

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7389575号
(P7389575)

(45)発行日 令和5年11月30日(2023.11.30)

(24)登録日 令和5年11月21日(2023.11.21)

(51)国際特許分類

F I

G 0 6 N 3/045(2023.01)

G 0 6 N 3/045

請求項の数 18 (全24頁)

(21)出願番号	特願2019-123945(P2019-123945)	(73)特許権者	000001007 キャノン株式会社 東京都大田区下丸子3丁目30番2号
(22)出願日	令和1年7月2日(2019.7.2)	(74)代理人	110003281 弁理士法人大塚国際特許事務所
(65)公開番号	特開2021-9622(P2021-9622A)	(72)発明者	舘 俊太 東京都大田区下丸子3丁目30番2号 キャノン株式会社内
(43)公開日	令和3年1月28日(2021.1.28)	(72)発明者	チン ソクイ 東京都大田区下丸子3丁目30番2号 キャノン株式会社内
審査請求日	令和4年7月1日(2022.7.1)	審査官	千葉 久博
		最終頁に続く	

(54)【発明の名称】 データ処理装置、データ処理方法、学習装置、学習方法、ニューラルネットワーク、及びプログラム

(57)【特許請求の範囲】

【請求項1】

ニューラルネットワークの第1の層と接続される第2の層の第3の部分の特徴量データを算出するために参照される前記第1の層の第1の部分と、前記第2の層の第4の部分の特徴量データを算出するために参照される前記第1の層の第2の部分と、をそれぞれ規定する結合パラメータを取得する取得手段と、

前記取得手段によって取得された前記結合パラメータに規定される前記第1の部分の特徴量データから前記第3の部分の特徴量データを算出する第1の算出処理と、前記結合パラメータに規定される前記第2の部分の特徴量データから前記第4の部分の特徴量データを算出する第2の算出処理とを並列に行う演算手段と、を備える

ことを特徴とするデータ処理装置。

【請求項2】

前記第1の層の前記第1の部分と前記第2の部分とは重複しないことを特徴とする、請求項1に記載のデータ処理装置。

【請求項3】

前記第1の層の前記第1の部分と前記第2の部分とが部分的に重複することを特徴とする、請求項1に記載のデータ処理装置。

【請求項4】

前記第2の層の前記第3の部分の大きさと前記第4の部分の大きさとが異なることを特徴とする、請求項1から3のいずれか1項に記載のデータ処理装置。

【請求項 5】

前記結合パラメータは、チャンネル単位で前記第 1 の層の前記第 1 の部分及び前記第 2 の部分を規定する

ことを特徴とする、請求項 1 から 4 のいずれか 1 項に記載のデータ処理装置。

【請求項 6】

前記結合パラメータは、複数のチャンネルを含むブロック単位で前記第 1 の層の前記第 1 の部分及び前記第 2 の部分を規定する

ことを特徴とする、請求項 5 に記載のデータ処理装置。

【請求項 7】

前記演算手段は、前記第 1 の層の前記第 1 の部分に含まれる全てのチャンネルの特徴量データを用いて、かつ前記第 1 の層の前記第 2 の部分に含まれるチャンネルの特徴量データを用いず、前記第 2 の層の第 3 の部分に含まれるそれぞれのチャンネルの特徴量データを算出する

ことを特徴とする、請求項 5 又は 6 に記載のデータ処理装置。

【請求項 8】

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又は再帰的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 から 7 のいずれか 1 項に記載のデータ処理装置。

【請求項 9】

前記取得手段は、前記ニューラルネットワークに入力される学習データに基づいて学習された結合パラメータを取得する

ことを特徴とする、請求項 1 から 8 のいずれか 1 項に記載のデータ処理装置。

【請求項 10】

学習データと、前記学習データに対する処理結果を示す教師データと、を取得する取得手段と、

ニューラルネットワークに前記学習データを入力することにより、前記学習データに対する処理結果を得るデータ処理手段であって、前記ニューラルネットワークの第 1 の層と接続される第 2 の層の第 3 の部分の特徴量データを算出するために参照される前記第 1 の層の第 1 の部分と、前記第 2 の層の第 4 の部分の特徴量データを算出するために参照される前記第 1 の層の第 2 の部分と、をそれぞれ規定する結合パラメータに規定される前記第 1 の部分の特徴量データから前記第 3 の部分の特徴量データを算出する第 1 の算出処理と、前記結合パラメータに規定される前記第 2 の部分の特徴量データから前記第 4 の部分の特徴量データを算出する第 2 の算出処理とを並列に行うデータ処理手段と、

前記学習データに対する処理結果と、前記教師データと、に基づいて、前記結合パラメータ及び前記ニューラルネットワークの階層間の重み係数の学習を行う学習手段と、を備える

ことを特徴とする学習装置。

【請求項 11】

前記学習手段は、第 1 の結合パラメータに従って前記データ処理手段が得た前記学習データに対する処理結果と、第 2 の結合パラメータに従って前記データ処理手段が得た前記学習データに対する処理結果と、に基づいて前記結合パラメータの学習を行う

ことを特徴とする、請求項 10 に記載の学習装置。

【請求項 12】

前記学習手段は、前記第 2 の層の前記第 1 の部分及び前記第 2 の部分の大きさを学習により決定する

ことを特徴とする、請求項 10 又は 11 に記載の学習装置。

【請求項 13】

前記取得手段は、さらに、学習により得られる前記結合パラメータに対する制約条件を指示するユーザ入力を取得し、

前記学習手段は、前記制約条件に従って前記結合パラメータの学習を行う

10

20

30

40

50

ことを特徴とする、請求項 10 から 12 のいずれか 1 項に記載の学習装置。

【請求項 14】

順に接続された前階層、第 1 の層、及び第 2 の層を有するニューラルネットワークであって、

前記第 1 の層は、前記前階層の一部のニューロンセットと結合した第 1 のニューロンセットと、前記前階層の一部とは異なる前記前階層の他の一部のニューロンセットと結合した第 2 のニューロンセットと、を有し、

前記第 2 の層は、結合パラメータにより規定される前記第 1 のニューロンセットの一部及び前記第 2 のニューロンセットの一部と結合した第 3 のニューロンセットと、前記結合パラメータにより規定され前記第 1 のニューロンセットの一部とは異なる前記第 1 のニューロンセットの他の一部及び前記第 2 のニューロンセットの一部とは異なる前記第 2 のニューロンセットの他の一部と結合した第 4 のニューロンセットと、を有し、

前記ニューラルネットワークは、前記結合パラメータを取得し、前記結合パラメータに規定される前記第 1 のニューロンセットの前記一部及び前記第 2 のニューロンセットの前記一部における特徴量データから前記第 3 のニューロンセットにおける特徴量データを算出する第 1 の算出処理と、前記結合パラメータに規定される前記第 1 のニューロンセットの前記他の一部及び前記第 2 のニューロンセットの前記他の一部における特徴量データから前記第 4 のニューロンセットにおける特徴量データを算出する第 2 の算出処理とを並列に行うよう、コンピュータを機能させることを特徴とするニューラルネットワーク。

【請求項 15】

ニューラルネットワークの第 1 の層と接続される第 2 の層の第 3 の部分の特徴量データを算出するために参照される前記第 1 の層の第 1 の部分と、前記第 2 の層の第 4 の部分の特徴量データを算出するために参照される前記第 1 の層の第 2 の部分と、をそれぞれ規定する結合パラメータを取得する取得工程と、

前記結合パラメータに規定される前記第 1 の部分の特徴量データから前記第 3 の部分の特徴量データを算出する第 1 の算出処理と、前記結合パラメータに規定される前記第 2 の部分の特徴量データから前記第 4 の部分の特徴量データを算出する第 2 の算出処理とを並列に行う演算工程と、

を含むことを特徴とするデータ処理方法。

【請求項 16】

学習データと、前記学習データに対する処理結果を示す教師データと、を取得する取得工程と、

ニューラルネットワークに前記学習データを入力することにより、前記学習データに対する処理結果を得るデータ処理工程であって、前記ニューラルネットワークの第 1 の層と接続される第 2 の層の第 3 の部分の特徴量データを算出するために参照される前記第 1 の層の第 1 の部分と、前記第 2 の層の第 4 の部分の特徴量データを算出するために参照される前記第 1 の層の第 2 の部分と、をそれぞれ規定する結合パラメータに規定される前記第 1 の部分の特徴量データから前記第 2 の層の前記第 3 の部分の特徴量データを算出する第 1 の算出処理と、前記結合パラメータに規定される前記第 2 の部分の特徴量データから前記第 4 の部分の特徴量データを算出する第 2 の算出処理とを並列に行うデータ処理工程と、

前記学習データに対する処理結果と、前記教師データと、に基づいて、前記結合パラメータ及び前記ニューラルネットワークの階層間の重み係数の学習を行う学習工程と、

を備えることを特徴とする学習方法。

【請求項 17】

コンピュータを、請求項 1 から 9 のいずれか 1 項に記載のデータ処理装置として機能させるためのプログラム。

【請求項 18】

コンピュータを、請求項 10 から 13 のいずれか 1 項に記載の学習装置として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、データ処理装置、データ処理方法、学習装置、学習方法、ニューラルネットワーク、及びプログラムに関し、特に、ニューラルネットワーク演算の処理コストの低減に関する。

【背景技術】

【0002】

畳み込みニューラルネットワーク（以下CNN）のようなニューラルネットワークは複数の層を有しており、各層は複数のニューロンで構成される。通常のニューラルネットワークにおいては、各ニューロンは前階層の全てのニューロンから入力信号を受け取り、後階層の全てのニューロンへと信号を送る。このため、ニューロンとニューロンとの間の重み係数の数は、ニューロンの数の2乗に比例する。

10

【0003】

実時間でニューラルネットワークを用いた認識処理を行うためには、十分な処理速度が必要とされる。また、様々な機器でニューラルネットワークを用いた処理を行うためには、処理に必要なメモリ量を削減することが求められることがある。これらの目的のために、階層間の重み係数の数を削減する手法が知られている。

【0004】

特許文献1は、ニューラルネットワークを並列化する方法を開示している。特許文献1では、同じ入力画像が入力される複数の並列ニューラルネットワークが用いられる。それぞれの並列ニューラルネットワークは、相互結合層と非相互結合層とを有している。1つの並列ニューラルネットワークの相互結合層からの出力は全ての並列ニューラルネットワークに入力されるが、非相互結合層からの出力は同じ並列ニューラルネットワークにしか入力されない。

20

【先行技術文献】

【特許文献】

【0005】

【文献】米国特許第9811775号明細書

【発明の概要】

【発明が解決しようとする課題】

30

【0006】

特許文献1のように非相互結合層を用いることにより、全ての層において各ニューロンが前階層及び後階層の全てのニューロンと接続している通常の畳み込みニューラルネットワークと比較して、重み係数の数を減らすことができる。

【0007】

一方で、ニューラルネットワークを用いた認識精度は、それぞれの並列ニューラルネットワークにより抽出された特徴を組み合わせることによって向上すると考えられる。このため、相互結合層をニューラルネットワークの後段に配置することにより、認識精度が向上するものと考えられる。ここで、CNNのような一般的なニューラルネットワークにおいては、後段の層ほどニューロンの数が多い傾向がある。このため、特許文献1の手法に従ってニューラルネットワークの後段に相互結合層を設け、前段に非相互結合層を設けても、重み係数の削減効果があまり大きくないという課題がある。

40

【0008】

本発明は、例えばニューラルネットワークに従う処理の速度を向上させる又は処理に必要なメモリ量を削減する等の目的で、ニューラルネットワークで用いられる重み係数の数を削減することを目的とする。

【課題を解決するための手段】

【0009】

本発明の目的を達成するために、一実施形態に係るデータ処理装置は以下の構成を備える。すなわち、

50

ニューラルネットワークの第 1 の層と接続される第 2 の層の第 3 の部分の特徴量データを算出するために参照される前記第 1 の層の第 1 の部分と、前記第 2 の層の第 4 の部分の特徴量データを算出するために参照される前記第 1 の層の第 2 の部分と、をそれぞれ規定する結合パラメータを取得する取得手段と、

前記取得手段によって取得された前記結合パラメータに規定される前記第 1 の部分の特徴量データから前記第 3 の部分の特徴量データを算出する第 1 の算出処理と、前記結合パラメータに規定される前記第 2 の部分の特徴量データから前記第 4 の部分の特徴量データを算出する第 2 の算出処理とを並列に行う演算手段と、

を備える。

【発明の効果】

10

【 0 0 1 0 】

例えばニューラルネットワークに従う処理の速度を向上させる又は処理に必要なメモリ量を削減する等の目的で、ニューラルネットワークで用いられる重み係数の数を削減することができる。

【図面の簡単な説明】

【 0 0 1 1 】

【図 1】一実施形態に係るデータ処理装置の基本機能を示す図。

【図 2】一実施形態に係るニューラルネットワークの構成を示す図。

【図 3】一実施形態に係るデータ処理方法のフローチャート。

【図 4】一実施形態に係るデータ処理装置のハードウェア構成例を示す図。

20

【図 5】結合設定部 105 の動作例を示す図。

【図 6】様々な結合パラメータの例を示す図。

【図 7】一実施形態に係るデータ処理装置の基本機能を示す図。

【図 8】一実施形態に係る学習方法のフローチャート。

【図 9】結合パラメータの変更方法の例を示す図。

【図 10】一実施形態に係るデータ処理装置による処理方法を説明する図。

【図 11】コンピュータの構成例を示す図。

【図 12】一実施形態に係るデータ処理方法のフローチャート。

【図 13】一実施形態に係るニューロングループの設定例を示す図。

【図 14】遺伝的アルゴリズムを用いた学習方法を説明する図。

30

【図 15】再帰的ニューラルネットワークへの適用例を説明する図。

【発明を実施するための形態】

【 0 0 1 2 】

以下、添付図面を参照して実施形態を詳しく説明する。なお、以下の実施形態は特許請求の範囲に係る発明を限定するものではない。実施形態には複数の特徴が記載されているが、これらの複数の特徴の全てが発明に必須のものとは限らず、また、複数の特徴は任意に組み合わせられてもよい。さらに、添付図面においては、同一若しくは同様の構成に同一の参照番号を付し、重複した説明は省略する。

【 0 0 1 3 】

[実施形態 1]

40

まず、図 1 を参照して、本実施形態に係る処理の概要を説明する。本実施形態に係るニューラルネットワークは、階層型のニューラルネットワークであり、第 1 の層 111 と、第 2 の層 112 と、を有している。また、第 1 の層 111 及び第 2 の層 112 のそれぞれは、第 1 のニューロングループ 101 と第 2 のニューロングループ 102 とを有している。第 1 のニューロングループ 101 は、ニューロン演算 106 a の対象となるニューロンのグループである。また、第 2 のニューロングループ 102 は、ニューロン演算 106 b の対象となるニューロンのグループである。ここで、第 1 の層 111 と第 2 の層 112 との間のニューロン演算 106 a、106 b は、同じニューロングループ内のニューロンの間のみを結合する。図 1 の例において、ニューロングループ 101、102 は、第 1 の層 111 と第 2 の層 112 の両方にまたがって存在している。以下で説明するように、ニュー

50

ーロン演算 1 0 6 a , 1 0 6 b は、それぞれ別の演算ユニット (1 7 a 又は 1 7 b) が行うことができる。

【 0 0 1 4 】

ニューロン演算 1 0 6 a , 1 0 6 b は、ニューラルネットワークに従う演算である。ニューロン演算 1 0 6 a , 1 0 6 b としては、畳み込み演算又は活性化関数の適用のような、ニューラルネットワークにおける一般的な処理であってもよい。具体的な演算の種類としては、A. Krizhevsky et al. "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012. (以下Krizhevskyと呼ぶ) 及びM.D. Zeiler, R. Fergus "Visualizing and Understanding Convolutional Networks", ECCV 2014. (以下Zeilerと呼ぶ) 等に挙げられているものを用いることができる。

10

【 0 0 1 5 】

第 1 の層 1 1 1 の第 1 のニューロングループ 1 0 1 には、入力特徴 1 1 a と、入力特徴の参照順序を示す参照入力特徴 1 2 a とが模式的に示されている。第 2 の層 1 1 2 の第 1 のニューロングループ 1 0 1 には、出力特徴 1 3 a が示されている。同様に、第 1 の層 1 1 1 の第 2 のニューロングループ 1 0 2 には入力特徴 1 1 b 及び参照入力特徴 1 2 b が、第 2 の層 1 1 2 の第 2 のニューロングループ 1 0 2 には出力特徴 1 3 b が、それぞれ示されている。図 1 の例において、入力特徴 1 1 a には複数の特徴チャンネル (又は特徴面) が含まれており、それぞれの特徴チャンネルが 1 つのニューロンに相当する。特徴量の演算を行う際には、複数の特徴チャンネルのそれぞれに異なるフィルタ (例えば 2 × 2 画素のフィルタなど) が適用される。入力特徴 1 1 b、参照入力特徴 1 2 a , 1 2 b、及び出力特徴 1 3 a , 1 3 b についても同様である。

20

【 0 0 1 6 】

結合設定部 1 0 5 は、第 1 の層 1 1 1 について、ニューロン演算 1 0 6 a で用いられる第 1 の部分と、ニューロン演算 1 0 6 b で用いられる第 2 の部分とを設定する。図 1 の例において、結合設定部 1 0 5 は、ニューロン演算 1 0 6 a , 1 0 6 b で用いられる特徴チャンネルを設定している。結合設定部 1 0 5 は、結合パラメータに従ってこのような設定を行う。例えば、結合設定部 1 0 5 は結合パラメータに従って入力特徴 1 1 a , 1 1 b の参照順序を変更することができる。結合設定部 1 0 5 は、入力特徴 1 1 a , 1 1 b の参照先アドレスを設定するレジスタ等によって実現することができる。図 1 においては、結合設定部 1 0 5 は、模式的に、特徴チャンネルの順序を入れ換えるユニットとして示されている。参照入力特徴 1 2 a , 1 2 b は、参照順序に従って並び替えられた入力特徴 1 1 a , 1 1 b を模式的に示しており、それぞれ、ニューロン演算 1 0 6 a で用いられる第 1 の部分、及びニューロン演算 1 0 6 b で用いられる第 2 の部分に相当する。

30

【 0 0 1 7 】

このように本実施形態においては、ニューロン演算 1 0 6 a により、第 1 の層 1 1 1 の参照入力特徴 1 2 a (第 1 の部分の特徴量データ) から、第 2 の層 1 1 2 の出力特徴 1 3 a (第 3 の部分の特徴量データ) が算出される。また、ニューロン演算 1 0 6 b により、第 1 の層 1 1 1 の参照入力特徴 1 2 b (第 2 の部分の特徴量データ) から、第 2 の層 1 1 2 の出力特徴 1 3 b (第 4 の部分の特徴量データ) が算出される。図 1 の例においては、参照入力特徴 1 2 a の全ての特徴チャンネル (ニューロン) と、出力特徴 1 3 a の全ての特徴チャンネル (ニューロン) と、が互いに結合されている。同様に、参照入力特徴 1 2 b の全ての特徴チャンネルと、出力特徴 1 3 b の全ての特徴チャンネルと、が互いに結合されている。このように、出力特徴 1 3 a に含まれるそれぞれのチャンネルの特徴量データは、参照入力特徴 1 2 a に含まれる全てのチャンネルの特徴量データを用いて、かつ参照入力特徴 1 2 b に含まれるチャンネルの特徴量データを用いずに、算出される。

40

【 0 0 1 8 】

結合設定部 1 0 5 が参照する結合パラメータは、結合パラメータ保持部 1 0 4 に格納されており、結合パラメータ保持部 1 0 4 から取得することができる。図 1 の例において、結合パラメータ決定部 1 0 3 は、結合パラメータ保持部 1 0 4 に格納されている結合パラメータを結合設定部 1 0 5 に供給することができる。この結合パラメータは、第 2 の層 1

50

1 2 の第 3 の部分 (出力特徴 1 3 a) の特徴量データを算出するために参照される第 1 の層 1 1 1 の第 1 の部分 (参照入力特徴 1 2 a) を規定する。また、この結合パラメータは、第 2 の層 1 1 2 の第 4 の部分 (出力特徴 1 3 b) の特徴量データを算出するために参照される第 1 の層 1 1 1 の第 2 の部分 (参照入力特徴 1 2 b) も規定する。この結合パラメータはさらに、第 2 の層 1 1 2 の第 3 の部分 (出力特徴 1 3 a) 及び第 4 の部分 (出力特徴 1 3 b) を規定してもよい。

【0 0 1 9】

一方で、本実施形態において、結合設定部 1 0 5 は入力特徴 1 1 a, 1 1 b の参照順序を変更し、又は入力特徴 1 1 a, 1 1 b を並び替えることで参照入力特徴 1 2 a, 1 2 b を決定している。したがって、第 2 の層 1 1 2 は、入力特徴 1 1 a (第 1 の層の第 1 のニューロンセット) の一部及び入力特徴 1 1 b (第 1 の層の第 2 のニューロンセット) の一部と結合している、出力特徴 1 3 a (第 2 の層の第 3 のニューロンセット) を有しているといえる。また、第 2 の層 1 1 2 は、入力特徴 1 1 a (第 1 の層の第 1 のニューロンセット) の一部及び入力特徴 1 1 b (第 1 の層の第 2 のニューロンセット) の一部と結合している、出力特徴 1 3 b (第 2 の層の第 4 のニューロンセット) を有しているといえる。ここで、入力特徴 1 1 a (第 1 の層の第 1 のニューロンセット) は、第 1 の層の前階層の第 1 のニューロンセットと結合していてもよい。また、入力特徴 1 1 b (第 1 の層 1 1 1 の第 2 のニューロンセット) は、第 1 の層 1 1 1 の前階層の第 2 のニューロンセットと結合していてもよい。これらの前階層、第 1 の層、及び第 2 の層は、ニューラルネットワークに含まれており、順に接続されている。本明細書において、ニューロンセットとは、1 つの階層にある複数のニューロンのサブセットのことを指す。

【0 0 2 0】

以上のように、本実施形態において、結合設定部 1 0 5 はニューロングループの間で入出力信号を交換している。図 1 の構成は、入力特徴 1 1 a と入力特徴 1 1 b との間で特徴量データを交換することにより、参照入力特徴 1 2 a 及び参照入力特徴 1 2 b を設定することと等価である。

【0 0 2 1】

図 1 には、ニューラルネットワークの 2 層分の構造が示されているが、図 2 に示されるような、より大きなニューラルネットワークに本実施形態に係る処理を適用することもできる。図 2 には、4 層の多層 CNN が示されており、層間のニューロン演算は畳み込み演算を含んでいる。図 2 にはまた、ニューラルネットワークにデータを入力するデータ入力部 1 0 0 が示されている。本実施形態に係るニューラルネットワークが処理対象とするデータの種類は特に限定されず、例えば音声、静止画像、動画、又は文章等であってもよい。図 2 の例では、入力データ 2 0 として 3 チャンネルのカラー静止画像が入力されている。以下では、チャンネルのことを c h と略すことがある。

【0 0 2 2】

図 2 にはさらに、ニューラルネットワークから出力された演算結果に基づく処理結果を出力する結果出力部 1 1 0 も示されている。処理結果は、ニューラルネットワークが処理するタスクに応じて形態が異なる。例えば、意味的領域分割を行う場合、処理結果は、入力画像中の各画素について、C 個のクラスそれぞれについての尤度を示す、C 個のチャンネルを有する尤度マップであってもよい。また、分類を行う場合、処理結果は、C 個の分類カテゴリーそれぞれについての尤度であってもよい。

【0 0 2 3】

図 2 には、第 1 の層 2 1 1 から第 4 の層 2 1 4 の間に、3 回のニューロン演算 1 0 6 - 1 a, b, 1 0 6 - 2 a, b, 及び 1 0 6 - 3 a, b と、2 回の結合設定部 1 0 5 による処理が示されている。図 2 には、ニューロングループ 1 0 1 とニューロングループ 1 0 2 の 2 つのニューロングループが存在し、ニューロングループ間のニューロン演算による結合は存在しない。また、このニューラルネットワークは、入力層から最終層までにわたって 2 つのニューロングループに分かれている。一方で、結合設定部 1 0 5 により、ニューロンの入出力の参照順序がニューロングループをまたいで変更されている。このため、2

つのニューロングループ間でニューロンの入出力信号を交換することができる。

【 0 0 2 4 】

次に、本実施形態に係るニューラルネットワークに従って処理を行うデータ処理装置について、図 3 及び図 4 を参照して説明する。図 4 は、ニューラルネットワークを用いて認識処理を行う、本実施形態に係るデータ処理装置である I P 等のハードウェア構成例を示す。図 4 は、ニューラルネットワークを用いて認識処理を行う際の、データ処理装置の動作フローを示す。なお、C N N の層間の畳み込み演算の重みパラメータは、すでに学習により得られているものとする。また、結合パラメータ保持部 1 0 4 はあらかじめ結合パラメータを格納している。結合パラメータの決定方法については後述する。

【 0 0 2 5 】

データ処理装置の動作が開始すると、ステップ S 3 0 1 で演算制御部 1 0 8 はデータ入力部 1 0 0 に制御信号を送る。制御信号に応じてデータ入力部 1 0 0 は入力データ 1 0 を受け取り、メモリ 1 0 7 に入力特徴として保存する。データ入力部 1 0 0 は、例えば、撮像装置等から入力データ 1 0 を受け取ることができる。以下の例において、入力データ 1 0 は多チャンネルの画像特徴であり、とりわけ 6 c h のマルチスペクトル画像であるものとする。

【 0 0 2 6 】

次に、ステップ S 3 0 2 ~ S 3 1 2 のループ処理において、情報処理装置は、メモリ 1 0 7 上の入力特徴に対して、ニューラルネットワークの第 1 層目の処理を行う。図 2 の例であれば、入力データ 2 0 を用いて、特徴データ 2 1 a , 2 1 b が算出される。

【 0 0 2 7 】

まずステップ S 3 0 3 で、結合設定部 1 0 5 は、結合パラメータ決定部 1 0 3 を介して、結合パラメータ保持部 1 0 4 から、現在処理中の階層についての結合関係を示す、結合パラメータを読み出す。本実施形態において、結合パラメータ保持部 1 0 4 は、結合パラメータを示すテーブルを格納している。

【 0 0 2 8 】

図 6 (A) は、このようなテーブルの一例を示す。また、図 5 は、図 6 (A) に示すテーブル 6 0 1 に従う結合設定部 1 0 5 の動作を示している。図 6 (A) に示すテーブル 6 0 1 は、チャンネル単位で、参照入力特徴 5 1 2 a 及び参照入力特徴 5 1 2 b を特定している。すなわち、テーブル 6 0 1 は、6 c h の入力特徴のそれぞれについて、参照順序を示す 6 つの値を有している。以下では、このようなテーブルのことを順序指定テーブルと呼ぶ。

【 0 0 2 9 】

後述するステップにおいて、演算部 1 7 a , 1 7 b は、図 5 に示されるように、順序指定テーブルが指定する順序に従って特徴チャンネルを参照し、演算処理を行う。図 6 (A) の例においては、第 1 のニューロングループ 1 0 1 については、入力特徴の [1 c h , 2 c h , 4 c h] の 3 つのチャンネルに対応する、参照入力特徴 5 1 2 a に対する畳み込み演算により、出力特徴 5 1 3 a が算出される。また、第 2 のニューロングループ 1 0 2 については、入力特徴の [3 c h , 5 c h , 6 c h] の 3 つのチャンネルに対応する、参照入力特徴 5 1 2 b に対する畳み込み演算により、出力特徴 5 1 3 b が算出される。すなわち、結合パラメータに従って、入力特徴 5 1 1 a , 5 1 1 b のうち、第 1 のニューロングループ 1 0 1 についてのニューロン演算 5 0 6 a に用いられる参照入力特徴 5 1 2 a が設定される。また、同様にニューロン演算 5 0 6 b に用いられる参照入力特徴 5 1 2 b も設定される。

【 0 0 3 0 】

ステップ S 3 0 4 ~ S 3 1 1 までの処理は、ニューロングループ 1 0 1 , 1 0 2 のそれぞれについて並列に行われる。以下ではニューロングループ 1 0 1 についての処理であるステップ S 3 0 4 a ~ S 3 1 1 a について説明するが、ニューロングループ 1 0 2 についての処理であるステップ S 3 0 4 b ~ S 3 1 1 b も同様に行われる。

【 0 0 3 1 】

ステップ S 3 0 5 a ~ S 3 1 0 a のループ処理により、出力チャンネルの 1 チャンネルごとに特徴データが算出される。このようなループ処理を所定回数 (N 回) 行うことにより、 N チャンネルの出力特徴が生成される。

【 0 0 3 2 】

ステップ S 3 0 6 a で演算部 1 7 a は、メモリ 1 0 7 から畳み込み演算の重み係数 (畳み込みカーネルとも呼ばれる) を読み込む。演算部 1 7 a は、読み込んだパラメータを演算部 1 7 中のレジスタ領域 (不図示) にセットする。

【 0 0 3 3 】

ステップ S 3 0 7 a で参照設定部 1 6 a は、畳み込み演算に使う入力特徴のアドレスを、演算部 1 7 a のレジスタ領域 (不図示) にセットする。参照設定部 1 6 a は、ステップ S 3 0 3 a で読み込まれた順序指定テーブルにより指示される順序で、メモリ 1 0 7 上の入力特徴の指定されたチャンネルが参照されるように、アドレスをセットする。

【 0 0 3 4 】

ステップ S 3 0 8 a で演算部 1 7 a は、ステップ S 3 0 7 a でセットされたアドレスに位置する入力特徴を、メモリ 1 0 7 から取得する。そして、演算部 1 7 a は、取得した入力特徴に対して、ステップ S 3 0 6 a でセットされた重み係数を用いた畳み込み演算を行う。

【 0 0 3 5 】

演算部 1 7 a は、下式 (1) に従って畳み込み演算を行うことができる。

$$F^{OUT}_j(x, y) = b_j + \sum_i x \cdot y \cdot W_{ij}(x, y) F^{IN}_{LUT(i)}(x + x, y + y) \dots \dots (1)$$

この式において、 F^{IN} 及び F^{OUT} は入力特徴及び出力特徴を表し、この例では、ともに縦 × 横 × チャンネル方向を有する 3 次元特徴データである。 F^{IN}_k は、入力特徴 F^{IN} の k 番目のチャンネルを表し、 F^{OUT}_j は出力特徴 F^{OUT} の j 番目のチャンネルを表す。 $LUT(i)$ は、畳み込み演算の対象となる i 番目の参照入力特徴のチャンネルの番号に対応する、入力特徴のチャンネルの番号を示す。例えば、図 5 の例では、ニューロン演算 5 0 6 a では、参照入力特徴 5 1 2 a の [1 c h , 2 c h , 3 c h] が用いられ、これらは入力特徴 5 1 1 a . 5 1 1 b の [1 c h , 2 c h , 4 c h] の 3 つのチャンネルに対応する。すなわち、 $i = [1 , 2 , 3]$ の時に、 $LUT(i)$ は [1 , 2 , 4] を表す。 b_j はバイアス項である。 W_{ij} は重み係数であり、前階層のニューロン i と後階層のニューロン j との間の結合重みを示す。 $x, y (x, y [- 1 , 0 , 1])$ は、畳み込み範囲を示す変数である。

【 0 0 3 6 】

また、演算部 1 7 a は、畳み込み演算の結果に対して、活性化演算をさらに行うことができる。演算部 1 7 a は、下式 (2) に従って活性化演算を行うことができる。

$$F^{OUT'}_j(x, y) = (F^{OUT}_j(x, y)) \dots \dots (2) \\ (x) = \text{Max}(0, x)$$

この式において、 (x) は活性化関数と呼ばれる非線形関数である。なお、演算部 1 7 a は、最大値プールと呼ばれる演算処理、及び全結合層の演算のような、CNNにおいて用いられるその他の演算を行うこともできる。

【 0 0 3 7 】

ステップ S 3 0 9 a で出力部 1 8 は、ステップ S 3 0 8 a で得られた 1 チャンネル分の出力特徴を、メモリ 1 0 7 上の所定の位置に保存する。以上の処理により、出力特徴の j 番目のチャンネルの特徴データである $F^{OUT'}_j$ がメモリ 1 0 7 上に生成される。

【 0 0 3 8 】

以上のステップ S 3 0 4 a ~ S 3 1 1 a の処理を、所定回数 (N 回) 繰り返すことで、第 1 のニューロングループ 1 0 1 についての出力特徴の N 個のチャンネルの特徴データがメモリ 1 0 7 上に生成される。並行して、以上のステップ S 3 0 4 a ~ S 3 1 1 a の処理

10

20

30

40

50

を、所定回数（N回）繰り返すことで、第2のニューロングループ102についての出力特徴のN個のチャンネルの特徴データがメモリ107上に生成される。これらの処理により、ニューラルネットワークの1階層分の演算処理が完了する。

【0039】

このようにしてステップS302～S312までのループを複数回繰り返すことで、ニューラルネットワークの全階層について同様の処理が行われる。ループの終了時点で、最終層の出力特徴がメモリ107上に生成されている。ステップS313で結果出力部110は、メモリ107上の最終層の出力特徴を出力する。こうして出力された出力特徴が、ニューラルネットワークによる入力データ10に対する処理結果である。こうして得られた処理結果に基づいて、分類処理などの様々な認識処理を行うことができる。結果出力部110は、処理結果を用いて分類処理などの認識処理を行い、認識結果を出力してもよい。こうして、ニューラルネットワークを用いた認識動作が終了する。

10

【0040】

結合パラメータのデータ形式は、図6（A）の順序指定テーブル601に限られない。例えば、結合パラメータは、図6（A）に示す結合テーブル603又は結合リスト602により表されてもよい。なお、図6（A）において、順序指定テーブル601、結合リスト602、及び結合テーブル603は、全て同じ結合関係を示している。結合テーブル603は、入力特徴のチャンネルと、出力特徴のチャンネルと、結合の有無を表す。記号は、チャンネル間の結合が存在することを示す。図6（A）では、入力特徴の[1ch, 2ch, 4ch]のそれぞれが、出力特徴の[1ch, 2ch, 3ch, 4ch]のそれぞれと結合していることを示している。図6（A）はまた、次階層のニューロングループが、1ch～4chと、5ch～8chの2グループに分かれていることも示している。結合リスト602は、出力特徴のそれぞれのチャンネルについて、結合している入力特徴のチャンネルの番号を示している。

20

【0041】

図6（A）の例では、入力特徴のチャンネルの順序を変更したものが、参照入力特徴に相当する。このため、入力特徴のそれぞれのチャンネルが1回ずつ参照され、畳み込み演算に用いられた。すなわち、結合設定部105は、結合パラメータに従って、第1の層について第1の部分及び第2の部分を設定するが、この場合第1の部分と第2の部分とは重複していない。一方で、図6（B）に示すように、入力特徴の1つのチャンネルが、次階層の複数のニューロングループにより参照されてもよい。すなわち、第1の部分と第2の部分とが部分的に重複していてもよい。図6（B）の例では、入力特徴の2ch及び5chが、第1のニューロングループ（出力特徴の1ch～4chを算出）と第2のニューロングループ（出力特徴の5ch～8chを算出）の双方から参照されている。図6（B）においても、順序指定テーブル611、結合リスト612、及び結合テーブル613は、全て同じ結合関係を示している。このように、入力特徴の参照方法の設定は、順序の変更のみに限られず、さまざまな結合関係を設定することができる。

30

【0042】

上記の例では、2つのニューロングループが存在していた。しかしながら、ニューロングループの数は2つに限定されない。各階層のニューロンが、任意のn個のグループに分けられていてもよい。

40

【0043】

また、上記の例では、入力データ20は6chのデータであったが、入力データのチャンネル数等の構成は特に限定されない。また、最初の階層においては、通常のCNNと同様に畳み込み処理を行ってもよく、2層目からニューロンを複数のニューロングループに分けてもよい。また、特徴データ20a, 20bが、それぞれ同一の入力データ20であってもよい。例えば、特徴データ20a, 20bの参照アドレスを、入力データ20が格納されている同一のアドレスにすることができる。こうして、第1層目から、同一の入力データ20が入力される2つのニューロングループについての並列処理を行ってもよい。これらの構成は、3chのRGB静止画像データのような、チャンネル数の少ない入力デ

50

ータに適用されてもよい。

【 0 0 4 4 】

本実施形態によれば、第 1 の層の第 1 の部分（例えば参照入力特徴 1 2 a）と、第 2 の層の第 3 の部分（例えば出力特徴 1 3 a）とが結合される。また、第 1 の層の第 2 の部分（例えば参照入力特徴 1 2 b）と、第 2 の層の第 4 の部分（例えば出力特徴 1 3 b）とが結合される。一方で、第 1 の層の第 1 の部分と、第 2 の層の第 4 の部分と、の間の結合が省略されるため、この結合に対応する畳み込み演算の演算量、及び重み係数を保持するメモリ量を削減することができる。その一方で、特許文献 1 のように非相互結合層を設ける場合とは異なり、ニューロングループ間での結合は維持されるため、認識精度の劣化が抑えられることが期待される。

10

【 0 0 4 5 】

さらに、本実施形態によれば、それぞれのニューロングループについてのニューロン演算を、同様の演算を用いて並列に行うことができる。例えば図 4 の例では、ニューロン演算 1 0 6 a のための構成（参照設定部 1 6 a、演算部 1 7 a、及び出力部 1 8 a）と、ニューロン演算 1 0 6 b のための構成（参照設定部 1 6 b、演算部 1 7 b、及び出力部 1 8 b）と、は同様である。すなわち、本実施形態に係るデータ処理装置は、並列に動作する演算部 1 7 a（第 1 の処理部）及び演算部 1 7 b（第 2 の処理部）を有している。演算部 1 7 a は、ニューロン演算 1 0 6 a により、参照入力特徴 1 2 a から出力特徴 1 3 a を算出することができる。また、演算部 1 7 b は、ニューロン演算 1 0 6 b により、参照入力特徴 1 2 b から出力特徴 1 3 b を算出することができる。

20

【 0 0 4 6 】

さらには、本実施形態によれば、異なる結合関係を用いたニューロン演算を、同様の構成を用いて実現することができる。例えば図 2 の例では、結合設定部 1 0 5 が読み出した結合パラメータを用いることで、第 3 の層 2 1 3 及び第 4 の層 2 1 4 を算出するためのニューロン演算を、異なる結合関係に従って行うことができる。また、第 2 の層 2 1 2 を算出するためのニューロン演算においては、特徴データ 2 1 a は特徴データ 2 0 a から、特徴データ 2 1 b は特徴データ 2 0 a からそれぞれ算出されている。結合パラメータは、このようにニューロングループ間での結合が存在しない結合関係も示すことができる。さらには、一部の階層についての結合パラメータが、前階層の全てのチャンネルと後階層の全てのチャンネルとが互いに結合されることを示してもよい。このような構成によれば、特許文献 1 のように相互結合層と非相互結合層とが混在する場合と比較して、演算ハードウェアの回路規模を小さくすることができ、並列処理の効率を向上できる。

30

【 0 0 4 7 】

図 2 の例では、2 つのニューロングループは排他的なグループであり、1 つのニューロンは 2 つのニューロングループのいずれか一方に含まれていた。このような構成は並列処理を容易にするが、ニューロングループの構成はこのような形態には限定されない。例えば、ニューロングループ同士がオーバーラップしていてもよい。例えば、図 6（C）に示す結合テーブル 6 2 3 は、出力特徴についてのニューロングループには、1 c h ~ 3 c h、3 c h ~ 5 c h、5 c h ~ 7 c h の 3 つのグループが含まれることを示している。この例では、3 c h 及び 5 c h は複数のニューロングループに含まれている。このように、様々な結合関係を示す結合パラメータを設計することができる。

40

【 0 0 4 8 】

本実施形態の適用対象は、特定の種別のニューラルネットワークには制限されない。例えば、本実施形態は、再帰的ニューラルネットワーク（Recursive Neural Network, Byeon et al. "Scene labeling with LSTM recurrent neural networks", CVPR 2015.（以下Byeon）を参照）又はオートエンコーダーのような様々な種類のニューラルネットワークに適用可能である。

【 0 0 4 9 】

[実施形態 2]

以下では、実施形態 1 で説明したニューラルネットワークの結合パラメータの決定方法

50

について説明する。実施形態 2 に係る学習装置は、ニューラルネットワークの学習を行うことができる。実施形態 2 に係る学習装置は、通常の機械学習の手法を用いて、ニューラルネットワークの重み係数の学習を行うことができる。さらに、実施形態 2 に係る学習装置は、結合パラメータの決定も行うことができる。

【0050】

本実施形態における学習方法を概説する。学習には、学習データと、学習データに対する処理結果を示す教師データと、が用いられる。実施形態 1 で説明したように、所定の結合パラメータに従って処理を行うニューラルネットワークに学習データを入力すると、学習データに対する処理結果が得られる。そして、結合パラメータと、学習データに対する処理結果と、教師データと、に基づいて、結合パラメータ及びニューラルネットワークの階層間の重み係数との学習が行われる。結合パラメータの具体的な学習方法としては、学習中に結合パラメータに摂動的な変化を与えることにより、ニューラルネットワークの性能が向上する結合パラメータを探索する方法が挙げられる。例えば、結合パラメータに変化を与えたことによりニューラルネットワーク全体の性能が向上した場合に、この変化を採用するステップを繰り返すことができる。

10

【0051】

以下、図 7 及び図 8 を参照して本実施形態について説明する。図 7 は、本実施形態に係る学習装置の機能構成例を示す。図 7 に示す学習装置は、教師信号提供部 121 及び制約指示部 122 を有することを除き、実施形態 1 と同様である。図 7 に示される第 1 のニューロングループ 101 及び第 2 のニューロングループ 102 も、図 4 に示される結合設定部 105、参照設定部 16a、16b、演算部 17a、17b、及び出力部 18a、18b を用いて実現できる。

20

【0052】

データ入力部 100 は、学習データを取得し、メモリ 107 を介してニューラルネットワークに入力する。また、教師信号提供部 121 は、学習データに対する処理結果を示す教師データを取得し、結果出力部 110 に入力する。教師データは、学習データに対して、ニューラルネットワークが出力する結果の目標値である。このような教師データは予め作成されていてもよい。以下では、教師データを示す信号のことを教師信号と呼ぶ。

【0053】

制約指示部 122 は、学習により得られるニューラルネットワークについての制約条件を指示する。制約条件に、例えば、ニューラルネットワークの階層数などを含むことができる。また、制約条件に、例えばそれぞれの階層におけるニューロングループの数のような、学習により得られる結合パラメータに対する制約条件を含むことができる。また、制約指示部 122 は、摂動の大きさを決めるパラメータのような、学習装置による学習時に用いられるパラメータを指示することができる。これらの制約条件及びパラメータは、ユーザが指定することができ、制約指示部 122 は、このような制約条件を指示するユーザ入力を取得することができる。ユーザによる指定がない項目については、制約指示部 122 はデフォルト値を用いることができる。

30

【0054】

図 8 は、本実施形態における学習処理のフローチャートの一例である。学習処理がスタートすると、ステップ S801 で制約指示部 122 は、制約情報及びパラメータを示すユーザ入力を取得する。

40

【0055】

ステップ S802 で演算制御部 108 は、ステップ S801 で取得したユーザ入力に従う構成を有するニューラルネットワークについて、全ての重み係数を初期化する。重み係数の初期値としては乱数値を用いることができる。また、結合パラメータ決定部 103 は、結合パラメータを初期化する。本実施形態においては、図 6 (A) に示すような順序指定テーブル 601 が、結合パラメータを示すために用いられる。結合パラメータの初期値としては、図 9 の順序指定テーブル 901 のような、順序が変更されていない状態を用いることができる。

50

【 0 0 5 6 】

ステップ S 8 0 3 でデータ入力部 1 0 0 は、複数の学習画像をメモリ 1 0 7 を介してニューラルネットワークへ入力する。また、教師信号提供部 1 2 1 は、複数の学習画像のそれぞれに対応する教師信号を結果出力部 1 1 0 に入力する。このように、ニューラルネットワークの学習のためには複数の学習画像を用いることができる。

【 0 0 5 7 】

ステップ S 8 0 4 では演算制御部 1 0 8 は、実施形態 1 と同様にニューラルネットワークを用いた処理を行う。すなわち、結合設定部 1 0 5 は、結合パラメータに従って、ニューロン演算において参照される入力特徴のチャンネルを設定する。そして、このような設定に従って、各層のニューロン演算が行われる。こうして、それぞれの学習画像に対するニューラルネットワークによる処理結果が得られる。

10

【 0 0 5 8 】

ニューラルネットワークを用いた処理が終了すると、結果出力部 1 1 0 は、処理結果を教師信号と比較する。そして、結果出力部 1 1 0 は、比較結果を示す誤差信号を生成し、ニューラルネットワークに誤差信号を入力することで、ニューラルネットワークの重み係数の学習を行うことができる。ニューラルネットワークの重み係数の学習には、公知の方法を用いることができる。例えば、誤差の算出方法は特に限定されず、Krizhevsky又はZeiler等に記載の方法を用いることができる。具体例としては、学習タスクに応じて、交差エントロピー等の損失関数を用いて誤差を算出することができる。また、こうして算出された誤差に基づいて、例えば誤差逆伝播法を用いることにより、ニューラルネットワークの重み係数の学習を行うことができる。

20

【 0 0 5 9 】

ステップ S 8 0 4 においては、複数の学習画像に対するニューラルネットワークを用いた処理と、処理結果と教師信号との比較に基づく重み係数の学習とを含む学習ステップが繰り返される。この学習ステップは、所定の回数、又は学習による誤差値の変化が所定の範囲以下に収束するまで、繰り返される。

【 0 0 6 0 】

ステップ S 8 0 5 で結合パラメータ決定部 1 0 3 は、ステップ S 8 0 4 終了時の誤差値を変数 $Loss_1$ として記憶する。また、結合パラメータ決定部 1 0 3 は、ステップ S 8 0 4 終了時のニューラルネットワークの重み係数を記憶する。

30

【 0 0 6 1 】

ステップ S 8 0 6 で結合パラメータ決定部 1 0 3 は結合パラメータを変更する。例えば、結合パラメータ決定部 1 0 3 は、第 1 の結合パラメータを第 2 の結合パラメータに変更することができる。図 9 (A) は、結合パラメータを保持する結合関係テーブルの変更例を示す。結合パラメータ決定部 1 0 3 は、例えば、ランダムに選ばれた階層についての結合関係テーブルを結合パラメータ保持部 1 0 4 から取得し、ランダムに選ばれたテーブル上の m 個の値を入れ替えることができる。図 9 (A) の例では、 $m = 2$ である。

【 0 0 6 2 】

ステップ S 8 0 7 で演算制御部 1 0 8 は、ステップ S 8 0 6 における変更後の結合パラメータに従ってニューラルネットワークを用いた処理を行う。そして結果出力部 1 1 0 は、ステップ S 8 0 4 と同様に、複数の学習画像に対するニューラルネットワークによる処理結果に基づいて、ニューラルネットワークの重み係数の学習を行う。ステップ S 8 0 6 における結合パラメータの変更により、誤差値は一時的に上昇するかもしれないが、ステップ S 8 0 7 においてステップ S 8 0 4 と同様に重みの更新を行うことで、変更後の結合関係に適した重み係数が得られる。

40

【 0 0 6 3 】

ステップ S 8 0 8 で結合パラメータ決定部 1 0 3 は、ステップ S 8 0 7 終了時の誤差値を $Loss_2$ として記憶する。

【 0 0 6 4 】

ステップ S 8 0 9 で結合パラメータ決定部 1 0 3 は、ステップ S 8 0 4 で第 1 の結合パ

50

ラメータに従って得られた処理結果と、ステップ S 8 0 7 で第 2 の結合パラメータに従って得られた処理結果と、に基づいて結合パラメータの学習を行う。この例において結合パラメータ決定部 1 0 3 は、L o s s 1 と L o s s 2 とを比較し、この比較に基づいて第 1 の結合パラメータ又は第 2 の結合パラメータを採用することができる。例えば、結合パラメータ決定部 1 0 3 は、下式 (3) に従って L o s s 1 と L o s s 2 とを比較することができる。

$$L o s s 2 < L o s s 1 + \dots\dots (3)$$

この式が真である場合、処理はステップ S 8 1 0 に進み、偽である場合、処理はステップ S 8 1 2 に進む。上式において、 α は摂動の許容度を定めるパラメータであり、ユーザが指定することができる。より大きい値の α は、結合関係をより頻繁に変更することを許容する。

10

【 0 0 6 5 】

ステップ S 8 1 0 において結合パラメータ決定部 1 0 3 は、ステップ S 8 0 6 における変更後の第 2 の結合パラメータを採用する。さらにステップ S 8 1 1 において結合パラメータ決定部 1 0 3 は α の値を小さくする。例えば、結合パラメータ決定部 1 0 3 は、 α の値に $(1 - \alpha)$ を乗算することができる。その後、処理はステップ S 8 0 5 に戻り、重み係数及び学習パラメータの学習が継続される。

【 0 0 6 6 】

ステップ S 8 1 2 において結合パラメータ決定部 1 0 3 は、ステップ S 8 0 6 における変更後の第 2 の結合パラメータを棄却し、変更前の第 1 の結合パラメータを採用する。また、結合パラメータ決定部 1 0 3 は、ニューラルネットワークの重み係数を、ステップ S 8 0 5 で記憶された、結合パラメータ変更以前の状態に戻す。ステップ S 8 1 3 において演算制御部 1 0 8 は、学習を終了するかどうかを判定する。例えば、これまでの学習の総ステップ数が所定数以上である場合に学習を終了することができる。学習を終了しない場合、処理はステップ S 8 0 6 に戻り、結合パラメータの学習が継続される。以上が、学習処理の流れである。

20

【 0 0 6 7 】

学習方法は上記の例には限定されない。例えば、図 8 の例ではランダムに選ばれた層の結合関係が変更されたが、入力に近い低次の階層から高次の階層へと順に結合関係の学習が行われてもよい。また、上記の例では結合パラメータの初期値として順序が変更されていない状態が用いられたが、所定の割合だけランダムに順序が入れ替えられた状態が初期値として用いられてもよい。この割合は、ユーザ入力により指定されていてもよい。

30

【 0 0 6 8 】

また、図 8 の例では、結合パラメータを変更し、変更を棄却するか採用するかが判定された。別の学習方法として、複数の結合パラメータを入れ替えながら同じステップ数の学習を行い、誤差が小さくなった結合パラメータを採用することもできる。

【 0 0 6 9 】

さらに、例えば図 9 (B) に示すように、結合を追加する摂動を結合パラメータに与えることもできる。図 9 (B) の結合リスト 9 1 1 は、入力特徴の 1 ~ 3 c h を参照する第 1 のニューロングループと、入力特徴の 4 ~ 6 c h を参照する第 2 のニューロングループが分かれていることを示している。一方で、変更後の結合リスト 9 1 2 は、第 1 及び第 2 のニューロングループが、それぞれ入力特徴の 4 c h 及び 2 c h をさらに参照することを示している。このような変更によればニューラルネットワークにおける結合数が増加する。この場合、ニューラルネットワークの規模の増加と、ニューラルネットワークの性能の向上とを比較し、規模の増加と比較して性能が向上している場合に変更を採用することができる。具体例としては、ステップ S 8 0 9 において、下式 (4) に従って L o s s 1 と L o s s 2 とを比較することができる。

40

$$- \log (L o s s 1 / L o s s 2) < \log (| W 1 | / | W 2 |) \dots\dots (4)$$

この式において、 $| W 1 |$ 及び $| W 2 |$ は、それぞれ結合パラメータの変更前及び変更後におけるニューラルネットワークの結合数 (又は重み係数の数) を示す。このように重

50

み係数の総数の変化を許容する場合、最終的な重み係数の上限を指示するユーザ入力に従って結合パラメータの学習を行ってもよい。

【0070】

さらに、蒸留学習を用いることもできる。具体的には、目的タスクについて、予め大規模なニューラルネットワークの学習を行うことができる。そして、それぞれ異なる結合パラメータに従うニューラルネットワークの候補を複数用意し、大規模なニューラルネットワークと近似する結果が得られるように、それぞれの候補の転移学習を行うことができる。その結果、もっとも誤差が小さくなった候補を選択することができる。転移学習の具体的な手法としては、例えば、Hinton, Geoffrey, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015).などに記載の方法を用いることができる。以上のように、機械学習による結合パラメータの学習方法は、特定の形態に限定されない。

10

【0071】

このように、機械学習によりニューロングループの結合関係を決定することにより、実施形態1に係るニューラルネットワークの認識精度を向上させることができる。また、上記のように、ニューラルネットワークの演算量（又は重み係数の数）を考慮しながら、ニューロングループの結合関係を決定することもできる。これらの手法によれば、実用的な演算速度及び精度を有するニューラルネットワークを得ることを容易にすることができる。

【0072】

[実施形態3]

20

実施形態1, 2では、ニューロングループごとに並列処理が行われた。一方でCNNにおいては、後段の階層が多くの特徴チャンネルを持つことが多い。したがって、1つのニューロングループにおけるニューロン演算の出力チャンネルの数も、後段の階層においては膨大になることがある。実施形態3では、1つのニューロングループについても並列処理が行われる。

【0073】

図10は、実施形態3におけるニューラルネットワークに従う処理の模式図である。図10には、2つのニューロングループ101, 102が存在し、それぞれ8chの入力特徴に対応する8chの出力特徴を算出する。一方で、図10には、4つの並列処理ユニット1001a~1001dが存在する。1つの並列処理ユニット1001a~1001dは、それぞれ、1つの処理サイクルにおいて、4chの入力特徴1011a~1011dに対する畳み込み演算を行うことにより、8chの出力特徴1013a~1013cを出力することができる。それぞれの並列処理ユニット1001a~1001dは、図4に示される結合設定部105、参照設定部16a、演算部17a、及び出力部18aと同様の構成を有することができる。

30

【0074】

図12は、実施形態3における、ニューラルネットワークの1階層についてのニューロン演算のフローチャートである。ステップS1201では、各並列処理ユニット1001a~1001dの動作が並列に開始される。ステップS1202において結合設定部105は、ステップS303と同様に結合パラメータを読み出す。ここで、結合パラメータは、各並列処理ユニット1001a~1001dのレジスタに予め保持されていてもよい。

40

【0075】

ステップS1203~S1208において、各並列処理ユニット1001a~1001dは、実施形態1のステップS305a~S310aと同様に畳み込み演算処理を行う。この結果として、各並列処理ユニット1001a~1001dは、メモリ上に8chの出力特徴1013a~1013dを格納する。ここで、各並列処理ユニット1001a~1001dの処理が終了する。

【0076】

ステップS1209において積算ユニット1021a及び積算ユニット1021bは、出力特徴1013a~1013b及び1013c~1013dを対応するニューロングル

50

ープごとに積算することで、出力特徴 1 0 1 4 a 及び出力特徴 1 0 1 4 b を生成する。こうして、ニューラルネットワークの 1 階層分の演算処理が終了する。

【 0 0 7 7 】

なお、結合設定部 1 0 5 は、実施形態 1 と同様に、ニューロングループ間でチャンネルを入れ替えることができる。すなわち、結合設定部 1 0 5 は、入力特徴 1 0 1 1 a ~ 1 0 1 1 b と、入力特徴 1 0 1 1 c ~ 1 0 1 1 d と、の間でチャンネルを入れ替えることにより、参照入力特徴 1 0 1 2 a ~ 1 0 1 2 b を設定することができる。

【 0 0 7 8 】

なお、結合設定部 1 0 5 は、特徴のチャンネル単位で結合関係を設定する代わりに、複数のチャンネルを含むブロック単位で結合関係を設定することもできる。すなわち、結合パラメータは、ブロック単位で参照入力特徴 (1 0 1 2 a ~ 1 0 1 2 d) を特定していてもよい。例えば、図 1 0 に示すように、並列処理ユニット 1 0 0 1 a ~ 1 0 0 1 d が 1 サイクルで処理する 4 c h を 1 つのブロックとしてもよい。例えば、入力特徴 1 0 1 1 b と 1 0 1 1 c の順序を入れ替えてもよい。この場合、参照入力特徴 1 0 1 2 a , 1 0 1 2 b , 1 0 1 2 c , 1 0 1 2 d は、それぞれ入力特徴 1 0 1 1 a , 1 0 1 1 c , 1 0 1 1 b , 1 0 1 1 d に対応する。このように、チャンネルのブロック単位で結合関係を設定することにより、チャンネルごとに結合関係を設定する場合と比較して、ニューラルネットワークに従う処理を行うデータ処理装置の回路規模を小さくすることができる。

【 0 0 7 9 】

[その他の実施形態]

ニューロングループの設定方法は、上記の例には限定されない。例えば、図 6 (A) では、複数のニューロンが、出力特徴の前半部分 (1 c h から 4 c h) に対応するニューロングループと、出力特徴の後半部分 (5 c h から 8 c h) に対応するニューロングループとに、ブロック状に分割されていた。しかしながら、例えば、複数のニューロンが、奇数番号のチャンネルに対応するニューロングループと、偶数番号のチャンネルに対応するニューロングループと、に分割されてもよい。このように、チャンネル番号についての巡回的な規則に従って、ニューロングループが構成されていてもよい。ハードウェア又はメモリ等の構成に合わせて、様々なニューロングループを設定することができる。

【 0 0 8 0 】

図 9 (C) はこのような一例を示す。図 9 の例では、入力特徴及び出力特徴のチャンネル数はそれぞれ 8 c h である。この例において、ニューロンは、入力特徴と出力特徴とのチャンネル間の結合の有無に従って、[1 c h , 5 c h] [2 c h , 6 c h] [3 c h , 7 c h] [4 c h , 8 c h] の 4 つのニューロングループに分かれており、グループ間で信号の交換は行われない。このような例において、結合パラメータの学習を行う際に、入力特徴のチャンネルの入れ替えを摂動として与えることにより、結合テーブル 9 2 1 を結合テーブル 9 2 2 へと変更することができる。結合テーブル 9 2 2 に従って処理を行うと、1 番目と 2 番目のグループ間で信号が交換されることになる。実施形態 2 と同様に、ニューラルネットワークの性能が向上すればこの入れ替えを採用することができ、演算量を増やさずにニューラルネットワークの性能が向上する結合関係を得ることができる。このように、ニューロングループがブロックに分割されていない場合にも、上述の各実施形態の方法は適用可能である。

【 0 0 8 1 】

また、入力特徴のチャンネルを入れ替える代わりに、出力特徴のチャンネルを入れ替えてもよい。例えば、結合テーブル 9 2 1 を結合テーブル 9 2 3 へと変更してもよい。このように、結合関係の変更方法としては、様々な方法を用いることができる。

【 0 0 8 2 】

また、べき集合に従ってニューロングループを設定することもできる。例えば、入力特徴が A , B , C , D の 4 チャンネルを有し、出力特徴が W , X , Y , Z の 4 チャンネルを有する場合について考える。この場合、出力特徴のそれぞれのチャンネルについての特徴データを算出するために、入力特徴のチャンネルのべき集合を用いるように、ニューロン

10

20

30

40

50

グループを設定することができる。例えば、WチャンネルはA B Cチャンネルから、XチャンネルはB C Dチャンネルが、YチャンネルはC D Aチャンネルから、ZチャンネルはD A Bチャンネルから、それぞれ算出することができる。また、1つの出力チャンネルに対して、ベキ集合の一部（例えばベキ集合 $\times \{A, B, C, D\}$ のうちの $4C_3$ の組み合わせ）を初期値として対応づけ、さらに結合パラメータの学習を行ってもよい。

【0083】

さらに、図13(A)に示すように、各階層についてニューロングループの数が異なっている。一般に、より前の階層においては多くの特徴を用いる必要性が低いことが多い。このような階層においてニューロングループの数を増やしても、性能に対する影響は低いことが期待される。

【0084】

それぞれのニューロングループのサイズは同じであってもよいし、異なっている。すなわち、第1のニューロングループが出力特徴の第3の部分の特徴量データを算出し、第2のニューロングループが出力特徴の第4の部分の特徴量データを算出する際に、第3の部分と第4の部分との間で大きさが異なっている。例えば、第3の部分と第4の部分との間で特徴量データのデータ量が異なっている。

【0085】

例えば図13(B)の例では、異なるチャンネル数の出力特徴を算出するように、複数のニューロングループが設定されている。すなわち、32chの入力特徴から32chの出力特徴を生成する層において、以下のようにニューロンを2つのグループに分けてもよい。すなわち、1つのニューロングループは24chの入力特徴から24chの出力特徴を算出し、もう1つのニューロングループは8chの入力特徴から8chの出力特徴を算出する。このような構成によれば、各ニューロングループで行われるニューロン演算の処理の回数は異なることになる。例えば、演算量は入力チャンネル数と出力チャンネル数との積により決まるため、この例では演算量は9対1になる。ハードウェアを用いてこのようなニューラルネットワークに従う演算を行う際には、演算量に応じた数の並列処理ユニットを設けることにより、その並列処理ユニットによる演算量を均等に行うことができる。

【0086】

図13(C)の例では、各ニューロングループで行われるニューロン演算の演算量が等しくなるように、ニューロングループの数がコントロールされている。この例では、第1の層1311における8chの入力特徴から、第2の層1312における16chの出力特徴が算出される。また、第2の層1312における16chの入力特徴から、第3の層1313における32chの出力特徴が算出される。この場合、第1の層1311からの入力特徴から第2の層1312への出力特徴を与える2つのニューロングループと、第2の層1312からの入力特徴から第3の層1313への出力特徴を与える4つのニューロングループと、を設けることができる。この場合、4chの入力特徴に対して、8chの出力特徴を算出する並列処理ユニットが、それぞれのニューロングループに対応する演算を行うことができる。この例では、後ろの階層においてチャンネル数が増えているが、チャンネル数に合わせてニューロングループの数も増えているので、それぞれのニューロングループに対応する演算の演算量は同一となる。このような構成によれば、並列処理を容易に行うことができる。

【0087】

このように、様々なニューロングループの設定方法を採用することができる。学習装置は、ニューロングループの数を学習により決定してもよいし、ニューロングループのサイズを学習により決定してもよい。すなわち、第1のニューロングループが出力特徴の第3の部分の特徴量データを算出し、第2のニューロングループが出力特徴の第4の部分の特徴量データを算出する構成において、第3の部分及び第4の部分の大きさを学習により決定することができる。

【0088】

10

20

30

40

50

ニューロングループの設定を学習により決定する方法は特に限定されない。例えば、適切なニューロングループの設定方法を機械学習により決定するために、前述の蒸留学習を用いることができる。また、実施形態2と同様に、ニューロングループの数をランダムに増減させることにより、結合パラメータの1つであるニューロングループの分割数に摂動的な変化を与えることができる。このような方法により、ニューラルネットワークの性能が向上する分割数を学習により決定することができる。

【0089】

別の学習方法として、遺伝的アルゴリズムを用いる方法が挙げられる。すなわち、それぞれ異なる結合パラメータに従うニューラルネットワークの学習を行い、性能のより高いニューラルネットワーク同士を組み合わせることで新たなニューラルネットワークを生成することができる。例えば、1つのニューラルネットワークの前半の階層と、他のニューラルネットワークの後半の階層と、を組み合わせることにより、新たなニューラルネットワークを生成することができる。

【0090】

図14は、遺伝的アルゴリズムの適用例を示す。図14には、世代1の候補として3つのニューラルネットワーク1401～1403が示されている。この中で、認識性能が低いと判定されたニューラルネットワーク1403が淘汰（削除）される。そして、性能が高いと判定されたニューラルネットワーク1401、1402の一部を組み合わせることにより、世代2の候補としてニューラルネットワーク1404～1406が生成されている。この場合には、式(4)などを用いることにより、ニューラルネットワークの規模を考慮してニューラルネットワークの評価を行うことができる。

【0091】

さらに、上述の実施形態に係る方法は、再帰的ニューラルネットワークのような様々な形態のニューラルネットワークに適用できる。図15(A)は、再帰的ニューラルネットワークの一つである、Long short term memory (LSTM) ネットワークの構成を示す。Byeonに開示されているように、LSTMにおいては、複数のLSTMニューロンを相互結合することにより、階層的なパターン情報の処理又は二次元情報の処理を行うことができる。それぞれのLSTMニューロンは、1つのサイクルにおいては独立して並列処理を行う。図15(A)において、ここで y_i^t は、 i 番目のLSTMニューロンへの再帰的入力であり、高次元のベクトルである。図15(B)では、LSTMニューロンの再帰的結合が展開されている。図15(B)は、本質的には図15(A)と同じ構成を表す。図15(B)に示される構成においては、階層 t の入力特徴から、階層 $t+1$ の出力特徴が生成されていると理解することができる。また、 y_i^t は i 番目のチャンネルの特徴量データであると理解することができる。

【0092】

図15(C)は、上述の実施形態に係る構成を、図15(B)に示されるLSTMネットワークに適用した様子を表す。結合設定部105は、各LSTMニューロンが生成した再帰的入力 y_i^t の一部を混交させ、入力 y_i^t を生成することができる。このような処理は、図1の例において結合設定部105が入力特徴11a、11bから参照入力特徴12を生成した処理と同様に行うことができる。こうして生成された入力 y_i^t は、LSTMニューロンへと再帰的に入力される。このような構成においては、これにより位相的に離れたLSTMニューロンからの信号にも基づいて、各LSTMニューロンの状態を変更することができる。このため、各LSTMニューロンによる並列処理を可能としながら、また演算処理の増加を抑えながら、図15(A)の構成と比較してより複雑な演算処理を行うことができる。

【0093】

本発明の一実施形態に係るデータ処理装置及び学習装置は、例えばネットワークを介して接続された複数の情報処理装置によって構成されていてもよい。

【0094】

上述の実施形態に係るデータ処理装置及び学習装置は、専用のハードウェアによって実

10

20

30

40

50

現される。一方で、データ入力部 1 0 0、結合パラメータ決定部 1 0 3、演算制御部 1 0 8、結果出力部 1 1 0、及び制約指示部 1 2 2 のような一部又は全部の処理部が、コンピュータにより実現されてもよい。

【 0 0 9 5 】

図 1 1 はコンピュータの基本構成を示す図である。図 1 1 においてプロセッサ 1 1 1 0 は、例えば CPU であり、コンピュータ全体の動作をコントロールする。メモリ 1 1 2 0 は、例えば RAM であり、プログラム及びデータ等を一時的に記憶する。コンピュータが読み取り可能な記憶媒体 1 1 3 0 は、例えばハードディスク又は CD - ROM 等であり、プログラム及びデータ等を長期的に記憶する。本実施形態においては、記憶媒体 1 1 3 0 が格納している、各部の機能を実現するプログラムが、メモリ 1 1 2 0 へと読み出される。そして、プロセッサ 1 1 1 0 が、メモリ 1 1 2 0 上のプログラムに従って動作することにより、各部の機能が実現される。図 1 1 において、入力インタフェース 1 1 4 0 は外部の装置から情報を取得するためのインタフェースである。また、出力インタフェース 1 1 5 0 は外部の装置へと情報を出力するためのインタフェースである。バス 1 1 6 0 は、上述の各部を接続し、データのやりとりを可能とする。

10

【 0 0 9 6 】

また、上述の実施形態で説明されたニューラルネットワークも、本発明の一実施形態である。ニューラルネットワークは、階層構造（例えば階層数及び各階層で行われる演算の種類）を示す情報と、階層間の重み係数を示す情報と、により規定することができる。上述の実施形態に係るニューラルネットワークは、さらに結合パラメータにより規定されている。このようなニューラルネットワークを規定する情報に従って、上述のデータ処理装置又はコンピュータは、入力データに対してニューラルネットワークに従う演算を行うことができる。したがって、このようなニューラルネットワークを規定する情報は、プログラムの一種に相当する。

20

【 0 0 9 7 】

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、ASIC）によっても実現可能である。

【 符号の説明 】

30

【 0 0 9 8 】

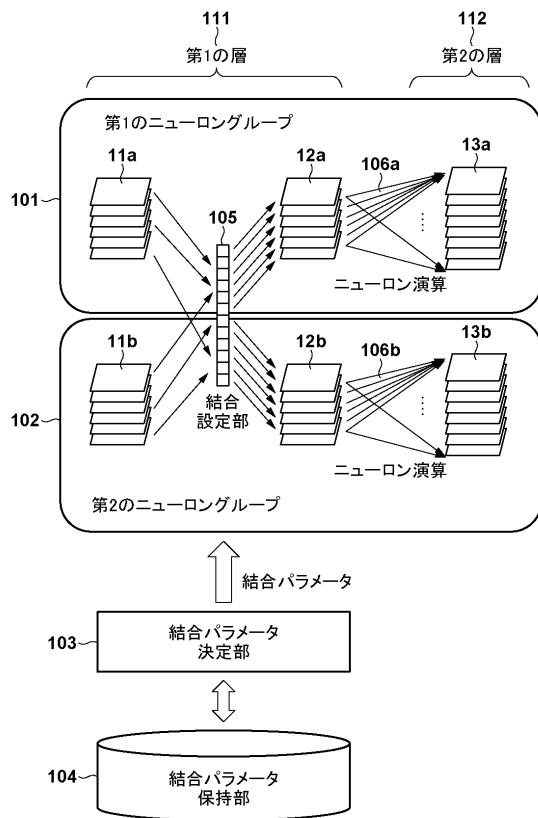
1 0 1 , 1 0 2 : ニューロングループ、1 0 3 : 結合パラメータ決定部、1 0 4 : 結合パラメータ保持部、1 0 5 : 結合設定部、1 0 6 a , 1 0 6 b : ニューロン演算、1 1 1 : 第 1 の層、1 1 2 : 第 2 の層

40

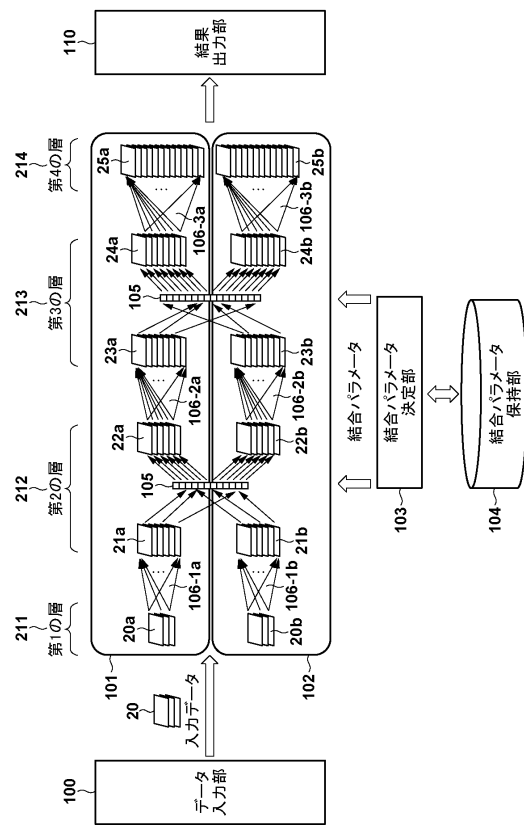
50

【図面】

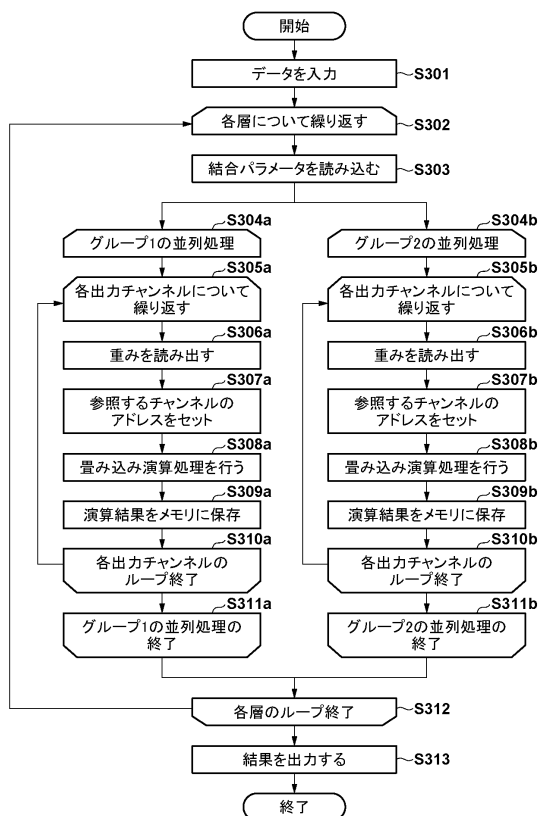
【図 1】



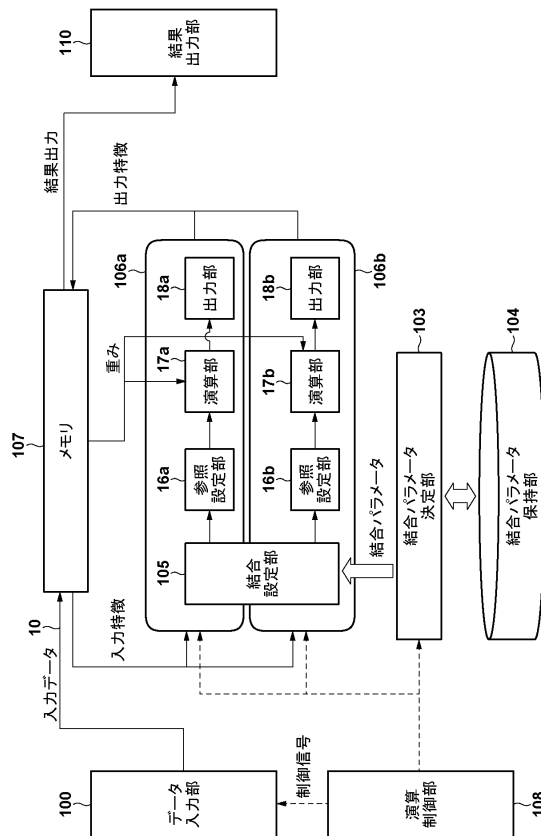
【図 2】



【図 3】



【図 4】



10

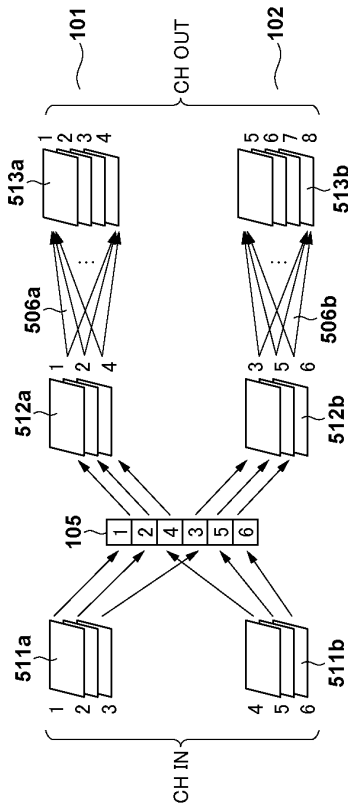
20

30

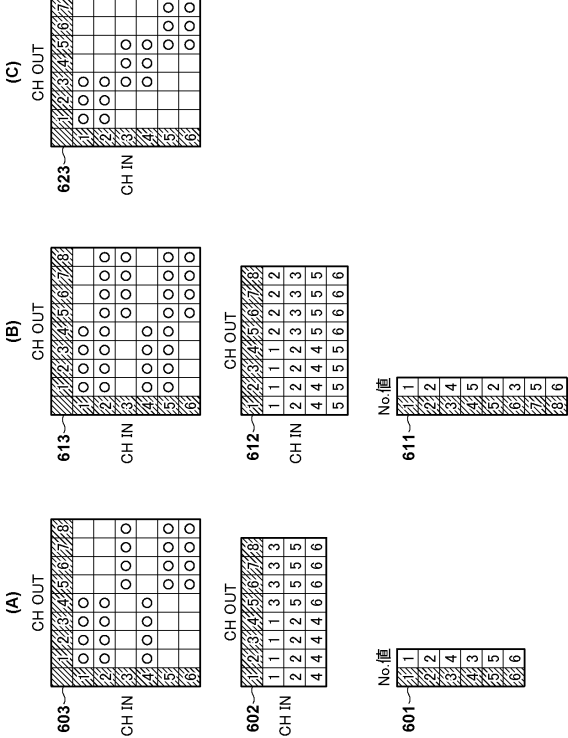
40

50

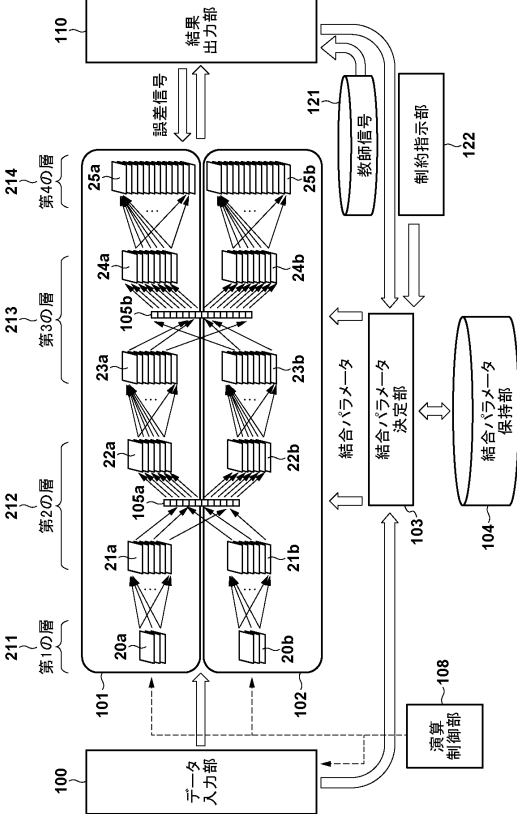
【図 5】



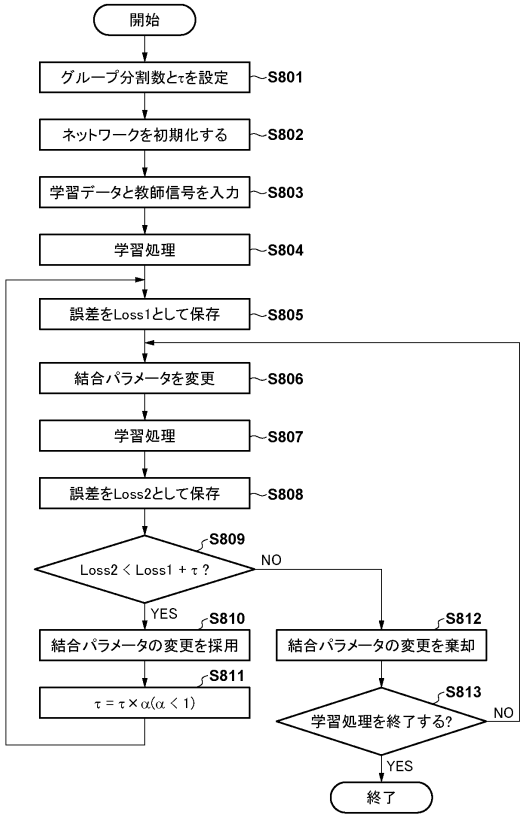
【図 6】



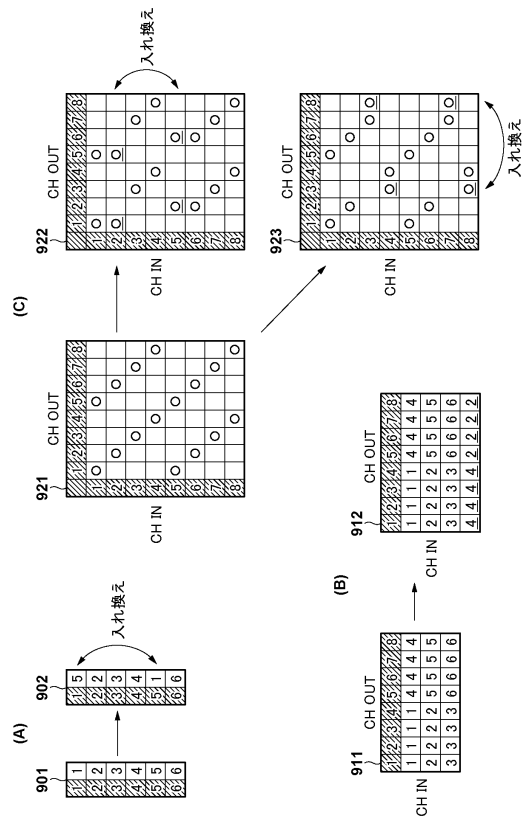
【図 7】



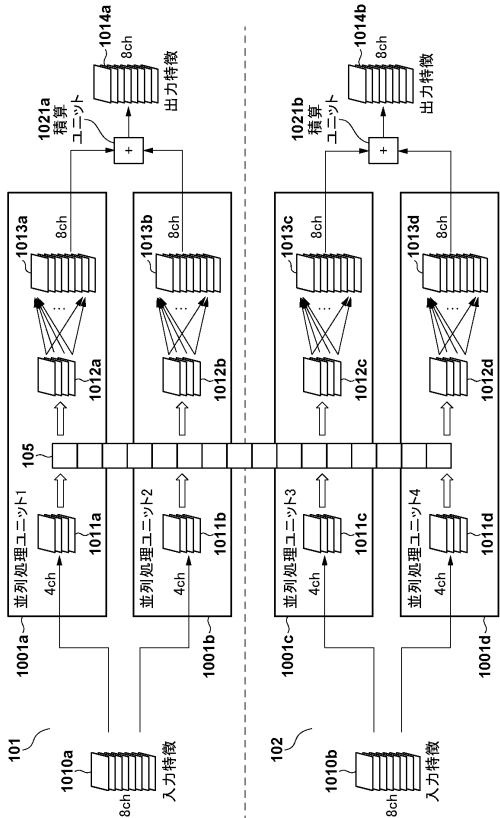
【図 8】



【図 9】



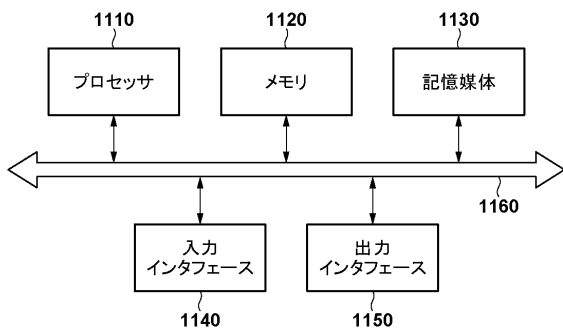
【図 10】



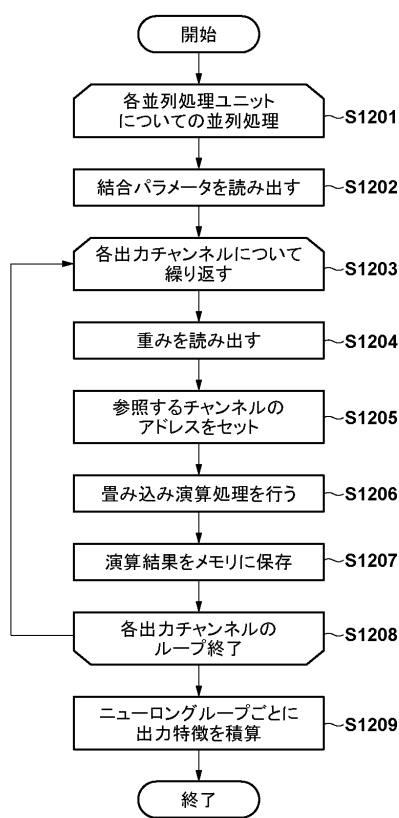
10

20

【図 11】



【図 12】

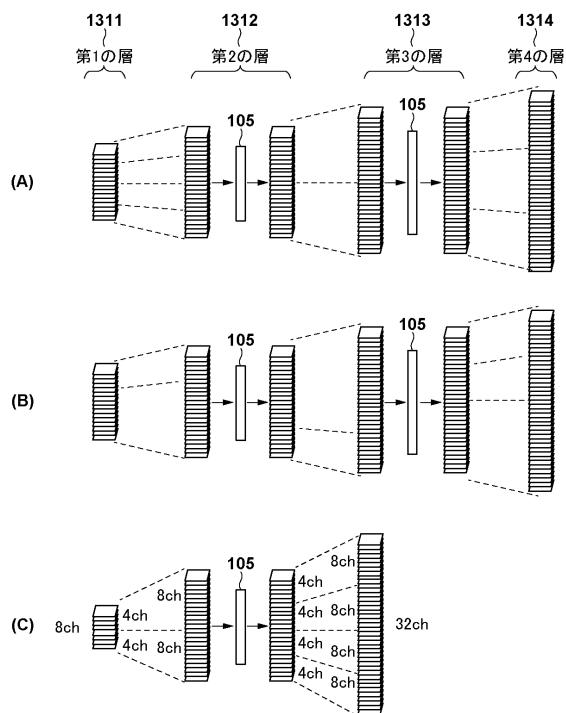


30

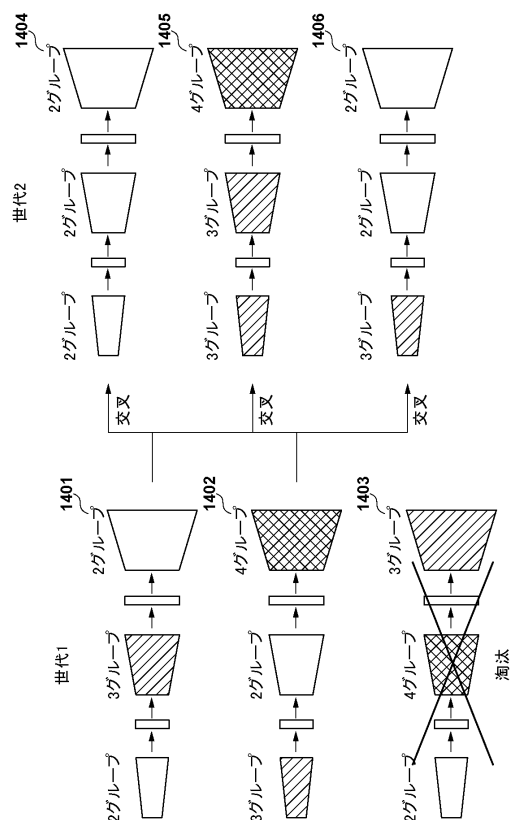
40

50

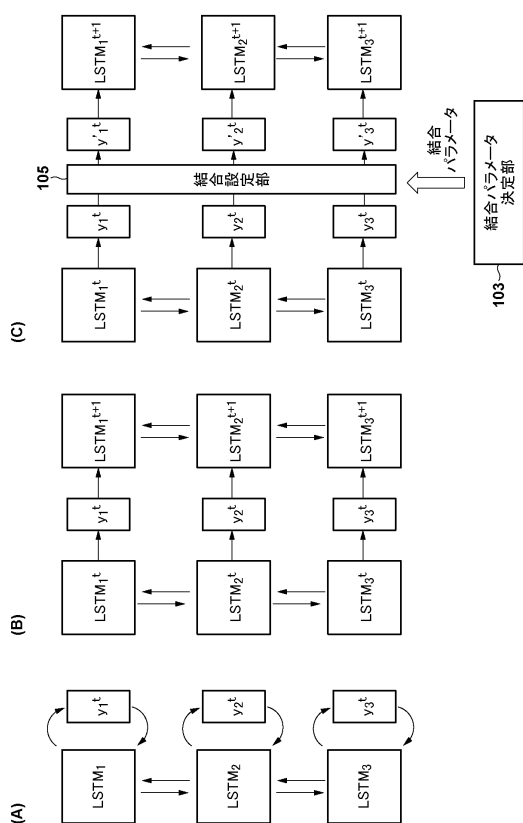
【 図 1 3 】



【圖 14】



【 図 1 5 】



フロントページの続き

- (56)参考文献 特開 2 0 1 5 - 2 1 0 7 4 7 (J P , A)
 特開 2 0 0 7 - 3 0 5 0 7 2 (J P , A)
 Pablo Barros, 外3名, "A multichannel convolutional neural network for hand posture recognition", Proceedings of the 24th International Conference on Artificial Neural Networks(ICANN 2014), 2014年09月30日
- (58)調査した分野 (Int.Cl., D B 名)
 G 0 6 N 3 / 0 4 5