを以下、OpenPose特徴マップと記す)とを使用してよい。但し、コンポーネントCM2-2は上記に限定されず、OpenPose以外のポーズ推定モデルにおいてポーズ推定タスクに固有な高次の特徴量を抽出する畳み込み層を使用してよい。

【0039】

コンポーネントCM2-3は、コンポーネントCM1から特徴量FM1を入力し、セマンティックセグメンテーション推定タスクに固有な高次の特徴量FM2-3を抽出するように構成されている。コンポーネントCM2-3は、1以上の畳み込み層を含んで構成してよい。例えば、コンポーネントCM2-3は、SegNetの構成要素であるデコーダを使用してよい。但し、コンポーネントCM2-3は上記に限定されず、SegNet以外のセマンティックセグメンテーション推定モデルにおいてセマンティックセグメンテーション推定タスクに固有な高次の特徴量を抽出する畳み込み層を使用してよい。

【0040】

コンポーネントCM3は、コンポーネントCM2-1、CM2-2、CM2-3から特徴量FM2-1、FM2-2、FM2-3を入力し、これら3つの特徴量FM2-1、FM2-2、FM2-3を結合(concatenate)して得られる特徴量FM3-1、FM3-2、FM3-3を生成するように構成されている。

【0041】

図5は、コンポーネントCM3の一例を示す構成図である。この例のコンポーネントCM3は、リサイズ部CM3-1と結合部CM3-2とリサイズ部CM3-3とを含んで構成されている。

【0042】

リサイズ部CM3-1は、特徴量FM2-1、FM2-2、FM2-3を結合し得るように、それらのサイズを合わせるように構成されている。リサイズ部CM3-1は、3つの特徴量のうちの何れか1つの特徴量を基準特徴量とし、基準特徴量のサイズに合わせて、残り2つの特徴量のサイズを変更する。例えば、特徴量FM2-1、FM2-2、FM2-3のサイズを、それぞれ、 $38 \times 38$ 、 $70 \times 70$ 、 $240 \times 320$ とし、基準特徴量を特徴量FM2-1とする。この場合、リサイズ部CM3-1は、特徴量FM2-2のサイズを $70 \times 70$ から $38 \times 38$ に変更した特徴量FM2-2'を生成して出力する。また、リサイズ部CM3-1は、特徴量FM2-3のサイズを $240 \times 320$ から $38 \times 38$ に変更した特徴量FM2-3'を生成して出力する。また、リサイズ部CM3-1は、特徴量FM2-1についてはサイズを変更せず、特徴量FM2-1そのものを特徴量FM2-1'として出力する。

【0043】

結合部CM3-2は、リサイズ部CM3-1から特徴量FM2-1'、FM2-2'、F

10

20

30

40

50

M 2 - 3 ' を入力し、これらを結合して得られる特徴量 F M 3 を生成して出力する。例えば、結合部 C M 3 - 2 は、それぞれが  $38 \times 38$  のサイズである特徴量 F M 2 - 1 '、F M 2 - 2 '、F M 2 - 3 ' を入力し、 $38 \times 38 \times 3$  のサイズの特徴量 F M 3 を生成して出力する。このように、特徴量の結合により、チャンネル数（次元数）が増加する。

【 0 0 4 4 】

リサイズ部 C M 3 - 3 は、結合部 C M 3 - 2 から特徴量 F M 3 を入力し、各タスクに応じたサイズに変更した特徴量 F M 3 - 1、F M 3 - 2、F M 3 - 3 を生成して出力する。例えば、コンポーネント C M 4 - 1、C M 4 - 2、C M 4 - 3 の入力サイズをそれぞれ、 $38 \times 38 \times 3$ 、 $70 \times 70 \times 3$ 、 $240 \times 320 \times 3$  とする。この場合、リサイズ部 C M 3 - 3 は、特徴量 F M 3 のサイズを  $38 \times 38 \times 3$  から  $70 \times 70 \times 3$  に変更した特徴量 F M 3 - 2 を生成して、コンポーネント C M 4 - 2 に出力する。また、リサイズ部 C M 3 - 3 は、特徴量 F M 3 のサイズを  $38 \times 38 \times 3$  から  $240 \times 320 \times 3$  に変更した特徴量 F M 3 - 3 を生成して、コンポーネント C M 4 - 3 に出力する。また、リサイズ部 C M 3 - 3 は、 $38 \times 38 \times 3$  のサイズの特徴量 F M 3 そのものを特徴量 F M 3 - 1 として、コンポーネント C M 4 - 1 に出力する。

10

【 0 0 4 5 】

図 6 は、コンポーネント C M 3 の他の一例を示す構成図である。この例のコンポーネント C M 3 は、3 つのサブコンポーネント C M 3 A、C M 3 B、C M 3 C を含んで構成されている。

【 0 0 4 6 】

サブコンポーネント C M 3 A は、特徴量 F M 2 - 1、F M 2 - 2、F M 2 - 3 から物体検出タスクのコンポーネント C M 4 - 1 のための特徴量 F M 3 - 1 を生成して出力するように構成されている。サブコンポーネント C M 3 A は、特徴量 F M 2 - 1 のサイズに合わせて特徴量 F M 2 - 2、F M 2 - 3 のサイズを変更して得られる特徴量 F M 2 - 2 '、F M 2 - 3 ' を生成して出力するリサイズ部 C M 3 A - 1 と、特徴量 F M 2 - 1、F M 2 - 2 '、F M 2 - 3 ' の 3 つを結合して得られる特徴量 F M 3 - 1 を生成して出力する結合部 C M 3 A - 2 とを含んで構成されている。例えば、特徴量 F M 2 - 1、F M 2 - 2、F M 2 - 3 のサイズを、それぞれ、 $38 \times 38$ 、 $70 \times 70$ 、 $240 \times 320$  とし、コンポーネント C M 4 - 1 の入力サイズを  $38 \times 38 \times 3$  とする。この場合、リサイズ部 C M 3 A - 1 は、特徴量 F M 2 - 2 のサイズを  $70 \times 70$  から  $38 \times 38$  に変更した特徴量 F M 2 - 2 ' を生成して出力し、特徴量 F M 2 - 3 のサイズを  $240 \times 320$  から  $38 \times 38$  に変更した特徴量 F M 2 - 3 ' を生成して出力する。結合部 C M 3 A - 2 は、同じ  $38 \times 38$  のサイズの特徴量 F M 2 - 1、F M 2 - 2 '、F M 2 - 3 ' を結合して、 $38 \times 38 \times 3$  のサイズの特徴量 F M 3 - 1 を生成して出力する。これにより、結合に伴う特徴量 F M 2 - 1 の劣化を抑えることができる。

20

30

【 0 0 4 7 】

サブコンポーネント C M 3 B は、特徴量 F M 2 - 1、F M 2 - 2、F M 2 - 3 からポーズ推定タスクのコンポーネント C M 4 - 2 のための特徴量 F M 3 - 2 を生成して出力するように構成されている。サブコンポーネント C M 3 B は、特徴量 F M 2 - 2 のサイズに合わせて特徴量 F M 2 - 1、F M 2 - 3 のサイズを変更して得られる特徴量 F M 2 - 1 '、F M 2 - 3 ' を生成して出力するリサイズ部 C M 3 B - 1 と、特徴量 F M 2 - 1 '、F M 2 - 2、F M 2 - 3 ' の 3 つを結合して得られる特徴量 F M 3 - 2 を生成して出力する結合部 C M 3 B - 2 とを含んで構成されている。例えば、特徴量 F M 2 - 1、F M 2 - 2、F M 2 - 3 のサイズを、それぞれ、 $38 \times 38$ 、 $70 \times 70$ 、 $240 \times 320$  とし、コンポーネント C M 4 - 2 の入力サイズを  $70 \times 70 \times 3$  とする。この場合、リサイズ部 C M 3 B - 1 は、特徴量 F M 2 - 1 のサイズを  $38 \times 38$  から  $70 \times 70$  に変更した特徴量 F M 2 - 1 ' を生成して出力し、特徴量 F M 2 - 3 のサイズを  $240 \times 320$  から  $70 \times 70$  に変更した特徴量 F M 2 - 3 ' を生成して出力する。結合部 C M 3 B - 2 は、同じ  $70 \times 70$  のサイズの特徴量 F M 2 - 1 '、F M 2 - 2、F M 2 - 3 ' を結合して、 $70 \times 70 \times 3$  のサイズの特徴量 F M 3 - 2 を生成して出力する。これにより、結合に伴う特徴量 F M 2 - 2 の劣

40

50

化を抑えることができる。

【0048】

サブコンポーネントCM3Cは、特徴量FM2-1、FM2-2、FM2-3からセマンティックセグメンテーション推定タスクのコンポーネントCM4-3のための特徴量FM3-3を生成して出力するように構成されている。サブコンポーネントCM3Cは、特徴量FM2-3のサイズに合わせて特徴量FM2-1、FM2-2のサイズを変更して得られる特徴量FM2-1'、FM2-2'を生成して出力するリサイズ部CM3C-1と、特徴量FM2-1'、FM2-2'、FM2-3の3つを結合して得られる特徴量FM3-3を生成して出力する結合部CM3C-2とを含んで構成されている。例えば、特徴量FM2-1、FM2-2、FM2-3のサイズを、それぞれ、 $38 \times 38$ 、 $70 \times 70$ 、 $240 \times 320$ とし、コンポーネントCM4-3の入力サイズを $240 \times 320 \times 3$ とする。この場合、リサイズ部CM3C-1は、特徴量FM2-1のサイズを $38 \times 38$ から $240 \times 240$ に変更した特徴量FM2-1'を生成して出力し、特徴量FM2-2のサイズを $70 \times 70$ から $240 \times 320$ に変更した特徴量FM2-2'を生成して出力する。結合部CM3C-2は、同じ $240 \times 320$ のサイズの特徴量FM2-1'、FM2-2'、FM2-3を結合して、 $240 \times 320 \times 3$ のサイズの特徴量FM3-3を生成して出力する。これにより、結合に伴う特徴量FM2-3の劣化を抑えることができる。

10

【0049】

再び図4を参照すると、コンポーネントCM4-1は、コンポーネントCM3から特徴量FM3-1を入力し、特徴量FM3-1から物体検出タスクの推定結果ER1を推定して出力するように構成されている。特徴量FM3-1は、物体検出タスクに固有な高次の特徴量FM2-1だけでなく、ポーズ推定タスクに固有な高次の特徴量FM2-2とセマンティックセグメンテーション推定に固有な高次の特徴量FM2-3とを含んでいる。そのため、コンポーネントCM4-1は、それら3つの高次の特徴量を考慮した学習および推定が可能になる。コンポーネントCM4-1は、例えば、SSDを構成する特別な畳み込み層につながる出力層(Detections: 8732 per Class, Non-Maximum Suppression)を使用してよい。

20

【0050】

ここで、コンポーネントCM4-1は、物体検出タスクに固有な高次の特徴量FM2-1の優先度合いを定める重みをそれ以外の第2の特徴量の優先度合いを定める重みより大きく設定してよい。例えば、コンポーネントCM4-1は、物体検出タスクに固有な高次の特徴量FM2-1の優先度合いを定める重みを0.5、それ以外の第2の特徴量の優先度合いを定める重みを0.25としてよい。このように特徴量FM2-1に相対的に大きな重みを与えることにより、3つの高次の特徴量を考慮した学習および推定を可能にしつつ、物体検出タスクに固有な高次の特徴量FM2-1の重要度を上げることができる。

30

【0051】

また、コンポーネントCM4-1は、入力された特徴量FM3-1に対して $1 \times 1$ の畳み込み(Channel-Wise Convolution)を行うことにより、高次の特徴量の次元数を、例えば、 $38 \times 38 \times 3$ から $38 \times 38 \times 1$ に削減してよい。これにより、コンポーネントCM4-1として、SSDなどの既存モデルにおける、高次の特徴量から推定結果を推定して出力するネットワーク部分をそのまま利用することができる。

40

【0052】

コンポーネントCM4-2は、コンポーネントCM3から特徴量FM3-2を入力し、特徴量FM3-2からポーズ推定タスクの推定結果ER2を推定して出力するように構成されている。特徴量FM3-2は、ポーズ推定タスクに固有な高次の特徴量FM2-2だけでなく、物体検出タスクに固有な高次の特徴量FM2-1とセマンティックセグメンテーション推定に固有な高次の特徴量FM2-3とが含まれている。そのため、コンポーネントCM4-2は、それら3つの高次の特徴量を考慮した学習および推定が可能になる。コンポーネントCM4-2は、例えば、OpenPoseの構成要素である、OpenPose特徴マップからポーズ推定結果を推定するネットワーク部分を用いてよい。

50

## 【 0 0 5 3 】

ここで、コンポーネントCM4 - 2は、ポーズ推定タスクに固有な高次の特徴量FM2 - 2の優先度合いを定める重みをそれ以外の第2の特徴量の優先度合いを定める重みより大きく設定してよい。例えば、コンポーネントCM4 - 2は、ポーズ推定タスクに固有な高次の特徴量FM2 - 2の優先度合いを定める重みを0.5、それ以外の第2の特徴量の優先度合いを定める重みを0.25としてよい。このように特徴量FM2 - 2に相対的に大きな重みを与えることにより、3つの高次の特徴量を考慮した学習および推定を可能にしつつ、ポーズ推定タスクに固有な高次の特徴量FM2 - 2の重要度を上げることができる。

## 【 0 0 5 4 】

また、コンポーネントCM4 - 2は、入力された特徴量FM3 - 2に対して1×1の畳み込み(Channel-Wise Convolution)を行うことにより、高次の特徴量の次元数を、例えば、70×70×3から70×70×1に削減してよい。これにより、コンポーネントCM4 - 2として、OpenPoseなどの既存モデルにおける、高次の特徴量から推定結果を推定して出力するネットワーク部分をそのまま利用することができる。

## 【 0 0 5 5 】

コンポーネントCM4 - 3は、コンポーネントCM3から特徴量FM3 - 3を入力し、特徴量FM3 - 3からセマンティックセグメンテーション推定タスクの推定結果ER3を推定して出力するように構成されている。特徴量FM3 - 3は、セマンティックセグメンテーション推定タスクに固有な高次の特徴量FM2 - 3だけでなく、物体検出タスクに固有な高次の特徴量FM2 - 1とポーズ推定に固有な高次の特徴量FM2 - 2とが含まれている。そのため、コンポーネントCM4 - 3は、それら3つの高次の特徴量を考慮した学習および推定が可能になる。コンポーネントCM4 - 3は、例えば、SegNetの構成要素であるソフトマックス層などを用いてよい。

## 【 0 0 5 6 】

ここで、コンポーネントCM4 - 3は、セマンティックセグメンテーション推定タスクに固有な高次の特徴量FM2 - 3の優先度合いを定める重みをそれ以外の第2の特徴量の優先度合いを定める重みより大きく設定してよい。例えば、コンポーネントCM4 - 3は、セマンティックセグメンテーション推定タスクに固有な高次の特徴量FM2 - 3の優先度合いを定める重みを0.5、それ以外の第2の特徴量の優先度合いを定める重みを0.25としてよい。このように特徴量FM2 - 3に相対的に大きな重みを与えることにより、3つの高次の特徴量を考慮した学習および推定を可能にしつつ、セマンティックセグメンテーション推定タスクに固有な高次の特徴量FM2 - 3の重要度を上げることができる。

## 【 0 0 5 7 】

また、コンポーネントCM4 - 3は、入力された特徴量FM3 - 3に対して1×1の畳み込み(Channel-Wise Convolution)を行うことにより、高次の特徴量の次元数を、例えば、240×320×3から240×320×1に削減してよい。これにより、コンポーネントCM4 - 3として、SegNetなどの既存モデルにおける、高次の特徴量から推定結果を推定して出力するネットワーク部分をそのまま利用することができる。

## 【 0 0 5 8 】

続いて、学習部162の詳細を説明する。

## 【 0 0 5 9 】

まず、モデル153の機械学習に用いられる訓練データについて説明する。

## 【 0 0 6 0 】

図7は、モデル153の機械学習に用いられる訓練データのリストの一例を示す。図7を参照すると、このリストには、合計n個の訓練データが登録されている。個々の訓練データは、訓練データを一意に識別するID、画像、物体検出ラベル、ポーズ推定ラベル、セマンティックセグメンテーション推定ラベルの各項目から構成されている。

10

20

30

40

50

## 【 0 0 6 1 】

画像の項目には、カメラ 1 8 で撮影されたフレーム画像が設定される。物体検出ラベルの項目には、ラベルの有無と、ラベルが有る場合には、ラベル情報である画像中に存在する人物などのクラスとその位置情報（矩形情報）が設定される。ポーズ推定ラベルの項目には、ラベルの有無と、ラベルが有る場合には、画像中に存在する関節の関節名（関節 ID）とその位置情報が設定される。セマンティックセグメンテーション推定ラベルの項目には、ラベルの有無と、ラベルが有る場合には、画像の各ピクセルのクラスが設定される。このように、訓練データ群の中には、3つのラベル（物体検出ラベル、ポーズ推定ラベル、セマンティックセグメンテーション推定ラベル）の項目の全てにラベル情報が設定されたもの以外に、一部のラベルの項目だけにラベル情報が設定されたものが含まれていてよい。

10

## 【 0 0 6 2 】

上述したような訓練データは、例えば、ユーザとの間の対話的処理によって作成されてよい。例えば、学習部 1 6 2 は、取得部 1 6 1 によって取得されたカメラ 1 8 の画像を画面表示部 1 4 に表示し、操作入力部 1 3 を通じてユーザから当該画像のラベル情報を受け付ける。そして、学習部 1 6 2 は、表示した画像と受け付けたラベル情報との組を1つの訓練データとして作成する。学習部 1 6 2 は、同様の方法により、必要十分な数の訓練データを作成する。但し、訓練データの作成方法は上記に限定されない。

## 【 0 0 6 3 】

次に、訓練データを用いて、学習部 1 6 2 がモデル 1 5 3 を学習する方法について説明する。

20

## 【 0 0 6 4 】

図 8 は、学習部 1 6 2 の学習処理の一例を示すフローチャートである。この例の学習処理は、図 4 に示される構成のモデル 1 5 3 を学習対象モデルとする。また、この例の学習処理は、モデル 1 5 3 全体を一気に学習するのではなく、学習するネットワーク部分を徐々に拡大しながら学習を行う。これにより、安定した学習が行える。具体的には、以下の4つの学習段階を経る。

## 【 0 0 6 5 】

## ( 1 ) 学習段階 1

学習段階 1 では、学習部 1 6 2 は、物体検出に係る深い層のネットワーク部分であるコンポーネント CM 2 - 1 および CM 4 - 1 だけを学習する。このとき、バックボーンであるコンポーネント CM - 1、ポーズ推定に係る深い層のネットワーク部分であるコンポーネント CM 2 - 2 および CM 4 - 2、並びに、セマンティックセグメンテーション推定に係る深い層のネットワーク部分であるコンポーネント CM 2 - 3 および CM 4 - 3 のパラメータは固定しておく。

30

## ( 2 ) 学習段階 2

学習段階 2 では、学習部 1 6 2 は、物体検出およびポーズ推定に係る深い層のネットワーク部分であるコンポーネント CM 2 - 1、CM 2 - 2、CM 4 - 1、および、CM 4 - 2 だけを学習する。このとき、バックボーンであるコンポーネント CM - 1、セマンティックセグメンテーション推定に係る深い層のネットワーク部分であるコンポーネント CM 2 - 3 および CM 4 - 3 のパラメータは固定しておく。

40

## ( 3 ) 学習段階 3

学習段階 3 では、学習部 1 6 2 は、全ての推論タスク、すなわち物体検出、ポーズ推定およびセマンティックセグメンテーション推定に係る深い層のネットワーク部分であるコンポーネント CM 2 - 1、CM 2 - 2、CM 2 - 3、CM 4 - 1、CM 4 - 2、および CM 4 - 3 だけを学習する。このとき、バックボーンであるコンポーネント CM - 1 のパラメータは固定しておく。

## ( 4 ) 学習段階 4

学習段階 4 では、学習部 1 6 2 は、モデル全体、すなわちバックボーンであるコンポーネント CM - 1、物体検出、ポーズ推定およびセマンティックセグメンテーション推定に

50

係る深い層のネットワーク部分であるコンポーネントCM2-1、CM2-2、CM2-3、CM4-1、CM4-2、およびCM4-3を学習する。

【0066】

図8を参照すると、学習部162は、モデル153の機械学習に用いられる訓練データ群から、各学習段階で使用される訓練データ群を作成する(ステップS21)。

【0067】

例えば、学習部162は、ステップS21において、図7で説明したような訓練データのリストから、学習段階3で使用される訓練データ群および学習段階4で使用される訓練データ群を、それぞれ必要な数だけ作成する。学習段階3および学習段階4では、3つのラベル(物体検出ラベル、ポーズ推定ラベル、セマンティックセグメンテーション推定ラベル)の項目の全てにラベル情報が設定された訓練データが必要である。そのため、学習部162は、そのような条件を満たす訓練データをリストから抽出することにより、学習段階3で使用される訓練データ群および学習段階4で使用される訓練データ群を作成する。

10

【0068】

また、学習部162は、ステップS21において、リストの残りの訓練データ群から、学習段階2で使用される訓練データ群を作成する。学習段階2では、物体検出ラベルおよびポーズ推定ラベルの項目にラベル情報が設定された訓練データ(セマンティックセグメンテーション推定ラベル情報の有無は不問である)が必要である。そのため、学習部162は、そのような条件を満たす訓練データをリストから抽出することにより、学習段階2で使用される訓練データ群を作成する。

20

【0069】

また、学習部162は、ステップS21において、リストの残りの訓練データ群から、学習段階1で使用される訓練データ群を作成する。学習段階1では、物体検出ラベルの項目にラベル情報が設定された訓練データ(ポーズ推定ラベル情報およびセマンティックセグメンテーション推定ラベル情報の有無は不問である)が必要である。そのため、学習部162は、そのような条件を満たす訓練データをリストから抽出することにより、学習段階1で使用される訓練データ群を作成する。

【0070】

次に、学習部162は、学習段階1、学習段階2、学習段階3、学習段階4の順に、それぞれ所定の終了条件が成立するまで、各段階の学習を行う(ステップS22~S25)。各段階の学習では、訓練データに含まれる画像をモデル153に入力したときにモデル153の出力として得られる推論タスクの推論結果と訓練データに含まれるラベル情報との誤差を、予め与えられた損失関数を用いて算出する。損失関数は、物体検出タスク、ポーズ推定タスク、セマンティックセグメンテーション推定タスク毎に存在する。物体検出タスクの損失関数をL1、ポーズ推定タスクの損失関数をL2、セマンティックセグメンテーション推定タスクの損失関数をL3とそれぞれ表記する。

30

【0071】

学習段階1では、損失関数L1で算出した損失を最小化するようにモデル153のコンポーネントCM2-1、CM4-1のパラメータを学習する。学習段階2では、損失関数L1で算出した損失と損失関数L2で算出した損失の総和(例えば重み付き和)を最小化するようにモデル153のコンポーネントCM2-1、CM2-2、CM4-1、CM4-2のパラメータを学習する。学習段階3では、損失関数L1で算出した損失と損失関数L2で算出した損失と損失関数L3で算出した損失の総和(例えば重み付き和)を最小化するようにモデル153のコンポーネントCM2-1、CM2-2、CM2-3、CM4-1、CM4-2、CM4-3のパラメータを学習する。学習段階4では、損失関数L1で算出した損失と損失関数L2で算出した損失と損失関数L3で算出した損失の総和(例えば重み付き和)を最小化するようにモデル153のコンポーネントCM1、CM2-1、CM2-2、CM2-3、CM4-1、CM4-2、CM4-3のパラメータを学習する。各学習では、例えば、勾配降下法と誤差逆伝搬法を用いてよい。

40

【0072】

50

以上、訓練データを用いてモデル153を学習する方法の例について説明した。しかし、本発明に適用可能な学習方法は以上の例に限定されない。例えば、次のような学習方法であってもよい。即ち、最初に、物体検出に係るコンポーネントCM2-1、CM4-1だけを学習する（他のコンポーネントCM1、CM2-2、CM2-3、CM4-2、CM4-3のパラメータは固定する）。次に、ポーズ推定に係るコンポーネントCM2-2、CM4-2だけを学習する（他のコンポーネントCM1、CM2-1、CM2-3、CM4-1、CM4-3のパラメータは固定する）。次に、セマンティックセグメンテーション推定に係るコンポーネントCM2-3、CM4-3だけを学習する（他のコンポーネントCM1、CM2-1、CM2-3、CM4-1、CM4-3のパラメータは固定する）。次に、全ての推論タスクに係るコンポーネントCM2-1~CM2-3、CM3-1~CM3-3）だけを学習する（コンポーネントCM1のパラメータは固定する。次に、モデル全体のコンポーネントCM1、CM2-1~CM2-3、CM4-1~CM4-3を学習する。

10

## 【0073】

以上説明したように、本実施形態に係る画像処理装置10によれば、複数のタスク間でタスク固有な高次の特徴量を相互に利用することができる。このため、複数のタスクのそれぞれにおいて、当該タスク固有の高次の特徴量と他のタスク固有の高次の特徴量とを考慮した学習および推定が可能となる。

## 【0074】

続いて、本実施形態の変形例について説明する。

20

## 【0075】

## &lt;変形例1&gt;

上記実施形態では、モデル153は、セマンティックセグメンテーション推定を行うように構成されていた。しかし、モデル153は、セマンティックセグメンテーション推定の代わりにインスタントセマンティックセグメンテーション推定を行うように構成されてよい。この場合、例えば図4に示されるマルチタスクモデル153のコンポーネントCM3とコンポーネントCM4-3との間に、特徴量FM3-3から物体検出を行うコンポーネントを追加し、コンポーネントCM4-3は物体検出された個々のクラスの矩形毎にピクセル単位でクラスを推定するように構成してよい。

## 【0076】

30

## &lt;変形例2&gt;

上記実施形態では、モデル153は、物体検出、ポーズ推定、セマンティックセグメンテーション推定の3つの推論タスクを行うように構成されていた。しかし、モデル153は、物体検出、ポーズ推定、および、セマンティックセグメンテーション推定のうちの何れか2つの推論タスクのみを行うように構成してよい。或いは、モデル153が行う推論タスクは、物体検出、ポーズ推定、セマンティックセグメンテーション推定に限定されず、それら以外のタスクであってもよい。

## 【0077】

## [第2の実施の形態]

図9は、本発明の第2の実施形態に係る画像処理システム20のブロック図である。図9を参照すると、画像処理システム20は、学習部21と学習済みモデル22とを備えている。

40

## 【0078】

学習部21は、画像から互いに相違する複数の推論タスクを行う学習済みモデル22を生成するように構成されている。学習部21は、例えば図1の学習部162と同様に構成することができるが、それに限定されない。

## 【0079】

学習済みモデル22は、上記画像から上記複数の推論タスクに共通な第1の特徴量を抽出する第1のコンポーネントと、上記推論タスクに対応して設けられ、上記第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出する第2のコンポーネントと、上記

50

推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成する第3のコンポーネントと、上記推論タスクに対応して設けられ、上記第3の特徴量から対応する推論タスクの推論結果を出力する第4のコンポーネントと、を含むように構成されている。

【0080】

上述のように構成された画像処理システム20は、以下のように動作する。即ち、学習部21は、画像から互いに相違する複数の推論タスクを行う学習済みモデル22を生成する。学習部21は、上記生成では、学習済みモデル22に、画像から複数の推論タスクに共通な第1の特徴量を抽出させ、次に、推論タスク毎に、第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出させ、次に、推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成させ、次に、推論タスク毎に、第3の特徴量から対応する推論タスクの推論結果を出力させる。

10

【0081】

以上のように構成され動作する画像処理システム20によれば、複数の推論タスク間でタスク固有な特徴量を相互に利用することができる。その理由は、画像処理システム20は、推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成し、第3の特徴量から対応する推論タスクの推論結果を出力するように構成されているためである。このため、複数の推論タスクのそれぞれにおいて、当該タスク固有の特徴量と他のタスク固有の特徴量とを考慮した学習および推定が可能となる。

【0082】

[第3の実施の形態]

20

図10は、本発明の第3の実施形態に係る画像処理システム30のブロック図である。図10を参照すると、画像処理システム30は、推定部31と学習済みモデル32とを備えている。

【0083】

推定部31は、学習済みモデル32を用いて、画像から互いに相違する複数の推論タスクの推論結果を出力するように構成されている。推定部31は、例えば図1の推定部163と同様に構成することができるが、それに限定されない。

【0084】

学習済みモデル32は、上記画像から上記複数の推論タスクに共通な第1の特徴量を抽出する第1のコンポーネントと、上記推論タスクに対応して設けられ、上記第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出する第2のコンポーネントと、上記推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成する第3のコンポーネントと、上記推論タスクに対応して設けられ、上記第3の特徴量から対応する推論タスクの推論結果を出力する第4のコンポーネントと、を含むように構成されている。

30

【0085】

上述のように構成された画像処理システム30は、以下のように動作する。即ち、推定部31は、学習済みモデル32を用いて、画像から互いに相違する複数の推論タスクの推論結果を推定して出力する。推定部31は、上記推定では、学習済みモデル32に、先ず、画像から複数の推論タスクに共通な第1の特徴量を抽出させ、次に、推論タスク毎に、第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出させ、次に、推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成させ、次に、推論タスク毎に、第3の特徴量から対応する推論タスクの推論結果を出力させる。

40

【0086】

以上のように構成され動作する画像処理システム30によれば、複数の推論タスク間でタスク固有な特徴量を相互に利用することができる。その理由は、画像処理システム30は、推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成し、第3の特徴量から対応する推論タスクの推論結果を出力するように構成されているためである。このため、複数の推論タスクのそれぞれにおいて、当該タスク固有の特徴量と他のタスク固有の特徴量とを考慮した学習および推定が可能となる。

【0087】

50

以上、上記各実施形態を参照して本発明を説明したが、本発明は、上述した実施形態に限定されるものではない。本発明の構成や詳細には、本発明の範囲内で当業者が理解しうる様々な変更をすることができる。

【産業上の利用可能性】

【0088】

本発明は、カメラ画像などの画像から、物体検出、ポーズ推定、セマンティックセグメンテーション推定などの複数の推論タスクを行う分野全般に利用できる。

【0089】

上記の実施形態の一部又は全部は、以下の付記のようにも記載され得るが、以下には限られない。

[付記1]

画像から互いに相違する複数の推論タスクを行う学習済みモデルを生成する学習部を含み、

前記学習済みモデルは、

前記画像から前記複数の推論タスクに共通な第1の特徴量を抽出する第1のコンポーネントと、

前記推論タスクに対応して設けられ、前記第1の特徴量から対応する推論タスクに固有な第2の特徴量を抽出する第2のコンポーネントと、

前記推論タスク毎に抽出された第2の特徴量を結合して第3の特徴量を生成する第3のコンポーネントと、

前記推論タスクに対応して設けられ、前記第3の特徴量から対応する推論タスクの推論結果を出力する第4のコンポーネントと、

を含む画像処理システム。

[付記2]

前記第3のコンポーネントは、前記複数の第2の特徴量のうちの1つを基準特徴量とし、前記基準特徴量以外の前記第2の特徴量のサイズを前記基準特徴量のサイズに合わせて変更し、前記サイズ変更後の前記基準特徴量以外の前記第2の特徴量と前記基準特徴量とを結合して前記第3の特徴量を生成し、前記推論タスク毎に、前記第3の特徴量のサイズを前記第4のコンポーネントの入力サイズに合わせて変更して前記第4のコンポーネントへ出力する、

付記1に記載の画像処理システム。

[付記3]

前記第3のコンポーネントは、前記推論タスクに対応するサブコンポーネントを含み、前記サブコンポーネントは、対応する前記推論タスクの前記第2の特徴量を基準特徴量とし、対応する前記推論タスク以外の前記第2の特徴量のサイズを前記基準特徴量のサイズに合わせて変更し、前記サイズ変更後の前記対応する前記推論タスク以外の前記第2の特徴量と前記基準特徴量とを結合して前記第3の特徴量を生成し、前記第4のコンポーネントへ出力する、

付記1に記載の画像処理システム。

[付記4]

前記学習部は、複数の学習段階に分けて前記学習済みモデルの学習を行い、

前記複数の学習段階は、少なくとも、

前記複数の推論タスクのうち何れか1つを学習対象タスクとし、前記学習対象タスク以外の推論タスクに係る前記第2のコンポーネントおよび前記第3のコンポーネントと前記第1のコンポーネントのパラメータを固定して、前記学習対象タスクに係る前記第2のコンポーネントおよび前記第3のコンポーネントのパラメータを学習する第1の学習段階と、

前記第1のコンポーネントのパラメータを固定して、前記全ての推論タスクに係る前記第2のコンポーネントおよび前記第3のコンポーネントのパラメータを学習する第2の学習段階と、を含む、

10

20

30

40

50

付記 1 乃至 3 の何れかに記載の画像処理システム。

[ 付記 5 ]

前記推論タスクに対応して設けられた第 4 のコンポーネントは、前記第 3 の特徴量を構成する複数の前記第 2 の特徴量のうち、対応する前記推論タスクの第 2 の特徴量の優先度合いを定める重みをそれ以外の第 2 の特徴量の優先度合いを定める重みより大きくする、付記 1 乃至 4 の何れかに記載の画像処理システム。

[ 付記 6 ]

前記推論タスクに対応して設けられた第 4 のコンポーネントは、入力された前記第 3 の特徴量に対して  $1 \times 1$  の畳み込みを行うことにより、前記第 3 の特徴量の次元数を削減する、

10

付記 1 乃至 5 の何れかに記載の画像処理システム。

[ 付記 7 ]

前記複数の推論タスクは、物体検出タスク、ポーズ推定タスク、セマンティックセグメンテーション推定タスクを含む、

付記 1 乃至 6 の何れかに記載の画像処理システム。

[ 付記 8 ]

前記学習済みモデルを用いて、画像から前記複数の推論タスクの推論結果を出力する推論部を、更に含む、

付記 1 乃至 7 の何れかに記載の画像処理システム。

[ 付記 9 ]

20

学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を出力する推論部を含み、

前記学習済みモデルは、

前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出する第 1 のコンポーネントと、

前記推論タスクに対応して設けられ、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出する第 2 のコンポーネントと、

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成する第 3 のコンポーネントと、

前記推論タスクに対応して設けられ、前記第 3 の特徴量から対応する推論タスクの推論結果を出力する第 4 のコンポーネントと、を含む画像処理システム。

30

[ 付記 10 ]

画像から互いに相違する複数の推論タスクを行う学習済みモデルを生成し、

前記生成では、前記学習済みモデルに、

前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、

前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、

前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、

40

画像処理方法。

[ 付記 11 ]

学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を推定して出力し、

前記推定では、前記学習済みモデルに、

前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、

前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、

50

前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、  
 画像処理方法。

[ 付記 1 2 ]

コンピュータに、画像から互いに相違する複数の推論タスクを行う学習済みモデルを生成する処理を行わせるためのプログラムであって、

前記生成では、前記学習済みモデルに、

前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、

前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、

10

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、

前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、

プログラムを記録したコンピュータ読み取り可能な記録媒体。

[ 付記 1 3 ]

コンピュータに、学習済みモデルを用いて、画像から互いに相違する複数の推論タスクの推論結果を推定して出力する処理を行わせるためのプログラムであって、

前記推定では、前記学習済みモデルに、

前記画像から前記複数の推論タスクに共通な第 1 の特徴量を抽出させ、

前記推論タスク毎に、前記第 1 の特徴量から対応する推論タスクに固有な第 2 の特徴量を抽出させ、

20

前記推論タスク毎に抽出された第 2 の特徴量を結合して第 3 の特徴量を生成させ、

前記推論タスク毎に、前記第 3 の特徴量から対応する推論タスクの推論結果を出力させる、

プログラムを記録したコンピュータ読み取り可能な記録媒体。

【符号の説明】

【 0 0 9 0 】

1 0 画像処理装置

1 1 カメラ I / F 部

1 2 通信 I / F 部

30

1 3 操作入力部

1 4 画面表示部

1 5 記憶部

1 6 演算処理部

1 7 画像サーバ

1 8 カメラ

1 5 1 プログラム

1 5 2 画像情報

1 5 3 モデル

1 5 4 推定結果情報

40

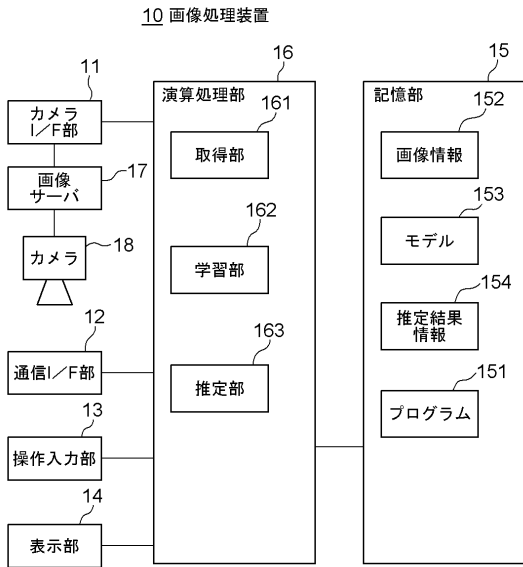
1 6 1 取得部

1 6 2 学習部

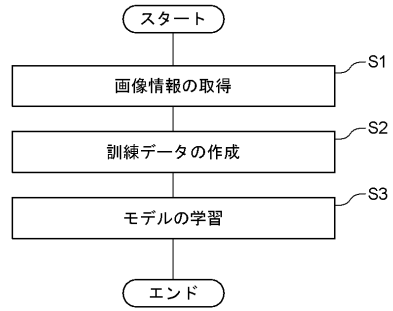
1 6 3 推定部

【図面】

【図 1】



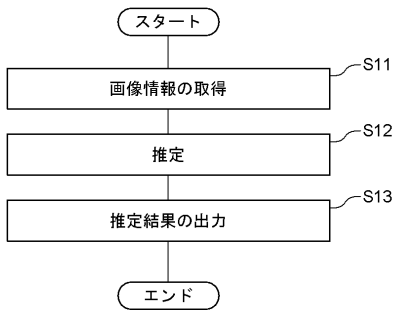
【図 2】



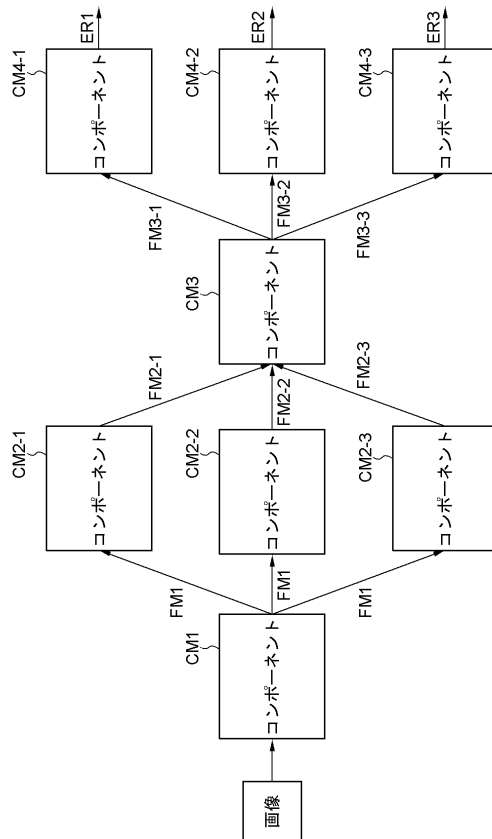
10

20

【図 3】



【図 4】

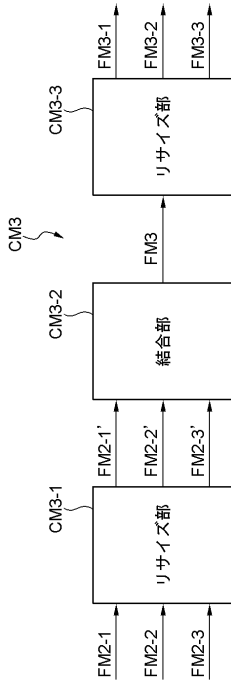


30

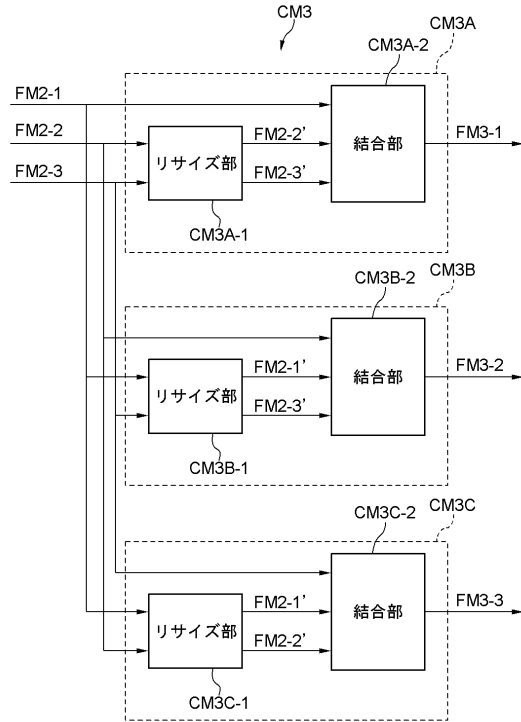
40

50

【図5】



【図6】



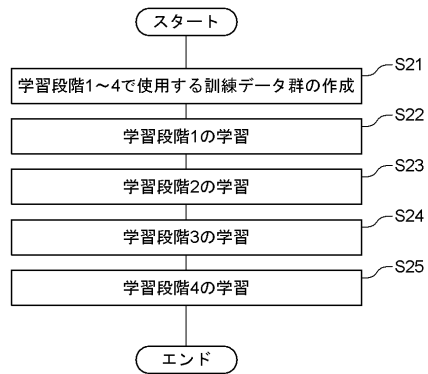
10

20

【図7】

ID	画像	物体検出ラベル	ポーズ検定ラベル	セグメンテーションラベル
1	G1	有	有	有
2	G2	有	有	無
3	G3	有	無	無
⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮
n	Gn	有	有	有

【図8】

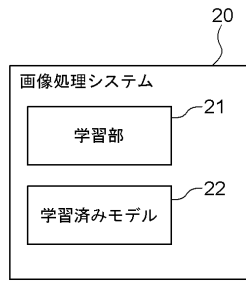


30

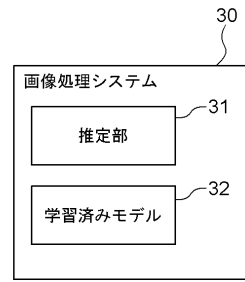
40

50

【図 9】



【図 10】



10

20

30

40

50

---

フロントページの続き

- (56)参考文献 特開 2021 - 021978 (JP, A)  
特開 2019 - 192009 (JP, A)
- (58)調査した分野 (Int.Cl., DB名)
- |      |               |
|------|---------------|
| G06T | 7/00 - 7/90   |
| G06V | 10/00 - 20/90 |
| G06N | 3/00 - 99/00  |