

【特許請求の範囲】**【請求項 1】**

認識対象の情報を含む認識対象データを入力として畳み込みニューラルネットワークの認識処理の結果を取得する認識処理部と、

前記認識対象データおよび / または前記畳み込みニューラルネットワークの中間層出力に対する注目領域を設定する注目領域設定部と、

前記注目領域中の前記認識対象データおよび / または前記中間層出力に対して前記認識処理よりも詳細な詳細認識処理を行う詳細認識部と、

前記詳細認識処理の結果と、前記中間層出力と、を統合処理する統合処理部と、

前記統合処理の結果を前記中間層出力として前記畳み込みニューラルネットワークに入力する中間入力処理部と、

前記認識処理の結果を出力する出力部と

を有することを特徴とする情報処理装置。

【請求項 2】

前記詳細認識部が有する認識器が前記畳み込みニューラルネットワークであることを特徴とする請求項 1 に記載の情報処理装置。

【請求項 3】

前記詳細認識部は少なくとも二つの認識器を備えることを特徴とする請求項 1 又は 2 に記載の情報処理装置。

【請求項 4】

前記注目領域設定部は、前記認識対象の一部に基づいて前記注目領域を設定し、

前記詳細認識部は、前記一部の種類ごとに認識器を備えることを特徴とする請求項 1 乃至 3 の何れか 1 項に記載の情報処理装置。

【請求項 5】

前記一部は少なくとも人体の部位に対応する情報を含むことを特徴とする請求項 4 に記載の情報処理装置。

【請求項 6】

前記認識対象データは画像を含み、

前記統合処理部は、複数の前記一部の前記画像における位置に基づいて前記統合処理を行うことを特徴とする請求項 4 又は 5 に記載の情報処理装置。

【請求項 7】

前記認識処理部は、前記認識対象データに含まれる物体または現象に関する異常な状態を検知することを特徴とする請求項 1 乃至 6 の何れか 1 項に記載の情報処理装置。

【請求項 8】

前記認識処理部および前記詳細認識部は、マルチタスク学習に基づいて前記認識処理を行うことを特徴とする請求項 1 乃至 7 の何れか 1 項に記載の情報処理装置。

【請求項 9】

前記統合処理部は、前記中間層出力に対して正規化処理を行うことを特徴とする請求項 1 乃至 8 の何れか 1 項に記載の情報処理装置。

【請求項 10】

前記一部は、あらかじめ決められた複数の候補の中から手動または自動で選択されることを特徴とする請求項 4 乃至 6 の何れか 1 項に記載の情報処理装置。

【請求項 11】

前記一部は、前記認識対象の一部を検出する検出器の検出結果に基づいて動的に決定されることを特徴とする請求項 4 乃至 6 の何れか 1 項に記載の情報処理装置。

【請求項 12】

前記注目領域は、少なくとも人体の部位の位置を検出する検出器の検出結果に基づいて決定される領域を含むことを特徴とする請求項 1 乃至 11 の何れか 1 項に記載の情報処理装置。

【請求項 13】

前記注目領域は、少なくとも物体の位置を検出する検出器の検出結果に基づいて決定される領域を含むことを特徴とする請求項 1 乃至 11 の何れか 1 項に記載の情報処理装置。

【請求項 14】

学習対象の情報を含む学習データを入力として畳み込みニューラルネットワークの認識処理の結果を取得する認識処理部と、

前記学習データおよび / または前記畳み込みニューラルネットワークの中間層出力に対する注目領域を設定する注目領域設定部と、

前記注目領域中の前記学習データおよび / または前記中間層出力に対して前記認識処理よりも詳細な詳細認識処理を行う詳細認識部と、

前記詳細認識処理の結果と、前記中間層出力と、を統合処理する統合処理部と、

前記統合処理の結果を前記中間層出力として前記畳み込みニューラルネットワークに入力する中間入力処理部と、

前記認識処理部、前記注目領域設定部、前記詳細認識部、前記統合処理部、前記中間入力処理部、のいずれか一つ以上に関する学習処理を行う学習部と

を備えることを特徴とする情報処理装置。

【請求項 15】

前記詳細認識部が有する認識器が前記畳み込みニューラルネットワークであることを特徴とする請求項 14 に記載の情報処理装置。

【請求項 16】

前記詳細認識部は少なくとも二つの認識器を備えることを特徴とする請求項 14 又は 15 に記載の情報処理装置。

【請求項 17】

前記注目領域設定部は、前記学習対象の一部に基づいて前記注目領域を設定し、

前記詳細認識部は、前記一部の種類ごとに認識器を備えることを特徴とする請求項 14 乃至 16 の何れか 1 項に記載の情報処理装置。

【請求項 18】

前記一部は少なくとも人体の部位に対応する情報を含むことを特徴とする請求項 17 に記載の情報処理装置。

【請求項 19】

前記学習データは画像を含み、

前記統合処理部は、複数の前記一部の前記画像における位置に基づいて前記統合処理を行うことを特徴とする請求項 17 又は 18 に記載の情報処理装置。

【請求項 20】

前記認識処理部は、前記学習データに含まれる物体または現象に関する異常な状態を検知することを特徴とする請求項 14 乃至 19 の何れか 1 項に記載の情報処理装置。

【請求項 21】

前記認識処理部および前記詳細認識部は、マルチタスク学習に基づいて前記認識処理を行うことを特徴とする請求項 14 乃至 20 の何れか 1 項に記載の情報処理装置。

【請求項 22】

前記統合処理部は、前記中間層出力に対して正規化処理を行うことを特徴とする請求項 14 乃至 21 の何れか 1 項に記載の情報処理装置。

【請求項 23】

前記一部は、あらかじめ決められた複数の候補の中から手動または自動で選択されることを特徴とする請求項 17 乃至 19 の何れか 1 項に記載の情報処理装置。

【請求項 24】

前記一部は、前記学習対象の一部を検出する検出器の検出結果に基づいて動的に決定されることを特徴とする請求項 17 乃至 19 の何れか 1 項に記載の情報処理装置。

【請求項 25】

前記注目領域は、少なくとも人体の部位の位置を検出する検出器の検出結果に基づいて決定される領域を含むことを特徴とする請求項 14 乃至 24 の何れか 1 項に記載の情報処

10

20

30

40

50

理装置。

【請求項 26】

前記注目領域は、少なくとも物体の位置を検出する検出器の検出結果に基づいて決定される領域を含むことを特徴とする請求項 14 乃至 24 の何れか 1 項に記載の情報処理装置。

【請求項 27】

前記統合処理は、少なくとも加法演算を含むことを特徴とする請求項 1 乃至 26 の何れか 1 項に記載の情報処理装置。

【請求項 28】

前記加法演算は、注目領域に基づいて部分てきに加法演算が適用されることを特徴とする請求項 27 記載の情報処理装置。

【請求項 29】

情報処理装置が行う情報処理方法であって、

前記情報処理装置の認識処理部が、認識対象の情報を含む認識対象データを入力として畳み込みニューラルネットワークの認識処理の結果を取得する工程と、

前記情報処理装置の注目領域設定部が、前記認識対象データおよび / または前記畳み込みニューラルネットワークの中間層出力に対する注目領域を設定する工程と、

前記情報処理装置の詳細認識部が、前記注目領域中の前記認識対象データおよび / または前記中間層出力に対して前記認識処理よりも詳細な詳細認識処理を行う工程と、

前記情報処理装置の統合処理部が、前記詳細認識処理の結果と、前記中間層出力と、を統合処理する工程と、

前記情報処理装置の中間入力処理部が、前記統合処理の結果を前記中間層出力として前記畳み込みニューラルネットワークに入力する工程と、

前記情報処理装置の出力部が、前記認識処理の結果を出力する工程と

を有することを特徴とする情報処理方法。

【請求項 30】

情報処理装置が行う情報処理方法であって、

前記情報処理装置の認識処理部が、学習対象の情報を含む学習データを入力として畳み込みニューラルネットワークの認識処理の結果を取得する工程と、

前記情報処理装置の注目領域設定部が、前記学習データおよび / または前記畳み込みニューラルネットワークの中間層出力に対する注目領域を設定する工程と、

前記情報処理装置の詳細認識部が、前記注目領域中の前記学習データおよび / または前記中間層出力に対して前記認識処理よりも詳細な詳細認識処理を行う工程と、

前記情報処理装置の統合処理部が、前記詳細認識処理の結果と、前記中間層出力と、を統合処理する工程と、

前記情報処理装置の中間入力処理部が、前記統合処理の結果を前記中間層出力として前記畳み込みニューラルネットワークに入力する工程と、

前記情報処理装置の学習部が、前記認識処理、前記注目領域の設定、前記詳細認識処理、前記統合処理、前記入力、のいずれか一つ以上に関する学習処理を行う工程と

を備えることを特徴とする情報処理方法。

【請求項 31】

コンピュータを、請求項 1 乃至 28 の何れか 1 項に記載の情報処理装置の各部として機能させるためのコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、認識対象を認識するための技術に関するものである。

【背景技術】

【0002】

映像等のデータから、例えば物体およびその状態を認識するために、学習データを用い

10

20

30

40

50

て認識器を学習し、該認識器を用いて該認識を行う装置や方法が知られている。このような方法の一つとして、Convolutional Neural Network (CNN) がある。

【0003】

CNNは、物体認識・行動認識・シーン認識など、様々な応用を目的として近年用いられている。行動認識については、例えば非特許文献1では、それぞれRGB画像とオプティカルフローとを入力する二つのストリームからなるCNNを用いて行動認識を行うアーキテクチャが提案されている。また詳細な行動を認識するために、非特許文献2では、人体パーツの局所領域ごとにCNNを用いて特徴量を抽出する方法に関する技術が提案されている。

10

【先行技術文献】

【非特許文献】

【0004】

【非特許文献1】Two-Stream Convolutional Networks for Action Recognition in Videos. K. Simonyan, A. Zisserman. NIPS, 2014

【非特許文献2】P-CNN: Pose-based CNN Features for Action Recognition. Guilhem Cheron, Ivan Laptev, and Cordelia Schmid. ICCV, 2015

【発明の概要】

【発明が解決しようとする課題】

【0005】

入力データの詳細な認識をより高精度に行うために、認識能力がさらに向上した認識システムが望まれる。本発明はこのような点に鑑みてなされたものであり、認識対象の認識精度を向上させるための技術を提供する。

20

【課題を解決するための手段】

【0006】

本発明の一様態は、認識対象の情報を含む認識対象データを入力として畳み込みニューラルネットワークの認識処理の結果を取得する認識処理部と、前記認識対象データおよび/または前記畳み込みニューラルネットワークの中間層出力に対する注目領域を設定する注目領域設定部と、前記注目領域中の前記認識対象データおよび/または前記中間層出力に対して前記認識処理よりも詳細な詳細認識処理を行う詳細認識部と、前記詳細認識処理の結果と、前記中間層出力と、を統合処理する統合処理部と、前記統合処理の結果を前記中間層出力として前記畳み込みニューラルネットワークに入力する中間入力処理部と、前記認識処理の結果を出力する出力部とを有することを特徴とする。

30

【発明の効果】

【0007】

本発明の構成により、認識対象の認識精度を向上させることができる。

【図面の簡単な説明】

【0008】

【図1】異常検知システムの構成例を示すブロック図である。

【図2】異常検知システムの構成例を示すブロック図である。

40

【図3】異常検知システムの構成例を示すブロック図である。

【図4】異常検知システムの構成例を示すブロック図である。

【図5】異常検知システムが行う処理のフローチャート。

【図6】ステップS501における処理の詳細を示すフローチャート。

【図7】異常検知システムが行う処理のフローチャート。

【図8】前段処理部11の構成例を示すブロック図。

【図9】前段処理部11が行う処理のフローチャート。

【図10】1フレームの画像の一例を示す図。

【図11】認識部12の構成例を示すブロック図。

【図12】認識部12が行う処理のフローチャート。

50

【図 1 3】CNN の一例を示す図。

【図 1 4】注目領域制御部 2 1 の構成例を示すブロック図。

【図 1 5】認識部 2 2 の構成例を示すブロック図。

【図 1 6】実行部 2 2 2 の構成例を示すブロック図。

【図 1 7】出力部 2 3 の構成例を示すブロック図。

【図 1 8】学習時の動作のフローチャート。

【図 1 9】第 1 の CNN 及び第 2 の CNN の構成とデータフローに関する模式図。

【図 2 0】パーツのデータの形成例を示す図。

【図 2 1】第 1 の CNN 及び第 2 の CNN の構成とデータフローに関する模式図。

【図 2 2】認識部 2 2 a の構成例を示すブロック図。

【図 2 3】コンピュータ装置のハードウェア構成例を示すブロック図。

【図 2 4】第 1 の CNN 及び第 2 の CNN の構成とデータフローに関する模式図。

【図 2 5】第 1 の CNN 及び第 2 の CNN の構成とデータフローに関する模式図。

【図 2 6】第 1 の CNN 及び第 2 の CNN 及び後段のモデルの構成とデータフローに関する模式図。

【図 2 7】第 1 の CNN 及び第 2 の CNN のデータフローに基づいて得られる中間処理結果の可視化例。

【図 2 8】第 1 の CNN 及び第 2 の CNN の構成とデータフローに関する模式図。

【発明を実施するための形態】

【0009】

以下、添付図面を参照し、本発明の実施形態について説明する。なお、以下説明する実施形態は、本発明を具体的に実施した場合の一例を示すもので、特許請求の範囲に記載した構成の具体的な実施例の 1 つである。

【0010】

[第 1 の実施形態]

本実施形態では、認識器としての Convolutional Neural Network (CNN) を学習する例を示す。なお、具体例は後述するが、本実施形態では、映像中の対象人物の行動に関する認識を行う CNN と、さらに詳細な認識を行う CNN (詳細 CNN) と、を用いる場合の構成及び動作の例を示す。

【0011】

本実施形態では認識処理を行う際に、通常の特徴量だけでなく、認識に寄与しうる注目領域の詳細な特徴量 (詳細特徴量) をも抽出し、該抽出した詳細特徴量を通常の特徴量に順次統合しながら認識処理を行う。具体的には、通常の特徴を認識する第 1 の CNN の中間層特徴マップに、詳細な特徴を認識する第 2 の CNN の特徴量 (詳細特徴量) を統合する。このとき、第 1 の CNN の中間層特徴マップ上の注目領域に対応する領域に対して詳細特徴量を統合することで、詳細特徴量の統合を実現することができる (ROI Embedding)。ここで「統合」とは、一つ以上の詳細特徴量を第 1 の CNN の中間層特徴マップに統合することを指し、例えば足し合わせや、代入、非線形変換等を通して、中間層特徴マップ上に詳細特徴量を反映させることを指す。「ROI Embedding」は、時々刻々と変化する注目領域の位置・大きさ・形等に対応して、第 1 の CNN の特徴量と詳細特徴量との統合処理のためのデータフローを変更し、全体のニューラルネットワーク (NN) の構造を変化させる、という特性がある。特に、本実施形態では、注目領域を人体の部位 (パーツ) として、各パーツの詳細な特徴量を抽出し、全身の特徴を認識する第 1 の CNN の中間層特徴マップ上に順次統合処理する例を示す。

【0012】

本実施形態では、監視カメラなどの撮像装置により撮影された映像中の物体や現象に係る異常を検知するためのシステムである異常検知システムを例にとり説明する。本実施形態に係る異常検知システムでは、監視対象を撮像装置で撮影し、該撮影により得られる映像に基づいて、監視対象に異常があるか否かを判定する。そして異常検知システムは、監視対象に異常があると判定した場合は、警備室等の監視センタに常駐する監視者に警告す

10

20

30

40

50

るための処理を行う。この監視対象には、例えば、一般家庭の屋内及び屋外、又は病院、駅などの公共施設が含まれる。

【0013】

本実施形態に係る異常検知システムの動作は、「学習時:Forward Stage」、「学習時:Backward Stage」、「学習時:Final Stage」、「検出時」、の4つのステージからなる。これら4つのステージのうち、最初の3ステージは「学習時」の動作を示しており、4つ目のステージは学習結果に基づいて上記の異常検知処理を行う際（「検出時」）の動作を示している。

【0014】

図1は、本実施形態に係る異常検知システムが有する構成のうち、CNNを学習させる動作ステージである「学習時:Forward Stage」に関連する構成の一例を示すブロック図である。認識装置10と詳細認識装置20との間は電子回路を介して接続されても良いし、記憶装置などの装置を介して接続されても良いし、ネットワーク（無線ネットワーク及び有線ネットワークのうち1つ以上により構成される）を介して接続されても良い。このネットワークには、例えば、携帯電話回線網やインターネットが適用できる。なお、以下に説明する装置や機能部についても同様に、装置間や機能部間は何なる接続形態で接続されても良い。

【0015】

図2は、本実施形態に係る異常検知システムが有する構成のうち、CNNを学習させる動作ステージである「学習時:Backward Stage」に関連する構成の一例を示すブロック図である。また、図3は、本実施形態に係る異常検知システムが有する構成のうち、CNNを学習させる動作ステージである「学習時:Final Stage」に関連する構成の一例を示すブロック図である。図3に示した構成では、図1、2に示した構成に異常判定装置30が加わっている。

【0016】

異常検知システムのCNNの学習時における動作について、図5のフローチャートに従って説明する。ステップS501では、認識装置10は、学習データを入力データとして第1のCNNの学習を行い、詳細認識装置20は、学習データと、認識装置10からの出力である「第1のCNNの出力」と、を入力データとして、第2のCNNの学習を行う。ここで、第1のCNNは、認識装置10が使用するCNNである。また、第2のCNNは、詳細認識装置20が使用するCNNであり、第1のCNNの出力等を受け取って詳細な認識処理を行うCNNで、詳しくは後述する。また、第1のCNNは、第2のCNNの（中間の）詳細認識処理結果をさらに受け取って認識処理を行うように動作するものとし、詳しくは後述する。なお、ステップS501の動作は、ステージ「学習時:Forward Stage」における動作と、ステージ「学習時:Backward Stage」における動作と、を含んでいる。各ステージでの詳細な動作については後述する。

【0017】

ステップS502では、異常判定装置30は、認識装置10から認識結果を受け取り、異常検知に係る閾値（以下、異常閾値と称する）を求める。ステップS502の動作は、ステージ「学習時:Final Stage」の動作に対応し、詳細な内容については後述する。

【0018】

ステップS503では、認識装置10による学習処理で得られた第1のCNNのモデル、詳細認識装置20による学習処理で得られた第2のCNNのモデル、異常判定装置30が求めた異常閾値、のそれぞれの保存処理が行われる。第1のCNNのモデル及び第2のCNNのモデルはそれぞれ、認識装置10及び詳細認識装置20によって、記憶部M1に保存され、異常閾値は、異常判定装置30によって記憶部M2に保存される。ここで、CNN（第1のCNN及び第2のCNNのそれぞれ）のモデルとは、CNNを規定するパラメータを保持するデータセットである。CNNのモデルは、CNNにおけるニューロン間の結合荷重を規定するデータ、CNNのネットワーク構造を規定するデータ、ニューロン応答関数を規定するデータ、などを含む。

10

20

30

40

50

【 0 0 1 9 】

図 4 は、本実施形態に係る異常検知システムが有する構成のうち、異常検出時の動作に係る動作ステージである「検出時」に関連する構成の一例を示すブロック図である。図 4 に示した構成では、図 3 に示した構成に、PC (Personal Computer) のディスプレイやタブレット PC、スマートフォン、フューチャーフォン等が適用可能な端末装置 4 0 が加わっている。

【 0 0 2 0 】

異常検知システムのステージ「検出時」における動作について、図 7 のフローチャートに従って説明する。ステップ S 7 0 1 では、認識装置 1 0 及び詳細認識装置 2 0 はそれぞれ、第 1 の CNN を用いた検出対象データの認識処理 (第 1 の CNN の認識処理)、第 2 の CNN を用いた検出対象データの認識処理 (第 2 の CNN の認識処理) を行う。このとき用いられる第 1 の CNN 及び第 2 の CNN は、上記の学習時に得られたモデルを用いるものとするが、異常検知システムの学習時に得られたモデルを用いず、その他のモデルを用いてもよい。

10

【 0 0 2 1 】

ステップ S 7 0 2 では、異常判定装置 3 0 は、検出対象データに対する第 1 の CNN の認識処理の結果に対して異常判定を行う。このとき、異常判定を行う際の閾値として、学習時に得られた異常閾値を用いるものとするが、異常検知システムの学習時に得られた閾値を用いず、その他の閾値 (例えば、ユーザが指定した閾値など) を用いてもよい。このとき、どのように閾値判定を行うかに関しては、後述する。

20

【 0 0 2 2 】

ステップ S 7 0 3 では、端末装置 4 0 は、異常判定装置 3 0 による異常判定の結果を受け取り、該受け取った判定結果に基づく表示処理を行う。この表示処理の詳細については後述する。

【 0 0 2 3 】

次に、各装置の動作に関して説明を行う。ここで示すのは、各装置の動作の大まかな内容と処理の流れの順番である。それぞれの装置を構成する各機能部の詳細な構成・動作に関しては、機能部の説明をする際にあらためて述べるものとする。

【 0 0 2 4 】

まず、学習時の動作に関する説明を行う。学習時の動作の流れはすでに図 5 を用いて述べたとおりであるから、ここでは、図 5 の説明を補完する形で図 6 に基づいて説明を行う。

30

【 0 0 2 5 】

ステップ S 5 0 1 では上記の通り、ステージ「学習時:Forward Stage」及びステージ「学習時:Backward Stage」の学習処理が実行される。ステップ S 5 0 1 の処理の詳細について、図 6 のフローチャートに従って説明する。

【 0 0 2 6 】

ステップ S 6 0 1 では、図 1 の構成により、ステージ「学習時:Forward Stage」の学習処理が実行される。まず前段処理部 1 1 は、記憶部 D 1 から学習データを受け取る。本実施形態では、記憶部 D 1 には学習データとして、監視カメラなどの撮像装置が撮影した人体を含む映像が保存されているものとするので、前段処理部 1 1 は、この映像を学習データとして記憶部 D 1 から取得する。そして前段処理部 1 1 は、映像中の人体の領域 (人体領域) を抽出する等の前処理を行ったのち、学習データを第 1 の CNN 用の学習データとして整形して認識部 1 2 に送る (詳細は後述する)。認識部 1 2 は、前段処理部 1 1 から受け取った学習データに基づいて、第 1 の CNN の認識処理を行う。第 1 の CNN の認識処理とは、第 1 の CNN によってデータを認識するための一連のプロセスのことを指す。なお、本実施形態に係る認識処理は、第 1 の CNN の認識処理と第 2 の CNN の認識処理とが相互に依存する構成になっており、具体例を以降で説明する。認識部 1 2 は、第 1 の CNN の認識処理において得られる該第 1 の CNN における規定の中間層の出力 (中間認識情報) を、注目領域制御部 2 1 に送る。また、前段処理部 1 1 は、認識部 1 2 が受け取

40

50

った学習データ（後述する教師データを含む）を、注目領域制御部 2 1 に送る。注目領域制御部 2 1 は、受け取った学習データ及び第 1 の CNN における規定の中間層の出力から、第 2 の CNN が認識する対象を設定し、該設定した対象を規定する詳細認識対象データ（詳細認識対象情報）を認識部 2 2 に送る。詳細認識対象データの具体的な作成方法については後述する。本実施形態では、詳細認識対象データとは、認識対象となる人物の行動認識を行うために着目すべき人体のパーツに関するデータを保持するものとする。注目領域制御部 2 1 は、認識対象となる人物のパーツの領域（注目領域）を規定する領域データ（注目領域情報）を、出力部 2 3 及び認識部 1 2 に送る。認識部 2 2 は、受け取った詳細認識対象データに対して認識処理を行い、該認識処理の結果（詳細認識情報）を詳細認識処理結果として出力部 2 3 に送る。出力部 2 3 は、認識部 2 2 による認識の結果を認識部 1 2 に出力するために、上記の受け取った領域データに基づいて、詳細認識の結果を認識部 1 2 に送る。すなわち、第 1 の CNN の規定の中間層の出力の領域に対して送り返す形で詳細認識情報を送信することになる（この点に関してはバリエーションがあり、詳細は後述する）。認識部 1 2 は、第 1 の CNN の中間層出力と、出力部 2 3 から受け取った詳細認識情報と、に基づいて、第 1 の CNN の認識処理を続行する。

10

【0027】

その際、さらに第 2 の CNN を利用して詳細な認識処理を行ってもよい。なお認識部 1 2 は、受け取った領域データを、以降のステージ「学習時:Backward Stage」で利用する。最後に認識部 1 2 は、あらかじめ定めた損失関数に基づいて、処理を行う。例えば、学習データに含まれる教師データに基づいて、二乗誤差によって誤差を算出する方法などがありえる。なお、二乗誤差ではなく、その他の公知な損失関数を用いてもよい。

20

【0028】

次に、ステップ S 6 0 2 では、図 2 の構成により、ステージ「学習時:Backward Stage」の学習処理が実行される。ここでは、認識部 1 2 による認識処理の結果に対する評価結果をもとにモデルの更新を行う。学習アルゴリズムには、例えば誤差逆伝播法を用いることができる。誤差逆伝播法は、評価結果の誤差をニューラルネットワークの上層から下層に伝播してゆく形で学習を行う手法であり、詳細は後述する。まず、認識部 1 2 は、誤差情報をもとにモデルの更新を行う。また、第 2 の CNN の出力が入力された第 1 の CNN の領域に関する誤差情報を、第 2 の CNN に対して伝播する。このとき、認識部 1 2 は、領域データに基づいて第 1 の CNN の注目領域を特定し、出力部 2 3 に該注目領域の誤差情報を送るものとする。なお、ここで述べる誤差とは、誤差逆伝播法による学習の過程において、第 1 の CNN の上層ネットワークから注目領域に逆伝播された誤差を示す。さらに出力部 2 3 は、認識部 2 2 に誤差情報を送る。このとき、第 2 の CNN のどの部分に誤差情報を送るかどうかを決定する必要がある。ここでは、注目領域に対して詳細情報を出力した第 2 の CNN の出力ニューロンに対して誤差情報を送るものとする。すなわち、ステージ「学習時:Forward Stage」の際に入力データが認識処理された経路の逆順をたどるように、誤差情報が伝播するものとする。認識部 2 2 は、受け取った誤差情報に基づいて第 2 の CNN のモデルを更新する（詳細は後述する）。また、受け取った誤差情報を、さらに認識部 1 2 に対して逆伝播し、認識部 1 2 は、受け取った誤差情報に基づいて第 1 の CNN のモデルを更新する。以上のように、誤差逆伝播法に基づいて第 1 の CNN 及び第 2 の CNN を更新することができる。

30

40

【0029】

次に、ステップ S 6 0 3 では、学習が完了したか否かの判定が行われる。この判定の基準については後述する。学習が完了していないと判定された場合は、ステップ S 6 0 1 に戻り、学習が完了したと判定された場合は、図 6 の動作は終了し、処理はステップ S 5 0 2 に移行する。

【0030】

図 5 に戻って、次にステップ S 5 0 2 では、ステージ「学習時:Final Stage」として、ステップ S 6 0 1 及びステップ S 6 0 2 において学習した第 1 の CNN 及び第 2 の CNN の認識結果に基づいて、異常検知をするための閾値が決定される。認識部 1 2 は、学習デ

50

ータに対する認識結果を異常判定部 3 1 に送る。異常判定部 3 1 は、認識部 1 2 の認識結果に基づいて、認識対象が異常かどうかを判定するための基準である閾値を決定する。このときの閾値の決定方法には、公知の既存の手法を用いればよい。例えば、第 1 の CNN が認識対象の行動を認識するように学習が行われており、認識対象が「転倒」したか否かを認識することができるとする。このとき、「転倒」クラスが異常と判定すべきクラスであるものとして、認識対象の行動が「転倒」クラスであることを示す確率が 50 % 以上であつたら、異常であると判定するとしてもよい。このとき「50 %」が異常閾値である。このようにあらかじめ固定の閾値を与える方法以外のアプローチとして、閾値を評価用データによって決定してもよい。具体的には、学習データセットとは異なる評価用データセット（例えば映像）を用意する。このとき、学習済みの第 1 の CNN を用いて各評価データに対する「転倒」の確率を求めることで、例えば真に「転倒」であるすべての評価データを「転倒」であると判定できる閾値を決定することができ、このような閾値を異常閾値として設定してもよい。その他の方法として、第 1 の CNN の特徴量に基づいて正常の範囲のモデル（正常モデル）をさらに学習し、学習されたモデル（正常モデル）に基づいて正常か否かを判定してもよい。このとき用いる第 1 の CNN の特徴量は、例えば最終層より一層前の層の出力（中間層出力）を特徴量として用いてもよい。また、このとき用いる正常モデルを学習する方法には如何なる方法を採用しても良く、例えば特開2014-203289号公報に記載の方法を用いることができる。ここで決定した閾値や、正常モデルなど、異常の判定に用いる特徴量は、記憶部 M 2 に格納される。

10

20

30

40

50

【0031】

次に、異常検知システムのステージ「検出時」における動作について、図 7 のフローチャートに従って説明する。ステップ S 7 0 1 では、第 1 の CNN の認識処理、第 2 の CNN の認識処理が行われる。先ず、前段処理部 1 1 は、記憶部 D 2 から検出対象データとして、監視カメラなどの撮像装置により撮影された人体を含む映像を取得する。そして前段処理部 1 1 は、学習時と同様に映像から人体領域を抽出するなどの前処理を行ったのち、検出対象データを第 1 の CNN 用の検出対象データとして整形（詳細は後述する）して、認識部 1 2 に送る。そして認識部 1 2 及び認識部 2 2 はそれぞれ、この検出対象データに基づいて第 1 の CNN の認識処理及び第 2 の CNN の認識処理を行う。この処理の詳細は、上記のステージ「学習時:Forward Stage」と同様であるため、その説明は省略する。

【0032】

次にステップ S 7 0 2 では、異常判定部 3 1 は、検出対象データに対する認識部 1 2 の認識結果に基づいて、異常かどうかの判定を行う。このとき異常判定部 3 1 は、上記のステージ「学習時:Final Stage」で記憶部 M 2 に格納した異常閾値（あるいは正常モデル）を受け取り、異常の判定に用いるものとする。

【0033】

次にステップ S 7 0 3 では、端末装置 4 0 は、異常判定部 3 1 による異常判定結果を受け取り、受け取った結果に基づく表示を表示部 4 1 に表示させるための表示処理を行う。表示部 4 1 は、検出対象が異常であることを示す異常判定結果が送られてきた場合、警告を示す表示を行ってもよい。例えば、監視カメラなどの撮像装置による撮影映像上で異常と判定された検出対象の箇所を強調表示してもよい。このとき、異常な箇所を強調表示するには、映像上の異常な箇所を特定する必要がある。そのためには、例えば入力される検出対象データに画面座標データを付与しておき、異常判定結果に応じてその座標データを利用してよい。

【0034】

なお、異常検知システムは、検出対象が異常である場合における警告動作として警告を示す表示を行うことに加えて若しくは代えて、他の処理を行うようにしても良い。例えば、異常検知システムにランプとサイレンが備わっている場合はサイレンにより警告音を鳴らすと共にランプを点滅させても良い。

【0035】

次に、次の映像が残っている場合には、処理はステップ S 7 0 4 を介してステップ S 7

01に戻る。一方、次の映像が残っていなければ、図7のフローチャートに従った処理は終了する。

【0036】

次に、上記の各装置が有する各機能部の構成及び動作について説明する。図8は、上記の前段処理部11が有する機能部のうち、学習時及び検出時に関連する構成の一例を示すブロック図である。また、図11は、上記の認識部12が有する構成のうち、学習時に関連する構成の一例を示すブロック図である。図14は、上記の注目領域制御部21が有する構成のうち、学習時に関連する構成の一例を示すブロック図である。図15は、上記の認識部22が有する構成のうち、学習時に関連する構成の一例を示すブロック図である。図16は、図15に示す認識部12の実行部222が有する構成において、学習時に関連する構成の一例を示すブロック図である。図17は、上記の出力部23が有する構成のうち、学習時に関連する構成の一例を示すブロック図である。

10

【0037】

次に、前段処理部11(図8)の学習時の動作について、図9のフローチャートに従って説明する。ステップS901では、抽出部111は、記憶部D1から学習データを読み込む。本実施形態では、この学習データは、監視カメラなどの撮像装置により撮影された人体を含む映像を含むとする。抽出部111は、映像から人体領域を抽出する。映像が動画像である場合には、各フレームの画像から人体領域を抽出することになる。ここで、動画像としての映像(監視映像)における、ある1フレームの画像の一例を図10に示す。図10の画像1001は、ある交差点における監視映像における1フレーム分の画像の一例であり、1002~1005は、該画像1001に含まれているオブジェクトである。1006, 1010, 1009, 1007はそれぞれ、オブジェクト1002~1005について抽出部111が抽出したBouding Boxであり、これらのBouding Boxに囲われた部分画像それぞれが人体領域に対応している。なお、ここで示したBouding Boxはあくまでも人体領域が抽出された際の具体例の一つであり、例えば後述する背景差分法によって撮像されたオブジェクトの輪郭に沿った小領域を抽出してもよい。1099は信号を示している。

20

【0038】

映像からこのような人体領域を抽出するための方法は複数存在し、如何なる方法を採用しても良い。このような方法には、例えば背景差分法、物体検出・追尾法、領域分割法の三つがある。

30

【0039】

抽出部111は、人体のように監視対象のオブジェクトがあらかじめ既知である場合は、ターゲットのオブジェクトのみを検出・追尾する目的に絞られた物体検出・追尾法が比較的適していると考えられる。物体検出・追尾法には、例えば以下の文献1に記載の方法がある。

【0040】

文献1 ... Real-Time Tracking via On-line Boosting, H. Grabner, M. Grabner and H. Bischof, Proceedings of the British Machine Conference, pages 6.1-6.10. BMV A Press, September 2006

40

さらに抽出部111は、学習データ(監視映像)にあらかじめ付与された教師データを利用して、人体領域に対して教師データを付与する。例えば、矩形の人体領域を抽出する場合は人体領域のBouding Boxを定義し、該Bouding Boxに教師データを付与することができる。また、例えば人体の輪郭に沿って人体領域を抽出する場合は、画像に対する人体領域のマスクを定義し、該マスクに教師データを付与することができる。教師データは対象をどのように分類すべきかを示すラベルである。ラベルの種類や、どのような対象にどのようなラベルを付与するか、ということは問題に依存するため、異常検知システムを使用または導入するユーザがあらかじめ決定し、学習データに対して教師データを付与しておくものとする。ここで、学習データにあらかじめ付与する教師データは、例えば撮像された被写体の画像上の領域に対して、ユーザが付与することができる。より具体的には、

50

例えば歩行者であるオブジェクト 1 0 0 3 の領域をユーザが手動で指定し、その領域に対して例えば歩行者を示すラベルを付与することができる。このとき、抽出部 1 1 1 は、抽出された人体領域と、付与された教師データの領域と、が重畳している場合、その人体領域に対して最も大きな面積の割合で重畳された教師データを付与してよい。なお、必ず上記のやり方で教師データを付与しなければならないわけではなく、例えばあらかじめ人体領域を抽出しておき、それぞれの人体領域に対して、ユーザが教師データを付与してもよい。なお、本実施形態では人体領域を抽出する例を示したが、抽出する領域は人体領域に限らず、他の対象物の領域を抽出するようにしても良い。また、領域を抽出せずに画像全体を認識処理の対象としてもよい。どのような領域を抽出し、認識処理の対象とするかは問題依存であるため、問題に応じて設定する必要がある。最後に、抽出部 1 1 1 は、抽出した人体領域（画像）を前処理部 1 1 3 と検出部 1 1 2 とに送る。なお、説明を簡略化するために、ステップ S 9 0 1 では、異常検知システムで用いられるすべての学習データを一度に受け取り、受け取った学習データから人体領域を抽出し、抽出されたすべての人体領域を以降の処理の対象とする。なお、学習データのデータサイズが非常に大きい場合は、すべてのデータを通信により一度に取得することは難しいため、その場合は一部ずつ学習データを取得するようにしてもよい。以降の説明では、特に注意書きを記載する場合（後述するMinibatch等に関して記述している場合等）を除いては、すべての学習データに対して一度に処理・通信を行うものとする。

10

【0041】

ステップ S 9 0 2 では、検出部 1 1 2 は、抽出部 1 1 1 から人体領域を受け取る。ここで受け取る人体領域の例を、図 1 0 下部に 1 0 1 1 として示す。人体領域 1 0 1 1 における 1 0 0 3 は、図 1 0 上部の画像 1 0 0 1 中のオブジェクト（人物）1 0 0 3 と同一の人物であることを示している。1 0 1 2 は頭部領域、1 0 1 3 は右手領域、1 0 1 4 は右足領域、1 0 1 5 は左手領域、1 0 1 6 は左足領域であり、それぞれのパーツ領域は、検出部 1 1 2 が検出したパーツ領域の一例である。ここで検出する対象のパーツ領域は、例えば認識すべき問題に応じてあらかじめ決めてもよい。例えば、現象の 1 つである万引き行為を検知するために、手を検出対象のパーツとしてもよい。ここでは、あらかじめ上記 5 つのパーツ領域を検出するものとする。パーツ領域を検出するための方法には如何なる方法を用いても良く、例えば以下の文献 2 に記載の方法がある。

20

【0042】

文献 2 ... Convolutional Pose Machines, Shih-En Wei, Varun Ramakrishna, Takeo Kanade and Yaser Sheikh, CVPR, 2016.

30

なお、ここでは、人体領域を抽出した後にパーツ領域の検出を行う例を示したが、これに限るものではなく、多数の人物のパーツ領域を特定したのち、人体領域を特定するようなプロセスを用いてもよい。

【0043】

図 9 に戻って、次にステップ S 9 0 3 では、検出部 1 1 2 は、ステップ S 9 0 2 において検出部 1 1 2 が人体領域から検出したパーツ領域を抽出し、該抽出したパーツ領域を前処理部 1 1 3 に送る。

【0044】

次にステップ S 9 0 4 では、前処理部 1 1 3 は、認識処理の前処理を行う。本実施形態では、各人体領域（画像）を規定の画像サイズ（例えば 2 2 4 画素 × 2 2 4 画素）に変形（縮小）したのち、平均画像を引く処理を行うものとする。ここで平均画像は、学習データにおける複数枚の画像のそれぞれの画素位置（x、y）の画素値の平均値を平均画像における画素位置（x、y）の画素値とすることで得られる画像である。また、前処理部 1 1 3 は、学習データ（学習画像）を水平方向に反転して水増しし、水増し前の画像と水増し後の画像とを、すべて学習画像として扱うこととする。

40

【0045】

次にステップ S 9 0 5 では、データ選択部 1 1 4 は、学習データの部分集合群（Minibatch集合）を作成する。Minibatchは、第 1 の CNN 及び第 2 の CNN のパラメータを S t

50

ochastic Gradient Descent (SGD、確率的勾配降下法)によって繰り返し最適化する際に用いる学習データの集合である。詳細は後述するが、本実施形態では、CNN等のパラメータを最適化する方法としてSGDを用いる例を示すこととする。なお、Minibatchのサイズは例えば本実施形態では、同じ教師データのラベルを持つ50個の学習画像をひとまとめにしたものとする。また、Minibatchは、学習データが重複なく属するように、かつすべての学習データが含まれるように作成するものとする。

【0046】

次にステップS906では、データ選択部114は、Minibatch集合からランダムに一つMinibatchを選択し、認識部12及び詳細認識装置20に送る。このとき選択するMinibatchは、その学習のIterationの中でまだ使用していないMinibatchを選択するものとする。ここでIterationとは、一般的にはエポックとも言われ、Minibatch集合のすべてのMinibatchを何回学習したかを表すものとする。

10

【0047】

次にステップS907では、データ選択部114は、学習処理の終了判定を行う。学習処理の終了判定の基準には様々な基準を適用することができ、例えば、Iterationが50以上になった場合には終了するとしてもよい。学習処理の終了条件が満たされていない場合には、処理はステップS906に戻り、学習処理の終了条件が満たされている場合には、図9のフローチャートに従った処理は終了する。

20

【0048】

次に、認識部12(図11)の学習時の動作について、図12のフローチャートに従って説明する。ステップS1201では、初期化部121は、第1のCNNのパラメータ(結合加重及びバイアス項)を初期化する。初期化にあたっては、第1のCNNのネットワーク構造をあらかじめ決めておく必要がある。ここで用いるネットワーク構造や初期パラメータは特定のパラメータに限らず、例えば以下の文献3と同じものを用いてよいし、独自に定義したネットワーク構造を用いてもよい。ここで図13に、本実施形態に適用可能なCNNの一例を示す。

【0049】

文献3 ... A. Krizhevsky et al. "ImageNet Classification with Deep Convolutional Neural Networks", Advances in Neural Information Processing Systems 25 (NIPS), 2012.

30

図13のニューラルネットワーク1320は、本実施形態に係る第1のCNNのネットワーク構造の一例を示したものである。図13に示す如く、ニューラルネットワーク1320は、入力層1301、convolution層1302、pooling層1303、convolution層1304、pooling層1305、Inner product層1306、Inner product層1307、出力層1308、を有している。また、2つの階層間の処理方法として、convolution処理1310、pooling処理1311、Inner Product処理1312が設定されている。それぞれの処理の具体的な内容は周知であり、例えば上記の文献3と同様であるため、ここでは省略する。

40

【0050】

convolution処理1310では、畳み込みフィルタを用いてデータ処理を実行し、pooling処理1311では、例えば、max poolingであれば局所的な最大値を出力する処理を行う。また、Inner Product処理1312では、内積処理を実行する。

【0051】

ここで、畳み込みフィルタは、学習後の各ニューロンの受容野に該当する結合加重(結合重み係数)の分布で定義される。また、図13では、convolution層およびpooling層には複数の特徴マップ(中間層特徴マップ)が存在し、入力層の画像上のピクセルに対応する位置には、複数のニューロンが存在する。例えば学習画像がRGB

50

形式である場合、RGBチャンネルに対応する3つのニューロンが存在する。また、撮像された映像の動き情報を持つOpticalFlow画像であれば、画像の横軸方向と縦軸方向とをそれぞれ表現する2種類のニューロンが存在することになる。また、複数の画像を同時に入力として用いる場合は、入力画像の数に対応する分だけ入力層のニューロンを増やすことで対応することが可能である。本実施形態では、標準的なRGB画像を対象とする例を示すものとする。

【0052】

なお、上記のステップS1201における初期化の方法にはバリエーションが存在する。初期化部121は、公知の手法で第1のCNNのパラメータを初期化することができる。初期化のバリエーションには大きく分けて、データを用いて初期化する場合と、データを用いないで初期化する場合と、がある。ここでは、データを用いないで初期化する場合の簡単な方法として、平均0分散1の正規分布からランダムにサンプリングした値を用いて重みパラメータを初期化し、さらにバイアス項パラメータはすべて0で初期化するものとする。なお、データを用いて初期化する場合は、例えば画像認識などを用途としてあらかじめ学習（プレトレーニング）した第1のCNNを、本実施形態に係る第1のCNNの初期モデル（初期化された結合重みパラメータなど）として用いてもよい。そして初期化部121は、上記のようにして初期化した第1のCNNのモデルを実行部122に送る。

10

【0053】

次にステップS1202では、実行部122は、変数Iterationの値を0に初期化する。次にステップS1203では、実行部122は、学習用のMinibatchを前段処理部11から取得する。

20

【0054】

次にステップS1204では実行部122は、第1のCNNの認識処理を実行する。このとき、第1のCNNの認識と第2のCNNの認識とを連携して行うための処理を実行する（詳細は後述する）。

【0055】

次にステップS1205ではCNN更新部123は、ステップS1204で得られた認識処理結果に基づいて、第1のCNNのモデルを更新するための処理を行うために、認識誤差および誤差情報の伝播を行う。本実施形態では、第1のCNNの学習を行うための方法として、誤差逆伝播法とSGDとを組み合わせた方法を用いる。誤差逆伝播法とSGDとを組み合わせた方法は周知であり、例えば、上記の文献3に記載の通り、Minibatchを選択し、第1のCNNのパラメータを逐次更新するという手順を繰り返す方法である。

30

【0056】

なお、誤差情報は認識処理のデータフローとは逆方向に伝播するのが一般的であり、ここでも第2のCNNに対して誤差情報が伝播するものとする。ここで本実施形態においては、誤差逆伝播処理のためのデータフローは動的に変わりうることに注意が必要である。これについての詳細な説明は後述する。

【0057】

次にステップS1206では、CNN更新部123は、ステップS1205で伝播した誤差情報を用いて、第1のCNNの更新を行う。次にステップS1207では、実行部122は、Minibatchをすべて学習に利用したか否かの判定を行う。すべてのMinibatchを学習に利用した場合は、ステップS1208に移行し、未だ学習に利用していないMinibatchが残っているのであれば、処理はステップS1203に移行する。

40

【0058】

ステップS1208では、実行部122は、変数Iterationの値に1を加算する。ステップS1209では、実行部122は、変数Iterationの値（Iteration回数）が予め設定した上限値に達しているか否かを判定する。この判定の結果、Iteration回数が予め設定した上限値に達している場合は、第1のCNNの学習を終了し、処理はステップS1210に移行する。一方、Iteration回数が予め設定した上限値に達していない場合は、処理はステップS1203に移行する。

50

【 0 0 5 9 】

なお多くの場合、NNの学習停止条件は、「Iteration回数が予め設定した上限値に達した」であったり、学習曲線の勾配を用いて自動的に決めたり、そのどちらかを採用する。本実施形態では、学習の停止条件として、「Iteration回数が予め設定した上限値に達している」を採用するが、これに限るものではない。本実施形態では、「予め設定した上限値」を「20000」とするが、この値に限らない。なお、Iteration回数が増えることによってNNの学習率を低下させる方法など、Iteration回数に基づく学習処理の工夫を導入してもよい。

【 0 0 6 0 】

次に、図14～図17の構成による学習時の動作について、図18のフローチャートに従って説明する。ステップS1801では、初期化部221は第2のCNNの初期化を行う。初期化の方法は公知の方法を用いることが可能であり、例えば上述したCNNの初期化方法を用いてもよい。ここで初期化したモデルを操作部222に送り、操作部222は受け取ったモデルをストリーム223が備える第2のCNNのモデルとして読み込む。

10

【 0 0 6 1 】

ステップS1802では、受信部211は、データ選択部114から詳細認識対象データを受け取り、さらに受信部211は、実行部122から第1のCNNの中間層出力を受け取る。受信部211は、受け取った詳細認識対象データを、実行部222の受信部222に送る。受信部222に送られた詳細認識対象データは、以降のステップで認識処理にかけられる。

20

【 0 0 6 2 】

ステップS1803では、設定部212は、前段処理部11から注目領域情報を受け取る。ステップS1804では、受信部222は、詳細認識対象データをストリーム223に送る。ストリーム223は、受け取った詳細認識対象データに対して、第2のCNNを用いて認識処理を行う。ここで、具体例を示すために、図19に第1のCNN及び第2のCNNの構成とデータフローに関する模式図を例示する。図19における例では、認識対象人物の行動を認識するために、第1のCNNは認識対象人物の全身に関わる行動、第2のCNNは認識対象人物の頭部に関する詳細な特徴をそれぞれ認識する役割があるとする。このとき第2のCNNは、認識対象人物の頭部領域の画像と、第1のCNNの認識した全身に関する中間層出力と、に基づいて、頭部に関わる動作を認識（特徴量を抽出）するものとする。ここで抽出された詳細な特徴量を、詳細認識処理結果として、第2のCNNは第1のCNNに送り、第1のCNNは受け取った詳細認識結果をさらなる認識処理を行うための入力データとして用いる構成をとる。このとき第2のCNNは、第1のCNNの特徴マップ上の頭部領域の位置を考慮して、詳細認識結果を送ってもよい。このとき、第2のCNNの出力する詳細認識結果は、頭部に関する詳細な特徴量である。この詳細な特徴量が、第1のCNNの持つ認識対象人物の全身に関する中間層特徴マップの頭部領域に送られ、統合されることで、第1のCNNの中間層特徴マップは頭部に関して詳細化されることが期待できる。このように詳細な特徴量を用いる構成によって、例えば人物の「キョロキョロ」を精度よく検知することが期待できる。1003、1011、1012は、それぞれ図10において示した認識対象人物の例、認識対象人物の全身領域の画像データ、認識対象人物の頭部領域の画像データ、を示している。1901は、頭部画像データと頭部領域の位置データを受け取る設定部212を表しており、この処理はステップS1803の処理を示している。1902は第2のCNNの認識処理を行う対象の頭部領域画像データ、1903は第1のCNNが認識処理を全身画像の例である。1904、1914はそれぞれストリーム223が備える第2のCNNの第一層、第二層の畳み込み処理を示している。1902は、頭部領域の画像データを第2のCNNの第一層に入力するためにサイズの変形を行った画像を示しており、第2のCNNの第一層に入力され、出力結果としての詳細認識処理結果1905が得られる。1905は第2のCNNの第二層に入力され、以下同様に第2のCNNの処理がなされる。ここまでが、ステップS

30

40

50

1804の処理である。なお、ここで説明した方法は一つの例であり、その他の方法を用いて第2のCNNと第1のCNNとを連携させてもよい。その他の例に関しては後述する。

【0063】

ステップS1805では、ストリーム2223は、変形部2224に詳細認識処理結果を送る。変形部2224は、受け取った詳細認識処理結果を変形することで、第1のCNNの中間層特徴マップに統合するための前段階としての処理を行う。このステップの処理の例は、図19において、詳細認識処理結果1905が変形部2224である1908に送られ、変形結果1909が得られるという工程として表されている。このとき、どのような情報に基づいてどのように変形するかどうかはバリエーションがあり、後述する。

10

【0064】

ステップS1806では、統合部2225は、変形部2224から変形された詳細認識処理結果を受け取り、これを入力データとして第1のCNNの中間層特徴マップに統合するための処理を行う。この際に、どのように統合処理を行うかどうかのバリエーションが複数考えられるため、詳細は後述する。このステップの処理の例を、図19において、変形結果1909と、第1のCNNの認識処理結果（詳細認識対象データ）である1907を受け取り、詳細認識処理結果1912を出力する統合部2225（図19では1910）として示す。

【0065】

ステップS1807で、指定部231は、設定部212から注目領域情報を受け取る。出力部232は、指定部231が指定する第1のCNNの注目領域に対して、詳細認識処理結果を送信する。ここで注目領域は、本実施形態においては認識対象人物の頭部領域である。第1のCNNの特徴マップ（ある層のある特徴マップ）上の認識対象人物の頭部領域の大きさは、認識対象人物の原画像上での頭部の大きさに対して、第1のCNNの処理等を経て変動しうる。例えば、認識対象人物の原画像上での頭部の大きさが既知であれば、第1のCNNのある特定の層の特徴マップ上で表現される頭部の大きさは計算可能であり、設定部212はそのように頭部の注目領域の大きさを求めてもよい。ここで得られる頭部の注目領域の大きさの変動は、先述の変形処理によって吸収される（詳細認識処理結果の特徴マップ上の大きさが、その時々々の注目領域の大きさにリサイズされる）。また、認識対象人物の頭部領域の位置は、人物の姿勢等によって変動しうる。この位置に関する変動は、設定部212が頭部領域の位置を指定することによって吸収される。以上のような変動は、時々刻々と認識対象人物の頭部領域が変化することで生じうる。上記の説明は、その変動を吸収するために、ニューラルネットワークの構造の一部が時々刻々と変化しうることを示している。例えば、頭部領域の位置が変わることによって詳細認識処理結果の送信先の領域が変動することは、ニューラルネットワークの結合先（ニューロン間の結合関係）が動的に変わることによって実現できる。ニューラルネットワークの結合先（ニューロン間の結合関係）を動的に変える構成は、ソフトウェアで実現してもよいし、ハードウェアで実現してもよい。

20

30

【0066】

例えば、以降で引用する文献6で用いられるVGG16というモデルでは、各層の特徴マップの大きさは、層のコンポーネントごとに定まっている（ここではプーリング層を挟むと層のコンポーネントが変わるという解釈のもと例示している）。具体的には、一段目から五段目までの層のコンポーネントにおいて、224x224、112x112、56x56、28x28、14x14というサイズの特徴マップとなる。また、入力画像サイズは224x224である。よって、各層のコンポーネントの特徴マップサイズと、入力画像のサイズとの比は、1、0.5、0.25、0.125、0.0625として定まる。例えばこの比に基づいて、ある特徴マップ上の注目領域の大きさを決定してもよい。

40

【0067】

例えばソフトウェアで実現する場合、上述の通り、統合部2225（図19では1910）が、詳細認識処理結果の送信先を認識時に変更してもよい。その場合、例えば認識時

50

において、同一のメモリ上に格納されたニューロン間でデータを送受信する際に、送受信先を決定するポイントを変更し、新たな送受信先となるニューロンを指定することが可能である。また、モデルの更新時には、認識時に用いたポイントを記憶しておき、誤差逆伝播法適用の際に誤差を伝播するデータフローを決定してもよい。その他の公知なソフトウェアによる方法によっても、ニューロン間の結合を変更してもよい。

【0068】

また、ハードウェアで実現する場合、例えばインターネットを介して通信機器間のデータ送受信先を切り替えることで、ニューロン間の結合を変更してもよい。例えばニューラルネットワークのモデルが複数の機器に分散して存在するとする。認識時において統合部2225（図19では1910）が、送受信先となる機器のアドレスを変更することでニューロン間のデータの送受信先変更を実現することが可能である。その他の公知なハードウェアによる方法によっても、ニューロン間の結合を変更してもよい。

【0069】

また、上述のハードウェア・ソフトウェアの実装を複合的に用いて、ニューロン間の結合を変更してもよい。また、ニューラルネットワーク中の結合重みを変更し、活性化状態を変化させることで、詳細認識処理結果の送信先を変更する処理を実現することができる。具体的には、詳細認識処理結果の送信先の候補に対して、ニューロン間の結合重みが存在するとする。送信先の候補とは、詳細認識処理結果を代入する第1のCNNの特徴マップ中の領域の候補とも言える。この送信先の候補に対して、詳細認識処理結果が実際に送信されるかが事前には決定されておらず、したがって詳細認識処理結果の送信元と受信先となるニューロン間の結合重みは0、バイアス項も0の値を持つとする。そして、認識時に送信先の候補が選ばれ、送信先が確定するとする。このとき、確定した送受信間関係に基づいて、統合部2225（図19では1910）が送受信の行われるニューロン間の結合重みを1にする。このように結合重みを変更されることで、詳細認識処理結果の転送が実現できる。これらの処理は、図19において以下のように例示される。

【0070】

1910から送られた詳細認識処理結果は、1911の指定部231から指定された頭部領域（1912）に送られる。送られた頭部に関する詳細認識処理結果は、ここでは、第1のCNNの中間層特徴マップである1907の頭部領域（1912）に対してそのまま上書き処理を行うものとする。なお詳細は統合方法の説明において例示するが、第1のCNNの中間層特徴マップと、第2のCNNの出力する詳細認識処理結果と、の特徴量の次元数はここでは同じであるとする。最後に、詳細認識処理結果が統合された第1のCNNの中間層出力は、実行部122第2層（図19の1913）に送られ、以降も同様の統合処理を行ってもよいし、行わなくてもよい。第1のCNNは、図13で示したように、最終的に認識結果を出力する。以降も同様に統合処理を行う場合は、第2のCNNのさらなる認識結果を得るために、1905をストリーム2223（頭部）第2層（図19の1914）に送ってもよい。そして、第2のCNNと第1のCNNの認識結果の統合を行う際に、ある一つの層の特徴マップの出力を用いてもよいし、さらに深い層で得られた第2のCNNの詳細認識処理結果を用いてもよい。なお、ここで述べる第2のCNNの詳細認識処理結果とは、第2のCNNの各層の中間層出力と、最終層出力とを含む。

【0071】

今回の例では、統合する対象となる第2のCNNの詳細認識処理結果と第1のCNNの中間層出力は、同一の層の情報であるものとする。一方で、認識結果の統合には異なる層の出力を用いてもよい。例えば、第2のCNNの第一層の詳細認識処理結果を、第1のCNNの第三層の中間層特徴マップに統合してもよい。また、第2のCNNの複数の異なる層の詳細認識処理結果を、第1のCNNの特定の層の中間層特徴マップに統合してもよい。この場合の統合方法については後述する。

【0072】

ステップS1808では、受信部233は、CNN更新部123から誤差情報を受け取る。このとき受け取る誤差情報は、ステップS1807で送信した詳細認識処理結果に関

10

20

30

40

50

する誤差情報であり、ステップS 1 8 0 7で用いた注目領域情報に基づいて受信される。具体的には、ステップS 1 2 0 5で得られた第1のCNNの誤差情報のうち、ステップS 1 8 0 7での送信先の領域に関わる誤差情報のみを受け取ることとなる。受信部2 3 3は、得られた誤差情報を更新部2 2 3に送る。更新部2 2 3は、受け取った誤差情報から、誤差逆伝播法に基づいて第2のCNN上の誤差を算出し、ここで得られた誤差情報を用いて第2のCNNのモデルを更新する。このように第2のCNNに逆伝播した誤差情報は、第2のCNNに対して認識対象情報を入力した第1のCNNの一部に再び伝播してもよく、第1のCNNと第2のCNNとでEnd-to-endな学習系を構築することが可能である。ここで更新された第2のCNNのモデルは、操作部2 2 2 1に送られる。操作部2 2 2 1は、ストリーム2 2 2 3および統合部2 2 2 5にモデルを送り、認識処理で用いるモデルの更新を行う。なお、ここでは変形部2 2 2 4が変形処理に学習パラメータを用いていないものとして、変形部2 2 2 4のモデルを更新しなかった。なお、変形処理になんらかの学習パラメータを用いている場合は、変形部2 2 2 4のモデルを更新してもよい。なお、変形処理に関する詳細な説明は後述する。

10

【0073】

そして、学習のためのMinibatchをすべて利用した場合には、処理はステップS 1 8 0 9を介してステップS 1 8 1 0に移行し、利用していないMinibatchが残っている場合には、処理はステップS 1 8 0 2に移行する。ここでの判定には、ステップS 1 2 0 7の判定結果をそのまま利用してもよい。

【0074】

20

学習が完了した場合には、処理はステップS 1 8 1 0を介してステップS 1 8 1 1に移行する。一方、学習が完了していなければ、処理はステップS 1 8 0 2に移行する。ここでの判定には、ステップS 1 2 0 9の判定結果をそのまま利用してもよい。

【0075】

ステップS 1 8 1 1では、操作部2 2 2 1は、第2のCNNのモデルを保存部2 2 4に送る。保存部2 2 4は、操作部2 2 2 1から受けた第2のCNNのモデルを記憶部M 1に格納する。

【0076】

以上、学習時に関する各部の詳細な構成および動作について説明を行った。次に、検出時に関する各部の詳細な構成および動作について説明を行う。前述のとおり、検出時には、図7のステップS 7 0 1において第1のCNN及び第2のCNNの認識処理を行う。ここで用いる各部の構成は、学習時と同じであってよい。すなわち、入力データを学習データではなく検出対象データとして、それ以外の部分を共通化して用いてよい。なお、モデルの更新などの学習処理は、検出時にも必要に応じて使い分けることができる。具体的には、検出対象データに教師データが付与されている場合、例えば学習時と同様の方法で学習を行うことによって、追加学習を行ってもよい。また、検出対象データに教師データが付与されていない場合、教師データに基づいて誤差を計算することはできないため、学習処理に関する一部の構成は用いなくてもよい。また、その場合、異常検知システムはユーザに対して教師データの付与を依頼してもよい。

30

【0077】

40

以降では、詳細認識対象データをどのように設定するかに関して、一例を説明する。本実施形態においては、認識対象人物の行動（特に、キョロキョロという行動を異常として定めたときの行動）を認識するための構成の例を示した。ここでは、第1のCNNの入力は認識対象人物の全身画像であり、第2のCNNの入力は認識対象人物の頭部画像であるとしている。このうち、認識対象人物の頭部を第2のCNNの詳細特徴抽出対象とする、ということが、詳細認識対象データを頭部と設定する、という定義に相当する。この定義は、“詳細認識対象データ＝頭部”という一対一の固定な関係を生じている。なお、詳細認識対象データを頭部以外に設定することも可能である。例えば、万引き行動を検知するために、“詳細認識対象データ＝商品に触れている手”などとしてもよい。この場合、例えば認識対象人物の手を検知することと、手が商品に触れていることを検知し、その手の

50

領域データを活用して詳細認識を行ってもよい。また、詳細認識対象データはユーザが指定してもよく、あらかじめシステム上に定義されていてもよい。

【 0 0 7 8 】

なお、詳細認識処理結果を統合するための一構成例として、本実施形態では、頭部の位置を考慮して第 2 の CNN の詳細認識処理結果を第 1 の CNN の中間層特徴マップに統合する例を示した。しかし、必要であれば、位置を考慮せず、第 2 の CNN の詳細認識処理結果を第 1 の CNN の中間層特徴マップに統合してもよい。その場合は、前出の変形処理を用いず、例えば第 2 の CNN と第 1 の CNN の各層の中間層出力を、例えば以下の文献 4 ～ 6 の方法で統合してもよい。この際、例えば第 2 の CNN の詳細認識結果を、第 1 の CNN の中間層特徴マップに統合する構成をとることができる。

10

【 0 0 7 9 】

文献 4 ... Spatiotemporal Residual Networks for Video Action Recognition. C. Feichtenhofer, A. Pinz, R. P. Wildes. Advances in Neural Information Processing Systems (NIPS), 2016.

文献 5 ... Spatiotemporal Multiplier Networks for Video Action Recognition. C. Feichtenhofer, A. Pinz, R. P. Wildes. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.

文献 6 ... Convolutional Two-Stream Network Fusion for Video Action Recognition. C. Feichtenhofer, A. Pinz, A. Zisserman. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.

20

また、頭部などの注目領域の位置を用いて、第 2 の CNN の詳細認識処理結果と第 1 の CNN の中間層特徴マップとを統合する方法として、その具体例を以下に述べる。なお以降で述べる方法は、これまで説明した「第 2 の CNN の出力する詳細認識処理結果を第 1 の CNN の特徴マップ上の領域に送信したうえで統合処理をする」という手順に必ずしも限らない点に注意が必要である。第一の候補としては、第 2 の CNN の出力のみを用いる場合であり、例えば以下の式 1 を第 1 の CNN の中間層特徴マップの注目領域に代入することができる。なお代入ではなく、足し合わせてもよい（加法演算）。

【 0 0 8 0 】

【 数 1 】

30

$$f_1(\text{Reshape}(A)) \cdots (1)$$

【 0 0 8 1 】

ここで、A は第 2 の CNN のある層の中間層出力ないし最終層出力（詳細認識処理結果）、Reshape は後述の変形処理関数、 f_1 は特徴変換のための関数であり、例えば 1×1 畳み込み処理などを導入してもよい。また、 f_1 を用いない場合は恒等写像を用いてもよい。第二の候補としては、第 2 の CNN の詳細認識処理結果と、第 1 の CNN の注目領域に関する中間認識情報を用いる場合であり、例えば以下の式 2 を第 1 の CNN の中間層特徴マップの注目領域に代入することができる。

40

【 0 0 8 2 】

【 数 2 】

$$f_2(\text{Reshape}(A), \text{Crop}(B)) \cdots (2)$$

【 0 0 8 3 】

ここで、B は第 1 の CNN のある層の中間層出力、Crop は注目領域（頭部パーツの領域

50

）の特徴マップの出力を抜き出す処理であり、例えば以下の文献 7 のROI poolingを用いてもよい。 f_2 は特徴変換のための関数であり、Concatenate処理を行ったのち、 1×1 畳み込み処理を導入してもよい。これによって、第 2 の CNN と第 1 の CNN の頭部に関する局所的な中間層出力を考慮して認識処理を行うことができる。なお、 f_2 は、例えば文献 8 に記載の residual networkを用いてもよい。

【 0 0 8 4 】

文献 7 ... Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. S. Ren, K. He, R. Girshick, J. Sun. Advances in Neural Information Processing Systems (NIPS), 2015.

文献 8 ... Identity Mappings in Deep Residual Networks. K. He, X. Zhang, S. Ren, J. Sun. ECCV, 2016 10

第三の候補としては、第 2 の CNN の詳細認識処理結果と、第 1 の CNN の中間層出力とを、一部の領域に限らずすべて用いる場合であり、例えば以下の式 3 を第 1 の CNN の中間層特徴マップに代入することができる。

【 0 0 8 5 】

【数 3】

$$f_3(\text{PartialAdd}(\text{Reshape}(A), B)) \cdots (3)$$

【 0 0 8 6 】

20

なお中間層特徴マップに代入するための構成は、ソフトウェアによる実現と、ハードウェアによる実現とがあり、図 18 のステップ S 1807 に係る説明で述べた方法によって実現可能なため、ここでは詳細な説明を省略する。ここで、PartialAdd は第一引数を第二引数の一部の領域に足し合わせる関数であり、ここでは B の注目領域に足し合わせる処理を担う。 f_3 は特徴変換のための関数であり、中間層特徴マップ全体の処理を行うことを特徴とする。そのため、より広い範囲の情報を捉えつつ、詳細認識処理結果との統合処理を行うことができるという特徴がある。第四の候補としては、第 2 の CNN の詳細認識処理結果を、第 1 の CNN の中間層特徴マップに基づいて入力する場合であり、例えば以下の式 4 を中間層特徴マップに代入することができる。

【 0 0 8 7 】

30

【数 4】

$$f_5 \left(\text{PartialAdd} \left(\text{Reshape} \left(A * \sigma(f_4(B)) \right), B \right) \right) \cdots (4)$$

【 0 0 8 8 】

なお中間層特徴マップに代入するための構成は、ソフトウェアによる実現と、ハードウェアによる実現とがあり、図 18 のステップ S 1807 に係る説明で述べた方法によって実現可能なため、ここでは詳細な説明を省略する。ここで f_4 、 f_5 は特徴変換の関数であり、 σ は sigmoid 関数、 $*$ は elementwise な掛け算処理である。すなわち、B の関数 f_4 の出力の大きさに応じて、A との掛け算の結果が変化し、これによって A をどの程度重要な情報として扱うかの決定を行うことができる。なお、関数 f_4 を A にも依存するようにしてもよく、その場合は第 2 の CNN と第 1 の CNN の認識処理に基づいて統合処理をコントロールするシステムとなる。

40

【 0 0 8 9 】

なお、本実施形態では第 2 の CNN の詳細認識処理結果を第 1 の CNN の中間層特徴マップに統合する例を示した。しかし、逆に第 1 の CNN の中間層特徴マップを第 2 の CNN の詳細認識処理結果に統合し、その統合結果に基づいて第 2 の CNN の処理を行ったのち、第 1 の CNN の中間層特徴マップに統合するような処理を採用してもよい。

【 0 0 9 0 】

なお、本実施形態では第 2 の CNN のある層の詳細認識処理結果を第 1 の CNN に統合

50

する例を示した。その他の例として、第2のCNNの複数の層の詳細認識処理結果を第1のCNNに統合するような構成をとってもよい。この場合は、例えば複数層からなる詳細認識処理結果の大きさをすでに述べたReshapeによって揃え、Concatenateし、 1×1 畳み込み処理をして特徴次元数をもとの次元数に低減する。これにより、本実施形態ですでに述べた通常の詳細認識処理結果として扱うことができるため、このようにしてもよい。

【0091】

なお、変形処理（上述のReshape関数）は、例えば、畳み込み処理に基づいて大きさを変更してもよいし、上記の非特許文献1で用いられるプーリング処理（max pooling, average poolingなど）を用いてダウンサンプリング処理を施してもよい。

【0092】

なお、例えばmax poolingを用いる際に、入力特徴マップを均一にカバーしたプーリングを行うという指針のもとで、プーリング処理を行う窓のサイズを決める一つの例として、入力特徴マップサイズを出力特徴マップサイズで割ってもよい（例えばその値を天井関数で丸めてよい）。また、そこで決めた値のもとで、 $(\text{入力マップサイズ} - \text{窓サイズ}) / (\text{出力特徴マップサイズ} - 1)$ という式によって、プーリング処理のストライドサイズを決めてもよい（その値を床関数で丸めてもよい）。なお出力特徴マップサイズは1以下である場合に、仮の値として2を用いることで、0割を回避することができる。また、窓のサイズはストライドの値以上の値であり、また入力特徴マップサイズと出力特徴マップサイズが同じ場合はストライドの値を1としてもよい。

【0093】

なお、本実施形態では一枚のRGB画像を入力する第1のCNNの例を示したが、他の入力形式であってもよい。例えば、以下の文献9で示されるTemporal Streamの入力層のように、複数の画像を入力として受け付ける入力層を持ってもよい。また、以下の文献9では、入力画像としてオプティカルフローを使う場合について例示されているが、本実施形態においても、このようにしてもよいし、さらに文献9に記載のTwo-stream構成にしてもよい。また、文献10のように、後段にLong-short term memoryなどのrecurrent neural network(RNN)を備える構成を用いてもよい。後段にLong-short term memoryなど時系列情報を処理する機能を持つコンポーネントを用いることで、時間情報の変化を捉えることが可能である。また、Long-short term memoryなどのRNNに畳み込み層を付与した公知のconvolutional long-short term memoryを用いてもよい。Long-short term memoryに畳み込み層を付与することで、後段の処理においても、本実施形態で用いる統合処理における第1・2のCNNの位置関係を考慮することができ、これによって認識精度が向上する場合がある。

【0094】

文献9 ... Two-Stream Convolutional Networks for Action Recognition in Videos . K. Simonyan, A. Zisserman. NIPS, 2014.

文献10 ... A Multi-Stream Bi-Directional Recurrent Neural Network for Fine-Grained Action Detection. B. Singh, et al. CVPR, 2016.

なお本実施形態では、第1のCNNと第2のCNNとが異なる入力画像（224画素×224画素の画像サイズにリサイズ済み）を受け付け、第2のCNNの詳細認識処理結果を第1のCNNの中間層特徴マップに統合するマルチストリームの構成例を示している。しかし、上記の構成ではなく、第2のCNNと第1のCNNとを連携する他の構成を用いてもよい。例えば、第2のCNNが第1のCNNの中間層出力を受け取る構成や、第2のCNNが入力画像を受け取らない構成や、入力画像をリサイズしない構成が考えられる。このようなアーキテクチャを組み合わせる例として上記すべての特徴を兼ねた構成の例を図21に示す。

【0095】

図21の2101は、入力画像としての人体領域1011と同じものであり、第1のCNNは人体領域1011と同じものである人体領域2101を受け取っている。図19の全身画像1903とは異なり、人体領域2101はリサイズされていない。実行部122

第1層2102は人体領域2101を受け取り、畳み込み処理などをかけ、中間層出力2103として出力する。そして、中間層出力2103と、その頭部領域の中間層出力2104とは第2のCNNに送出される。具体的には、中間層出力2104を第2のCNNの入力として畳み込み処理などをストリーム2223第1層2106で行い、その出力である詳細認識処理結果と中間層出力2103とを統合部2225第1層2105で統合処理する。そして、その処理結果を詳細認識処理結果として、中間層出力2103に上書きする。そして、その結果を実行部122第2層2107に入力し、以降の処理を繰り返す。最後に、認識処理対象のクラス数の特徴次元数を出力する畳み込み処理を行い、特徴マップ全体を包含するMax pooling処理を行い、その結果をSoftmax処理にかけることによって、認識結果を得ることができる。以上の統合処理の方法や、ニューラルネットワークのアーキテクチャなどは、(任意)に変更してもよい。これによって、図19とは異なる構成で、第2のCNNと第1のCNNとを組み合わせることが可能である。

10

【0096】

なお、本実施形態における説明では、例として詳細認識処理結果を第1のCNNに統合する場合を示した。しかし、第1のCNNではなく、Long-short term memory、ニューラルネットワーク、確率モデル、識別器、Convolutional long-short term memoryなどその他のモデルに統合するような構成であってもよい。このとき、第1のモデルがCNNではなく、その他のモデルであってもよいし、CNNではあるが一部のみその他のモデルであってもよい(第2のCNNも同様である)。

20

【0097】

なお、本実施形態においてROI embeddingを行う場合の例を図24に示す。図24の2701は時系列を表している。2702および2703はそれぞれある時刻におけるある人物の画像および右手周辺の画像であり、それぞれ第2・第1のCNNの中間層特徴マップを模擬的に示している。また、2704は第2のCNNの中間層特徴マップをプーリングする際のプーリング範囲の例を示しており、2705は第1の中間層特徴マップ2703上の当該人物の右手周辺(注目)領域として、2702を統合処理2706によって統合することを示す。2707、2708はそれぞれ2702、2703と同様であるが、異なるタイミングで取得された特徴量である。2710は2706と領域のサイズおよび位置が異なることが分かる。この場合、2709のように、プーリング範囲を変更することで対応することが可能である。なお、2703などに示される画像は文献12に記載のデータベースに収録された画像を使用している。

30

【0098】

文献12 ... A Database for Fine Grained Activity Detection of Cooking Activities. Marcus Rohrbach and et al., Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on.

なお、本実施形態において示される統合処理の効果や、中間処理結果などを可視化することができ、例えば図27のように示してもよい。3001、3002はそれぞれ第1・2のCNNの中間層特徴マップの模式図であり、3003の示す注目領域に対して3002を統合する3004があるとすると、このとき、例えば中間層特徴マップ(注目領域3003などの局所領域でもよいし、全体でもよい)の活性具合(ノルムの大きさや数値の大きさ)や、その他の統計量などで示される数値が変化したことなどによって、統合処理の影響の大きさを測ることができる。このとき、上述の数値を例えばマップ上に示すことで、上記影響の大きさを可視化することができる。図27では、3005の部分の値の大きさが上がっていることを例示している。なお、値の大きさを測る際には特定の次元の特徴マップを用いてもよいし、次元方向に総和などの処理を施して得た特徴マップを用いてもよい。

40

【0099】

なお、本実施形態においては一例として第1・2のCNNの学習手順などを示したが、例えば以下のような手順で学習をしてもよい。まず、第1・2のCNNを独立に学習する。このとき、各CNNは例えば前出のVGG16を用いてもよい。次に、統合処理を行うレイ

50

ヤーを決定し、第1のCNNを第2のCNNと統合したアーキテクチャのもとで、前の学習処理で得た各CNNのパラメータを再利用する形で学習を行う（統合処理を行うレイヤーは何層あってもよく、例えば前述のVGG16の定義として4_2層および5_2層のReLU処理の結果を用いることができる）。このとき、第1・2のCNNだけでなく、統合部分の学習処理を行ってよい。次に、Long-short term memoryの学習を行う。例えば上述の第1のCNNの4_2層および5_2層において第2のCNNと統合処理を行い、その結果をconcatenate処理したうえでベクトルとして扱い、Long-short term memoryの入力として用いてもよい。

【0100】

以上説明したように、本実施形態で示した一形態は、データ中の認識対象人物やその状態を認識するCNNと詳細CNNとがあり、CNNの中間層出力と詳細CNNの詳細認識処理結果とを統合するための学習を行う。これにより、認識に寄与する詳細な特徴量を抽出し、認識対象を精度よく認識することができ、最終的な認識結果（異常検知結果）の精度が向上しうる。

【0101】

[第2の実施形態]

本実施形態を含む以下の各実施形態では、第1の実施形態との差分について説明し、以下で特に触れないか居切は第1の実施形態と同様であるものとする。第1の実施形態では、データ中の物体やその状態を認識する畳み込みニューラルネットワーク（CNN）である第1のCNNと、さらに詳細な特徴を認識するCNNである第2のCNNと、があり、学習に基づいてそれらの中間の認識処理結果を統合して認識結果を得る例を示した。その際の構成の第一の例は、第1のCNNと第2のCNNとがそれぞれ一つあり、組として動作するTwo-streamの構成であった。また、その他の例として、第1のCNNと第2のCNNの組が複数ある場合の構成について例示した。これらはいずれも第1のCNNと第2のCNNとが組として一対一の対応関係にあるという特徴がある。

【0102】

本実施形態では、一つの第1のCNNに対して複数の第2のCNNが存在し、複数の第2のCNNが異なる詳細認識処理結果を抽出する場合の構成および動作について例示する。このとき、複数の第2のCNNの詳細認識処理結果をいかに統合するか、ということに関して、第1の実施形態とは異なる構成および動作について説明する。なお、本実施形態で示す構成は第1の実施形態で例示した構成と大部分が同一であり、一部の構成と動作が異なる。本実施形態において第1の実施形態と構成や動作が大きく異なるのは、先述のように、一つの第1のCNNに対して複数の第2のCNNが存在する点である。

【0103】

本実施形態においては、一つの構成例として、第1の実施形態の構成に加えて手の詳細な特徴量を抽出する第2のCNNを導入する場合の構成および動作に関して例示する。すなわち、認識対象人物の全身を認識する第1のCNNと、頭部を詳細に認識する第2のCNN（頭部）と、右手を詳細に認識する第2のCNN（右手）と、左手を詳細に認識する第2のCNN（左手）と、があるとする。これによって、例えば認識対象人物が万引き行動を行う際に、対象人物の全身と、頭部と、右手と、左手と、の詳細な特徴量を抽出・統合することで、万引き行動を比較的精度よく検知しうる。

【0104】

本実施形態において、認識対象人物の全身と、頭部と、右手と、左手と、の領域を抽出する方法は、第1の実施形態と同一の方法を用いるとする。ここで、図10の下部における1011、1012、1013、1015はそれぞれ認識対象人物の全身、頭部、右手、左手、が検出された領域を例示する。このうち、全身および頭部領域中の特徴量は図19の構成によって認識処理をかけることができ、図19は第1のCNNと第2のCNNとが一対一の関係である場合の例を示している。このとき、図19の頭部領域に関する処理をそれぞれ右手領域や左手領域に読み替えることで、「全身および右手領域」や、「全身および左手領域」の認識処理を行う構成になりうる。すなわち、複数の第2のCNNを用

10

20

30

40

50

意し、それぞれが頭部・右手・左手に関する詳細認識処理結果を抽出する機能を持つものとして行うことができる。このとき重要になるのは、上記の複数の詳細認識処理結果をどのように統合するか、という点である。例えば、図19の1909において、頭部・右手・左手に関する詳細認識処理結果の変形結果をそれぞれ得たとする。これら複数の変形結果と、第1のCNNの中間層出力1907と、を統合処理する機能が必要である。これらの特徴量を統合する方法としては、例えば各パーツの位置を考慮しない方法として、以下の式5を第1のCNNの中間層特徴マップの注目領域に代入することができる。

【0105】

【数5】

$$f_6(\text{Concat}(B, A, C, D)) \cdots (5)$$

10

【0106】

なお中間層特徴マップに代入するための構成は、ソフトウェアによる実現と、ハードウェアによる実現とがあり、図18のステップS1807に係る説明で述べた方法によって実現可能なため、ここでは詳細な説明を省略する。ここでB、A、C、Dはそれぞれ全身・頭部・右手・左手の詳細認識処理結果であり、ConcatはConcatenateの略で複数の入力の特徴次元方向に結合する処理であり、 f_6 は特徴変換処理を表す。例えば f_6 が 1×1 の畳み込み処理であり、特徴次元数が相対的に小さくなる場合、式5はB、A、C、Dの入力を次元圧縮する処理を行うことになる。式5に関して、パーツの位置を考慮する場合は、例えば以下の式6を第1のCNNの中間層特徴マップの注目領域に代入することができる。

20

【0107】

【数6】

$$f_7(\text{Concat}(B, \text{Pad}(\text{Reshape}(A)), \text{Pad}(\text{Reshape}(C)), \text{Pad}(\text{Reshape}(D)))) \cdots (6)$$

【0108】

なお中間層特徴マップに代入するための構成は、ソフトウェアによる実現と、ハードウェアによる実現とがあり、図18のステップS1807に係る説明で述べた方法によって実現可能なため、ここでは詳細な説明を省略する。ここでPadは領域を考慮した0埋め処理である。すなわち、第1のCNNの特徴マップ上における各パーツの領域の大きさと位置とを考慮した0埋め処理であり、その具体的な処理の例を図20に示す。

30

【0109】

2301、2302、2303はそれぞれ、頭部、右手、左手の詳細認識処理結果である。これらをReshape処理によって、統合先の第1のCNNの特徴マップ上の大きさに変換したものが、2311、2312、2313である。さらにこれらを、統合先の第1のCNNの特徴マップ上のパーツ位置に配置し、不要な領域を0埋め(Pad処理)した詳細認識処理結果が2321、2322、2323である。これによって、第1のCNNの特徴マップ上での領域の大きさと位置とを考慮した詳細認識処理結果を得ることができる。これらの出力とBをConcatし、特徴変換する処理が式6である。 f_7 は特徴変換処理を表す。

40

【0110】

なおConcat処理ではなく、第1の実施形態のPartialAdd処理を用いても、パーツ領域の大きさと場所とを考慮した処理を導入することが可能である。また、例えば第1の実施形態に記載のその他の処理を導入して、複数の詳細認識処理結果の統合処理を導入してもよいし、その他の方法を用いてもよい。例えば、以下の式7を第1のCNNの中間層特徴マップに代入することができる。

【0111】

【数 7】

$$f_8 \left(\text{PartialAdd}(B, \text{Reshape}(A), \text{Reshape}(C), \text{Reshape}(D)) \right) \cdots (7)$$

【0112】

なお中間層特徴マップに代入するための構成は、ソフトウェアによる実現と、ハードウェアによる実現とがあり、図18のステップS1807に係る説明で述べた方法によって実現可能なため、ここでは詳細な説明を省略する。基本的に式7の結果は第1のCNNの中間特徴マップ全体に代入すればよい。ここでPartialAddは第1の実施形態で説明したパーツの場所を考慮した統合処理関数であり、 f_8 は特徴変換処理を表す。式7のPartialAddは第1の実施形態とは異なり、複数パーツに対応したPartialAdd関数である。その処理内容は基本的に第1の実施形態と同じであり、変数A, C, Dがそれぞれ持つ頭部・右手・左手の詳細な特徴量を、B中の領域に対して、足し合わせる処理を示す。

【0113】

なお、複数の領域が重複する場所に関しては、例えば遮蔽されたパーツの領域に関する詳細認識処理結果は無視する（統合処理を行わない）という処理を導入してもよい。この際、パーツが遮蔽されていることを判定する方法として、例えばパーツの奥行に関する位置を取得して用いることで判定してもよい。

【0114】

なお、それぞれの詳細認識処理結果を抽出するか否かを判定し、不要な詳細認識処理結果は用いないという機能を導入してもよい。例えば遮蔽されているパーツや、パーツ領域の検出結果の信頼度が低いパーツを検出し、その場合に詳細認識処理結果を抽出しない、あるいは詳細認識処理を行わない、というように処理の変更をすることができる。

【0115】

なお、第1のCNNの中間層出力と第2のCNNの詳細認識処理結果とを統合することで、値のスケールに変化が起きるのを避けるために、正規化処理を導入してもよい。例えば統合処理を行った第1のCNNと第2のCNNの数を用いてもよい。具体的には、第1のCNNの中間層特徴マップのある領域には、もとの中間層出力と、2つの第2のCNNの詳細認識処理結果とが足し算によって統合された結果が代入されたとする。すなわち、3つの種類の値の和が得られている。

【0116】

なお中間層特徴マップに代入するための構成は、ソフトウェアによる実現と、ハードウェアによる実現とがあり、図18のステップS1807に係る説明で述べた方法によって実現可能なため、ここでは詳細な説明を省略する。ここで、該当領域部分の値を3で割ることで、正規化を行ってもよい。また、統合処理を行った第1のCNNと第2のCNNの数は陽には用いず、例えば平均値を用いて正規化を行ってもよい。具体的には、第1のCNNの各特徴マップの縦方向および/または横方向の平均値を用いて、特徴マップの値から引くことで、特徴マップの値のセンタリングを行ってもよい。このとき中心化行列を用いると、平易にセンタリングを行うことができる。いま、中間層特徴マップAの縦方向に関する平均を $\text{mean}(A, 1)$ とし、Aの中心化行列をHとすると、以下の式8, 9からHA及びAHを求めることができる。

【0117】

【数 8】

$$HA = A - \text{mean}(A, 1) \cdot \cdot \cdot \cdot (8)$$

$$AH = A - \text{mean}(A, 2) \cdot \cdot \cdot \cdot (9)$$

10

【0118】

正規化としてこれらの処理のいずれかを用いてもよいし、両方を導入したHAHを用いてもよい。なお中心化行列は $H = I - 1/n \cdot 1_n \cdot 1_n^T$ として求めることができる。ここで、 n は対称行列 A の一辺の大きさであり、 I は $n \times n$ のidentity matrix、 1_n は n サイズの1ベクトル、 T は転置処理を示す。これらの式は微分可能であり、誤差逆伝播法に用いることができる。なお、特徴マップが対称行列でない場合も同様にセンタリング処理を用いることができる。その他の正規化方法として、例えば以下の文献11の方法を用いてもよいし、それ以外の方法を用いてもよい。

【0119】

文献11 ... Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. S. Ioffe, C. Szegedy, JMLR. 20

以上のような構成においても、第1の実施形態で示したように、例えば誤差逆伝播法を用いて第1のCNNと複数の第2のCNNとを学習することが可能である。なお、第1の実施形態において図21に関して示したような構成を用いて、複数の第2のCNNの導入を図ってもよいし、その他の構成によって第1のCNNと複数の第2のCNNとの統合処理を行ってもよい。

【0120】

なお、本実施形態では第1のCNNが一つであり、第2のCNNが複数ある場合の例を示したが、第1のCNNが複数存在する構成であってもよい。その構成における一つの例としては、次のようなものがある。RGB画像を入力データとする第1のCNNのストリーム、複数枚のオプティカルフローを入力データとするストリームの第1のCNNのストリーム、を用意し、それぞれの第1のCNNに関する詳細認識処理結果を抽出する第2のCNNを複数用意してもよい。それ以外の構成としても、例えば第1の実施形態で示した構成を基本として、複数の第2のCNNを導入する形をとってもよい。

30

【0121】

なお、本実施形態では複数の第2のCNNをそれぞれ別個に保持する場合の例を示したが、第2のCNNのモデルを一部ないし全部共通化してもよい。例えば、一般的な行動認識処理をするにあたっては、認識対象人物の右手の動作も左手の動作も意味合いは同じであり、それぞれの特徴を抽出するための構成は右手と左手とで共通であってもよい場合がある。そこで、認識対象人物の右手に関する詳細な特徴を抽出する第2のCNNと、左手に関する第2のCNNと、では、共通のモデルを用いることとしてもよい。また、必要であれば、複数の第2のCNNのモデルをすべて共通化してもよい。これにより、大幅なモデルサイズ削減が期待できる。

40

【0122】

なお、本実施形態では詳細認識に用いる認識対象人物のパーツ（頭部や、手など）をあらかじめ詳細認識対象として設定し、それらをすべて認識処理に用いる場合の例を示した。このうち、認識に寄与する詳細認識対象を選択し、選択された詳細認識対象のみを認識処理に用いてもよい。例えば、式(4)は認識に寄与する詳細認識対象の特徴量を検出時に動的に重みづけする機能を持つ。その他の方法としては、例えば学習時に選択処理を行う場合がありえる。例えば、クロスバリデーションによって、どの詳細認識処理結果を用

50

いると認識処理が向上するかを確かめることが可能な場合は、これを用いてもよい。また、どの詳細認識処理結果を用いるかユーザに決定させてもよいし、必要であればユーザに新たな詳細認識対象を追加させてもよい。例えば、不審者を検知する不審行動検知のユースケースにおいては、人体のどの部分に不審さが現れるか、ということは現場の監視員が知悉するところである。そこで、監視員に認識対象人物のどのパーツに着目すべきかどうかという観点において、詳細認識対象をリコメンドしてもらうことが考えられる。検出時の実運用において新たな詳細認識対象が追加された場合は、必要であれば、はじめから学習処理をやり直してもよいし、使用していたモデルを初期値として学習処理をやり直してもよい。新たな詳細認識対象を検出するための一構成として、可能であれば上記の文献2の検出結果を用いてもよいし、必要であれば専用の検出器を導入してもよい。

10

【0123】

なお、本実施形態では詳細認識対象の種類が固定であるとしたが、動的に変動してもよい。例えば、認識対象の状態を認識するために、認識対象に付属する物体または物品の種類が重要である場合がありえる。具体的には、認識対象人物の属性を認識するために、当該人物の身に着けている品の種類が重要であるとする。しかしながら、検出時にどのような人物が現れるか事前には分からない場合、現れた人物が身に着けているものに応じて詳細に見るべき品の種類が変化しうる。これが、詳細認識対象の種類が動的に変動する場合に対応するとする。このような場合、物体/物品らしいものを検知することができる検出器を用いて、詳細認識対象の検知器として利用し、検知された詳細認識対象を認識処理に用いる構成を用いてもよい。例えば、上記の文献7の方法を用いて、認識対象人物が持つ物品/物体らしいものを検出し、これを詳細認識する構成が考えられる。その際、例えば検知された物品/物体がネックレスであるなら、ネックレス専用の第2のCNNを用いて詳細認識処理をするように、物品/物体のおおまかな種類に応じた詳細認識処理の構成が考えられる。この際、どのような種類か不明な物品/物体であれば、種別不明品を対象として第2のCNNを用いてもよい。また、物品/物体の種類によらずに、同一の第2のCNNを用いて認識処理をすることが可能である。

20

【0124】

なお、本実施形態では、認識処理を行う例として認識対象人物に関する異常な行動を認識する処理や、一般的な行動を認識する例を示したが、これらは一例であって、その他の認識処理を目的とした場合にも上述の方法を用いてもよい。例えば、静止画像中のシーンを認識するタスクにおいて、詳細認識処理を行う第2のCNNと、全体の荒い認識処理を行う第1のCNNと、を用い、それらを統合する構成を用いてもよい。より具体的には、先述した認識対象人物の全身の画像を入力とする第1のCNNと、認識対象人物のパーツに関する認識結果を出力する第2のCNNと、に関する構成と同様に、シーンを認識すべき対象である全体の静止画像を入力とする第1のCNNと、全体の画像の一部に含まれる物体の詳細な特徴量を抽出する第2のCNNと、があり、これらの中間層出力と認識処理結果とを統合したうえで、シーンを認識する構成が考えられる。この際、例えば物体を検知する方法としては、上記の文献7の方法を用いてもよい。これらの具体的な処理の方法は、本実施形態ですでに述べた認識対象人物の行動認識に関する構成・動作例によって説明ができるため、詳細な処理の説明を割愛する。

30

40

【0125】

なお、統合処理の一例として、図25を用いてもよい。2801は第1のCNNの1層が中間層特徴マップを出力する例を示し(例えば1層におけるReLU処理後の中間層特徴マップ)、2802~2804はそれぞれ第2のCNN(注視領域は異なってよい)の1層が中間層特徴マップを出力する例を示している。このとき、第1のCNNの中間層特徴マップに、第2のCNNの中間層特徴マップを統合するとする。まず2805において、中間層特徴マップ2801を受け取る(2814を伝達して送信される)。ここでは、中間層特徴マップ2801に対して中間層特徴マップ2802をROI embeddingする。次に、2808で非線形処理をかける。このとき非線形処理として、例えばBN-ReLU-Convなど、文献4などで記載の方法を用いてもよい(なにも処理をしなくてもよい。例えば、ここで

50

は2801の出力を足し合わせる場合を例示したが、0で埋められた仮の特徴マップにROI embeddingしてもよいし、そのあとに非線形処理をかけなくてもよい)。2803、2804は2802と同様であり、2806、2807は2805と同様であり、2809、2810は2808と同様であり、2815、2816は2814と同様である。また、3つ以上の第2のCNNが存在する場合を想定して、2818として例示した。これらの出力(中間層特徴マップの統合処理における中間状態)と、さらに2801の中間層特徴マップを、2811で統合する。このときの統合は、例えば足し算を用いてもよい(ここでは次元数は同じとする。次元数が異なる場合は別記参照のこと)。上記の足し算された結果が、2812で非線形処理にかけられる(ここでは非線形処理をかける例を示したが、かけなくてもよいし、必要であれば線形処理でもよい)。その後、2817でReLU処理にかけられ、2813で中間層特徴マップ2801とさらに足しあわされたうえで、次の第1のCNNの中間層への入力として用いられる。なお、上記の構成は、文献4で示されるようなResidual Networkをマルチパーツストリームに拡張した例として考えることができる。

10

20

30

40

50

【0126】

なお、全体の統合処理アーキテクチャの例として、図26を示す。2901~2903は、それぞれ認識対象人物、当該人物の左手、当該人物の右手の画像を示し、それぞれ第1のCNN2904(全身)、第2のCNN2905(左手)・2906(右手)に入力されるものとする。このとき、2908によってCNNストリーム(左手)の中間層特徴マップが第1のCNN(全身)の中間層の統合処理部に送られることで統合され、また2905によって同様に第2のCNNストリーム(右手)の中間層特徴マップが第1のCNNに送られ統合される。図26では、このような処理を複数の層で行うことができることを例示している。これらの統合結果が、2910によってRNN2907に送られる。RNNはLong-short term memoryなど、その他の公知のモデルでもよい。RNNを導入することで、このように、本実施形態におけるアーキテクチャをより長期の時間方向に対して拡張することができる。なお、2911は第2のCNNストリームがさらにあってもよいことを示している。

【0127】

なお、全体の統合処理アーキテクチャの例として、図28を示す。図28の第1のCNN3120は、本実施形態に係る第1のCNNのネットワーク構造の一例を示したものである。図28に示す如く、第1のCNN3120は、入力層3101、convolution1層3102、pooling1層3103、convolution2層3104、pooling2層3105、Innerproduct1層3106、Innerproduct2層3107、出力層3108、を有している。また、2つの階層間の処理方法として、convolution処理3110、pooling処理3111、InnerProduct処理3112が設定されている。それぞれの処理の具体的な内容は周知であり、例えば図13と同様であるため、ここでは省略する。また、第2のCNN3220は構成としては第1のCNNとここでは同じであるとし(違ってもよい)、3202~3212はそれぞれ3102~3112に対応しているとする。このとき、第2のCNNの1層目の出力(convolution1層3202がかけられた後や、さらにReLUをかけた後でもよい)を統合処理する例を3301に示す。3301は第2のCNNの中間層特徴マップから第1のCNNの中間層特徴マップに送信され、統合処理が行われる例である。

【0128】

以上説明したように、本実施形態で示した一形態は、データ中の認識対象人物やその状態を認識するCNNと詳細CNNとがあり、CNNの中間層出力と詳細CNNの認識処理結果とを統合するための学習を行う。このとき、一つのCNNに対して複数の詳細CNNがあり、複数の認識処理結果を統合する処理を行う。これにより、認識対象を精度よく認識するための複数の詳細な特徴量に基づいて認識処理を行うことができ、最終的な認識結果(異常検知結果)の精度が向上しうる。

【 0 1 2 9 】

[第 3 の 実 施 形 態]

第 1 , 2 の 実 施 形 態 で は 、 第 1 の C N N の 中 間 層 出 力 と 、 第 2 の C N N の 詳 細 認 識 処 理 結 果 と 、 を 統 合 す る こ と に よ っ て 、 単 一 の 目 的 に お け る 学 習 ・ 認 識 処 理 を 行 う 例 を 示 し た 。

【 0 1 3 0 】

本 実 施 形 態 で は 、 上 記 の 例 に 加 え て 、 複 数 の 目 的 関 数 を 用 い て 、 複 数 の 目 的 に お け る 学 習 ・ 認 識 処 理 を 行 う 例 を 示 す 。 具 体 的 に は 、 学 習 時 に 第 1 の C N N と 第 2 の C N N と で 異 な る 目 的 関 数 を 持 つ マ ル チ タ ス ク 学 習 を 行 う 場 合 の 例 を 示 す 。

【 0 1 3 1 】

な お 、 本 実 施 形 態 で 示 す 構 成 は 第 1 , 2 の 実 施 形 態 で 例 示 し た 構 成 と 大 部 分 が 同 一 で あ り 、 一 部 の 構 成 と 動 作 が 異 な る 。 本 実 施 形 態 に お い て 第 1 , 2 の 実 施 形 態 と 構 成 や 動 作 が 大 き く 異 な る の は 、 先 述 の よ う に 、 第 1 の C N N と 一 つ 以 上 の 第 2 の C N N と が あ る 構 成 に お い て 、 マ ル チ タ ス ク 学 習 を 行 う 点 で あ る 。

【 0 1 3 2 】

本 実 施 形 態 に お い て は 、 一 つ の 構 成 例 と し て 、 第 1 の 実 施 形 態 の 構 成 に お い て マ ル チ タ ス ク 学 習 を 行 う 場 合 の 構 成 お よ び 動 作 に 関 し て 例 示 す る 。 す な わ ち 、 認 識 対 象 人 物 の 全 身 の 特 徴 量 を 抽 出 す る 第 1 の C N N と 、 頭 部 の 詳 細 な 特 徴 量 を 抽 出 す る 第 2 の C N N と 、 が あ り 、 第 2 の C N N は 第 1 の C N N と は 異 な る 目 的 関 数 に 基 づ い て 学 習 処 理 を 行 う と す る 。 こ れ に よ っ て 、 例 え ば 全 身 の 行 動 と 頭 部 の 行 動 と で 、 異 な る 意 味 を 持 つ 行 動 を 別 々 に 識 別 処 理 す る た め の 学 習 が 可 能 に な り 、 ま た 、 そ れ ぞ れ の タ ス ク に 必 要 な 特 徴 量 を 内 部 的 に 統 合 す る 処 理 が な さ れ る こ と で 、 最 終 の 認 識 処 理 の 精 度 が 向 上 す る こ と が 期 待 で き る 。

【 0 1 3 3 】

全 身 の 認 識 処 理 を 行 う 第 1 の C N N の 目 的 関 数 と 、 頭 部 の 詳 細 認 識 処 理 を 行 う 第 2 の C N N の 目 的 関 数 と 、 を 以 下 に 式 1 0 , 1 1 と し て 示 す 。

【 0 1 3 4 】

【 数 9 】

$$\min Loss_{Body} = \Sigma (t_{Body} - y_{Body})^2 \cdot \cdot \cdot \cdot (10)$$

$$\min Loss_{Head} = \Sigma (t_{Head} - y_{Head})^2 \cdot \cdot \cdot \cdot (11)$$

【 0 1 3 5 】

こ こ で 、 $Loss_{Body}$ 、 $Loss_{Head}$ は それ ぞ れ 、 認 識 対 象 人 物 の 全 身 に 関 す る 認 識 処 理 を 行 う 第 1 の C N N の ロ ス 関 数 、 認 識 対 象 人 物 の 頭 部 に 関 す る 詳 細 認 識 処 理 を 行 う 第 2 の C N N の ロ ス 関 数 で あ る 。 t_{Body} 、 t_{Head} 、 y_{Body} 、 y_{Head} は それ ぞ れ 前 出 の 第 1 の C N N の 教 師 デ ー タ 、 第 2 の C N N の 教 師 デ ー タ 、 第 1 の C N N の 認 識 結 果 、 第 2 の C N N の 詳 細 認 識 処 理 結 果 を 表 す 。 な お 、 式 (1 0) 、 (1 1) の は Minibatch に 含 ま れ る デ ー タ 分 だ け 総 和 を 得 て 、 平 均 値 を 得 る 処 理 を 行 う 関 数 と す る 。 こ れ ら の ロ ス 関 数 は 平 均 二 乗 誤 差 と 呼 ば れ る 。 な お 、 平 均 二 乗 誤 差 で は な く 、 そ の 他 の ロ ス 関 数 を 用 い て も よ い 。 こ の と き 問 題 に な る の は 、 上 記 の 教 師 デ ー タ を ど の よ う な も の に す る か と 、 認 識 結 果 を ど の よ う に 得 る か (ネ ッ ト ワ ー ク ア ー キ テ ク チ ャ の 問 題) 、 と い う こ と で あ る 。

【 0 1 3 6 】

第 1 の C N N の 教 師 デ ー タ お よ び 認 識 結 果 は 、 第 1 の 実 施 形 態 で 示 し た 方 法 に よ っ て 得 る こ と が で き る 。 第 2 の C N N の 教 師 デ ー タ お よ び 認 識 結 果 は 、 第 1 の 実 施 形 態 で は 用 い な か っ た の で 、 こ こ で 図 1 お よ び 図 2 2 に 基 づ い て 新 た に 説 明 を 行 う 。

【 0 1 3 7 】

まず、第2のCNNの教師データは、第1のCNNの教師データと同様に図1の記憶部D1から読み込まれ、前段処理部11から注目領域制御部21に送信され、注目領域制御部21から認識部22aに送信されるものとする。ここで認識部22aは、第1の実施形態で用いる図1の認識部22の代わりに本実施形態で用いる機能であり、その構成の例を図22に示す。

【0138】

実行部222aは、実行部222とは異なり、注目領域制御部21から送られてきた入力データに対する詳細認識処理結果（具体例：ラベルの推定結果）と、第2のCNNの教師データとを、更新部223aに送る。更新部223aは、受け取った詳細認識処理結果と、第2のCNNの教師データとを用いて、誤差を算出する。ここで用いる誤差算出方法は、第1の実施形態で示した方法を用いてよいし、式11で例示した方法を用いてもよい。これらの誤差算出方法によって得られた誤差情報は、第1の実施形態で説明した方法によって学習に用いることができる。

【0139】

もし、式10, 11の2種の教師データが完全に同一のものであれば、式10, 11は異なる目的関数ではなく、同一の目的を持つ二つの目的関数となる。そのように学習処理を行ってもよいが、本実施形態では前述のとおり、異なる目的関数を持つ場合について説明を行うために、2種の教師データが異なる値を持ちうるという仮定のもとで説明を行う。例えば、第1のCNNは認識対象人物の状態が不審であるか否かを判定し、第2のCNNは認識対象人物が「キョロキョロ」行動をしたか否かを判定する、とする。このとき、式10, 11は異なる目的のもとで認識結果を評価するものの、「不審さ」と「キョロキョロ行動」は比較的相関の高い状態であると考えることができる。すなわち、式10, 11に基づくマルチタスク学習は、二つの目的に共通して有用な特徴を抽出するための機構を提供することが可能である。

【0140】

なお、上記のような教師データの組だけでなく、その他の教師データを用いて学習処理を行ってもよい。どのような教師データの組を用いるかは、解くべき問題に依存するため、検出の目的に応じて適宜変えてもよい。

【0141】

マルチタスク学習のためのネットワークアーキテクチャの構成方法として、例えば以下のような方法を用いてもよい。図19の上部は頭部に関する第2のCNNのストリーム、下部は全身に関する第1のCNNのストリームに関する構成の例を表している。このとき、上部のストリームの出力が、式11に基づいて評価され、下部のストリームの出力が、式10に基づいて評価されるとする。このような場合、第1のCNNと第2のCNNとで、それぞれ式10, 11の異なる目的関数を用いる構成になり、すなわちマルチタスク学習の動作が可能な構成になる。なお、検出時には、目的に応じて、第1のCNNと第2のCNNの出力を用いてよい。具体的には、検出時に認識対象人物の「不審さ」を出力すべきであるなら第1のCNNの出力を用いることにしてよく、また「キョロキョロ行動」をしたか否かを出力する必要があるれば、第2のCNNの出力を用いることにしてもよく、また必要であれば両方の出力結果を用いてもよい。なお、上述の方法以外のマルチタスク学習の構成例を用いて、第1のCNNと第2のCNNのマルチタスク学習を行ってもよい。

【0142】

以上説明したように、本実施形態で示した一形態は、データ中の認識対象人物やその状態を認識するCNNと詳細CNNとがあり、CNNの中間層出力と詳細CNNの詳細認識処理結果とを統合するための学習を行う。このとき、CNNと（複数の）詳細CNNとで、それぞれ異なる目的関数を用いる場合について示した。これにより、認識対象を精度よく認識するための複数の詳細な特徴量に基づいて認識処理を行うことができ、最終的な認識結果の精度が向上しうる。

【0143】

[第4の実施形態]

10

20

30

40

50

図 1 ~ 4 , 8 , 1 1 , 1 4 ~ 1 7 , 2 2 に示した各機能部はハードウェアで実装しても良いが、一部をソフトウェア（コンピュータプログラム）で実装しても良い。例えば、記憶部として説明した機能部をメモリで実装し、それ以外の機能部をソフトウェアで実装しても良い。後者の場合、このようなソフトウェアを実行可能なコンピュータ装置は上記の異常検知システムに適用可能である。

【 0 1 4 4 】

異常検知システムに適用可能なコンピュータ装置のハードウェア構成例について、図 2 3 のブロック図を用いて説明する。なお、図 2 3 に示したハードウェア構成は、以上説明した異常検知システムに適用可能なコンピュータ装置のハードウェア構成の一例に過ぎず、その構成は適宜変更しても構わない。

10

【 0 1 4 5 】

なお、以上説明した異常検知システムは、認識装置 1 0 、詳細認識装置 2 0 、端末装置 4 0 、それ以外の記憶部等の各装置を別個の装置としても良いし、一部の装置を一体化させても良い。然るに図 2 3 に示した構成を有するコンピュータ装置は、認識装置 1 0 、詳細認識装置 2 0 、端末装置 4 0 のそれぞれ若しくはその一部に適用しても良いし、異常検知システムを構成する装置のうち一部を一体化させた装置に適用しても良い。なお、以上説明した異常検知システムをどのような装置で構成するのかについては、適宜変形例が考えられる。

【 0 1 4 6 】

C P U 2 6 0 1 は、R A M 2 6 0 2 や R O M 2 6 0 3 に格納されているコンピュータプログラムやデータを用いて処理を実行する。これにより C P U 2 6 0 1 は、コンピュータ装置全体の動作制御を行うと共に、該コンピュータ装置を適用した装置が行うものとして上述した各処理を実行若しくは制御する。

20

【 0 1 4 7 】

R A M 2 6 0 2 は、R O M 2 6 0 3 や外部記憶装置 2 6 0 6 からロードされたコンピュータプログラムやデータ、I / F 2 6 0 7 を介して外部から受信したデータ、等を格納するためのエリアを有する。さらに R A M 2 6 0 2 は、C P U 2 6 0 1 が各種の処理を実行する際に用いるワークエリアを有する。このように R A M 2 6 0 2 は、各種のエリアを適宜提供することができる。R O M 2 6 0 3 には、コンピュータ装置の基本プログラムや設定データなど、書換不要のコンピュータプログラムやデータが保存されている。

30

【 0 1 4 8 】

操作部 2 6 0 4 は、キーボードやマウスなどのユーザインターフェースにより構成されており、ユーザが操作することで各種の指示を C P U 2 6 0 1 に対して入力することができる。

【 0 1 4 9 】

表示部 2 6 0 5 は、C R T や液晶画面などにより構成されており、C P U 2 6 0 1 による処理結果を画像や文字などでもって表示することができる。なお、操作部 2 6 0 4 と表示部 2 6 0 5 とを一体化させてタッチパネル画面を構成しても良い。

【 0 1 5 0 】

外部記憶装置 2 6 0 6 は、ハードディスクドライブ装置等の大容量情報記憶装置である。外部記憶装置 2 6 0 6 には、O S（オペレーティングシステム）や、本コンピュータ装置を適用した装置が行うものとして上述した各処理を C P U 2 6 0 1 に実行させるためのコンピュータプログラムやデータが保存されている。外部記憶装置 2 6 0 6 に保存されているコンピュータプログラムには、本コンピュータ装置を適用した装置の各機能部の機能を C P U 2 6 0 1 に実行させるためのコンピュータプログラムが含まれている。また、外部記憶装置 2 6 0 6 に保存されているデータには、上記の説明において既知の情報として説明したもの（パラメータや閾値、関数など）が含まれている。外部記憶装置 2 6 0 6 に保存されているコンピュータプログラムやデータは、C P U 2 6 0 1 による制御に従って適宜 R A M 2 6 0 2 にロードされ、C P U 2 6 0 1 による処理対象となる。

40

【 0 1 5 1 】

50

I / F 2 6 0 7 は、外部の機器との間のデータ通信を行うためのインターフェースとして機能するものである。C P U 2 6 0 1、R A M 2 6 0 2、R O M 2 6 0 3、操作部 2 6 0 4、表示部 2 6 0 5、外部記憶装置 2 6 0 6、I / F 2 6 0 7 は何れもバス 2 6 0 8 に接続されている。

【 0 1 5 2 】

なお、以上説明した各実施形態の一部若しくは全部を適宜組み合わせ使用しても構わないし、以上説明した各実施形態の一部若しくは全部を選択的に使用しても構わない。また、以上の説明において使用した数値や構成や処理順は説明上一例としてあげたものであり、以上の説明に限らない。

【 0 1 5 3 】

なお、以上説明した異常検知システムは、次のような各部を有する構成を有する情報処理装置の一例として説明したものである。

【 0 1 5 4 】

- ・ 認識対象の情報を含む認識対象データを入力として畳み込みニューラルネットワークの認識処理の結果を取得する認識処理部
- ・ 認識対象データおよび / または畳み込みニューラルネットワークの中間層出力に対する注目領域を設定する注目領域設定部
- ・ 注目領域中の認識対象データおよび / または中間層出力に対して認識処理よりも詳細な詳細認識処理を行う詳細認識部
- ・ 詳細認識処理の結果と、中間層出力と、を統合処理する統合処理部
- ・ 統合処理の結果を中間層出力として畳み込みニューラルネットワークに入力する中間入力処理部
- ・ 認識処理の結果を出力する出力部

また、以上説明した異常検知システムは、次のような各部を有する構成を有する情報処理装置の一例として説明したものである。

【 0 1 5 5 】

- ・ 学習対象の情報を含む学習データを入力として畳み込みニューラルネットワークの認識処理の結果を取得する認識処理部
- ・ 学習データおよび / または畳み込みニューラルネットワークの中間層出力に対する注目領域を設定する注目領域設定部
- ・ 注目領域中の学習データおよび / または中間層出力に対して認識処理よりも詳細な詳細認識処理を行う詳細認識部
- ・ 詳細認識処理の結果と、中間層出力と、を統合処理する統合処理部
- ・ 統合処理の結果を中間層出力として畳み込みニューラルネットワークに入力する中間入力処理部
- ・ 認識処理部、注目領域設定部、詳細認識部、統合処理部、中間入力処理部、のいずれか一つ以上に関する学習処理を行う学習部

(その他の実施例)

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、A S I C) によっても実現可能である。

【符号の説明】

【 0 1 5 6 】

1 0 : 認証装置 2 0 : 詳細認識装置 D 1 : 記憶部 M 1 : 記憶部

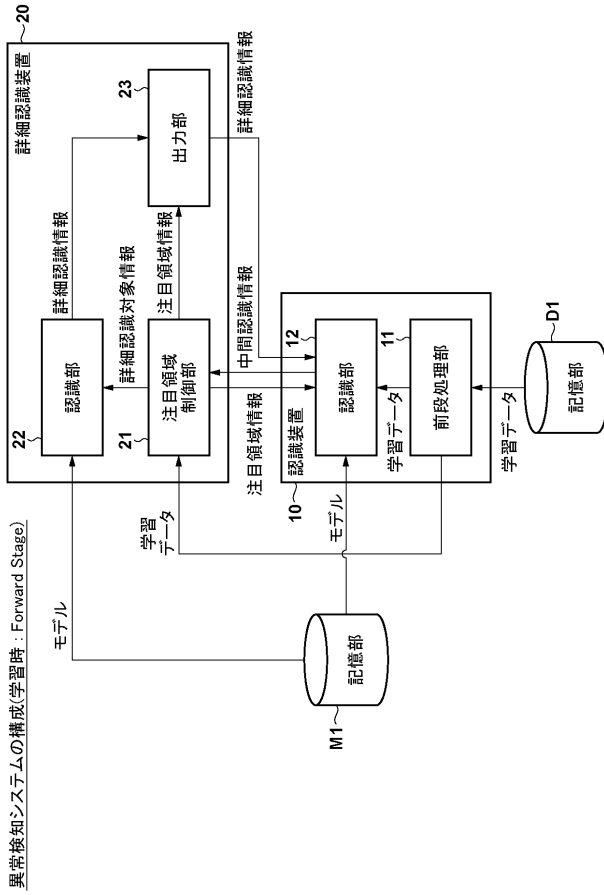
10

20

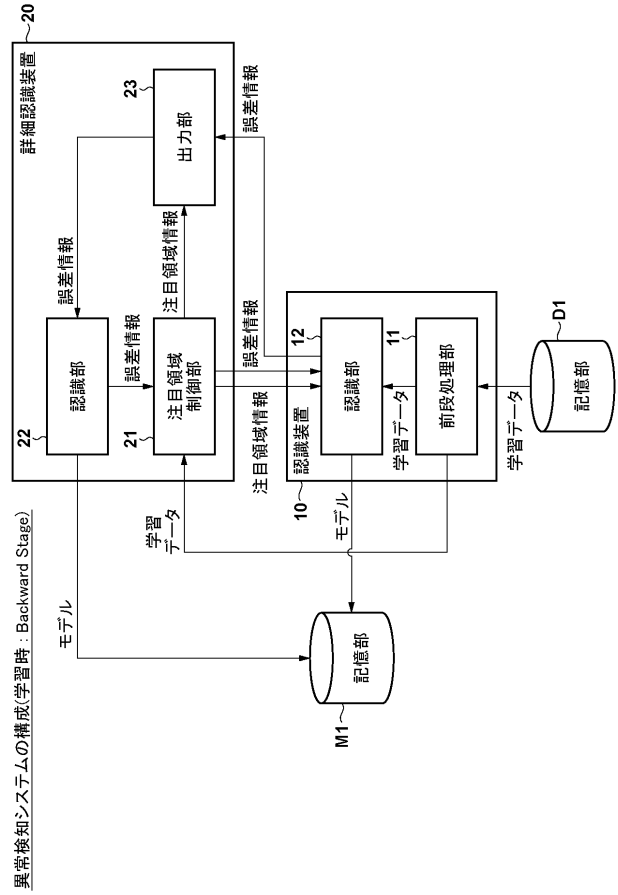
30

40

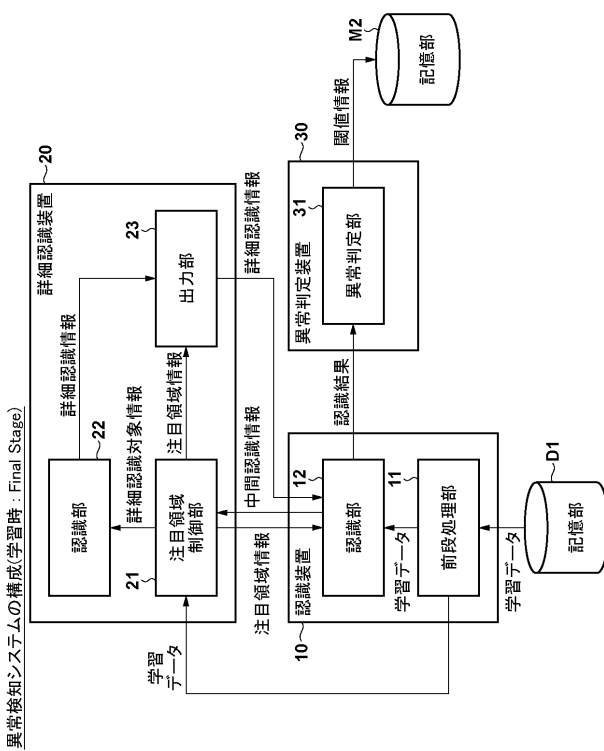
【 図 1 】



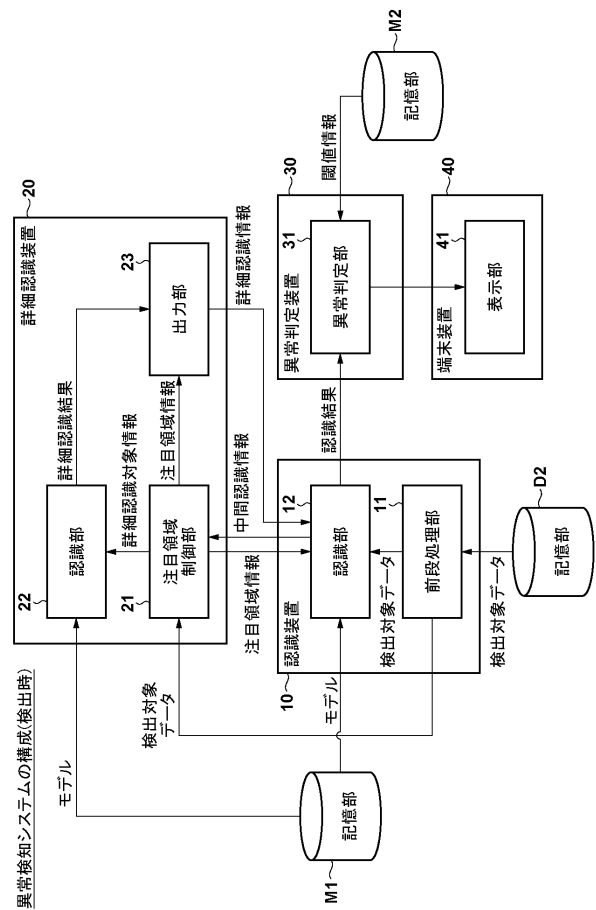
【 図 2 】



【 図 3 】

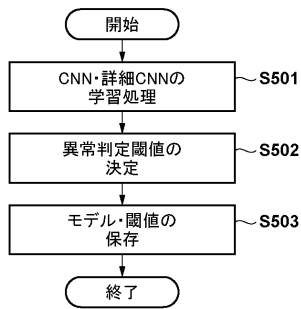


【 図 4 】



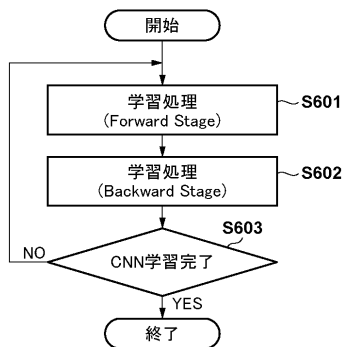
【図 5】

異常検知システムの動作(学習時)



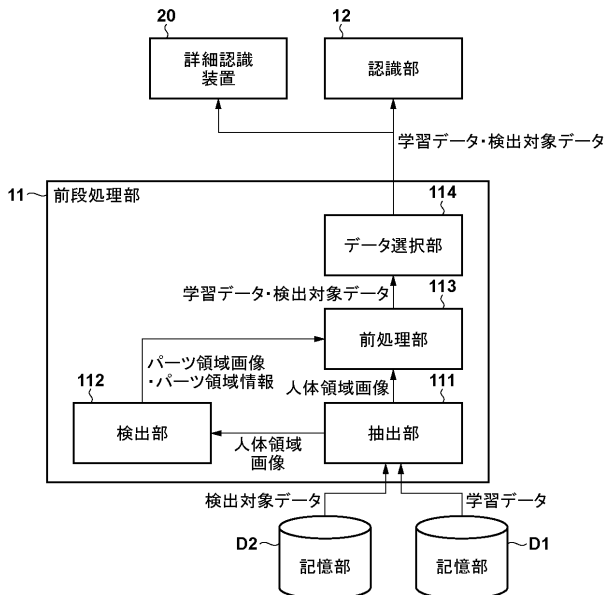
【図 6】

ステップS501の詳細動作



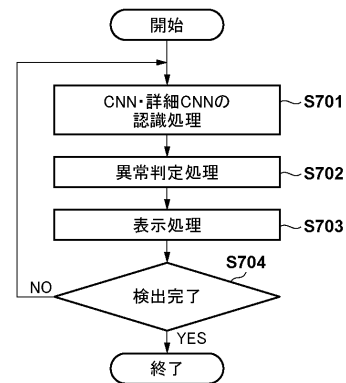
【図 8】

前段処理部11の構成(学習時・検出時併記)



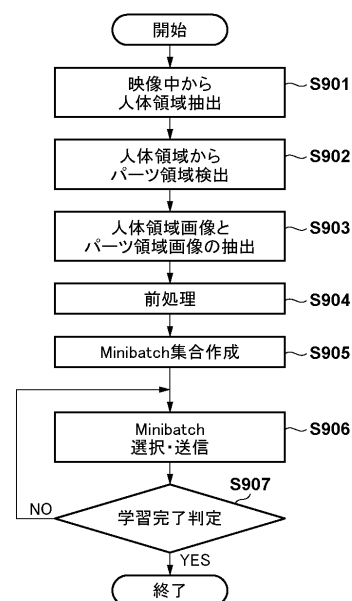
【図 7】

異常検知システムの動作(検出時)

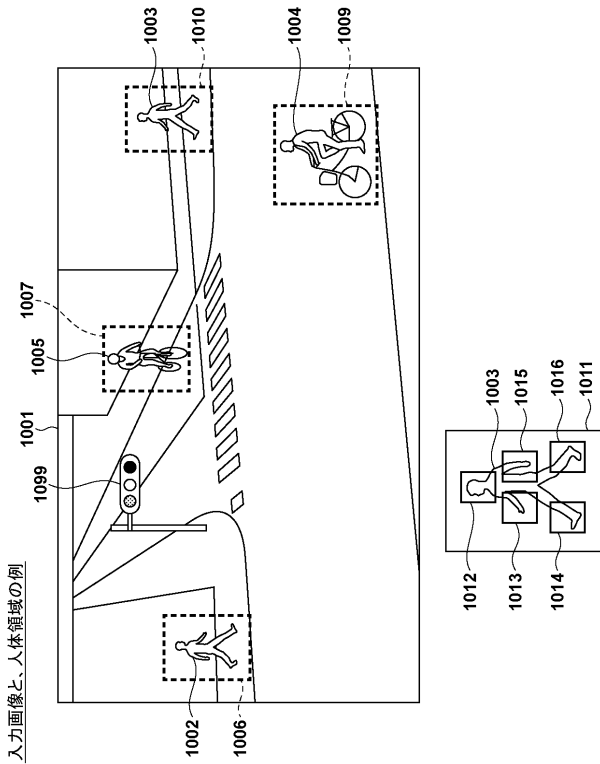


【図 9】

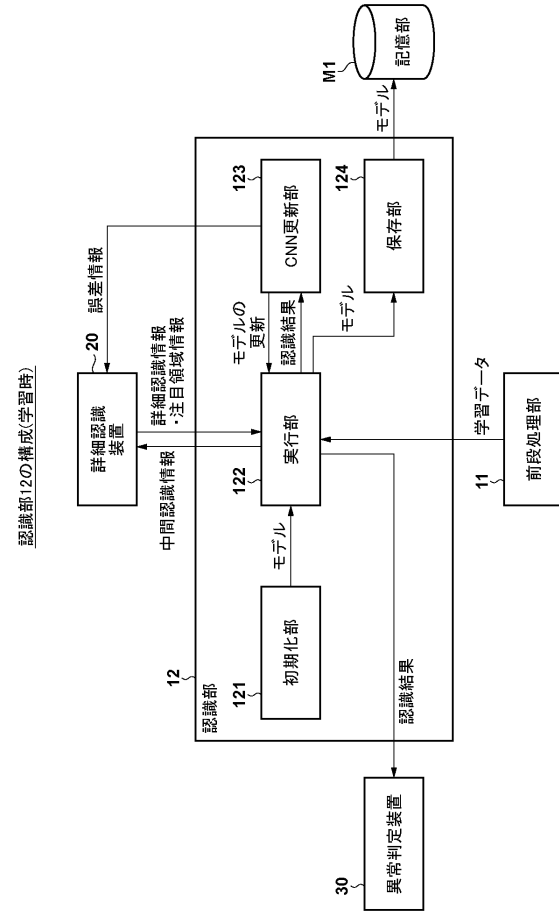
前段処理部11の動作(学習時)



【図 10】

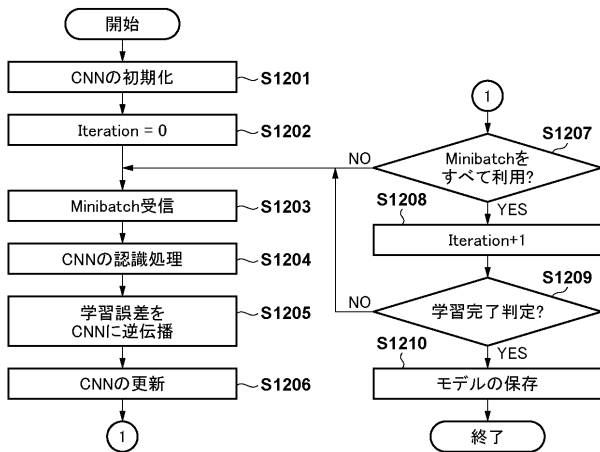


【図 11】

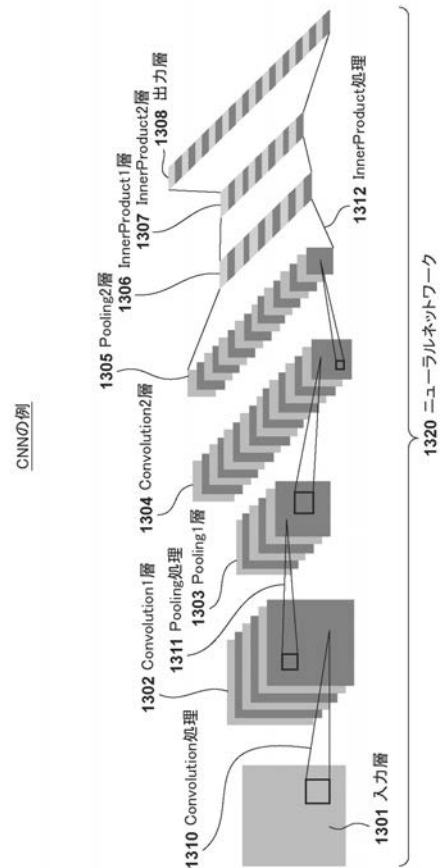


【図 12】

認識部12の動作(学習時)

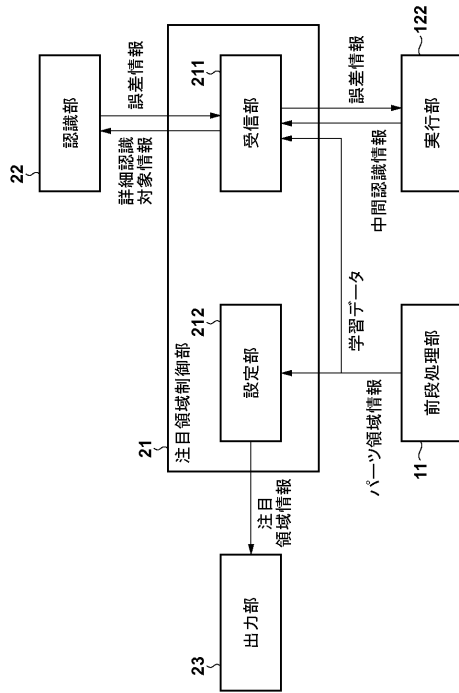


【図 13】



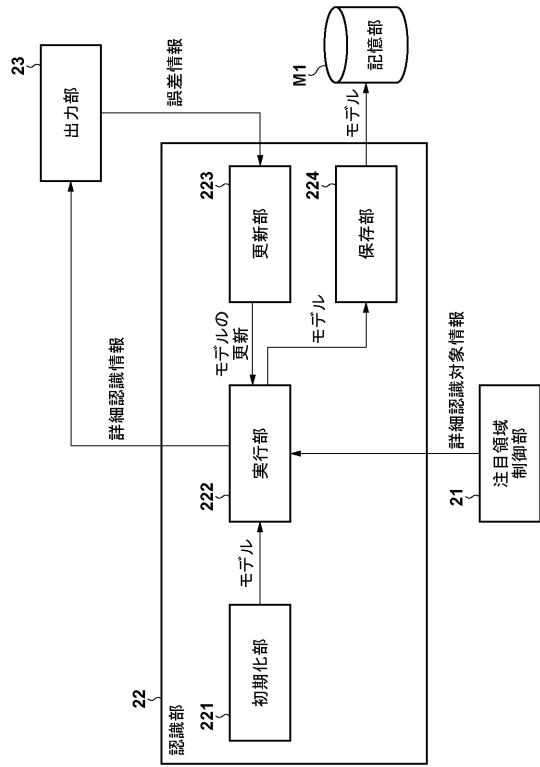
【図 14】

注目領域制御部21の構成(学習時)



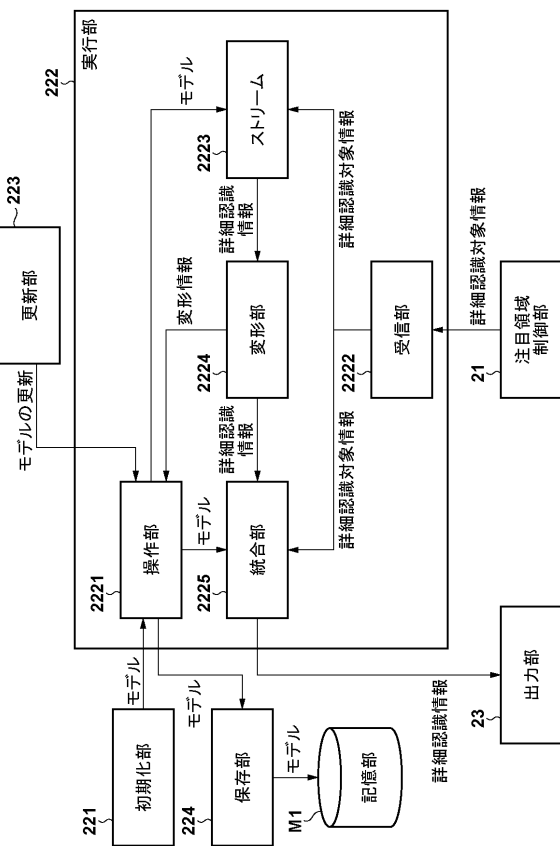
【図 15】

認識部22の構成(学習時)



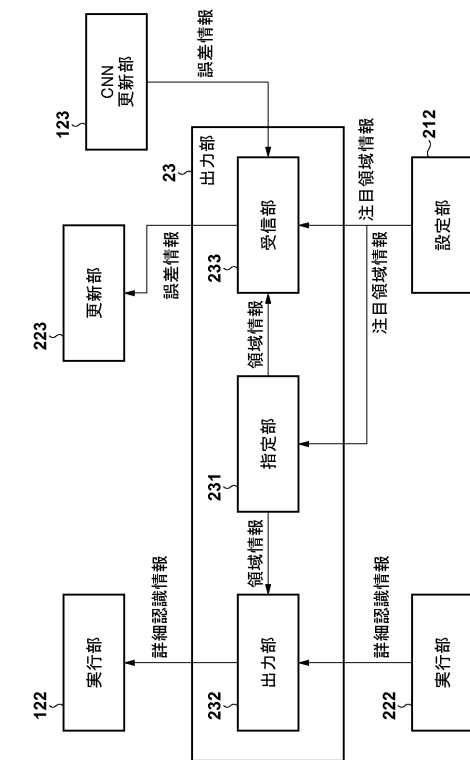
【図 16】

実行部222の構成(学習時)



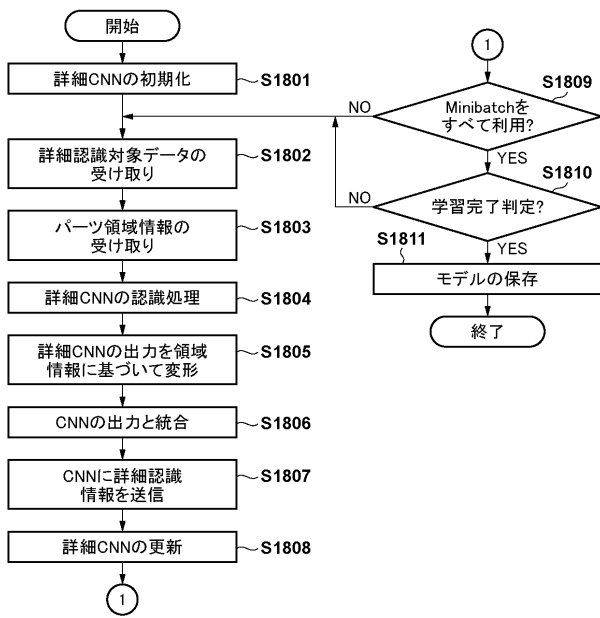
【図 17】

出力部23の構成(学習時)

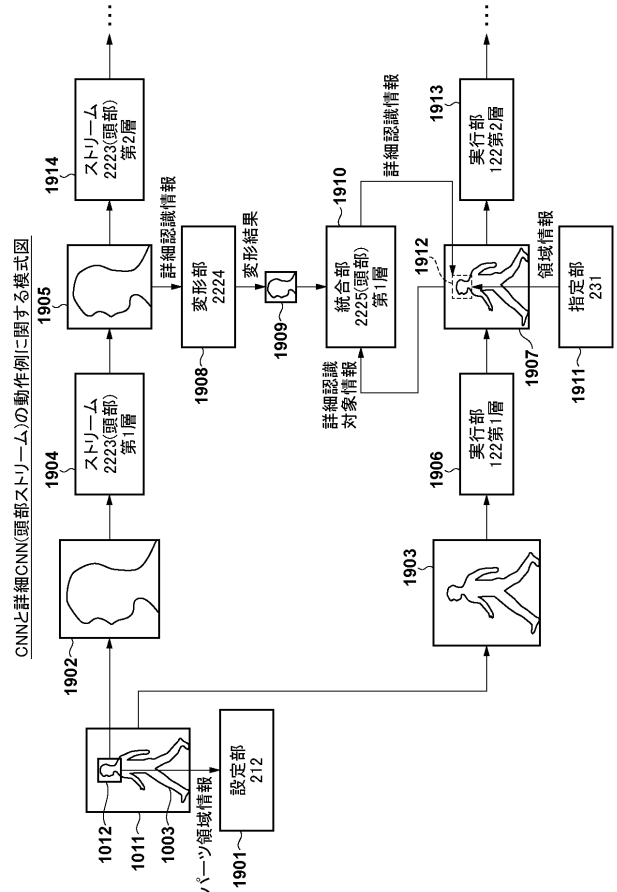


【図 18】

詳細CNN(学習時)の動作

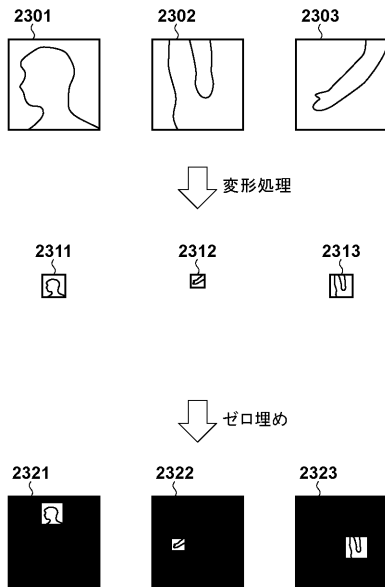


【図 19】



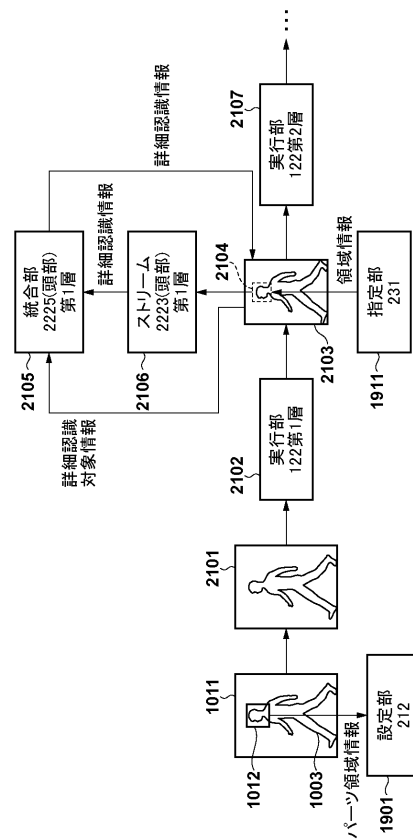
【図 20】

ゼロ埋めによる複数パーツデータ成形

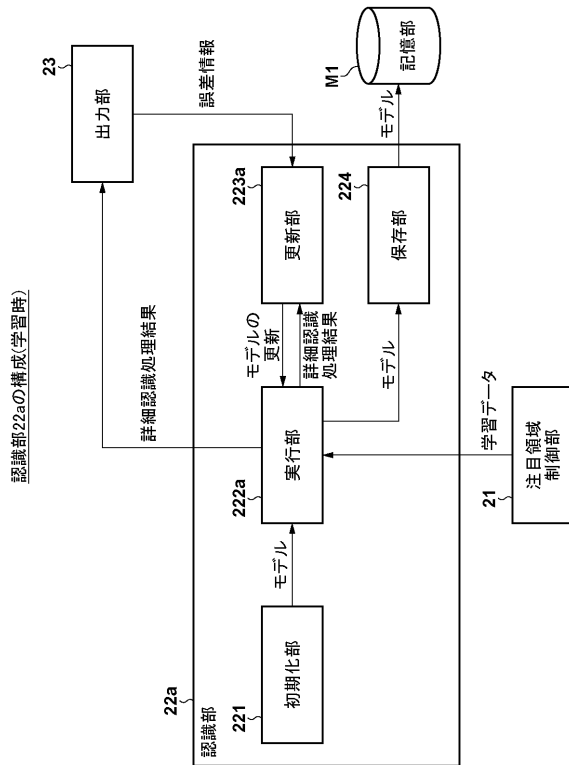


【図 21】

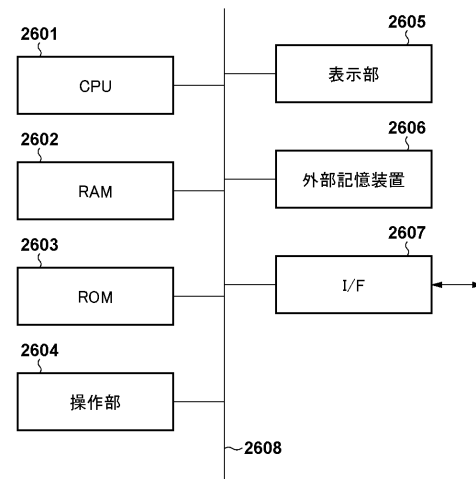
CNNと詳細CNN(頭部ストリーム)の動作例に関する模式図(ver.2)



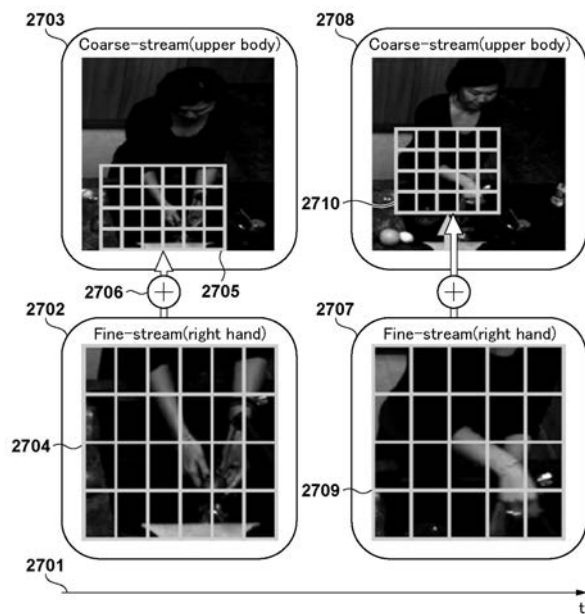
【図 2 2】



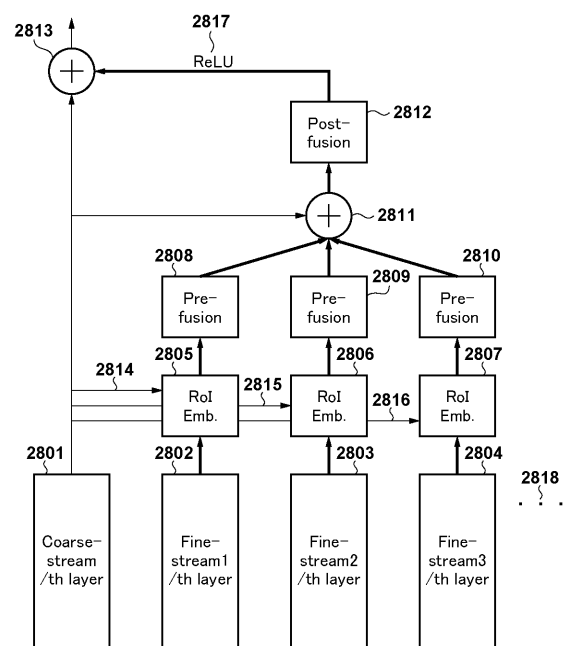
【図 2 3】



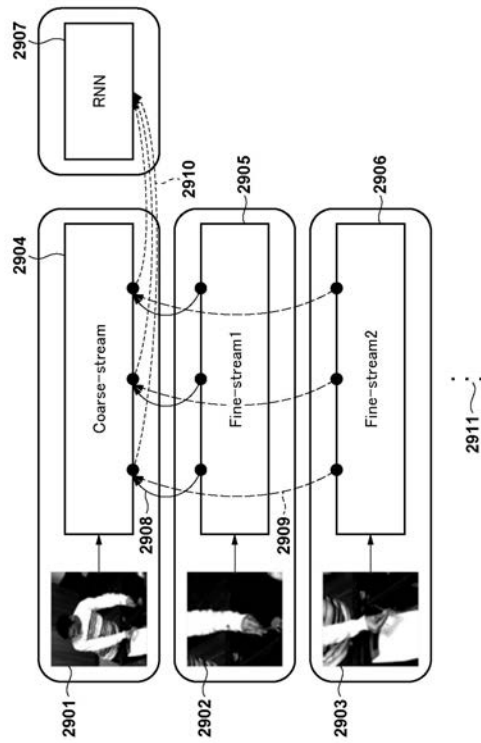
【図 2 4】



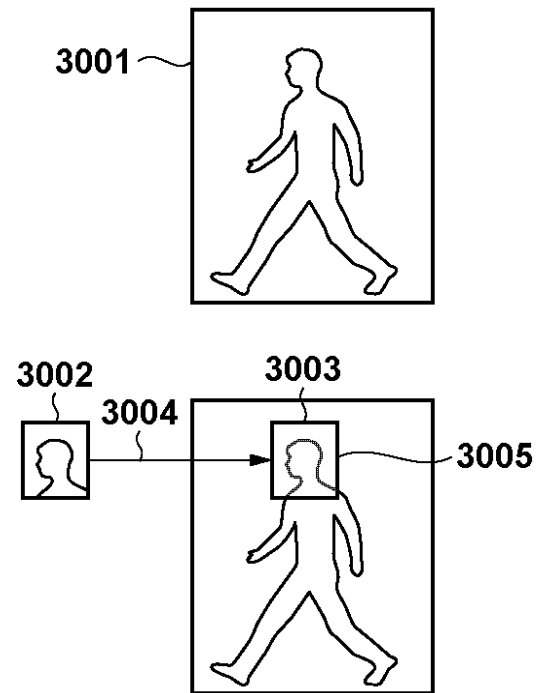
【図 2 5】



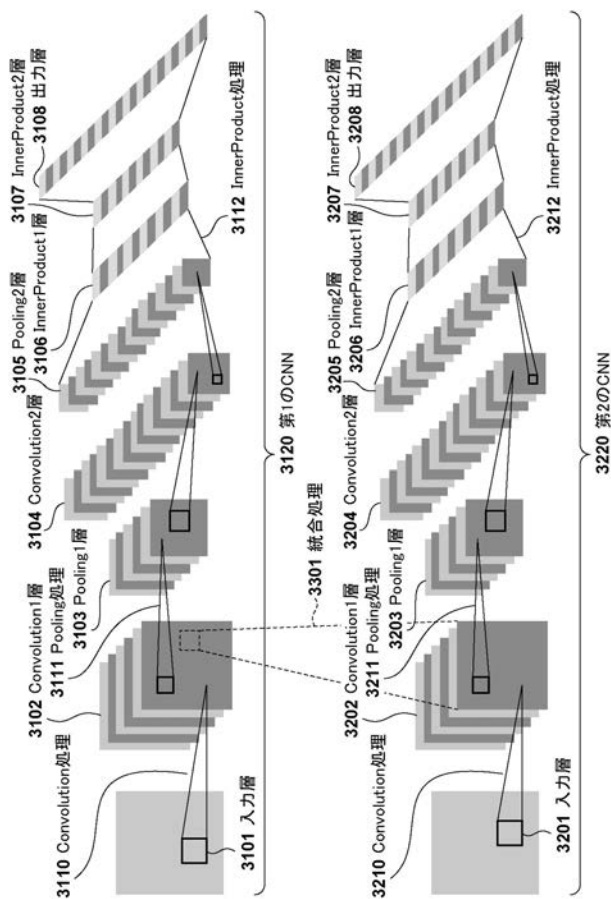
【図 26】



【図 27】



【図 28】



フロントページの続き

(72)発明者 斎藤 侑輝

東京都大田区下丸子3丁目30番2号 キヤノン株式会社内

(72)発明者 小森 康弘

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 BA02 BA18 FA16 FA69 GA51 HA11 JA11 KA04 KA15