

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7297468号
(P7297468)

(45)発行日 令和5年6月26日(2023.6.26)

(24)登録日 令和5年6月16日(2023.6.16)

(51)国際特許分類		F I	
G 0 6 N	3/08 (2023.01)	G 0 6 N	3/08
G 0 6 N	3/0464(2023.01)	G 0 6 N	3/0464
G 0 6 T	1/40 (2006.01)	G 0 6 T	1/40

請求項の数 10 (全17頁)

(21)出願番号	特願2019-36705(P2019-36705)	(73)特許権者	000001007 キヤノン株式会社 東京都大田区下丸子3丁目30番2号
(22)出願日	平成31年2月28日(2019.2.28)	(74)代理人	100126240 弁理士 阿部 琢磨
(65)公開番号	特開2020-140546(P2020-140546 A)	(74)代理人	100124442 弁理士 黒岩 創吾
(43)公開日	令和2年9月3日(2020.9.3)	(72)発明者	脇野 しおり 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
審査請求日	令和4年2月28日(2022.2.28)	審査官	山本 俊介

最終頁に続く

(54)【発明の名称】 データ処理装置及びその方法

(57)【特許請求の範囲】

【請求項1】

入力データに対し、ニューラルネットワークに対応する階層的な演算処理を実行するデータ処理装置であって、

前記ニューラルネットワークに対応する階層的な演算処理を複数の処理単位に分割した処理単位ごとに処理に用いる、前記演算処理の重み係数と動作パラメタとを含む制御データをセットとした複数の制御データセットを保持する保持手段と、

前記保持手段から前記複数の制御データセットを順に転送する転送手段と、

前記転送手段から転送された1つの制御データセットを用いて、当該制御データセットに対応する前記分割した処理単位の演算処理を行う演算処理手段とを備えることを特徴とするデータ処理装置。

10

【請求項2】

前記処理単位は、前記ニューラルネットワークの階層単位であることを特徴とする請求項1記載のデータ処理装置。

【請求項3】

前記処理単位は、前記ニューラルネットワークの特徴面単位であることを特徴とする請求項1記載のデータ処理装置。

【請求項4】

前記処理単位は、該処理単位の処理に必要な重み係数のメモリ量が、前記演算処理手段の参照するメモリにおける重み係数の格納部のメモリサイズ以下となるように定められて

20

いることを特徴とする請求項 1 記載のデータ処理装置。

【請求項 5】

前記演算処理手段は、演算処理の実行と並列に、次の演算処理に使用する重み係数を重み係数格納部へ格納することを特徴とする請求項 1 記載のデータ処理装置。

【請求項 6】

前記演算処理手段は、演算の実行と並列に次の処理に使用する動作パラメタをレジスタへ設定することを特徴とする請求項 1 記載のデータ処理装置。

【請求項 7】

前記処理単位は、該処理単位の処理に必要な特徴面を保持するためのメモリ量が、前記演算処理手段の参照するメモリにおける特徴面の格納部のメモリサイズ以下となるように定められていることを特徴とする請求項 1 記載のデータ処理装置。

10

【請求項 8】

前記転送手段は、複数のデータセットを連続して転送することを特徴とする請求項 1 記載のデータ処理装置。

【請求項 9】

前記動作パラメタは、前記演算処理において参照する特徴面の数、参照する特徴面のサイズ、前記演算処理に用いるフィルタカーネルのサイズの少なくとも 1 つを含むことを特徴とする請求項 1 記載のデータ処理装置。

【請求項 10】

入力データに対し、ニューラルネットワークに対応する階層的な演算処理を実行するデータ処理方法であって、

20

前記ニューラルネットワークに対応する階層的な演算処理を複数の処理単位に分割した処理単位ごとに処理に用いる、前記演算処理の重み係数と動作パラメタとを含む制御データをセットとした複数の制御データセットを保持手段に保持する保持工程と、

前記保持手段から前記複数の制御データセットを順に転送する転送工程と、

前記転送工程において転送された 1 つの制御データセットを用いて、当該制御データセットに対応する前記分割した処理単位の演算処理を行う演算処理工程とを備えることを特徴とするデータ処理方法。

【発明の詳細な説明】

【技術分野】

30

【0001】

本発明は、ニューラルネットワークに対応する階層的な演算処理を実行するデータ処理装置及びその方法に関するものである。

【背景技術】

【0002】

Convolutional Neural Networks (以下CNNと略記する) に代表される階層的な演算手法が認識対象の変動に対して頑健なパターン認識を可能にする手法として注目されている。例えば、非特許文献 1 では様々な応用例・実装例が開示されている。

【0003】

40

図 9 は簡単な CNN の例を示すネットワーク構成図である。901 は入力層であり、画像データに対して CNN 処理を行う場合、所定サイズの画像データに相当する。903a ~ 903d は第 1 階層 (908) の特徴面を、905a ~ 905d は第 2 階層 (909) の特徴面、907 は第 3 階層 910 の特徴面を示す。

【0004】

特徴面とは、所定の特徴抽出演算 (コンボリューション演算及び非線形処理) の処理結果に相当するデータ面である。特徴面は上位階層で所定の対象を認識するための特徴抽出結果に相当し、画像データに対する処理結果であるため、処理結果も面で表す。903a ~ 903d は 901 に対応するコンボリューション演算と非線形処理により生成されるものである。例えば、903a は 9021a に模式的に示す 2 次元のコンボリューション演

50

算と演算結果の非線形変換により算出する。

【 0 0 0 5 】

例えば、カーネル（係数マトリクス）サイズが $columnSize \times rowSize$ の2次元のコンボリューション演算は、以下に示すような積和演算により処理する。

【 0 0 0 6 】

【 数 1 】

$$output(x,y) = \sum_{row=-rowSize/2}^{rowSize/2} \sum_{column=-columnSize/2}^{columnSize/2} input(x+column,y+row) \times weight(column,row) \dots \quad (1)$$

10

【 0 0 0 7 】

ここで、

$input(x,y)$: 座標 (x,y) での参照画素値

$output(x,y)$: 座標 (x,y) での演算結果

$weight(column,row)$: $output(x,y)$ の演算に使用する重み係数

$columnSize, rowSize$: コンボリューションカーネルサイズ

である。

【 0 0 0 8 】

20

CNN処理では、複数のコンボリューションカーネルを画素単位で走査しながら積和演算を繰り返し、最終的な積和結果を非線形変換する事で特徴面を算出する。なお、903aを算出する場合は前階層との結合数が1であるため、コンボリューションカーネルは1つである。ここで、9021b、9021c、9021dはそれぞれ特徴面903b、903c、903dを算出する際に使用されるコンボリューションカーネルである。

【 0 0 0 9 】

図10は、特徴面905aを算出する場合の例を説明する図である。特徴面905aは4つの前階層(908)の特徴面903a~dと結合している。特徴面905aのデータを算出する場合、特徴面903aに対しては9041aで模式的に示すカーネルを用いたフィルタ演算を行い、結果を累積加算器1002に保持する。

30

【 0 0 1 0 】

同様に特徴面903b、903c、904cに対しては、それぞれカーネル9042a、9043a、9044aで示すカーネルのコンボリューション演算を行い(1001)、その結果を累積加算器1002に蓄積する。4種類のコンボリューション演算の終了後、ロジスティック関数や双曲正接関数(tanh関数)を利用した非線形変換処理を行う(1003)。

【 0 0 1 1 】

以上の処理を画像全体に対して1画素ずつ走査しながら処理する事で、特徴面905aを算出する。同様に、特徴面905bは、前階層(908)の特徴面に対して9041b/9042b/9043b/9044bで示す3つのコンボリューション演算を用いて算出する。更に特徴面907は前階層(909)の特徴面905a~dに対して9061、9062、9063、9064で示す4つのコンボリューション演算を用いて算出する。なお、各カーネル係数はパーセプトロン学習やバックプロパゲーション学習等の一般的な手法を用いて予め学習により決定されているものとする。

40

【 0 0 1 2 】

近年のディープネットと呼ばれる大規模な結合ネットワークでは、1つの特徴面を生成するために参照するデータ面の数、すなわち、結合数が増えており、それに伴い演算に必要な重み係数が増えている。さらに、認識精度の向上のため、階層数も増えており、階層毎の動作パラメータサイズも増えている。

【 0 0 1 3 】

50

CNN演算を行うCNN処理ハードウェアを、組み込みシステムに実装してネットワーク処理をする場合、CNN処理ハードウェアは、階層毎に、入力データと重み係数との演算を行う。そして演算結果を次の階層の入力データとし、次の階層の重み係数との演算を行うことを繰り返し、最終的なパターン認識結果を得る。

【0014】

例えば、特許文献1では、CNN処理に必要な各階層の重み係数や各階層の情報を内部のメモリへ転送し、内部のメモリから演算部へ効率よく転送しながら演算を行う。

【先行技術文献】

【特許文献】

【0015】

【文献】特表2018-521374号公報

【非特許文献】

【0016】

【文献】Yann LeCun, Koray Kavukcuoglu and Clement Farabet: Convolutional Networks and Applications in Vision, Proc. International Symposium on Circuits and Systems (ISCAS'10), IEEE, 2010.

【発明の概要】

【発明が解決しようとする課題】

【0017】

ハードウェア組み込みシステムに実装する場合、このように内部に大容量メモリを備えて全階層処理に必要な重み係数や動作パラメタを格納し、演算に使用することはコストの観点から望ましくない。そのため全階層の重み係数や動作パラメタは、安価なROMやDRAM等の外部メモリに格納し、1つの階層的な処理に使用する重み係数や動作パラメタのみを内部メモリと内部レジスタに格納し、使用する。

【0018】

このとき、各階層の処理ごとに、外部メモリからCNN処理ハードウェアへ、データ量の大きい重み係数や動作パラメタを転送してから処理を開始する。組み込みシステムでは一般的に、外部メモリや内部バスはCNN処理ハードウェア以外の機能ブロックも共有して使用するため、CNN処理ハードウェアがデータリードのリクエストを出してから、リードデータが到着するまでのメモリレイテンシが大きくなる。そのため各階層の切り替え時に処理の準備に時間がかかってしまう。

【0019】

一方、様々なニューラルネットワークを高速に処理する要求があり、高速で柔軟な処理制御が必要となっている。CNN処理ハードウェアを制御する制御プロセッサがニューラルネットワークに応じて処理制御を行うと、処理対象のニューラルネットワーク変更時、ファームウェアも変更する必要があり、ファームウェアの開発コストが増加してしまう。

【課題を解決するための手段】

【0020】

本発明の1態様によれば、入力データに対し、ニューラルネットワークに対応する階層的な演算処理を実行するデータ処理装置に、前記ニューラルネットワークに対応する階層的な演算処理を複数の処理単位に分割した処理単位ごとに処理に用いる、前記演算処理の重み係数と動作パラメタを含む制御データをセットとした複数の制御データセットを保持する保持手段と、前記保持手段から前記複数の制御データセットを順に転送する転送手段と、前記転送手段から転送された1つの制御データセットを用いて、当該制御データセットに対応する前記分割した処理単位の演算処理を行う演算処理手段とを備える。

【発明の効果】

【0021】

本発明によれば、1つの層のCNN処理中に、次階層のCNN処理に必要な重み係数と

10

20

30

40

50

動作パラメタを転送することで、各層の切り替え時にデータ転送時間を隠ぺいし、階層的な演算処理の処理性能を向上させることができる。

【0022】

さらに、処理対象のニューラルネットワークを変更するとき、処理単位に分割した動作パラメタと重み係数のデータセットを生成し、読み込み、実行することで、ニューラルネットワークに応じて制御プロセッサの処理制御を変更することなくニューラルネットワークの変更が可能となる。

【図面の簡単な説明】

【0023】

【図1】実施形態に関するCNN演算処理部のブロック図。

10

【図2】データ分配部101が受信するデータフォーマットを示す図。

【図3】動作設定レジスタ群を示す図。

【図4】CNN演算処理部601の特徴面生成処理のフローチャート。

【図5】CNN演算処理部601を動作させるためのデータセットを示す図。

【図6】本発明に関するパターン認識装置を利用した画像処理システムの構成例を示す図。

【図7】ネットワークの処理に必要なメモリ量と処理単位を示す図。

【図8】画像処理システムのシーケンス図。

【図9】CNNの例を示すネットワーク構成図。

【図10】特徴面905aを算出する場合の例を説明する図。

【発明を実施するための形態】

20

【0024】

以下、図面を参照しながら本発明の好適な実施形態について詳細に説明する。

【0025】

(実施形態1)

図1～図8を参照して本発明の実施形態1を説明する。

【0026】

図6は本発明に関するパターン認識装置を利用した画像処理システムの構成例を示すものである。当該システムは画像データから特定の物体の領域を検出する機能を有する。図6において600は画像入力部であり、光学系、CCD(Charge-Coupled Devices)又はCMOS(Complimentary Metal Oxide Semiconductor)センサー等の光電変換デバイスを有する。また画像入力部600は、センサーを制御するドライバー回路/A/Dコンバーター/各種画像補正を司る信号処理回路/フレームバッファ等を備える。

30

【0027】

601は、本実施形態に関する演算装置を含むCNN演算処理部である。605はDMAC(Direct Memory Access Controller)であり、画像バス602上の各処理部とCPUバス609間のデータ転送を司るデータ転送部である。603はブリッジであり、画像バス602とCPUバス609のブリッジ機能を提供する。604は前処理部であり、パターン認識処理を効果的に行うための各種前処理を行う。具体的には色変換処理/コントラスト補正処理等の画像データ変換処理をハードウェアで処理する。606はCPUであり、画像処理システム全体の動作を制御する。

40

【0028】

607はROM(Read Only Memory)であり、CPU606の動作を規定する命令や動作パラメタデータ、及び、CNN演算処理部601を動作させるための制御データセットを格納するデータ保持部である。制御データセットは、処理するニューラルネットワークに応じた重み係数と動作パラメタを含む(以下、単にデータセットとも表現する)。データセットはDMAC605を介してパターン認識処理部に入力される。608はCPU606の動作に必要なメモリであり、RAMが使用される。

【0029】

画像入力部600から入力された画像データは前処理部604で処理され、一旦RAM

50

608に格納する。その後、RAM608からDMAC605を介してCNN演算処理部601に入力する。CNN演算処理部601は前記データセットに従い、入力された前処理後の画像データに対して画素単位で所定の判別処理を行い、入力画像中の所定の物体の領域を判定する。判定結果はDMAC605を介してRAM608に格納する。

【0030】

図1に、CNN演算処理部601のブロック図を示す。各ブロックの機能について説明する。101はDMAC605から送信されるデータを受信し、内部のブロックへ送信するデータ分配部である。データのヘッダに付加された行き先を示す指示子と転送長を参照し、該当するブロックへデータを送信する。データの送受信は標準的な二線ハンドシェイクプロトコル(Valid/Ready)を使って行われるものとする。

10

【0031】

図2に、データ分配部101が受信するデータフォーマットを示し、データ分配部101の動作について説明する。受信データ幅は32ビットとする。201はヘッダ、202はデータボディである。ヘッダの下位8ビットの「DST」フィールドは行き先を示し、上位24ビットの「LENGTH」フィールドはデータボディの転送長(バイト数)を示す。予め内部のブロックにユニークなIDを付加し、DSTフィールドの値と一致したIDのブロックへLENGTH分のデータボディを送信する。例えば、図2に示すデータを受信すると、最初のデータは、ヘッダはDST=0、LENGTH=Nであるので、IDが0の内部ブロックへNバイトのデータボディを送信する。続くデータは、ヘッダはDST=1、LENGTH=Mであるので、IDが1の内部ブロックへMバイトのデータボディを送信する。なお、行き先を示すIDがアドレスであっても良い。

20

【0032】

102はCNN演算処理部601の全体を制御する制御部であり、内部に動作設定レジスタ群を備える。動作設定レジスタ群は複数の動作設定レジスタセットを含み、1つの特徴面生成処理を行うための情報が当該動作設定レジスタセット毎に保持される。制御部102の動作設定レジスタ群は8つの動作設定レジスタセットを備えるものとする。設定する値は、動作パラメタとしてROM607に配置されており、制御部102は、DMAC605、データ分配部101経由で受信し、内部の動作設定レジスタセットに設定する。制御部102のIDは「0」とする。

【0033】

設定したレジスタ値は、後述する重み係数格納部103、特徴面格納部104、畳み込み演算部105へ、処理制御信号として配る。処理制御信号はレジスタ値の他に処理開始を指示する信号を含み、1つ以上の動作設定レジスタセットが設定されたときに処理開始指示をする。

30

【0034】

図3は動作設定レジスタセットを示す図である。「最終層指定」レジスタは、本動作設定レジスタセットにより生成する特徴面が最終層か否かを指定するレジスタであり、1の場合、特徴面の生成処理を終了すると検出処理結果は外部へ出力する。「参照特徴面の数」レジスタは、生成特徴面と接続する前階層の特徴面数を指定するレジスタであり、例えば図9に示す特徴面905a~dを演算する場合「4」が設定される。「非線形変換」レジスタは、非線形変換の有無を指定するレジスタであり、1が設定されている場合、非線形変換処理を行い、0が設定されている場合非線形変換処理を行わない。

40

【0035】

「演算結果格納先ポインタ」レジスタは、生成特徴面の演算結果を保持する後述の特徴面格納部104の先頭ポインタを示すアドレスを指定するレジスタであり、当該ポインタ値を先頭ポインタとして演算結果をラスタスキャン順に格納する。「フィルタカーネルの水平サイズ」及び「フィルタカーネルの垂直サイズ」レジスタは生成特徴面の演算に使用するフィルタカーネルのサイズを指定するレジスタである。

【0036】

「重み係数格納先ポインタ」レジスタは生成特徴面の演算に使用する重み係数を保持す

50

る、後述の重み係数格納部 103 の格納先アドレスを示すレジスタである。例えば、重み係数データは「参照特徴面の数」レジスタと同じ数の係数の組を有し、「重み係数格納先ポインタ」レジスタで指定されるアドレスからラスタスキャン順に格納されている。すなわち、「フィルタカーネルの水平サイズ」×「フィルタカーネルの垂直サイズ」×「参照特徴面の数」の数の係数データが重み係数格納部 103 に格納されている。

【0037】

「参照特徴面の水平サイズ」及び「参照特徴面の垂直サイズ」レジスタはそれぞれ参照特徴面の水平方向画素数及び垂直方向ライン数を示すレジスタである。「参照特徴面格納先ポインタ」レジスタは、参照特徴面を保持する後述の特徴面格納部 104 の先頭ポインタを示すアドレスを指定するレジスタであり、該当ポインタ値を先頭ポインタとし参照特徴面がラスタスキャン順に格納されている。すなわち、「参照特徴面の水平サイズ」×「参照特徴面の垂直サイズ」×「参照特徴面の数」の数の特徴面データが特徴面格納部 104 に格納されている。

10

【0038】

このような複数のレジスタが各特徴面単位に設けられている。生成特徴面の「参照特徴面格納ポインタ」レジスタの内容が前階層結合対象特徴面の「演算結果格納先ポインタ」である場合、前階層の特徴面と生成特徴面が結合されていることになる。従って、ここでのレジスタ設定だけで任意の階層的結合関係を特徴面単位に構築することが可能である。ROM 607 に配置される動作パラメタはこのようなレジスタセットが生成特徴面の数分連続するデータ構造であり、制御部 102 は、8 セットある動作設定レジスタセットに順に設定する。8 セットすべてに設定すると、二線ハンドシェイク信号の Ready をネゲートし、続く動作パラメタデータを受信しない。先に設定した動作設定レジスタセットの処理が完了と、処理が完了した動作レジスタセットは解放され、次の動作パラメタデータを受信し、設定する。

20

【0039】

103 は、重み係数データを格納する重み係数格納部であり、メモリとメモリ制御部で構成される。メモリの読み出しアドレスを示すリードポインタ、書き込みアドレスを示すライトポインタを備え、メモリをリングバッファとして使用する。通常のリングバッファであればメモリの空き容量は、ライトポインタとリードポインタにより管理されるが、後述する図 4 のフローチャートによると、同一の重み係数が複数回読み出されて使用される。最後の読み出しが終わったときにその重み係数のメモリ領域を解放し、ライトポインタとの差分でメモリの空き容量を管理する。メモリに空き容量がないときは、二線ハンドシェイク信号の Ready をネゲートし、データを受信しない。重み係数格納部 103 の ID は「1」とする。

30

【0040】

制御部 102 より処理開始指示を受けると、レジスタ「重み係数格納先ポインタ」、「フィルタカーネルの水平サイズ」、「フィルタカーネルの垂直サイズ」、「参照特徴面の数」の値を参照する。リードポインタを「重み係数格納先ポインタ」レジスタ値に更新し、「フィルタカーネルの水平サイズ」×「フィルタカーネルの垂直サイズ」×「参照特徴面の数」の重み係数を読み出し、後述の畳み込み演算部 105 へ送信する。また、ライトポインタは初期化時にメモリの先頭アドレスに更新する。データ分配部 101 より重み係数を受信するとライトポインタの示すアドレスへデータを格納し、ライトポインタをインクリメントする。ライトポインタがメモリの最終アドレスまで到達するとメモリの先頭アドレスに戻す。

40

【0041】

104 は、特徴面データを格納する特徴面格納部であり、メモリとメモリ制御部で構成される。メモリの読み出しアドレスを示すリードポインタと、書き込みアドレスを示すライトポインタを備える。特徴面格納部 104 の ID は「2」とする。

【0042】

特徴面格納部 104 は、制御部 102 より処理開始指示を受ける。すると、レジスタ「

50

参照特徴面格納先ポインタ」、「参照特徴面の水平サイズ」、「参照特徴面の垂直サイズ」、「参照特徴面の数」、「フィルタカーネルの水平サイズ」、「フィルタカーネルの垂直サイズ」の値を参照する。そして、式(1)に示した参照画素値を読み出し、後述の畳み込み演算部105へ送信する。具体的には、リードポインタを「参照特徴面格納先ポインタ」レジスタ値に更新する。「フィルタカーネルの水平サイズ」($= \text{columnSize}$)「フィルタカーネルの垂直サイズ」($= \text{rowSize}$)で1枚目の参照特徴面座標(0, 0)を決定する。そしてこの座標での参照画素値($-\text{columnSize}/2$ 、 $-\text{rowSize}/2$) \sim ($\text{columnSize}/2$ 、 $\text{rowSize}/2$)を読み出し、送信する。続いて2枚目の参照特徴面座標(0, 0)での参照画素値、3枚目の参照特徴面(0, 0)における参照画素値 \dots の順に「参照特徴面の数」読み出し、送信する。その後、参照特徴面を水平方向に1画素ずつずらして参照特徴面(1, 0)における参照画素値($1 - \text{columnSize}/2$ 、 $-\text{rowSize}/2$) \sim ($1 + \text{columnSize}/2$ 、 $\text{rowSize}/2$)を同様に「参照特徴面の数」分送信する。このようにして「参照特徴面の水平サイズ」分水平方向にずらして参照特徴面を送信する。次は垂直方向に1ラインずつずらした参照特徴面(0, 1)での参照画素値($-\text{columnSize}/2$ 、 $1 - \text{rowSize}/2$) \sim ($\text{columnSize}/2$ 、 $1 + \text{rowSize}/2$)を同様に送信する。「参照特徴面の垂直サイズ」分ライン方向にずらして送信すると処理を完了する。

10

【0043】

また、特徴面格納部104は、制御部102より処理開始指示を受けると、参照特徴面の読み出しと送信を行うのと同時に、後述の畳み込み演算部105から演算結果である生成特徴面を受信し、メモリに書き込む。この処理はレジスタ「最終層指定」、「演算結果格納先ポインタ」、「参照特徴面の水平サイズ」、「参照特徴面の垂直サイズ」の値を参照する。具体的には、ライトポインタを「演算結果格納先ポインタ」レジスタ値に更新する。「最終層指定」レジスタが0のとき、生成特徴面を受信し、ライトポインタのアドレスにデータを格納する。ライトポインタはインクリメントする。「参照特徴面の水平サイズ」 \times 「参照特徴面の垂直サイズ」の演算結果を受信すると処理が完了する。「最終層指定」レジスタが1のときは生成特徴面を受信は行わない。

20

【0044】

また、ライトポインタは、初期化時にメモリの先頭アドレスに更新する。制御データ分配部101より画像データを受信するとライトポインタが示すアドレスにデータを格納し、ライトポインタをインクリメントする。入力層処理時のみ画像データを受信する。

30

【0045】

105は、畳み込み演算を行う畳み込み演算部である。制御部102より処理開始指示を受けると、重み係数格納部103と特徴面格納部104からデータを受信し、畳み込み演算を行う。「フィルタカーネル水平サイズ」 \times 「フィルタカーネル垂直サイズ」の重み係数と参照画素値を受信すると畳み込み演算を行い、演算結果を内部の累積加算器に保持し、演算結果を累積する。「参照特徴面の数」分の演算結果を累積加算した結果を、レジスタ「非線形変換」の値が1のとき、非線形変換処理をし、演算結果を特徴面格納部104へ送信する。レジスタ「非線形変換」の値が0のとき、非線形変換処理は行わず演算結果を特徴面格納部104へ送信する。また、「最終層指定」が1のときは演算結果を外部へ送信し、DMAC605を介してRAM608へ格納する。

40

【0046】

CNN演算処理部601による特徴面生成処理のフローチャートを図4に示す。CNN演算処理部601は、制御部102の動作設定レジスタ群のうち、1つ以上のレジスタセットが設定されるとフローチャートの処理を開始する。S401にて、制御部102は、重み格納部103、特徴面格納部104、畳み込み演算部105へ処理開始指示を出す。S402は最も外側のループで「参照特徴面の垂直サイズ」のループ、S403は「参照特徴面の水平サイズ」のループである。つまり、参照特徴面をラスタ順にスキャンし、特徴面をラスタ順に生成する。

50

【 0 0 4 7 】

続くループ S 4 0 4 は「参照特徴面の数」、S 4 0 5 は「フィルタカーネル垂直サイズ」、S 4 0 6 は「フィルタカーネル水平サイズ」である。つまり、参照特徴面毎に、参照画素範囲の参照画素値が、ラスト順に処理される。S 4 0 7 は重み係数格納部 1 0 3 が参照特徴面の重み係数を読み出すステップであり、重み係数の座標は、S 4 0 4、S 4 0 5 のループ変数により決定する。S 4 0 8 は特徴面格納部 1 0 4 が参照特徴面の画素値を読み出すステップであり、読み出す画素値の座標は S 4 0 2、S 4 0 3、S 4 0 5、S 4 0 6 のループ変数により決定する。

【 0 0 4 8 】

S 4 0 9 は、畳み込み演算部 1 0 5 が S 4 0 7、S 4 0 8 にて読み出した画素値と重み係数を畳み込み演算し、その結果を累積するステップである。畳み込み演算部 1 0 5 は、S 4 0 4 のループが終了すると、S 4 1 0 にて「非線形変換」レジスタを参照し、1 であると、S 4 1 1 へ遷移し、非線形変換処理を行う。レジスタ値が 0 であるとき、S 4 1 1 をスキップする。S 4 1 2 は畳み込み演算部 1 0 5 が「最終層指定」レジスタを参照し、1 であるとき、S 4 1 3 にて、検出処理結果を外部へ送信する。レジスタ値が 0 であるとき、S 4 1 4 にて演算結果を特徴面格納部 1 0 4 へ送信する。特徴面格納部 1 0 4 はライトポインタに従い演算結果を格納する。

【 0 0 4 9 】

1 枚の特徴面生成処理が完了すると、フローチャートの処理は完了となる。処理済の重み係数格納メモリ領域と、レジスタセットを解放する。

【 0 0 5 0 】

図 5 に、RAM 6 0 8 に格納された、前処理部 6 0 4 で処理された処理対象の画像データと、ROM 6 0 7 に格納された CNN 演算処理部 6 0 1 を動作させるためのデータセットの一例を示す。図 7 に示すネットワーク構成を処理するためのデータセットである。画像データを入力として、第 1 階層の特徴面 9 0 3 a ~ d を順に生成し、続いて生成した第 1 階層 9 0 8 の特徴面を入力として、第 2 階層の特徴面 9 0 5 a ~ d を順に生成する。引き続き第 2 階層の特徴面 9 1 0 を入力として、第 3 階層の特徴面 9 0 7 の特徴面を生成する一連の処理をするためのデータセットである。

【 0 0 5 1 】

5 0 1 は RAM 6 0 8 に格納された画像データ、5 0 2 ~ 5 0 7 は ROM 6 0 7 に格納されたデータセットであり、DMAC 6 0 5 は、5 0 1 ~ 5 0 7 の順にデータを CNN 演算処理部 6 0 1 へ送信する。5 0 1 は処理対象の画像データであり、ヘッダは、行き先を示す指示子 DST = 2、LENGTH = 画像データの転送長である。つまり、5 0 1 のデータ部は分配部 1 0 1 を経由して特徴面格納部へと送信される。

【 0 0 5 2 】

5 0 2 は第 1 階層の重み係数であり、ヘッダは、行き先を示す指示子 DST = 1、LENGTH = 重み係数の転送長である。つまり、5 0 2 のデータ部は分配部 1 0 1 を経由して重み係数格納部 1 0 3 へと送信される。データは、9 0 2 1 a ~ d の重み係数を含む。重み係数は演算に使用する順に並べるものとする。

【 0 0 5 3 】

5 0 3 は第 1 階層 9 0 8 の特徴面を生成するための動作パラメタであり、ヘッダは、行き先を示す指示子 DST = 0、LENGTH = レジスタセットの転送長である。つまり、5 0 3 のデータ部（動作パラメタ）は分配部 1 0 1 を経由して制御部 1 0 2 へと送信される。動作パラメタは、図 3 に示す第 1 階層の特徴面 9 0 3 a ~ d を生成するために動作設定レジスタセットに設定する値である。レジスタ「参照特徴面の水平サイズ」の設定値は、入力層 9 0 1 の水平サイズ、レジスタ「参照特徴面の垂直サイズ」の設定値は入力層 9 0 1 の垂直サイズ、レジスタ「参照特徴面格納先ポインタ」の設定値は特徴面格納部メモリの先頭アドレスとなる。動作パラメタは処理する順、ここでは 9 0 3 a ~ d の順に並べるものとする。

【 0 0 5 4 】

10

20

30

40

50

504、505は第2階層の重み係数9041a~d、9042a~d、9043a~d、9044a~dと、第2階層の特徴面905a~dを生成するための動作パラメタである。重み係数は演算に使用する順、9041a~9044a、9041b~9044b、9041c~9044c、9041d~9044d、に並べるものとする。動作パラメタに含まれる、レジスタ「参照特徴面の水平サイズ」、「参照面特徴面の垂直サイズ」の設定値は第1階層908の特徴面の水平サイズ、垂直サイズである。また、レジスタ「参照特徴面格納先ポイント」の設定値は、503の中で指定した演算結果格納先ポイントのアドレスである。

【0055】

506、507は、第3階層の重み係数9061~9064と第3階層の特徴面907を生成するための動作パラメタである。動作パラメタに含まれる、レジスタ「最終層指定」の設定値は1である。レジスタ「参照特徴面の水平サイズ」、「参照面特徴面の垂直サイズ」の設定値は第2階層909の特徴面の水平サイズ、垂直サイズであり、レジスタ「参照特徴面格納先ポイント」の設定値は、505の中で指定した演算結果格納先ポイントのアドレスである。ここでは、503、505、507の動作パラメタ値そのものを、制御部102がレジスタセット群に設定するが、制御部102が動作パラメタから値を算出し、レジスタセットへ設定してもよい。

10

【0056】

図5に示すように、CNN演算処理部601を動作させるためのデータセットは、ある処理単位の重み係数と動作パラメタ（生成するレジスタセットの設定値）のペアで構成する。処理対象のネットワークの情報と図1のハードウェアの構成に基づいて処理単位を決定する。処理単位の決定に作用するハードウェアの構成とは、具体的には、重み係数格納部103のメモリサイズや特徴面格納部104のメモリサイズである。単位当たりの処理に必要な重み係数のメモリ量が、重み係数格納部103のメモリサイズ以下になるよう決定する。且つ、単位当たりの処理に必要な特徴面のメモリ量が、特徴面格納部104のメモリサイズ以下になるよう決定する。つまり、ハードウェアで処理可能なように決定する。

20

【0057】

各階層を処理するときに必要な重み係数を保持するために必要なメモリ量 W_1 、特徴面を保持するために必要なメモリ量 S_1 は、

$$W_1 = W_x(1, f, f') \times W_y(1, f, f') \times (F_t \times F_{t-1})$$

30

$$S_1 = (I_x \times I_y) \times (F_t + F_{t-1})$$

1：生成対象となる階層番号（1，2，...）

f：生成特徴面番号

f'：参照特徴面番号

$W_x(1, f, f')$ ：カーネルサイズ水平方向

$W_y(1, f, f')$ ：カーネルサイズ垂直方向

F_t ：生成特徴面数

F_{t-1} ：参照特徴面数

I_x, I_y ：入力画像サイズ（水平方向，垂直方向）

として表すことができる。

40

【0058】

図9のネットワークを階層単位で処理するときに必要なメモリ量を図7の701に示す。図9のネットワークの入力画像サイズは 32×32 とする。第1階層の参照特徴面数は1、生成特徴面数は4である。カーネルフィルタのサイズを 3×3 とすると、4枚の特徴面を生成するために使用する重み係数9021a~dのサイズは、係数マトリクスの1要素あたりのデータを1バイトとすると、 $3 \times 3 \times 4 = 36$ バイトである。特徴面のデータ量は、 $32 \times 32 \times 5 = 5120$ バイトである。

【0059】

第2階層のカーネルフィルタのサイズを 5×5 とすると、4枚の特徴面を生成するために使用する重み係数9041a~d、9042a~d、9043a~d、9044a~

50

dのサイズは同様に算出すると400バイトとなる。また、特徴面のデータ量は8192バイトである。第3階層のカーネルフィルタのサイズを7×7とすると、4枚の特徴面を生成するために使用する重み係数9061～9064のサイズは同様に算出すると196バイト、特徴面のデータ量5120バイトである。

【0060】

ここでは重み係数格納部103のメモリサイズを436バイト、特徴面格納部104のメモリサイズが16Kバイトとする。各階層の処理に必要な重み係数のメモリ量が、重み係数格納部103のメモリサイズ以下、必要な特徴面のメモリ量が、特徴面格納部104のメモリサイズ以下であるので、処理単位を階層単位とし、図5に示すデータセットとなる。処理単位の決め方の他の例を述べる。図9のネットワークの第2階層の特徴面数が8の場合、階層単位で処理するときに必要なメモリ量は702に示す通り、第2階層の重み係数は800バイト、特徴面のデータ量は12288バイトである。重み係数格納部103のメモリサイズを超えるため、階層単位で処理できないので1つの階層を複数回に分けて処理する。703のように第2階層の処理を、重み係数格納部103に格納可能な2分割とし、4つの処理とする。データセットは4セットとなる。

10

【0061】

さらに他の例として、特徴面格納部104のメモリサイズが8Kバイトのハードウェアで、702に示すネットワークを処理する場合について述べる。第2階層の特徴面のデータ量は12288バイトであり、特徴面格納部104のメモリサイズを超える。この場合、入力画像を分割して複数回処理する。704は、32×32を、32×16の2ブロックに分割して、特徴面格納部104に格納可能な処理単位とした。データセットは6セットとなる。なお、データセットを3セットとし、同一のデータセットを2回使用してもよい。また、全階層の重み係数が格納可能であれば、2回目は重み係数を転送しなくてもよい。また、全階層の動作パラメタが格納可能であれば、2回目は動作パラメタを転送しなくてもよい。

20

【0062】

図6の画像処理システムにて、画像データ1枚を図5のデータセットにより処理する動作について、図8のシーケンス図を用いて説明する。図8はCPU606、DMAC605、制御部102、重み係数格納部103、特徴面格納部104、畳み込み演算部105のそれぞれの動作と相互作用を時系列に示している。

30

【0063】

画像入力部600が画像データを受信すると前処理部604が処理を行い、画像データをRAM608に格納する。1枚分の画像データ処理が完了するとCPU606に通知する(801)。CPU606は、CNN演算処理部601へ初期化を指示し、重み係数格納部103、特徴面格納部104は、それぞれのライトポイントをメモリの先頭アドレスに更新する(802)。CPU606は、RAM608の画像データ領域の直前のアドレスに、501に示す画像データのLENGTHとDST(=2)を付加する(803)。

【0064】

その後、DMAC605にデータ転送元に501の先頭アドレス、転送先にCNN演算処理部601、転送長に501のサイズを設定し、起動すると、DMAC605はRAM608からCNN演算処理部601へ501領域を転送する(804)。図8に図示しないデータ分配部101は、ヘッダのDSTが2であるため、特徴面格納部104へ送信し、特徴面格納部104はライトポイントが示すメモリの先頭アドレスから画像データを格納する(805)。DMAC605は全データの転送完了後、CPU606へ完了を通知する。

40

【0065】

CPU606は、完了通知を受けると、DMAC605のデータ転送元に502の先頭アドレス、転送先にCNN演算処理部601、転送長に502～507の合計サイズを設定する。そして、起動すると、DMAC605は、ROM607からCNN演算処理部601へ502領域を転送する(806)。データ分配部101はヘッダのDSTが1であ

50

るため、重み係数格納部 103 へ送信し、重み係数格納部 103 はメモリの空き領域があるときはデータを受信し、メモリのライトポイントの示すメモリの先頭アドレスから順に重み係数を格納する(807)。DMAC605は502領域の転送完了後、引き続き503領域を転送する(808)。データ分配部101はヘッダのDSTが0であるため、制御部102へ送信し、制御部102は第1階層の特徴面903a~dの4つのレジスタセットを動作設定レジスタへ設定する(809)。

【0066】

動作設定レジスタ1セット以上に値が設定されたため、図4のフローチャートに従い特徴面生成処理を開始する。最初に、第1階層の特徴面903aを生成するレジスタセットを処理する。制御部102は重み格納部103、特徴面格納部104、畳み込み演算部105へ処理開始指示を出す(S401)。重み格納部103、特徴面格納部104、畳み込み演算部105は、903aのレジスタセットの各レジスタ値を参照する。それにより、S407-S409を「フィルタカーネル水平サイズ」×「フィルタカーネル垂直サイズ」×「参照特徴面の数」回、繰り返し実行する(810)。

10

【0067】

さらに、903aのレジスタセットの各レジスタ値を参照し、810、S410、S414を「参照特徴面の水平サイズ」×「参照特徴面の垂直サイズ」回、繰り返し実行し、第1階層の特徴面903aの生成が完了する(811)。完了時に、処理済の重み係数格納メモリ領域と、レジスタセットが解放される。第1階層の特徴面903aの生成処理(811)と並列に、DMAC605は503領域の転送完了後、引き続き504領域を転送する(812)。

20

【0068】

データ分配部101はヘッダのDSTが1であるため、重み係数格納部103へ送信する。重み係数格納部103は第1階層の重み係数502が格納されており、まだメモリの空き領域があるため、データを受信しメモリに格納する(813)。DMAC605は504領域の転送完了後、引き続き505領域を転送する(814)。データ分配部101はヘッダのDSTが0であるため、制御部102へ送信する。制御部102は第1階層の4つのレジスタセット503が設定されており、まだレジスタセットの空きがあるため、データを受信し、第2階層の特徴面905a~dの4つのレジスタセットを動作設定レジスタへ設定する(815)。

30

【0069】

DMAC605は、505領域の転送完了後、引き続き506領域を転送する(816)。データ分配部101はヘッダのDSTが1であるため、重み係数格納部103へ送信する。重み係数格納部103は第1階層の重み係数502と第2階層の重み係数504が格納されており、メモリに空き領域がないため、二線ハンドシェイク信号を制御し、データを受信しない(817)。第1階層の特徴面903a生成完了後に、処理済の重み係数9021a格納メモリ領域が解放された後に、受信する(818)。

【0070】

第1階層の特徴面903a生成完了後、特徴面903aのレジスタセットも解放する。909にて4セット設定しており、特徴面903b~dのレジスタセットについても同様に処理する。第1階層の最後の特徴面903dの生成完了時、第2階層の重み係数は(813)にて既に格納済みである。動作設定レジスタも(815)にて既に設定済みである。従って制御部102はすぐに処理開始を指示し、第2階層特徴面905aの生成を開始する(819)。

40

【0071】

DMAC605は、506領域の転送完了後、引き続き507領域を転送する(820)。データ分配部101はヘッダのDSTが0であるため、制御部102へ送信する。制御部102は第1階層の4つのレジスタセットは解放されており、第2階層の4つのレジスタセット505が設定されている状態である。レジスタセットの空きがあるため、データを受信し、第3階層の特徴面907の1つのレジスタセットを動作設定レジスタへ設定す

50

る(821)。DMAC605は502～507全データの転送完了後、CPU606へ完了を通知する。CPU606は完了通知を受けると、DMAC605のデータ転送元にCNN演算処理部601、転送先にRAM608、転送長に検出処理結果のサイズを設定し、起動する。

【0072】

一方、CNN演算処理部601は引き続き特徴面生成処理をしており、第3階層の特徴面907の生成では、「最終層指定」レジスタは1であるため、畳み込み演算部105は検出処理結果を外部へ出力する。出力された検出処理結果は、DMAC605がRAM608へ転送する(822)。

【0073】

以上、説明したように本実施形態の処理装置によれば、各階層の特徴面生成完了時に次階層の処理に必要な重み係数は既に重み係数格納部に格納されており、また、次階層のレジスタセットのうち少なくとも1セット以上は動作設定レジスタに設定されている。そのため階層の切り替え時にデータ準備の時間を省き、すぐに演算処理が開始可能であるため、パターン認識装置の処理性能を向上させることができる。

【0074】

また、処理対象のニューラルネットを変更するときは、703、704に示すようにハードウェア処理装置で実行可能な処理単位にし、対応するデータセットにより動作させる。

【0075】

CPU606がDMAC605に設定するデータセットの転送長が異なるだけで、制御手順は図8と変わらない。

【0076】

本実施形態の処理装置によれば、処理対象のニューラルネットワークを変更するとき、ハードウェア処理装置で実行可能な処理単位に分割したデータセットに基づいて動作する。これにより、ニューラルネットワークに応じて制御プロセッサの処理制御を変更することなくニューラルネットワークの変更が可能となる。

【0077】

(その他の実施形態)

上記実施形態では特徴抽出処理としてCNN演算処理の場合について説明したが、これに限るわけではなく、パーセプトロン等他の様々な階層的な処理に適用可能である。

【0078】

また上記実施形態では、畳み込み演算をハードウェアで構成する場合について説明したが、本発明はこれに限るわけではない。全ての処理を汎用のプロセッサでソフトウェア処理する場合にも適用可能である。

【符号の説明】

【0079】

- 101 データ分配部
- 102 制御部
- 103 重み係数格納部
- 104 特徴面格納部
- 105 畳み込み演算部
- 601 CNN演算処理部

10

20

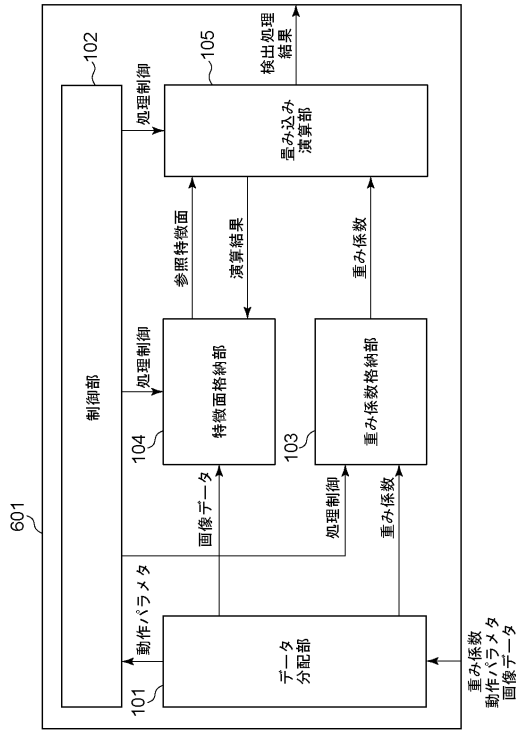
30

40

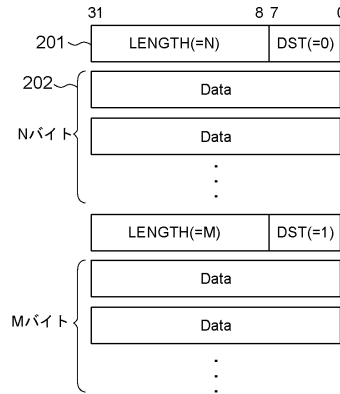
50

【図面】

【図 1】



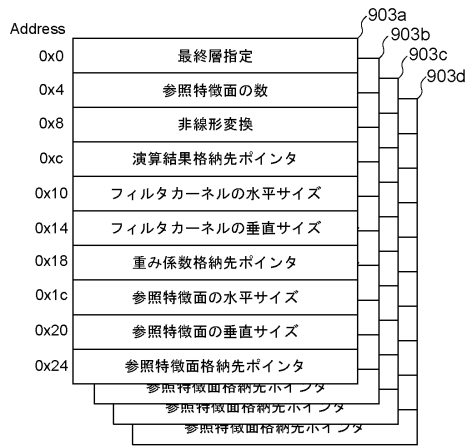
【図 2】



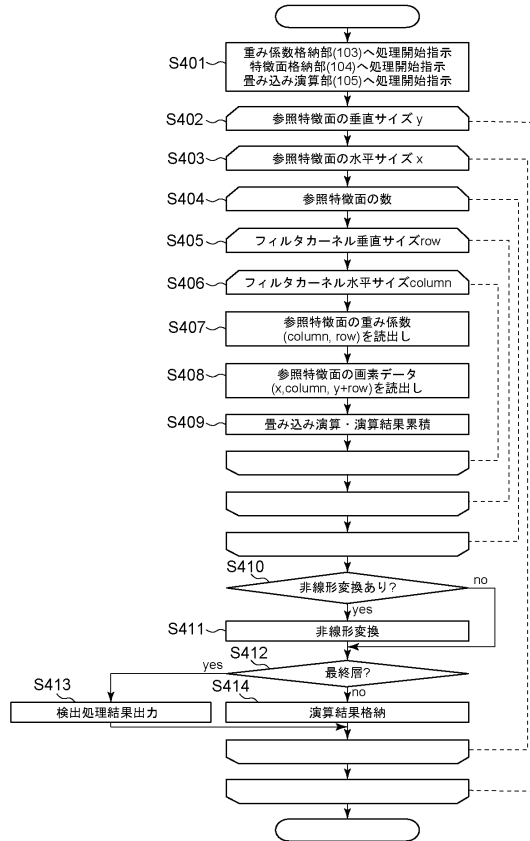
10

20

【図 3】



【図 4】

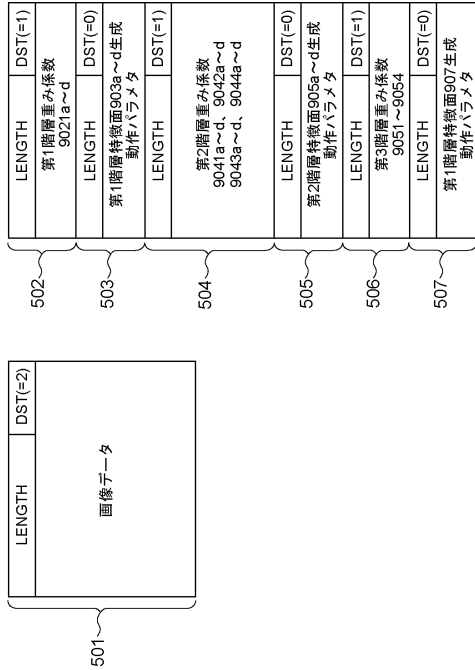


30

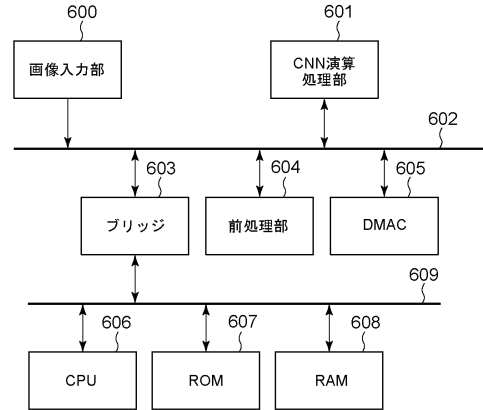
40

50

【図5】



【図6】



10

20

【図7】

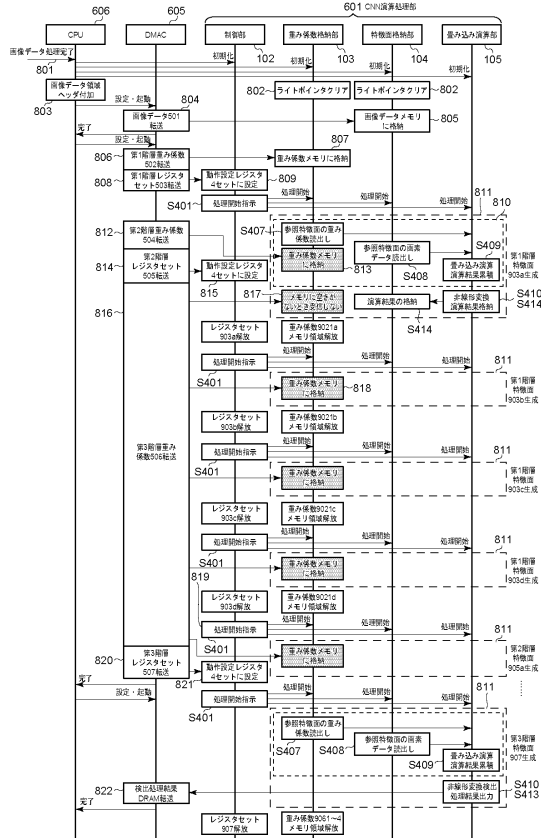
階層	参照特徴面数	カーネルフィルタサイズ	生成特徴面数	重み係数 (byte)	特徴面メモリ (byte)	処理番号
第1階層	1	3×3	4	36	5120	1
第2階層	4	5×5	4	400	8192	2
第3階層	4	7×7	1	196	5120	3

階層	参照特徴面数	カーネルフィルタサイズ	生成特徴面数	重み係数 (byte)	特徴面メモリ (byte)	処理番号
第1階層	1	3×3	4	36	5120	1
第2階層	4	5×5	8	800	12288	2
第3階層	4	7×7	1	196	5120	3

階層	参照特徴面数	カーネルフィルタサイズ	生成特徴面数	重み係数 (byte)	特徴面メモリ (byte)	処理番号
第1階層	1	3×3	4	36	5120	1
第2階層	4	5×5	4	400	12288	2
第3階層	4	7×7	1	196	5120	3

階層	参照特徴面数	カーネルフィルタサイズ	生成特徴面数	重み係数 (byte)	特徴面メモリ (byte)	処理番号
第1階層#BLK0	1	3×3	4	36	2560	1
第2階層#BLK0	4	5×5	8	800	6144	2
第3階層#BLK0	4	7×7	1	196	2560	3
第1階層#BLK1	1	3×3	4	36	2560	4
第2階層#BLK1	4	5×5	8	800	6144	5
第3階層#BLK1	4	7×7	1	196	2560	6

【図8】

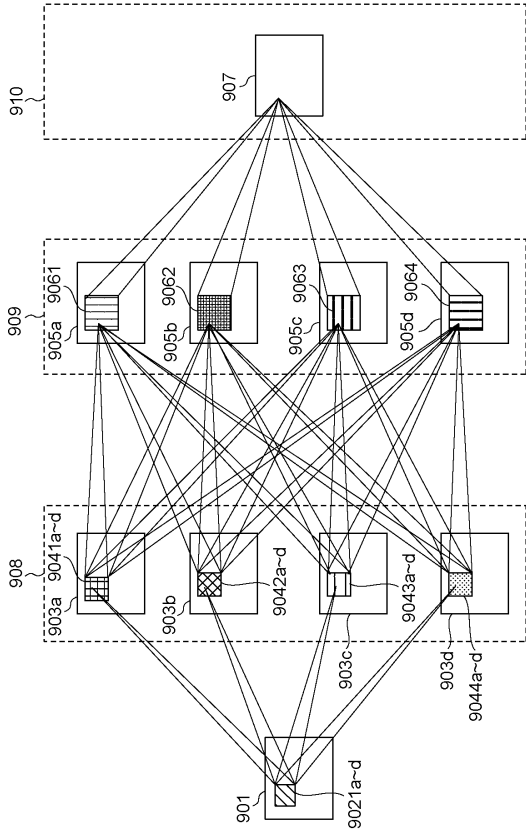


30

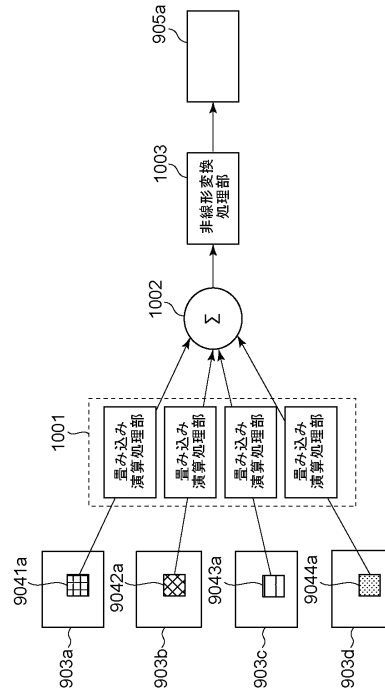
40

50

【図 9】



【図 10】



10

20

30

40

50

フロントページの続き

- (56)参考文献 特開2009-080693(JP,A)
特開2018-147182(JP,A)
特開2018-073102(JP,A)
特開2017-004142(JP,A)
- (58)調査した分野 (Int.Cl., DB名)
G06N 3/00-99/00
G06T 1/40