

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2018-190332

(P2018-190332A)

(43) 公開日 平成30年11月29日(2018.11.29)

(51) Int.Cl. F I テーマコード (参考)
G06T 7/00 (2017.01) G06T 7/00 350C 5L096
G06N 3/08 (2006.01) G06N 3/08

審査請求 未請求 請求項の数 18 O L (全 20 頁)

(21) 出願番号	特願2017-94694 (P2017-94694)	(71) 出願人	000001007
(22) 出願日	平成29年5月11日 (2017.5.11)		
		(74) 代理人	キヤノン株式会社 東京都大田区下丸子3丁目30番2号 100090273 弁理士 國分 孝悦
		(72) 発明者	角田 敬正 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	真継 優和 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
		Fターム(参考)	5L096 AA06 HA09 HA11 JA11 KA04

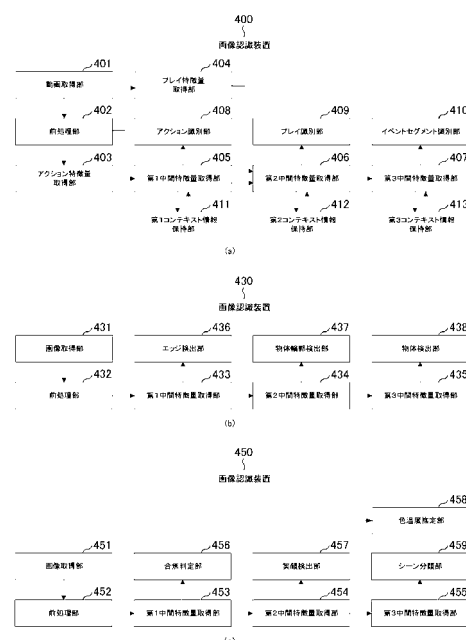
(54) 【発明の名称】 画像認識装置および学習装置

(57) 【要約】

【課題】DNNにおいて、スケールが互いに異なる複数の認識タスクを有し、1つの中間層から識別層を分岐させたネットワークで認識処理を行うことができるようにする。

【解決手段】ニューラルネットワークを用いたネットワーク構造により対象の認識を行う画像認識装置であって、入力画像から第1の識別を行うための第1の中間特徴量を取得する第1の取得手段と、前記第1の中間特徴量に基づいて前記第1の識別を行う第1の識別手段と、前記第1の中間特徴量から前記第1の識別よりもスケールの大きい第2の識別を行うための第2の中間特徴量を取得する第2の取得手段と、前記第2の中間特徴量に基づいて前記第2の識別を行う第2の識別手段と、を有する。

【選択図】図4



【特許請求の範囲】**【請求項 1】**

ニューラルネットワークを用いたネットワーク構造により対象の認識を行う画像認識装置であって、

入力画像から第 1 の識別を行うための第 1 の中間特徴量を取得する第 1 の取得手段と、
前記第 1 の中間特徴量に基づいて前記第 1 の識別を行う第 1 の識別手段と、

前記第 1 の中間特徴量から前記第 1 の識別よりもスケールの大きい第 2 の識別を行うための第 2 の中間特徴量を取得する第 2 の取得手段と、

前記第 2 の中間特徴量に基づいて前記第 2 の識別を行う第 2 の識別手段と、

を有することを特徴とする画像認識装置。

10

【請求項 2】

前記第 1 の識別手段は、前記ニューラルネットワークの中間層から分枝したネットワークによって実現されることを特徴とする請求項 1 に記載の画像認識装置。

【請求項 3】

前記スケールは、時間的尺度、空間的尺度、あるいは問題の複雑さの尺度の何れかであることを特徴とする請求項 1 又は 2 に記載の画像認識装置。

【請求項 4】

前記スケールは時間的尺度であり、前記ニューラルネットワークは再帰的ニューラルネットワーク (RNN) であることを特徴とする請求項 3 に記載の画像認識装置。

【請求項 5】

前記スケールは空間的尺度であり、前記ニューラルネットワークは畳み込みニューラルネットワーク (CNN) であることを特徴とする請求項 3 に記載の画像認識装置。

20

【請求項 6】

前記第 2 の識別手段は、さらに前記第 1 の中間特徴量に基づいて前記第 2 の識別を行うことを特徴とする請求項 5 に記載の画像認識装置。

【請求項 7】

前記入力画像は、静止画又は動画であることを特徴とする請求項 1 ~ 6 の何れか 1 項に記載の画像認識装置。

【請求項 8】

請求項 1 ~ 7 の何れか 1 項に記載の画像認識装置の識別器を学習する学習装置であって

30

、
前記第 1 の取得手段および前記第 2 の取得手段をそれぞれ構成する中間層と前記第 1 の識別手段および前記第 2 の識別手段をそれぞれ構成する識別層とを初期化する初期化手段と、

前記初期化手段によって初期化された中間層および識別層を、学習データを用い学習する学習手段と、

前記学習手段による学習の結果に応じて、中間層もしくは識別層を新たに追加する追加手段と、

を有することを特徴とする学習装置。

【請求項 9】

40

前記初期化手段は、前記中間層および前記識別層の数を最小限に構成したネットワーク構造を設定することを特徴とする請求項 8 に記載の学習装置。

【請求項 10】

前記追加手段は、前記学習手段による学習の結果に応じて、前記中間層の後段に新たな中間層を追加、もしくは前記中間層の最終層に新たな識別層を追加することを特徴とする請求項 8 又は 9 に記載の学習装置。

【請求項 11】

前記追加手段は、新たに中間層を追加する場合に前記識別層の下方に追加することを特徴とする請求項 10 に記載の学習装置。

【請求項 12】

50

前記学習手段による学習の結果を、検証用データを用いて評価する評価手段をさらに有し、

前記評価手段による評価の結果、精度が閾値未満である場合に、前記追加手段は、新たに中間層を追加することを特徴とする請求項 8 ～ 11 の何れか 1 項に記載の学習装置。

【請求項 13】

前記学習手段は、前記第 1 の識別手段および前記第 2 の識別手段による識別のそれぞれに対応する正解を含む学習データを用いて学習することを特徴とする請求項 8 ～ 12 の何れか 1 項に記載の学習装置。

【請求項 14】

前記学習手段は、確率的勾配降下法を用いて学習を行うことを特徴とする請求項 8 ～ 13 の何れか 1 項に記載の学習装置。

【請求項 15】

ニューラルネットワークを用いたネットワーク構造により対象の認識を行う画像認識装置の制御方法であって、

入力画像から第 1 の識別を行うための第 1 の中間特徴量を取得する第 1 の取得工程と、
前記第 1 の中間特徴量に基づいて前記第 1 の識別を行う第 1 の識別工程と、

前記第 1 の中間特徴量から前記第 1 の識別よりもスケールの大きい第 2 の識別を行うための第 2 の中間特徴量を取得する第 2 の取得工程と、

前記第 2 の中間特徴量に基づいて前記第 2 の識別を行う第 2 の識別工程と、

を有することを特徴とする画像認識装置の制御方法。

【請求項 16】

請求項 1 ～ 7 の何れか 1 項に記載の画像認識装置の識別器を学習する学習装置の制御方法であって、

前記第 1 の取得手段および前記第 2 の取得手段をそれぞれ構成する中間層と前記第 1 の識別手段および前記第 2 の識別手段をそれぞれ構成する識別層とを初期化する初期化工程と、

前記初期化工程において初期化された中間層および識別層を、学習データを用い学習する学習工程と、

前記学習工程による学習の結果に応じて、中間層もしくは識別層を新たに追加する追加工程と、

を有することを特徴とする学習装置の制御方法。

【請求項 17】

請求項 1 ～ 7 の何れか 1 項に記載の画像認識装置の各手段としてコンピュータを機能させるためのプログラム。

【請求項 18】

請求項 8 ～ 14 の何れか 1 項に記載の学習装置の各手段としてコンピュータを機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、特に、スケールの異なる複数の認識タスクを実行するために用いて好適な画像認識装置、学習装置、画像認識装置の制御方法、学習装置の制御方法およびプログラムに関する。

【背景技術】

【0002】

近年、画像や動画から物体などの対象を認識する技術が多様に研究されている。例えば、画像中の線分を抽出するエッジ検出や、類似する色やテクスチャの領域を切り出す領域分割、人物の目や口等のパーツを検出するパーツ検出、顔や人体、物体などを検出する物体認識、画像中の環境や状況を認識するシーン認識がある。ここで、認識処理の目的を認識タスクと呼ぶ。これらの認識対象は多様な尺度を持ち、個別の認識タスクを実行する認

10

20

30

40

50

識装置は、それぞれのスケールで定義された対象を、適切な特徴量に基づいて識別する。

【 0 0 0 3 】

近年、層の数が多くて深い多層ニューラルネットワーク (D N N : Deep Neural Network) を用いた画像認識の研究が活発である。例えば非特許文献 1 には、畳み込みニューラルネットワーク (C N N : Convolutional Neural Network) に正則化やデータ増し等の様々な工夫を加えた技術が開示されている。この技術は、従来の課題とされていた過学習等を克服し、物体検出タスクにおいて非常に高い性能を出しており、注目されている。

【 0 0 0 4 】

また、非特許文献 2 には、C N N に関する詳細な調査が開示されている。このような C N N では、入力層に近い低次の中間層で線分や単純図形、テクスチャ等の小さいスケールのプリミティブな特徴が捉えられている。また、識別層に近い高次の中間層で物体の形状やより詳細な特徴等、より大きいスケールの認識が行われる。

【 0 0 0 5 】

C N N では主に画像をターゲットとするが、一方で、再帰的ニューラルネットワーク (R N N : Recurrent Neural Networks) を用い、音声やテキスト等の時系列情報を認識する研究も盛んである。非特許文献 3 には、音声から音素の系列を認識するタスクに数種類の R N N を適用し、精度を比較評価している技術が開示されている。

【 先行技術文献 】

【 非特許文献 】

【 0 0 0 6 】

【 非特許文献 1 】 Krizhevsky, A., Sutskever, I., & Hinton, G. E., "Imagenet classification with deep convolutional neural networks.", In Advances in neural information processing systems (pp.1097-1105), 2012.

【 非特許文献 2 】 Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba, "Object detectors emerge in deep scene CNNs", ICLR2015

【 非特許文献 3 】 Alex Graves, Abdel-rahman Mohamed and Geoffrey Hinton, "Speech recognition with deep recurrent neural networks", Computing Research Repository 2013

【 非特許文献 4 】 J Donahue, Y Jia, O Vinyals, J Hoffman, N Zhang, E Tzeng, T Darrell, T Darrell, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", arXiv2013

【 非特許文献 5 】 Xi Li, Liming Zhao, Lina Wei, MingHsuan Yang, Fei Wu, Yueting Zhuang, Haibin Ling, Jingdong Wang, "DeepSaliency: Multi-Task Deep Neural Network Model for Salient Object Detection", IEEE Transactions on Image Processing, 2015

【 非特許文献 6 】 Karen Simonyan, Andrew Zisserman, "Two-Stream Convolutional Networks for Action Recognition in Videos", NIPS2015

【 非特許文献 7 】 Khurram Soomro, Amir Roshan Zamir, Mubarak Shah, "UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild", CRCV-TR-12-01, 2012

【 非特許文献 8 】 Li Shen, Teck Wee Chua, Karianto Leman, "Shadow optimization from structured deep edge detection", CVPR2015

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 7 】

以上の研究により、認識タスクが定義するカテゴリの空間的および時間的スケールなどの様々な尺度や粒度を考えた場合、小さい粒度の認識タスクでは低階層のニューラルネットワーク (N N) で十分である。一方、大きい粒度の認識タスクでは高階層の N N が望ましい。

【 0 0 0 8 】

一方、C N N の中間層の特徴量を他の認識タスクに転用する研究も多く行われ (例えば非特許文献 4 参照) 、認識タスクに依存する適切な層の存在は知られている。また、識別

10

20

30

40

50

層を分枝させ複数の認識タスクを１つのネットワークで実現する研究（例えば非特許文献５参照）も行われてきた。

【０００９】

しかし、前述したように、認識タスクによって時間的や空間的などのスケールは様々であり、さらには場合によっては、中間層の特徴量を他の認識タスクに転用し、１つの中間層から識別層を分枝させることが好ましい場合もある。このような場合に１つのネットワーク構造を実現するのが困難であった。

【００１０】

本発明は前述の問題点に鑑み、ＤＮＮにおいて、スケールが互いに異なる複数の認識タスクを有し、１つの中間層から識別層を分枝させたネットワークで認識処理を行うことができるようにすることを目的としている。

10

【課題を解決するための手段】

【００１１】

本発明に係る画像認識装置は、ニューラルネットワークを用いたネットワーク構造により対象の認識を行う画像認識装置であって、入力画像から第１の識別を行うための第１の中間特徴量を取得する第１の取得手段と、前記第１の中間特徴量に基づいて前記第１の識別を行う第１の識別手段と、前記第１の中間特徴量から前記第１の識別よりもスケールの大きい第２の識別を行うための第２の中間特徴量を取得する第２の取得手段と、前記第２の中間特徴量に基づいて前記第２の識別を行う第２の識別手段と、を有することを特徴とする。

20

【発明の効果】

【００１２】

本発明によれば、ＤＮＮにおいて、スケールが互いに異なる複数の認識タスクを有し、１つの中間層から識別層を分枝させたネットワークで認識処理を行うことができる。

【図面の簡単な説明】

【００１３】

【図１】動画の１フレームおよびカメラ配置の例を示す図である。

【図２】時間的スケールが互いに異なる各ラベルの関係を説明するための図である。

【図３】各実施形態に係る画像認識装置のハードウェア構成例を示すブロック図である。

【図４】各実施形態に係る画像認識装置の機能構成例を示すブロック図である。

30

【図５】各実施形態における画像認識装置による処理手順の一例を示すフローチャートである。

【図６】第１の実施形態に係るネットワーク構造の例を示す図である。

【図７】各実施形態に係る学習装置の機能構成例を示すブロック図である。

【図８】各実施形態における学習装置による処理手順の一例を示すフローチャートである。

。

【図９】最も浅いネットワークに構成した具体的な初期構造の例を示す図である。

【図１０】異なる識別層からの由来に基づく勾配を説明するための図である。

【図１１】認識タスクの精度が不十分である場合の対応を説明するための図である。

【図１２】識別層を追加する場合の処理を説明するための図である。

40

【図１３】ＣＮＮのネットワーク構造の例を示す図である。

【図１４】ＣＮＮにおける特徴抽出ユニットと識別ユニットとを説明するための図である。

。

【図１５】第２の実施形態におけるＣＮＮのネットワーク構造の例を示す図である。

【図１６】第３の実施形態におけるＣＮＮのネットワーク構造の例を示す図である。

【発明を実施するための形態】

【００１４】

（第１の実施形態）

以下、本発明の第１の実施形態について、図面を参照しながら説明する。

例えば、サッカーやラグビーなどのチームスポーツ競技の未編集動画から、試合内容に

50

即したハイライトビデオを自動作成するには、競技特有のプレイや、それらの連鎖である重要なイベントの時間的セグメントを認識することが重要である。さらに競技特有のプレイを詳細に認識するには、映像全体から得られる特徴と、各プレイヤーのプリミティブな行動とが、重要な手がかりとなる。すなわち、試合内容に即したハイライトビデオの自動生成を実現するために、複数の認識タスクを実行することが重要である。つまり、プレイヤーの行動認識、競技特有のプレイ認識、重要イベントの時間的セグメントの認識（イベントセグメント認識）という、少なくとも3つの時間的尺度の異なる階層的な認識タスクを実行することが重要である。

【0015】

各認識タスクは、上位のタスクほど時間的スケールの大きいタスクで、下位のタスクほど時間的スケールの小さいタスクであり、上位は下位を包含するため、1つのフレームワークで同時に実行することが効率的であるといえる。そこで、本実施形態では、例えばフットサル動画を題材に、RNNで認識時の処理として多層BLSTM（Bidirectional Long Short-Term Memory）で異なる時間的尺度のマルチタスクを行う例について説明する。そして、学習時の処理として、多層BLSTMのネットワーク構造を最適化する方法について述べる。

10

【0016】

図1(a)は、本実施形態において取得されるフットサル動画の1フレーム（静止画）101の例を示す図である。図1(a)において、動画の1フレーム101中には複数人のプレイヤー102とボール104とが存在する。また、図1(b)に示すカメラ配置例151のように、コート152の周囲に複数台のカメラ153が配置され、全てのプレイヤーは何れかのカメラでキャプチャされるものとする。

20

【0017】

図2は、動画（静止画の時系列）201に対する、時間的スケールが互いに異なる「イベントラベル」、「プレイラベル」、「アクションラベル」の関係を説明するための図である。アクションラベル205は、最も時間的スケールの小さいラベルで、スポーツ競技に依存せず汎用的なアクションを表現するラベルである。競技にはN人のプレイヤー211～21Nが存在し、アクションラベル205はプレイヤー毎に設定される。プレイラベル204はアクションラベル205より時間的スケールが大きく、ボールを中心として定義される競技固有のプレイを表現するラベルである。イベントラベル203は最も時間的スケールの大きいラベルで、ゴールシーンやシュートシーンなどの重要シーンの時間的セグメントを表現するラベルである。以上のラベルは、画面やプレイヤーに対して常に1つのラベルが設定され、複数存在しないものとする。なお、以下の表1には、アクションラベル、プレイラベル、イベントラベルの例を示す。

30

【0018】

【表 1】

アクションラベル	プレイラベル	イベントラベル
キック	パス	シュート
レシーブ	シュート	ゴール
ドリブル	ゴール	コーナーキック
トラップ	ドリブル	フリーキック
キープ	クリア	
キャッチ	コーナーキック	
カット	フリーキック	
ラン	パントキック	
ターン	インターセプト	
ステップ	キーパーセーブ	
スライディング		

10

20

【 0 0 1 9 】

図 3 は、本実施形態に係る画像認識装置 3 0 0 のハードウェア構成例を示すブロック図である。

図 3 において、画像認識装置 3 0 0 は、CPU 3 0 1 と、ROM 3 0 2 と、RAM 3 0 3 と、HDD 3 0 4 と、表示部 3 0 5 と、入力部 3 0 6 と、ネットワーク I / F 部 3 0 7 とを有している。CPU 3 0 1 は、ROM 3 0 2 に記憶された制御プログラムを読み出して各種処理を実行する。RAM 3 0 3 は、CPU 3 0 1 の主メモリー、ワークエリア等の一時記憶領域として用いられる。HDD 3 0 4 は、各種データや各種プログラム等を記憶する。表示部 3 0 5 は、各種情報を表示する。入力部 3 0 6 は、キーボードやマウスを有し、ユーザによる各種操作を受け付ける。

30

ネットワーク I / F 部 3 0 7 は、ネットワークを介して画像形成装置等の外部装置との通信処理を行う。また、他の例としては、ネットワーク I / F 部 3 0 7 は、無線により外部装置との通信を行ってもよい。

【 0 0 2 0 】

なお、後述する画像認識装置の機能や処理は、CPU 3 0 1 が ROM 3 0 2 又は HDD 3 0 4 に格納されているプログラムを読み出し、このプログラムを実行することにより実現されるものである。また、他の例としては、CPU 3 0 1 は、ROM 3 0 2 等に替えて、SD カード等の記録媒体に格納されているプログラムを読み出してもよい。

【 0 0 2 1 】

40

図 4 (a) は、本実施形態に係る画像認識装置 4 0 0 の機能構成例を示すブロック図である。

本実施形態の画像認識装置 4 0 0 は、動画取得部 4 0 1、前処理部 4 0 2、アクション特徴量取得部 4 0 3、及びプレイ特徴量取得部 4 0 4 を有する。さらに、第 1 中間特徴量取得部 4 0 5、第 2 中間特徴量取得部 4 0 6、第 3 中間特徴量取得部 4 0 7、アクション識別部 4 0 8、プレイ識別部 4 0 9、及びイベントセグメント識別部 4 1 0 を有する。さらに本実施形態における画像認識装置 4 0 0 は、記憶手段として、第 1 コンテキスト情報保持部 4 1 1、第 2 コンテキスト情報保持部 4 1 2、第 3 コンテキスト情報保持部 4 1 3 を有する。画像認識装置 4 0 0 が有するこれらの各機能の詳細については、図 5 (a) 等を用いて後述する。

50

【 0 0 2 2 】

図 5 (a) は、本実施形態における画像認識装置 4 0 0 による処理手順の一例を示すフローチャートである。

まず、S 5 0 1 において、動画取得部 4 0 1 は、複数のフレームからなる静止画像系列を取得する。例えば、図 1 (a) で例示した識別対象の動画の一部を動画取得部 4 0 1 が取得する。

【 0 0 2 3 】

本実施形態では、動画としてフットサルの試合映像を対象とし、最終的な動画の長さはフットサルの試合時間に準ずる。また、動画のフレームレートは、秒間 3 0 フレーム程度の標準的な速度を想定する。本工程では、動画全体の中から、連続する数フレームから数 1 0 フレームの予め設定した長さ（例えば試合時間）の静止画の系列を取得する。ここでは、説明を簡略化するため 6 0 フレームの静止画を取得するものとする。この 6 0 フレームは、外部装置に予め記憶されていた動画から取得しても良いが、リアルタイムにカメラで撮影した動画から取得してもよい。前者の場合、動画取得部 4 0 1 は、外部装置に記憶された試合映像から所定のフレーム数の静止画系列を取得し、後者の場合、カメラが出力する静止画系列から所定のフレーム数を取得する。また、図 1 (b) に示したように、複数台のカメラで撮影されたことを前提とし、本工程ではそのカメラ台数分の静止画系列を取得する。

【 0 0 2 4 】

次に、S 5 0 2 においては、前処理部 4 0 2 は、S 5 0 1 で取得した静止画系列に対して前処理を行う。まず、前処理部 4 0 2 は、カメラ台数分の静止画系列に映る各プレイヤーとボールとをトラッキングし、プレイヤー位置およびプレイヤー領域、ボール位置の情報を取得する。ここで取得するプレイヤー領域は、図 1 (a) に示したプレイヤー領域 1 0 3 に該当する領域である。

【 0 0 2 5 】

次に、前処理部 4 0 2 は、画面全体と各プレイヤー領域とに関して、図 1 (b) のように配置された複数台のカメラからよりよい配置のカメラを選択する。本実施形態では、画面全体に関しては最初のフレームでボールをより手前で撮影しているカメラ、プレイヤー領域に関しては最初のフレームで領域の隠れがなくプレイヤーがより大きく映っているカメラを選択する。次に、前処理部 4 0 2 は、画面全体および各プレイヤー領域の静止画系列の最初の 1 フレームに対し、色成分の正規化を行い、さらに静止画系列からオプティカルフローを取得するための処理を行う。オプティカルフローは、ここではフレーム間のブロックマッチング法を用いて取得する。

【 0 0 2 6 】

次に、S 5 0 3 においては、アクション特徴量取得部 4 0 3 は、プレイヤー毎のアクションを識別するためのアクション特徴量を取得する。ここでアクション特徴量とは、動画取得部 4 0 1 で取得した静止画系列から計算される、少なくともアピアランスとモーションとの 2 つに関する特徴量である。非特許文献 6 には、アピアランスおよびモーションに対応する 2 つの CNN を用いて、数秒から数分程度の動画で構成される公開データセット（非特許文献 7 参照）、UCF - 1 0 1 等に対するアクション分類が開示されている。

【 0 0 2 7 】

モーションに関する CNN には、オプティカルフローのベクトル場を鉛直方向成分と垂直方向成分とに分けた 2 つのマップを 1 単位として、それらを積み上げた複数パターン（1、5、1 0 等）を入力としている。本実施形態では、1 0 単位分のオプティカルフローを入力とし、UCF - 1 0 1 のアクション分類タスクで学習した CNN を用い、Softmax 識別層の前の全結合層の出力をモーションに関する特徴量として用いる。

【 0 0 2 8 】

一方、アピアランスに関する特徴量は、非特許文献 4 に記載された方法と同様に、物体検出タスクの公開データセットである ILSVRC を用いて学習した CNN の Softmax 識別層の前の全結合層の出力を用いる。予め学習した CNN の中間出力を別の認識タ

10

20

30

40

50

スクに転用するこのような利用方法は、非特許文献 4 等で記載の通り一般的な方法であるので、これ以上の詳細な説明は割愛する。

【0029】

次に、S504においては、プレイ特徴量取得部404は、正規化した画面全体の静止画とオプティカルフローとから、プレイを識別するためのプレイ特徴量を取得する。ここでは、アクション特徴量取得部403が行った処理と同様の方法でアピアランスおよびモーションに関する特徴量をプレイ特徴量として取得する。

【0030】

次に、画面全体から取得したプレイ特徴量と、プレイヤー毎のアクション特徴量とから、アクション、プレイ、イベントセグメントの3つを識別する。第1中間特徴量取得部405、第2中間特徴量取得部406、および第3中間特徴量取得部407は、後述する学習時の処理で最適化された層数を持つBLSTMに対応する。

10

【0031】

図6は、アクション、プレイ、イベントセグメントに関してBLSTM等のユニットでグラフ表現した模式図である。図6に示すように、プレイヤー毎のアクション特徴量601は、それぞれ独立に複数層で構成されるBLSTM605に入力され、アクション識別層606を介して、プレイヤー毎のアクションが識別される。そして画面全体から取得したプレイ特徴量604は、複数層のBLSTM605を介して、プレイヤー毎のアクション特徴量601と統合される。そして複数層のBLSTM605とプレイ識別層607とを介してプレイラベルが識別される。プレイ特徴量604およびアクション特徴量601は、さらに複数層のBLSTM605とイベントセグメント識別層608とを介し、最終層でイベントラベルが識別される。ここで、アクション識別層606およびプレイ識別層607は、ソフトマックス識別器と中間層とを1層ずつ持つマルチレイヤーの識別層であり、イベントセグメント識別層608はソフトマックス識別器で構成されている識別層である。

20

【0032】

まず、S505において、第1中間特徴量取得部405は、アクション特徴量の時系列に関する特徴量を取得する。第1中間特徴量取得部405は、BLSTMによって構成され、第1コンテキスト情報保持部411には、前ステップの第1中間特徴量が保存されている。第1中間特徴量取得部405は、アクション特徴量取得部403で取得したアクション特徴量と前ステップの第1中間特徴量とから現ステップの第1中間特徴量を計算する。また、このBLSTMには、後述する学習時の処理で記載した方法に従って学習したBLSTMを用いる。

30

【0033】

次に、S506において、アクション識別部408は、現ステップの第1中間特徴量に、ソフトマックス識別器を用いることで、アクションラベルの識別スコアを計算する。このソフトマックス識別器の学習方法に関しては後述する。

【0034】

S507においては、第2中間特徴量取得部406は、プレイ特徴量の時系列に関する特徴量を取得する。第2中間特徴量取得部406は、BLSTMによって構成され、第2コンテキスト情報保持部412には、前ステップの第2中間特徴量が保存されている。第2中間特徴量取得部406は、プレイ特徴量取得部404で取得したプレイ特徴量と、現ステップの第1中間特徴量と、前ステップの第2中間特徴量とを演算して現ステップの第2中間特徴量を計算する。また、このBLSTMには、後述する学習時の処理で記載した方法に従って学習したBLSTMを用いる。

40

【0035】

次に、S508において、プレイ識別部409は、現ステップの第2中間特徴量に、ソフトマックス識別器を用いることで、プレイラベルの識別スコアを計算する。

【0036】

S509においては、第3中間特徴量取得部407は、現ステップの第3中間特徴量を

50

取得する。第3中間特徴量取得部407は、BLSTMによって構成され、第3コンテキスト情報保持部413には、前ステップの第3中間特徴量が保存されている。第3中間特徴量取得部407は、現ステップの第2中間特徴量と前ステップの第3中間特徴量とを演算して現ステップの第3中間特徴量を計算する。また、このBLSTMには、後述する学習時の処理で記載した方法に従って学習したBLSTMを用いる。

【0037】

次に、S510において、イベントセグメント識別部410は、現ステップの第3中間特徴量に、ソフトマックス識別器を用いることで、イベントラベルの識別スコアを計算する。

【0038】

次に、本実施形態におけるS505～S510で利用するBLSTM（中間層）および識別層の学習方法と、図6に示したネットワーク構造の最適化方法とについて説明する。

【0039】

図7は、本実施形態における学習装置700の機能構成例を示すブロック図である。なお、本実施形態に係る学習装置700のハードウェア構成については、基本的には図3に示した構成と同様の構成である。

図7に示すように、学習装置700は、ネットワーク構造初期化部701、ネットワークパラメータ最適化部702、精度評価部703、層追加部704、およびタスク追加部705を有する。さらに学習装置700は、記憶手段として、学習用データ保持部706、検証用データ保持部707、およびネットワークパラメータ保持部708を有する。

【0040】

図8は、本実施形態における学習装置700が行う処理手順の一例を示すフローチャートである。ここで各工程の概要及び図7に示した各構成の機能について説明する。

S801においては、ネットワーク構造初期化部701は、認識タスク毎の認識層とネットワーク全体の構造とを初期化する。本実施形態では、まず認識タスクとして、最もスケールの小さいアクション認識と次にスケールの小さいプレイ認識とを、最も浅いネットワークに構成した構造を初期構造とする。図9には、最も浅いネットワークに構成した具体的な初期構造の例を示す。図6に示した構造と比較して、アクション特徴量901、プレイ特徴量904、BLSTM905、アクション識別層906およびプレイ識別層907に示すように、中間層及び識別層の数を最小限にすることにより、最も浅いネットワークを構成している。

【0041】

次に、S802において、ネットワークパラメータ最適化部702は、学習用データ保持部706に保持されている学習用データを用い、ネットワークを構成するBLSTMおよび識別層のパラメータの最適化（学習）を行う。ここで、最適化するパラメータのうち、プレイヤー毎のアクション識別層およびそれに接続されるBLSTMは、プレイヤーが変わっても同じ定義のラベルを識別し、プレイヤーの見えは試合毎に様々なバリエーションがあり得る。そこで、プレイヤー毎のBLSTMおよびアクション識別層において重みは共有するものとする。

【0042】

ここで、学習用データには、動画を構成する静止画の時系列に対応させて、イベントラベル、プレイラベル、プレイヤー毎のアクションラベルの正解ラベルが用意されている。各BLSTMおよび識別層のパラメータを最適化する際には確率的勾配降下法を用いる。この場合、図10に示すように、BLSTM1001において、プレイ識別層1003に由来する勾配1005と、アクション識別層1002に由来する勾配1004とが混合する。このようなBLSTMでは、複数の勾配を単純に平均しても良いが、検証用データを用意して誤差率を計量し、その誤差率で勾配を混合させても良い。

【0043】

誤差逆伝搬法で計算される当該BLSTMの全体の勾配を E とすると、アクション識別層に由来する勾配を $E^{(1)}$ 、プレイ識別層に由来する勾配を $E^{(2)}$ とした場合に、B

10

20

30

40

50

LSSTMの全体の勾配は以下の式(1)により算出される。

$$E = E^{(1)} + E^{(2)} \cdots (1)$$

【0044】

ここで、 $E^{(1)}$ はアクション識別の検証用データに関する誤差率に比例する値で、 $E^{(2)}$ はプレイ識別の検証用データに関する誤差率に比例する値である。 $E^{(1)}$ 、 $E^{(2)}$ は足して1となるように正規化されている。

【0045】

また、パラメータの最適化には確率的勾配降下法を適用するため、学習データはミニバッチに分割するが、検証用データも同様にミニバッチに分割し、ミニバッチ毎に誤差率を計算しても良い。こうすることで、ミニバッチ毎に混合比が変わるため、ミニバッチ毎に勾配を変える確率的勾配降下法と同様の効果、すなわちより良い局所解に降下する可能性が増すと考えられる。

10

【0046】

また、アクションの識別はプレイヤー毎に行うため、正解ラベルはN人のプレイヤー分あり、以下の式(2)によりN人のプレイヤー分の勾配が得られる。

$$E^{(1)} = E_1^{(1)} + E_2^{(1)} + \cdots + E_N^{(1)} \cdots (2)$$

BLSTMの全体の勾配を算出する際に、式(1)を用いたが、代わりに以下の式(3)を用いてもよい。この場合、アクション識別の勾配は、プレイ識別の勾配に対しN人分の重みが加わっているため、それを相殺するためにNでアクション識別の勾配を割っても良い。

20

$$E = (1/N) \cdot E^{(1)} + E^{(2)} \cdots (3)$$

以上の処理により学習されたパラメータは、ネットワークパラメータ保持部708に記憶される。

【0047】

次にS803において、精度評価部703は、検証用データ保持部707に記憶されている検証用データを用い、最適化されたパラメータに基づく各識別タスクの精度を得る。すなわち、S801およびS802を経た現段階では、アクション識別およびプレイ識別の2つのタスクに関して、精度を計算する。本工程で用いる検証用データは、S802での勾配の混合比の決定に用いる検証用データと同じものでも良いし、別の検証用データを用意しても良い。

30

【0048】

次に、S804において、精度評価部703は、S803で得た精度が閾値以上であるか否かを判定する。精度が閾値未満で精度が不十分な場合はS805に分岐し、そうでない場合はS806に分岐する。

【0049】

S805においては、層追加部704は、タスク毎に既に存在するBLSTMの後段にBLSTMを追加する。アクション識別、プレイ識別ともに精度が十分でない場合には図11(a)に示すように、アクション識別層1106の下方と、プレイ識別層1105の下方とにBLSTM1101、1102を追加する。また、アクション識別の精度のみが十分でない場合は、図11(b)に示すように、アクション識別層の下方のみにBLSTM1103を追加する。一方、プレイ識別の精度のみが十分でない場合には、図11(c)に示すように、プレイ識別層の下方のみにBLSTM1104を追加する。そして、S802に戻る。

40

【0050】

一方、識別タスクの精度がともに十分である場合、その識別タスクの識別層の分枝位置が決定する。そして、S806において、精度評価部703は、全てのタスクの精度が十分であるか否かを判定する。この判定の結果、精度が不十分であるタスクがある場合は、S807に分岐する。

【0051】

S807においては、タスク追加部705は、タスクを追加する。そして、S805を

50

経て、追加するタスクの B L S T M と識別層とを追加することになる。図 1 1 (a) が現在の状態である場合、図 1 2 に示すように、S 8 0 7 では最終層にイベントセグメント識別層 1 2 0 1 が追加され、イベントセグメント識別層 1 2 0 1 の下方には S 8 0 5 で B L S T M 1 2 0 2 が追加される。

【 0 0 5 2 】

以降、S 8 0 2 および S 8 0 3 を再度行い、最終的に全てのタスクで精度が十分になるまで、B L S T M の追加と識別層の分枝位置の探索とを繰り返す。

【 0 0 5 3 】

以上のように本実施形態によれば、アクション識別、プレイ識別、イベントセグメント識別という時間的スケールの異なるマルチタスクを行う R N N のニューラルネットワークを実現できる。これにより、スポーツ動画の入力に対してアクション、プレイ、イベントセグメントの各ラベルを同時に推定することができる。

10

【 0 0 5 4 】

(第 2 の実施形態)

本実施形態では、空間的尺度の異なるマルチタスクを行う C N N を例に説明する。非特許文献 2 に記載の技術では、1 0 0 0 クラスの物体検出タスクで学習した 7 層の C N N において各層が何に反応しているかをクラウドソーシングで調査している。受容野に対応する画像をワーカーに分類させると、1 層目が色や単純図形、2 層目がテクスチャ、3 層目が領域や表面、4 層目がオブジェクトパーツ、5 層目がオブジェクトに反応しているとする答えが多かった。このことは最終的に 1 0 0 0 クラスの物体検出を実現するために、低層層で暗黙的に上記の性質が獲得されていることを示している。最終層で物体検出を学習すると同時に、低次の層で単純図形やテクスチャ、領域分割等のプリミティブな画像認識を教師あり学習の枠組みで積極的に学習することで、1 つのネットワークで空間的スケールの異なるマルチタスクを行う C N N が学習できる。

20

【 0 0 5 5 】

そこで本実施形態では、コンピュータビジョンにとって特に重要で、かつ教師あり学習の仕組みで学習しやすい複数のタスクに関し、ニューラルネットワークの構造を決定して学習し、識別する方法を説明する。具体的には、「エッジ検出」、「物体輪郭検出」、「物体検出」の 3 つのタスクに関し、これを行うニューラルネットワークの構造を決定して学習し、識別する方法を説明する。

30

【 0 0 5 6 】

まず、本実施形態において学習される C N N を用いて画像を識別する際の処理について説明する。C N N は畳みこみ演算を多く行うニューラルネットワークであり、非特許文献 1 や非特許文献 2 に開示されている。つまり、C N N では、畳み込み層 (C o n v (Convolutional Layer)) と非線形処理 (R e L U (Rectified linear unit) や P o o l (Max pooling) など) との組み合わせで特徴抽出層が実現される。そのあと、C N N は全結合層 (F C (Fully-connected Layer)) を経て画像分類結果 (各クラスに対する尤度) を出力する。

【 0 0 5 7 】

図 1 3 に、非特許文献 1 に開示されているネットワーク構造の例を示す。画像 I m g 1 3 0 1 を入力すると、C o n v + P o o l 1 3 0 2、C o n v + P o o l 1 3 0 3、C o n v 1 3 0 4、C o n v 1 3 0 5、C o n v + P o o l 1 3 0 6 を順番に適用して特徴抽出が行われる。その後、識別層にて F C 1 3 0 7、F C 1 3 0 8、S o f t m a x 1 3 0 9 を行い、マルチクラスのカテゴリ尤度を出力している。なお、全ての C o n v および F C の後には非線形処理 (R e L U) を行うが、図 1 3 に示す例では簡略化のために省略している。また、画像 I m g 1 3 0 1 は、C N N に入力する際、所定サイズに画像をリサイズしてクロップし、正規化するなどの前処理が行われるのが一般的である。

40

【 0 0 5 8 】

図 4 (b) は、本実施形態の画像認識装置 4 3 0 の機能構成例を示すブロック図である。なお、本実施形態に係る画像認識装置 4 3 0 のハードウェア構成については図 3 と同様

50

である。

本実施形態の画像認識装置 430 は、画像取得部 431、前処理部 432、第 1 中間特徴量取得部 433、第 2 中間特徴量取得部 434、第 3 中間特徴量取得部 435、エッジ検出部 436、物体輪郭検出部 437、および物体検出部 438 を有する。画像認識装置 430 が有するこれらの各機能の詳細については、図 5 (b) 等を用いて後述する。

【0059】

図 5 (b) は、本実施形態における画像認識装置 430 による処理手順の一例を示すフローチャートである。

まず、S531 において、画像取得部 431 は、マルチタスクの認識を行う対象である静止画を取得する。静止画は外部の記憶装置から取得しても良いし、スチルカメラ等の撮像デバイスで新たに取得しても良い。

【0060】

次に S532 においては、前処理部 432 は、前工程で取得した入力画像に対し前処理を行う。ここでは入力画像に対して色成分を正規化して 256×256 程度の大きさにリサイズし、続いて 224×224 程度の大きさにクロップする前処理を行う。

【0061】

次に前処理した入力画像から、第 1 中間特徴量取得部 433、第 2 中間特徴量取得部 434、第 3 中間特徴量取得部 435 の 3 つの中間特徴量取得部がそれぞれ特徴量を取得する。そして、エッジ検出部 436、物体輪郭検出部 437、および物体検出部 438 が、それぞれエッジ抽出、物体輪郭抽出、物体検出の 3 つを行う。3 つの各中間特徴量取得部は、後述する学習時の処理で最適化された構造を持つ CNN である。

【0062】

ここで、スケールの異なるマルチタスクを行う CNN を実現するために、図 14 に示したような基本ユニットを考える。図 14 (a) ~ 図 14 (c) は、それぞれ Conv と Pool とが組み合わされた特徴抽出層のユニット (特徴抽出ユニット) の例を示している。より低層側では、図 14 (a) に示す特徴抽出ユニットを用い、中層側では図 14 (b) に示す特徴抽出ユニットを用い、高層側では、図 14 (c) に示す特徴抽出ユニットを用いるものとする。

【0063】

また、エッジ抽出および物体輪郭抽出では、ピクセル毎の 2 値分類を行うので、識別層として図 14 (d) に示すような FC と Logistic Regression とを組み合わせた識別ユニットを用いる。FC は通常のように特徴抽出ユニット最終層の全チャンネル全レスポンスマップと接続せずに、マップ上の 4×4 の領域をチャンネル方向に全結合しマップ上の各点で 2 値判定を行う。そして、最終アウトプットである物体検出では、マルチクラスの分類を行うので、図 14 (e) に示すような $FC \times 2 + Softmax$ となる識別ユニットを用いるものとする。

【0064】

図 15 は、図 4 (b) の各構成および入力画像のネットワークを、図 14 のユニットを用いて表現した図である。第 1 中間特徴量取得部 433 は図 14 (a) の $Conv \times 2 + Pool$ の特徴抽出ユニットで実現される。そして、第 2 中間特徴量取得部 434 は図 14 (b) の $Conv \times 3 + Pool$ の特徴抽出ユニット、第 3 中間特徴量取得部 435 は図 14 (c) の $Conv \times 4 + Pool$ の特徴抽出ユニットで実現される。また、前述の通り、エッジ検出部 436 および物体輪郭検出部 437 は、それぞれ図 14 (d) の Logistic Regression の識別ユニットで実現される。そして、物体検出部 438 は、図 14 (e) の $Softmax$ の識別ユニットで実現される。

【0065】

図 15 に示すように、入力画像 1501 から、特徴抽出ユニット 1502 を介して、Logistic Regression の識別ユニット 1505 でエッジ検出が行われる。そして、特徴抽出ユニット 1502 のもう一方の出力により、特徴抽出ユニット 1503 を介し、Logistic Regression の識別ユニット 1506 で物体輪郭抽出が行われる。さらに、特徴抽出ユニ

10

20

30

40

50

ット1503の出力により、特徴抽出ユニット1504を介し、Softmaxの識別ユニット1507で物体検出が行われる。

【0066】

S533においては、第1中間特徴量取得部433は、前処理が行われた入力画像からエッジ検出に係る第1中間特徴量を取得する。そして、S534において、エッジ検出部436は、第1中間特徴量から、画像中のエッジを検出する。

【0067】

次に、S535において、第2中間特徴量取得部434は、第1中間特徴量に基づき、物体の輪郭に係る第2中間特徴量を取得する。そして、S536において、物体輪郭検出部437は、第2中間特徴量から、物体の輪郭を検出する。

10

【0068】

次に、S537において、第3中間特徴量取得部435は、第2の中間特徴量に基づき、物体検出に係る第3中間特徴量を取得する。そして、S538において、物体検出部438は、第3中間特徴量から、画像中の物体を検出する。

【0069】

以上のように、S533、S535、S537およびS538における処理は、一般的なCNNにおける認識時の処理と同様であるので、詳細な説明は割愛する。そして、S534およびS536の処理は、画素毎の2値分類を行う処理であるが、ここでは非特許文献8に記載の方法を参考にして実現する。すなわち、特徴抽出ユニットの最終層のレスポンスマップ上の4×4の領域を、チャネル方向に25次元のFCと全結合させる。そして、レスポンスマップ上の4×4領域に対応する受容野において、中央に位置する5×5ピクセルの全点でエッジの有無を2値判定するという方法である。

20

【0070】

以上の処理により、入力画像に対し、エッジ抽出、物体輪郭抽出、物体検出の各識別タスクを1つのネットワークで同時に実現できる。次に、本実施形態におけるS533～S538で利用する特徴抽出ユニットおよび識別ユニットの学習方法と図15に示したネットワーク構造の探索方法とについて説明する。なお、本実施形態における学習装置の構成は、第1の実施形態と同様であり、本実施形態における学習の処理手順は、基本的には図8と同様である。以下、第1の実施形態との相違について説明する。

【0071】

30

まず、S801において、ネットワーク構造初期化部701は、CNNのネットワーク構造を初期化する。次のS802においては、ネットワークパラメータ最適化部702は、エッジ検出、物体輪郭検出（および物体検出）のネットワークに組み込まれたタスクの識別ユニットと特徴抽出ユニットとのネットワークパラメータを最適化する。S803では、精度評価部703は、エッジ検出、物体輪郭検出（および物体検出）の各タスクの評価用データを用い、タスク毎に精度を評価する。S805では、層追加部704は、精度が不十分であるタスクに応じて、図14(a)～図14(c)で例示した特徴抽出ユニットを追加する。S807では、タスク追加部705は、図14(e)に示す物体検出用の識別ユニットを追加する。

【0072】

40

ここで、学習用データは、第1の実施形態と同様に、エッジ検出、物体輪郭検出、物体検出の3つのタスクの正解ラベルが各学習用画像に付けられた物を用いても良い。一方、エッジ検出、物体輪郭検出、物体検出のタスク毎に別々に学習用画像と正解ラベルとを用意しても良い。前者の場合、S802では、第1の実施形態と同様に勾配を混合させ、パラメータの最適化を実行すれば良い。後者の場合、S802において、非特許文献5に開示された方法と同様にタスク毎にロス関数を切り替え、最適化すれば良い。

【0073】

また、後者の場合、同じタスクの公開データセットを利用しても良い。例えば物体検出ではLarge Scale Visual Recognition Challenge (ILSVRC)、物体輪郭検出では、Berkeley Segmentation Dataset (BSD500)である。また、Microsoft（登録商標）

50

COCOでは、物体ラベルとその物体領域との両方のアノテーションが付与されているため、物体検出タスクと物体輪郭検出タスクとの両方を同時に学習させるためにこのデータセットを用いても良い。また、エッジ検出では、正解ラベルを人間が与えて学習用データを作成しても良いが、Sobel法やCanny法のような既存アルゴリズムの出力を正解ラベルとして利用し、それらを模倣するようにネットワークを学習させても良い。

【0074】

以上のように本実施形態によれば、エッジ検出、物体輪郭検出、物体検出という空間的スケールの異なるマルチタスクを行うCNNのニューラルネットワークを実現できる。これにより、画像に対しエッジ検出、物体輪郭検出、イベントセグメント検出の各タスクが一度に処理できるようになる。

【0075】

(第3の実施形態)

本実施形態では、問題の複雑さの尺度が異なるマルチタスクによりカメラの撮像制御および現像用の画像認識を実現するCNNを例に説明する。本実施形態で扱う識別タスクは、具体的には合焦判定、笑顔検出、色温度推定、シーン分類とする。

【0076】

ここで、合焦判定は、カメラのコントラスト方式のオートフォーカス制御で行う合焦判定であり、低次の特徴で実現できる。笑顔検出は、一般に口角の形に基づいて行われるため、中間程度のスケールの特徴が有効であると考えられる。色温度推定は、画像撮影環境における光源の色味による被写体の色変化を軽減させるホワイトバランス補正に利用される。色温度は無彩色領域における色の偏りを元に推定されるが、プリミティブな特徴である色と、画像中の無彩色領域を見つけるより高次の特徴とが複合的に有効なことが考えられるので、複数の中間層の出力を元に識別する識別ユニットを考える。シーン分類では、「ポートレート」、「記念撮影」、「風景」、「夜景」、「市街地」、「マクロ」等のセンサ感度や絞り・シャッター速度等の制御パラメータを決める際の参考になる、撮影シーンの分類を行う。これは画像全体の情報が必要になるため最終層で実行する。以上から各タスクの空間的スケールは、合焦判定、笑顔判定、シーン分類の順に大きくなり、色温度推定は複数のスケールの特徴量を用いる。

【0077】

図4(c)は、本実施形態の画像認識装置450の機能構成例を示すブロック図である。なお、本実施形態に係る画像認識装置450のハードウェア構成については図3と同様である。

本実施形態の画像認識装置450は、画像取得部451、前処理部452、第1中間特徴量取得部453、第2中間特徴量取得部454、第3中間特徴量取得部455、合焦判定部456、笑顔検出部457、色温度推定部458、シーン分類部459を有する。画像認識装置450が有するこれらの機能の詳細については、図5(c)等を用いて後述する。

【0078】

図5(c)は、本実施形態における画像認識装置450による処理手順の一例を示すフローチャートである。本実施形態では、第2の実施形態と同様に図14に示したような特徴抽出ユニットおよび識別ユニットを採用し、基本的な処理の流れは第2の実施形態と同様であることから、図16を参照しながら第2の実施形態との差異を中心に説明する。

【0079】

まず、S551において、画像取得部451は、撮像デバイスから画像を取得する。次に、S552においては、前処理部452は、S551で取得した画像に対する前処理を実行する。この処理は第2の実施形態で説明したS532と同様である。

【0080】

次に、S553においては、第1中間特徴量取得部453は、前処理がなされた入力画像から第1中間特徴量を取得する。ここで、第1中間特徴量取得部453は第2の実施形態と同様に図16に示すようなConv×2+Poolの特徴抽出ユニット1602によ

10

20

30

40

50

り実現され、合焦判定に係る特徴量を入力画像 1601 から取得する。そして、S554 においては、合焦判定部 456 は、第 1 中間特徴量に基づいて画像が合焦か否かを判定する。この判定は 2 値分類であるが、第 2 の実施形態でのエッジ検出と異なり、特徴抽出ユニット 1602 のレスポンスマップ全体と識別ユニットの FC 層とが全結合され、識別ユニット 1605 で画像全体に対し 2 値の合焦か非合焦かが判定される。

【0081】

次に、S555 においては、第 2 中間特徴量取得部 454 は、第 1 中間特徴量に基づき、第 2 中間特徴量を取得する。ここで、第 2 中間特徴量取得部 454 は第 2 の実施形態と同様に $\text{Conv} \times 3 + \text{Pool}$ の特徴抽出ユニット 1603 により実現され、笑顔検出に係る特徴量を取得する。そして、S556 においては、笑顔検出部 457 は、第 2 実施形態での物体輪郭検出と同様に、2 つめの特徴抽出ユニット 1603 のレスポンスマップの各点で、識別ユニット 1606 により笑顔の有無を 2 値判定する。

【0082】

次に、S557 においては、第 3 中間特徴量取得部 455 は、第 2 中間特徴量に基づき、第 3 中間特徴量を取得する。ここで、第 3 中間特徴量取得部 455 は第 2 の実施形態と同様に $\text{Conv} \times 4 + \text{Pool}$ の特徴抽出ユニット 1604 により実現され、撮影シーンに係る特徴量を取得する。

【0083】

S558 においては、色温度推定部 458 は、第 1 中間特徴量および第 3 中間特徴量に基づき色温度の推定を行う。本工程に係る識別ユニット 1608 は複数の層から入力を受け取るため、FC 層が 1 つめの特徴抽出層と 3 つめの特徴抽出層と全結合している。また色温度は連続値であるため、分類ではなく回帰を行い、最終層は線形写像とする。

【0084】

次に S559 においては、シーン分類部 459 は、第 3 中間特徴量から画像の撮影シーンを分類する。この処理では、第 2 の実施形態における物体検出と同様に、画像全体に対してマルチクラスの分類を行う。そのため最終層は Softmax の識別ユニット 1607 を用いる。

【0085】

以上の認識時の処理により、図 16 に示すようなネットワーク構造によって、入力画像に対し、合焦判定、笑顔検出、色温度推定、シーン分類の 4 つのタスクを、1 つのネットワークの処理で同時に実現できる。次に、本実施形態における S553 ~ S559 で利用する特徴抽出ユニットおよび識別ユニットの学習方法と、図 16 に示したネットワーク構造の探索方法に関し、第 2 の実施形態と差異の部分について説明する。

【0086】

ネットワーク構造の探索では、まず複数の中間層の出力を用いる色温度推定は除き、合焦判定、笑顔検出、シーン分類の 3 つのタスクについて、第 2 の実施形態と同様の方法を用いる。その後、色温度推定の識別層のパラメータのみを学習パラメータとして、色温度推定がより高精度になる中間層を探索する。その際は、ネットワーク構造の探索の結果、全体で n 個の中間ユニットからなる構造になった場合、 n 個中の 2 つの組み合わせを探索すれば 2 つの中間層の入力を用いるパターンの全ての組み合わせを網羅できる。また、色温度推定は回帰の問題となるため、識別ユニットの最終層は線形写像になる。そして 2 乗誤差のロス関数を用いてパラメータの最適化を行う。その他の学習方法に関しては第 2 の実施形態と同様の処理で実現できるため説明を省略する。

【0087】

以上のように本実施形態によれば、合焦判定、笑顔検出、色温度推定、シーン分類というカメラ制御時に有用なマルチタスクを行う CNN のニューラルネットワークを実現できる。これにより、認識時の処理によって各タスクが一度に処理できるようになる。

【0088】

(その他の実施形態)

本発明は、上述の実施形態の 1 以上の機能を実現するプログラムを、ネットワーク又は

10

20

30

40

50

記憶媒体を介してシステム又は装置に供給し、そのシステム又は装置のコンピュータにおける１つ以上のプロセッサがプログラムを読み出し実行する処理でも実現可能である。また、１以上の機能を実現する回路（例えば、ＡＳＩＣ）によっても実現可能である。

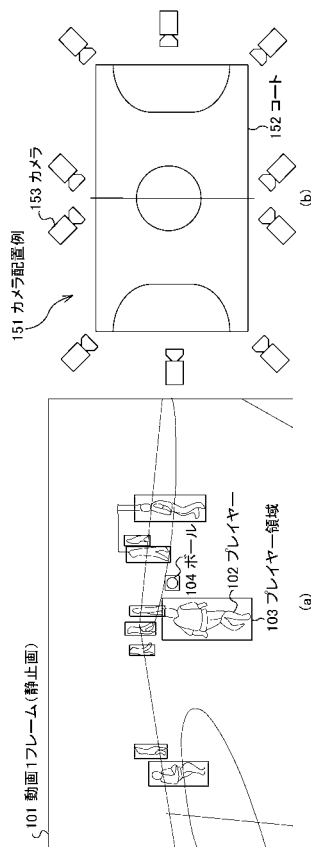
【符号の説明】

【００８９】

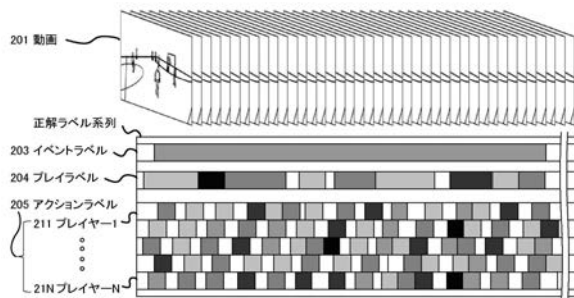
- ４０５ 第１中間特徴量取得部
- ４０６ 第２中間特徴量取得部
- ４０７ 第３中間特徴量取得部
- ４０８ アクション識別部
- ４０９ プレイ識別部
- ４１０ イベントセグメント識別部

10

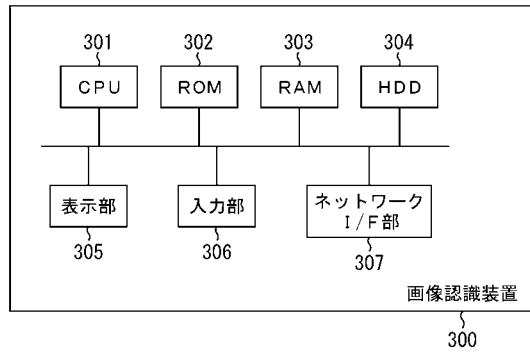
【図１】



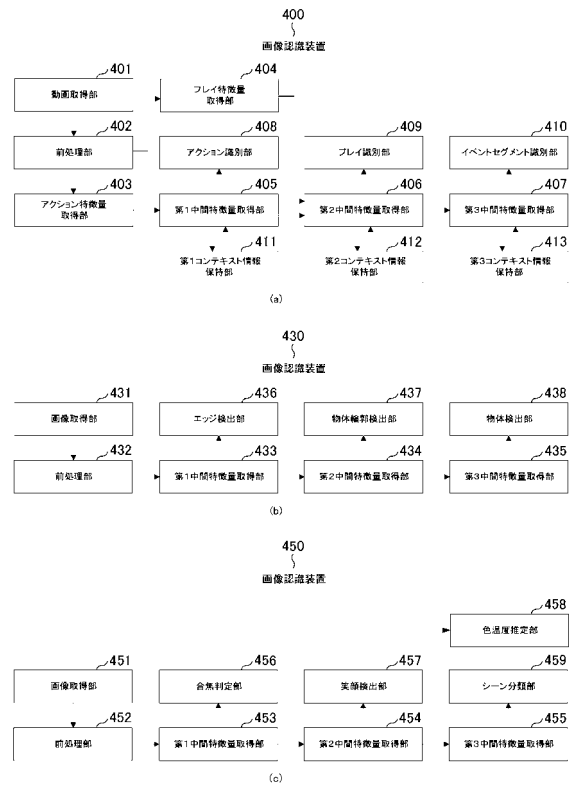
【図２】



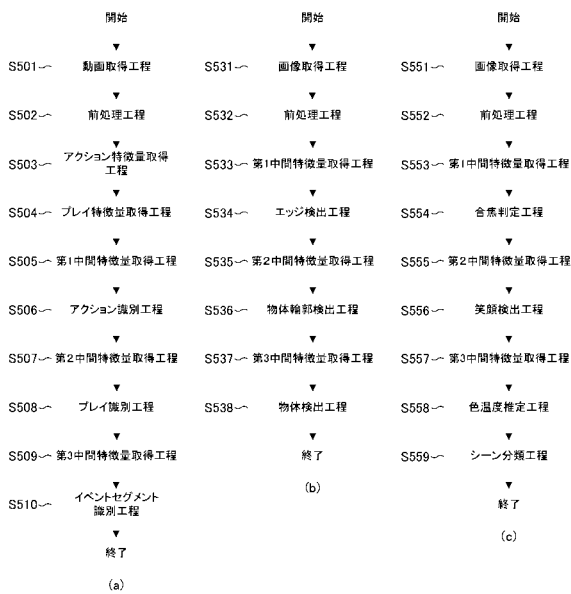
【図 3】



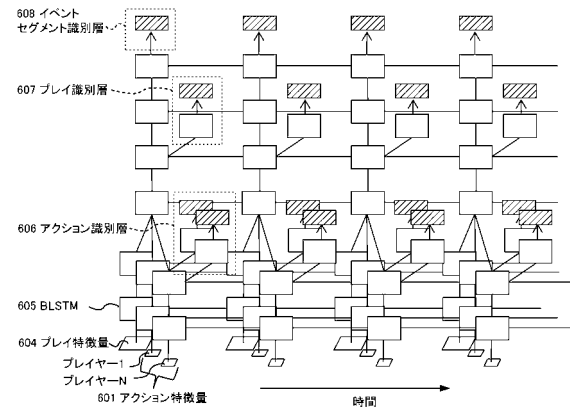
【図 4】



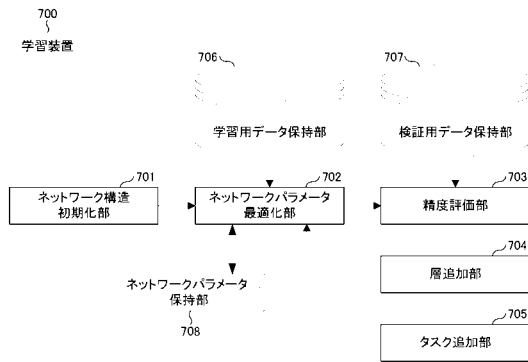
【図 5】



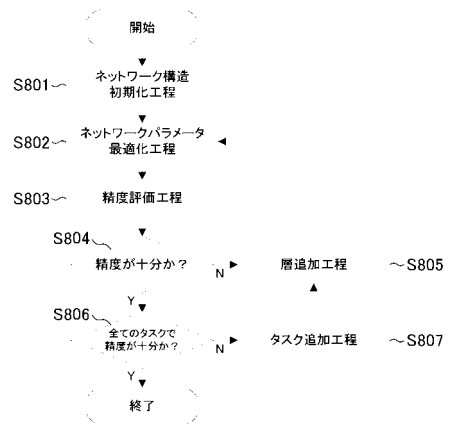
【図 6】



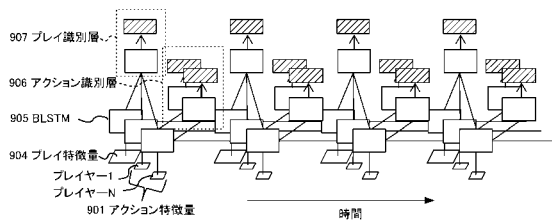
【図 7】



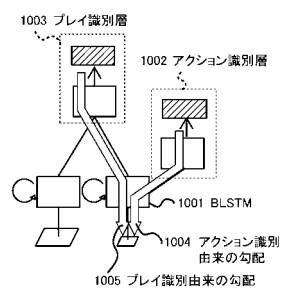
【図 8】



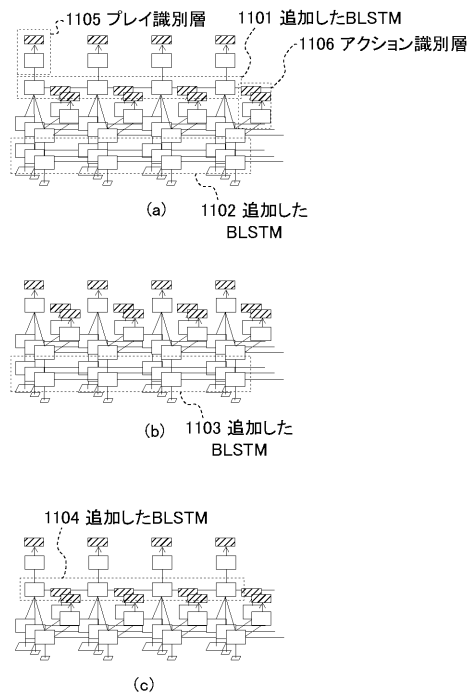
【図 9】



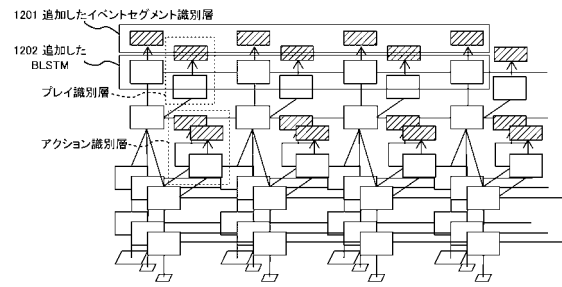
【図 10】



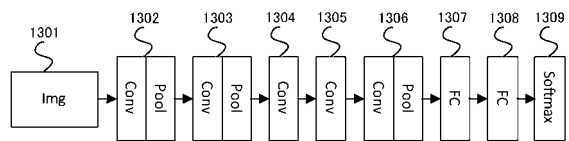
【図 1 1】



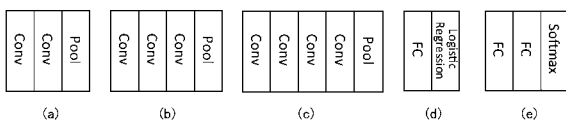
【図 1 2】



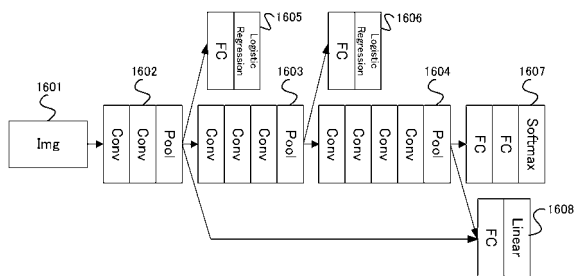
【図 1 3】



【図 1 4】



【図 1 6】



【図 1 5】

