

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5368687号
(P5368687)

(45) 発行日 平成25年12月18日(2013.12.18)

(24) 登録日 平成25年9月20日(2013.9.20)

(51) Int.Cl.		F I	
G06N	3/063	(2006.01)	G06N 3/063
G06N	3/00	(2006.01)	G06N 3/00 560C
G06T	7/00	(2006.01)	G06T 7/00 350C

請求項の数 13 (全 37 頁)

(21) 出願番号	特願2007-250063 (P2007-250063)	(73) 特許権者	000001007
(22) 出願日	平成19年9月26日 (2007.9.26)		キヤノン株式会社
(65) 公開番号	特開2009-80693 (P2009-80693A)		東京都大田区下丸子3丁目30番2号
(43) 公開日	平成21年4月16日 (2009.4.16)	(74) 代理人	100076428
審査請求日	平成22年9月27日 (2010.9.27)		弁理士 大塚 康德
		(74) 代理人	100112508
			弁理士 高柳 司郎
		(74) 代理人	100115071
			弁理士 大塚 康弘
		(74) 代理人	100116894
			弁理士 木村 秀二
		(74) 代理人	100130409
			弁理士 下山 治
		(74) 代理人	100134175
			弁理士 永川 行光

最終頁に続く

(54) 【発明の名称】 演算処理装置および方法

(57) 【特許請求の範囲】

【請求項1】

入力データに対して演算を行い演算結果データを生成する、複数の論理的な処理ノードを接続したネットワークでネットワーク演算を実行する演算処理装置であって、

前記ネットワークを構成する複数の処理ノードの各々が演算結果データを保持するためのバッファ用の記憶領域をメモリに割り当てるために、1つの処理ノードに対して演算結果データの一部を保持するバッファを割り当てるバンドバッファ方式と、1つの処理ノードに対して前記入力データに対する演算結果データの全てを保持するバッファを割り当てるページバッファ方式とを含む複数種類のバッファ割り当て方法のそれぞれについて、前記ネットワーク演算に必要な前記記憶領域のメモリ量を、当該ネットワークの構成に基づいて算出する算出手段と、

前記バッファとして割り当て可能なメモリ量を取得する取得手段と、

前記複数種類のバッファ割り当て方法のうち、前記算出手段で算出されたメモリ量が前記取得手段で取得されたメモリ量以下になるバッファ割り当て方法を選択する選択手段と

前記選択手段で選択されたバッファ割り当て方法に応じた実行順で、前記ネットワーク演算における各処理ノードによる演算を実行させる実行手段とを備えることを特徴とする演算処理装置。

【請求項2】

前記ネットワークは、複数の論理的な処理ノードを階層的に接続した階層型ネットワー

クであることを特徴とする請求項 1 に記載の演算処理装置。

【請求項 3】

前記選択手段は、前記算出手段で算出されたメモリ量が前記取得手段で取得されたメモリ量以下になるバッファ割り当て方法が複数存在する場合には、予め定めた優先順位に基づいてそれらのバッファ割り当て方法の中から 1 つを選択することを特徴とする請求項 1 に記載の演算処理装置。

【請求項 4】

前記算出手段は、前記ネットワーク演算を実行する全ての処理ノードの各々のバッファのサイズを当該処理ノードの後段に接続される処理ノードが必要とするデータ量に設定し、それらバッファのサイズを合計することにより、前記バンドバッファ方式を用いた場合の前記ネットワーク演算に必要なメモリ量を算出することを特徴とする請求項 1 に記載の演算処理装置。

10

【請求項 5】

前記算出手段は、前記ネットワークの構成に存在する連続する 2 つの階層のそれぞれの組について、2 つの階層に属する全ての処理ノードが前記入力データに対する演算結果データの全てを保持した場合に必要なバッファの合計サイズを計算し、それら合計サイズの内の最大のサイズを、前記ページバッファ方式を用いた場合の前記ネットワーク演算に必要なメモリ量とすることを特徴とする請求項 1 に記載の演算処理装置。

【請求項 6】

前記実行手段は、前記バンドバッファ方式が選択された場合には、前記複数の処理ノードの各々を予め定められた処理単位で実行させ、

20

前記複数の処理ノードの各々に割り当てられたバッファを、前記処理単位の演算結果の量に対応したメモリ領域を単位として循環させながら演算結果を書き込むリングバッファとして制御することを特徴とする請求項 1 乃至 5 のいずれか 1 項に記載の演算処理装置。

【請求項 7】

前記実行手段は、前記ページバッファ方式が選択された場合には、1 つの階層に属する全ての処理ノードの出力が生成された後に次階層に属する処理ノードの演算処理を開始するように制御し、

N 番目の階層に属する全ての処理ノードの出力が生成された後、N - 1 番目の階層に属する全ての処理ノードが使用していたバッファ領域を開放し、N + 1 番目以降の階層に属する処理ノードのバッファ領域に割り当てておくことを特徴とする請求項 1 乃至 6 のいずれか 1 項に記載の演算処理装置。

30

【請求項 8】

入力データに対して演算を行い演算結果データを生成する、複数の論理的な処理ノードを接続したネットワークでネットワーク演算を実行する演算処理装置であって、

前記ネットワークを構成する複数の処理ノードの各々が演算結果データを保持するためのバッファ用の記憶領域をメモリに割り当てる複数種類のバッファ割り当て方法のそれぞれについて、前記ネットワーク演算に必要な前記記憶領域のメモリ量を、当該ネットワークの構成に基づいて算出する算出手段と、

前記バッファとして割り当て可能なメモリ量を取得する取得手段と、

40

前記複数種類のバッファ割り当て方法のうち、前記算出手段で算出されたメモリ量が前記取得手段で取得されたメモリ量以下になるバッファ割り当て方法を選択する選択手段と

、前記選択手段で選択されたバッファ割り当て方法に応じた実行順で、前記ネットワーク演算における各処理ノードによる演算を実行させる実行手段と、

前記算出手段で算出されたメモリ量が前記取得手段で取得されたメモリ量以下になるバッファ割り当て方法が存在しない場合に、前記入力データを分割する分割手段とを備え、

前記選択手段は、分割された入力データの各々について前記ネットワーク演算に用いるべきバッファ割り当て方法を選択することを特徴とする演算処理装置。

【請求項 9】

50

前記処理ノードが行う演算はコンボリユーション演算であり、前記ネットワークは、コンボリユーションアル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乃至8のいずれか1項に記載の演算処理装置。

【請求項10】

入力データに対して演算を行い演算結果データを生成する、複数の論理的な処理ノードを接続したネットワークでネットワーク演算を実行する演算処理方法であって、

前記ネットワークを構成する複数の処理ノードの各々が演算結果データを保持するためのバッファ用の記憶領域をメモリに割り当てるために、1つの処理ノードに対して演算結果データの一部を保持するバッファを割り当てるバンドバッファ方式と、1つの処理ノードに対して前記入力データに対する演算結果データの全てを保持するバッファを割り当てるページバッファ方式とを含む複数種類のバッファ割り当て方法のそれぞれについて、前記ネットワーク演算に必要な前記記憶領域のメモリ量を、当該ネットワークの構成に基づいて算出する算出工程と、

10

前記バッファとして割り当て可能なメモリ量を取得する取得工程と、

前記複数種類のバッファ割り当て方法のうち、前記算出工程で算出されたメモリ量が前記取得工程で取得されたメモリ量以下になるバッファ割り当て方法のうちの1つを選択する選択工程と、

前記選択工程で選択されたバッファ割り当て方法に応じた実行順で、前記ネットワーク演算における各処理ノードによる演算を実行させる実行工程とを備えることを特徴とする演算処理方法。

20

【請求項11】

入力データに対して演算を行い演算結果データを生成する、複数の論理的な処理ノードを接続したネットワークでネットワーク演算を実行する演算処理方法であって、

前記ネットワークを構成する複数の処理ノードの各々が演算結果データを保持するためのバッファ用の記憶領域をメモリに割り当てる複数種類のバッファ割り当て方法のそれぞれについて、前記ネットワーク演算に必要な前記記憶領域のメモリ量を、当該ネットワークの構成に基づいて算出する算出工程と、

前記バッファとして割り当て可能なメモリ量を取得する取得工程と、

前記複数種類のバッファ割り当て方法のうち、前記算出工程で算出されたメモリ量が前記取得工程で取得されたメモリ量以下になるバッファ割り当て方法を選択する選択工程と

30

、前記選択工程で選択されたバッファ割り当て方法に応じた実行順で、前記ネットワーク演算における各処理ノードによる演算を実行させる実行工程と、

前記算出工程で算出されたメモリ量が前記取得工程で取得されたメモリ量以下になるバッファ割り当て方法が存在しない場合に、前記入力データを分割する分割工程とを備え、

前記選択工程では、分割された入力データの各々について前記ネットワーク演算に用いるべきバッファ割り当て方法を選択することを特徴とする演算処理方法。

【請求項12】

請求項10または11に記載の演算処理方法をコンピュータに実行させるためのコンピュータプログラム。

40

【請求項13】

請求項10または11に記載の演算処理方法をコンピュータに実行させるためのコンピュータプログラムを格納したコンピュータ読み取り可能な記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、例えばパターン識別装置、パターン識別システムや階層的なフィルタ演算処理装置等に適用可能なネットワーク演算を実行する演算処理装置および方法に関するものである。

【背景技術】

50

【 0 0 0 2 】

パターン識別システムや予測システム、制御システム等への応用として、ニューラルネットワークを利用した信号処理装置が広く利用されている。一般に、ニューラルネットワークはマイクロプロセッサ上で動作するソフトウェアとして実現される事が多く、パーソナルコンピュータやワークステーション等のアプリケーションソフトウェアとして提供されている。

【 0 0 0 3 】

図 2 は、一般的な階層結合型ニューラルネットワークを利用する画像処理装置の概念的な構成例を示すブロック図である。21 は検出対象のデータであり、例えば、ラスタスキャンされた画像データを示す。22 は画像データ 21 中から所定の物体を検出する演算ユニットであり、図示の例では 3 階層のニューラルネットワークで構成されている。23 は演算結果に相当する出力データ面である。演算ユニット 22 は、画像データ 21 中の所定の画像領域 24 を走査参照しながら処理を行う事で、画像中に存在する検出対象を検出する。出力データ面 23 は、検出対象の画像データ 21 と同じサイズのデータ面である。出力データ面 23 には、画像データ 21 の全ての領域を走査しながら処理した演算ユニット 22 の検出出力が走査順に格納される。演算ユニット 22 は、対象物が検出された位置で大きな値を出力する事から、当該出力データ面 23 を走査する事で対象物の画像面内の位置を把握する事ができる。

【 0 0 0 4 】

演算ユニット 22 において、25、26、27 は夫々ニューラルネットワークの階層を示し、各階層に所定の数のニューロン 28 が存在する。第 1 階層 25 は参照画像の画素数と同じ数のニューロン(ノード) 28 を有する。各ニューロンは所定の重み係数でフィードフォワード結合する。

【 0 0 0 5 】

図 3 は 1 つのニューロン 28 の構成例を示すブロック図である。in₁~in_n は入力値であり、第 2 階層以降では前階層のニューロンの出力値である。累積加算器 32 は、当該入力値と学習によって得られた係数 w₁~w_n を乗じた結果を累積加算する。非線形変換処理部 33 は、累積加算器 32 の累積加算結果をロジスティック関数や双曲正接関数(tanh 関数)等により非線形変換し、その変換結果を検出結果 out として出力する。なお、階層型ニューラルネットワークにおいて、夫々のニューロンに必要な重み係数 w₁~w_n は一般的に知られているバックプロパゲーション等の学習アルゴリズムを使用して、検出する対象物毎に決定されているものである。

【 0 0 0 6 】

このような階層結合型ニューラルネットワークを組み込み機器等へ高性能かつ安価に実装する事を目的として、アナログハードウェアやデジタルハードウェアで実現する手法が提案されている。

【 0 0 0 7 】

特許文献 1 では、単階層のアナログニューラルネットワークハードウェアを時分割多重利用する事で多階層化を実現する階層構造ニューラルネットのアーキテクチャが開示されている。又、特許文献 2 ではデジタルハードウェアにより実現する方法が開示されている。

【 0 0 0 8 】

一方、ニューラルネットワークの中でも、Convolutional Neural Networks と呼ばれるニューラルネットワークを用いた演算手法は識別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知られている。以下、Convolutional Neural Networks は CNN と略記する。例えば、特許文献 3 及び特許文献 4 では、CNN 演算を画像中の対象物識別や検出に適用した例が提案されている。

【 0 0 0 9 】

図 4 は簡単な CNN の例を示す論理的なネットワーク構成図である。ここでは、第 1 階層 406 の特徴数が 3、第 2 階層 410 の特徴数が 2、第 3 階層 411 の特徴数が 1 の 3

10

20

30

40

50

層CNNの例が示されている。401は画像データであり、ラスタスキャンされた画像データに相当する。403a~403cは第1階層406の特徴面を示す。特徴面とは、所定の特徴抽出フィルタ(コンボリューション演算の累積和及び非線形処理)で前階層のデータを走査しながら演算した結果を示す画像データ面である。特徴面はラスタスキャンされた画像データに対する検出結果であるため面で表す。特徴面403a~403cは401から夫々対応する特徴抽出フィルタにより生成される。例えば、特徴面403a~403cは、夫々模式的にコンボリューションカーネル404a~404cに対応する2次元のコンボリューションフィルタ演算とその演算結果の非線形変換により生成される。なお、402はコンボリューション演算に必要な参照画像領域を示す。

【0010】

例えば、カーネルサイズ(水平方向の長さ)と垂直方向の高さが11×11のコンボリューションフィルタ演算は以下に示すような積和演算により処理する。

【0011】

【数1】

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

... (1)

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):座標(x+column,y+row)での重み係数、

columnSize=11,rowSize=11:フィルタカーネルサイズ(フィルタタップ数)である

。

【0012】

404a~404cは夫々異なる係数のコンボリューションフィルタカーネルである。また、特徴面によってコンボリューションカーネル404a~404cのサイズも異なる

。

【0013】

CNN演算では、複数のフィルタ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生成する。特徴面403aを算出する場合は、前階層との結合数が1であるため、1つのコンボリューションカーネル404aが用いられる。一方、特徴面407a及び407bを計算する場合、前階層との結合数が3であるため、夫々409a~409c及び409d~409fに相当する3つのコンボリューションフィルタの演算結果を累積加算する。つまり、特徴面407aは、コンボリューションカーネル409a~409cの出力を累積加算し、最後に非線形変換処理する事によって得られる。

【0014】

ここで、409a~409fは何れも異なるフィルタ係数のコンボリューションカーネルである。また、コンボリューションカーネル409a~409cと409d~409fは、図4に示されるように、それぞれ異なるカーネルサイズを有する。コンボリューションフィルタの累積加算及び非線形変換処理の基本的構成は図3に示すニューロンの構成と同様である。即ち、コンボリューションカーネルの係数が重み係数w₁~w_nに相当する。特徴面407a、407b、408の様に複数の前階層の特徴面と結合される場合、複数のコンボリューションカーネルの演算結果が累積加算器32で蓄積される事になる。即ち、総結合数はコンボリューションカーネルサイズ×前階層の特徴数に相当する。

【0015】

図5はCNN演算における図形検出処理の一例を説明する図である。51a~51cは第1階層406の特徴抽出対象を模式的に示す図であり、それぞれ水平方向のエッジ及び斜め方向のエッジを抽出する様に学習されたコンボリューションカーネルである。52a、52bは複数の第1階層における特徴抽出結果とその空間的な配置関係から、第2階層

10

20

30

40

50

4 1 0で抽出される図形である。5 3は最終的に抽出される図形を示している。5 3は複数の第2階層4 1 0における特徴抽出結果とその空間配置関係から、第3階層4 1 1で抽出される図形である。コンボリューションカーネルの各フィルタ係数は特徴毎に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。物体の検出や認識等においては、10×10以上の大きなサイズのフィルタカーネルを使用する事が多い。また、一般的に特徴毎にコンボリューションカーネルのサイズは異なる。

【0016】

このように、CNN演算では、特徴抽出毎に画像面単位で結果を保持しながら階層的に結合する事で、プリミティブな特徴とその空間的な配置関係に基づく頑健なパターン検出を実現する。

10

【特許文献1】特許第2679730号明細書

【特許文献2】特開平3-55658号公報

【特許文献3】特開平10-021406号公報

【特許文献4】特開2002-358500号公報

【発明の開示】

【発明が解決しようとする課題】

【0017】

図2で説明した様に、一般的な階層型ニューラルネットワークを利用した画像中の物体検出装置において、演算処理に必要なメモリサイズは、入出力画像バッファを除くと、各ニューロン出力を保持するためのバッファメモリがあれば十分である。即ち、ニューロン数と等価な数の所定ビット数のメモリがあれば所望の演算処理を実行できる。

20

【0018】

一方、CNN演算の場合、前階層の複数の特徴抽出結果の空間的配置に基づいて特徴抽出を行うため、各階層間で所定サイズのデータバッファが必要になる。例えば、図4に示すCNN演算構成例の場合、入出力画像バッファを除くと画像サイズ×5個の特徴面バッファメモリを用意している(特徴面403a~403c、特徴面407a~407b)。このため、一般的な階層型ニューラルネットに比べ処理に必要なメモリサイズが増大する。

【0019】

特許文献3及び特許文献4に開示されている手法も、特徴抽出結果を画像面で保持する手法であり、処理に必要なメモリサイズが一般的な階層型ニューラルネットワークによる方式に比べて大きい。

30

【0020】

このため、特に、ハードウェアにより実現する場合、LSIの内部にサイズの大きいRAM(Random Access Memory)を用意する必要があり、回路規模が増大する。ソフトウェアにより実現する場合であっても、組み込み機器に実装する場合、システムに必要なメモリ量が増大する事で同様にコストが上昇する。すなわち、演算に使用可能なメモリ量は、システムにかけることのできるコストによって定まる有限な値となる。

【0021】

一方、メモリの増大を避ける手法として、入力するデータを領域分割して投入する方法が利用されている。しかしながら、参照領域が広い演算を階層的に処理する場合、分割投入するデータを広い範囲でオーバーラップさせる必要があるため、結果的に処理対象領域が増加してしまい、処理効率及び処理速度が低下する。

40

【0022】

本発明はこのような問題点を解決するためになされたものであり、CNN演算等のネットワーク構造で接続された複数の処理ノードによって行われる演算処理を、限られたメモリ量で効率良く実現する事を目的とする。

【課題を解決するための手段】

【0023】

上記課題を解決するために、本発明の一態様による演算処理装置は以下の構成を備える

50

。すなわち、

入力データに対して演算を行い演算結果データを生成する、複数の論理的な処理ノードを接続したネットワークでネットワーク演算を実行する演算処理装置であって、

前記ネットワークを構成する複数の処理ノードの各々が演算結果データを保持するためのバッファ用の記憶領域をメモリに割り当てるために、1つの処理ノードに対して演算結果データの一部を保持するバッファを割り当てるバンドバッファ方式と、1つの処理ノードに対して前記入力データに対する演算結果データの全てを保持するバッファを割り当てるページバッファ方式とを含む複数種類のバッファ割り当て方法のそれぞれについて、前記ネットワーク演算に必要となる前記記憶領域のメモリ量を、当該ネットワークの構成に基づいて算出する算出手段と、

前記バッファとして割り当て可能なメモリ量を取得する取得手段と、

前記複数種類のバッファ割り当て方法のうち、前記算出手段で算出されたメモリ量が前記取得手段で取得されたメモリ量以下になるバッファ割り当て方法を選択する選択手段と

前記選択手段で選択されたバッファ割り当て方法に応じた実行順で、前記ネットワーク演算における各処理ノードによる演算を実行させる実行手段とを備える。

【0024】

また上記課題を解決するための、本発明の他の態様による演算処理方法は、

入力データに対して演算を行い演算結果データを生成する、複数の論理的な処理ノードを接続したネットワークでネットワーク演算を実行する演算処理方法であって、

前記ネットワークを構成する複数の処理ノードの各々が演算結果データを保持するためのバッファ用の記憶領域をメモリに割り当てるために、1つの処理ノードに対して演算結果データの一部を保持するバッファを割り当てるバンドバッファ方式と、1つの処理ノードに対して前記入力データに対する演算結果データの全てを保持するバッファを割り当てるページバッファ方式とを含む複数種類のバッファ割り当て方法のそれぞれについて、前記ネットワーク演算に必要となる前記記憶領域のメモリ量を、当該ネットワークの構成に基づいて算出する算出工程と、

前記バッファとして割り当て可能なメモリ量を取得する取得工程と、

前記複数種類のバッファ割り当て方法のうち、前記算出工程で算出されたメモリ量が前記取得工程で取得されたメモリ量以下になるバッファ割り当て方法のうちの1つを選択する選択工程と、

前記選択工程で選択されたバッファ割り当て方法に応じた実行順で、前記ネットワーク演算における各処理ノードによる演算を実行させる実行工程とを備える。

【発明の効果】

【0025】

本発明によれば、CNN演算等のネットワーク構造で接続された複数の処理ノードによって行われる演算処理を、限られたメモリ量で効率良く実現することができる。中間演算結果を介して演算部を接続するネットワーク型演算器による演算処理を、限られたメモリ量において最適に処理することが可能となる。すなわち、同じ構成のハードウェアで、より多様なネットワーク型演算処理を実行することが可能となる。

【発明を実施するための最良の形態】

【0026】

以下、本発明の好適な実施形態について、添付の図面を用いて説明する。

【0027】

<第1実施形態>

図6は第1実施形態に関する階層的演算処理回路を具備したパターン検出装置の構成例である。当該装置は画像データ中の特定の物体(画像パターン)を検出する機能を有する画像処理装置である。図6において61は画像入力部であり、光学系、CCD(Charge-Coupled Devices)又はCMOS(Complimentary Metal Oxide Semiconductor)センサ等の光電変換デバイスを有する。また、画像入力部61は、光電変換間デバイスを制御するド

10

20

30

40

50

ライバ回路 / A Dコンバータ / 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具備する。62は前処理部であり、検出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、前処理部62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をハードウェアで処理する。63はCNN処理部であり、本実施形態による階層的演算処理回路を含み、特徴検出処理部として機能する。なお、CNN処理部63の詳細は図1を用いて後述する。66はDMAC (DirectMemory Access Controller) であり、画像バス64上の各処理部間のデータ転送、及び、画像バス64上のデバイスとCPUバス67上のRAM70間のデータ転送を司る。65はブリッジであり、画像バス64とCPUバス67のブリッジ機能を提供する。68はCPUであり、本装置全体の動作を制御するものである。69はROM (ReadOnly Memory) であり、CPU68の動作を規定する命令や各種演算に必要なパラメータデータを格納する。例えば、CNN処理部63の動作に必要な重み係数、ネットワーク結合情報、シーケンス情報等もROM69に格納されている。70はCPU68の動作に必要なメモリ (RAM : RandomAccess Memory) である。RAM70はDRAM (Dynamic RAM) 等の比較的容量の大きいメモリで構成される。CPU68はブリッジ65を介して画像バス64上の各種処理部にアクセスする事が可能である。画像バス64とCPUバス67を分離する事により、ハードウェアによる画像入力部61、前処理部62、CNN処理部63の各処理部の動作とCPU68の動作を同時に並列実行させることができる。

【0028】

図6の階層的演算処理回路を具備したパターン検出装置は、例えば図14に示されるようなConvolutionalNeural Networks (以下CNNと略記する) のような階層的な演算を行うのに用いられる。図14において処理ノードとは、コンボリューション演算の対象画像とコンボリューションカーネルからコンボリューション演算結果を得る処理を行うブロックを指す。なお、図14では便宜上第0処理ノードを設けたが、通常第0処理ノードでは特になにも処理は行われず、入力画像が第1~第3処理ノードへ入力される。例えば、図14の第4処理ノードでは、第1~第3処理ノードの出力に対し、それぞれ係数の異なるコンボリューションカーネルを適用してコンボリューション演算を行う。そして、それぞれのコンボリューション演算の結果を加算し、その加算結果に非線形変換を行って第4処理ノードの演算結果を得ている。

【0029】

CNN処理部63に図14に示されるCNNを適用する場合、演算処理部を処理ノード間で時分割に使用することで、各処理ノードで規定された演算を実行する。例えば、まず第1処理ノードで規定された演算を行い、その後第2処理ノードで規定された演算を行う、というようにCNNの演算が実行されていく。つまり、CNNを構成する処理ノードは複数存在し、論理的なネットワークを構成するが、処理ノードで規定された演算を実行する演算処理部は物理的に1つしか存在しない。

【0030】

図1はCNN処理部63の詳細を説明する図である。図1において、101は演算部であり、所定のデータ群に対してコンボリューション演算と非線形処理を実行する。図12に演算部101の一例を示す。図12において、1201は乗算器であり、係数選択部1204がカーネル選択信号に従って重み係数記憶部1205より選択し、出力する重み係数と、カーネル選択信号と同期して入力される入力データとを乗じる。1202は累積加算器であり、乗算器1201の出力を所定の期間累積加算する。1203は非線形変換処理であり、ロジスティック関数やtanh関数を用いて累積加算結果を非線形変換する。非線形変換は、例えば、各入力値に対して所定の関数値を列挙する関数テーブルで実現される。1205は重み係数記憶部であり、検出対象と処理ノードに応じた複数の重み係数データが格納されている。重み係数記憶部1205は、例えばRAM等により構成される。1204は係数選択部であり、ネットワーク構成管理部108が指示するカーネル選択信号に従って、対応する重み係数を順次記憶部から読み出す。

【0031】

10

20

30

40

50

図1に戻り、102はワークメモリ(以下、単にメモリという)であり、入力画像/中間層の特徴抽出結果/最終検出結果等を格納する。本実施形態では、コンボリユーション演算を高速に実行するため、メモリ102として、高速にランダムアクセスが可能なSRAM(Static RAM)を使用している。

【0032】

103はメモリアクセス制御部であり、メモリ102に対するアクセス、アドレスの生成、リード/ライト信号制御及びデータバスの方向制御等を司る。メモリアクセス制御部103はリングバッファ設定部104-1~104-nの出力に従ってメモリ102にアクセスする。

【0033】

104-1~104-nは複数のリングバッファ設定部であり、それぞれ各処理ノードがメモリ102をリングバッファとして利用するために用いられる。以下、リングバッファ設定部104-1~104-nの任意の1つを指す場合は、リングバッファ設定部104と記載する。リングバッファ設定部104は、CNN演算の論理的な各処理ノード毎に一つずつ用意される。各リングバッファ設定部104は、リングバッファのサイズを指定するリングサイズ設定部106、リングバッファの動作状況を保持するリングカウンタ105及びメモリ102上の物理アドレスを決定するためのオフセットアドレス設定部107を具備する。リングバッファ設定部104の出力はセレクタ1121、1122で選択されて、メモリアクセス制御部103へ提供される。この構成により、メモリ102には、ネットワークを構成する複数の処理ノードの各々に対応して、演算結果データを保持するための中間バッファ用の記憶領域が割り当てられることになる。その詳細は、後述する。

【0034】

108はネットワーク構成管理部であり、1つの演算部101を利用して論理的な階層ネットワーク処理を実現するために、各動作を制御する。論理的な階層結合関係を指定する構成情報は、ネットワーク構成情報設定部110内にテーブルデータ(以下、構成情報テーブルという)として保持される。構成情報テーブルはレジスタやRAMで構成される。ネットワーク構成管理部108は、構成情報テーブルに従ってメモリアクセス制御部103や演算部101の動作を順次制御することにより、後述する所定の単位演算をベースとした階層ネットワークの演算処理を実現する。

【0035】

シーケンス制御部109は、シーケンス情報設定部111に記されたシーケンス情報に従って、各処理ノードによる単位演算の実行順を制御する。本実施形態では、1ライン単位の出力を得る演算処理を所定の演算処理単位(単位演算)としている。すなわちライン単位で論理的な処理ノードを切り替えながら時分割で処理を実行することで階層型ネットワーク演算が遂行される。シーケンス情報設定部111はシーケンス情報を保持するRAM等により構成される。

【0036】

113はCPUバスアクセス制御部であり、CPU68がCNN処理部63内の各種レジスタやメモリにアクセスするためのバスインターフェースである。例えば、

- ・リングサイズ設定部106のリングバッファサイズ、
- ・ネットワーク構成情報設定部110の構成情報テーブル、
- ・シーケンス情報設定部111のシーケンス情報、
- ・演算部101の重み係数記憶部1205内の重み係数データ、

等の各種設定データは、当該インターフェースを介してCPU68から書き込むことができる。

【0037】

ここで、図2を用いて本実施形態における所定の単位演算について説明する。先に述べたとおり、本実施形態での所定の単位演算とは、演算部101を用いて行われるライン単位のコンボリユーション演算処理である。ただし、図22では、簡単のため、一つの処理ノードの演算出力画像(または、ネットワークへの入力画像)を演算対象画像としてコ

10

20

30

40

50

ンボリューション演算を行う場合が示されており、非線形変換も省略されている。

【 0 0 3 8 】

図 2 2 の (a) において、2 2 0 1 は演算対象画像 (参照画像) を表している。演算対象画像 2 2 0 1 において、模式的に示す最小一升が、ラストスキャン順で示された入力画像又は前階層の処理ノードでの演算結果画像である演算対象画像の画素 (input(x, y)、x : 水平方向位置、y : 垂直方向位置) を示す。

2 2 0 2 は演算結果画像を表し、模式的に示す最小一升が、ラストスキャン順の演算結果画素 (output(x,y)、x : 水平方向位置、y : 垂直方向位置) を示すものとする。

【 0 0 3 9 】

演算対象画像 2 2 0 1 内の太線で囲まれた領域 2 2 0 3 は、output(6, 7)位置のコンボリューション演算を処理する場合の参照画像の領域を示す。領域 2 2 0 3 では、コンボリューションカーネルのサイズが水平方向「 1 1」、垂直方向「 1 3」の場合が示されている。

10

【 0 0 4 0 】

演算結果画像 2 2 0 2 の太線で囲まれた領域 2 2 0 4 は、演算対象画像 2 2 0 1 に対して単位演算 (水平方向 1 行分の演算) を行った場合の結果領域を示す。ここで、領域 2 2 0 4 内の格子状の網掛け領域 2 2 0 6 は、コンボリューションカーネルのサイズに依存して発生する周辺領域 (演算が行われない領域) の画素である。つまり、output(5, 7)の位置の演算を行うべき参照画像領域は、領域 2 2 0 3 を左に 1 画素分ずらしたものとなる。しかしながら、そのような領域は演算対象画像 2 2 0 1 (参照領域) からはみ出してしまうため、一部の参照画素が存在しないことになる。なお、階層的処理においてこの周辺領域 (無効領域) をどう扱うか (削除するか、デフォルト値を埋め込むか等) は、本発明において本質的でないので、ここでは例えば、デフォルト値を埋め込むとする。尚、2 2 0 4 より上のラインについても同様に無効領域となる。

20

【 0 0 4 1 】

図 2 2 から明らかのように、1 ラインの単位演算を行うには、演算対象画像 2 2 0 1 の必要領域として、少なくとも領域 2 2 0 5 が必要となる。領域 2 2 0 5 は、図 2 2 において網掛け領域として示されており、水平方向サイズは演算対象画像 2 2 0 1 と同じサイズ、垂直方向サイズはコンボリューションカーネルの垂直方向サイズが必要となる。説明の都合上、この領域を単位演算対象画像領域 2 2 0 5 と呼ぶ。領域 2 2 0 4 で示されるような単位演算を、単位演算対象画像領域 2 2 0 5 をずらしながら行うことで、演算対象画像 2 2 0 1 の全領域にわたってコンボリューション演算を行うことができる。例えば、図 2 2 の (b) には、1 画素下にずらした単位演算対象画像領域に対して単位演算を行った場合を示している。この時、ある単位演算を実行できるか否かは、その単位演算の単位演算対象画像領域 2 2 0 5 ' の画素データが、前階層の処理ノードによって演算され、その結果が出力されているか否かに依存する。もちろん、複数の参照画像を入力として演算に用いる処理ノードの場合は、全ての参照画像についての単位演算対象画像領域の画素データが出力されている必要がある。

30

【 0 0 4 2 】

図 7 は、図 4 で説明した CNN ネットワークに対し、本実施形態を適用した場合の動作の一例を説明する図である。

40

【 0 0 4 3 】

図 7 において、7 0 1 は入力層であり、所定サイズの検出対象画像データが入力される。7 0 3 a ~ 7 0 3 c は第 1 階層 7 0 6 の演算出力を格納するメモリ領域を示す。すなわち、入力層 7 0 1 の入力面に対するコンボリューション演算 7 0 4 a ~ 7 0 4 c 及び非線形変換の結果が、メモリ領域 7 0 3 a ~ 7 0 3 c にそれぞれ格納される。第 1 階層の演算結果である特徴面は論理的には入力層 7 0 1 と同じサイズのデータ面となる。しかしながら、ここでは所定高さのリングバッファとして機能するメモリ領域 7 0 3 a、7 0 3 b に特徴面が格納される。このリングバッファは、幅が入力画像と同じであり、ライン単位で循環するバンドバッファである。

50

【 0 0 4 4 】

図 1 3 は、本実施形態によるリングバッファの動作を模式的に説明する図である。ここでは説明のためリングバッファの高さ（循環数）を 6 とする。また、ここでは、入力画像 1 3 0 0 の 8 ライン分の画像データが、L 1 ~ L 8 としてラスタスキャン順に入力された場合に、6 ライン分のリングバッファにどのように保持され、参照されるかを説明する。

【 0 0 4 5 】

ここでリングバッファに付随するリングカウンタは 0 ~ 5 の値を循環する。また、リングカウンタの初期値は 5 であり、1 ライン分のデータが投入されるときに 1 インクリメントされるものとする。ただし、リングバッファの循環数と同じ値になると、リングカウンタのカウント値は 0 に戻る。例えば、本リングバッファでは循環数は 6 であるので、カウンタ値は 5 の次は 0 に戻ることになる。

10

【 0 0 4 6 】

1 3 0 1 はリングバッファに入力画像 1 3 0 0 の先頭から 6 ライン分のデータ（L 1 ~ L 6）がフルに充填された状態を表し、リングカウンタの値は「5」となっている。次のラインを格納するとき、リングカウンタはインクリメントされ 0 に戻り、リングバッファの先頭行に L 7 が充填される。すなわちリングカウンタの値は、最新のラインを格納したリングバッファ中の行を示す（0 基準）。この状態を、図 1 3 の 1 3 0 2 に示す。

【 0 0 4 7 】

1 3 0 2 の状態では、リングバッファから L 2 ~ L 7 のバンドを参照する事が可能となり、その開始行はリングカウンタの値 + 1 の行である。更に次のライン L 8 を格納する場合は、1 3 0 3 に示すように 2 行目位置に L 8 が充填され、リングカウンタの値は 1 となる。この場合、L 3 ~ L 8 を参照する事が可能となり、先頭行はやはりリングカウンタの値 + 1 の行となっていることが分かる。

20

【 0 0 4 8 】

尚、リングバッファの循環数を入力画像データのライン数に一致させると、そのバッファにおいては、1 ページ分の処理中において前のラインが上書きされなくなる。すなわち、リングバッファは、バンドバッファとしてだけでなく、ページバッファとしても機能させることができる。

【 0 0 4 9 】

ここで、各処理ノードのリングバッファに最低限必要な高さ（循環数）は、次階層の処理ノードが単位演算処理を行う際の、単位演算対象画像領域（図 2 2 の 2 2 0 5）の高さに一致する。すなわち、各リングバッファの高さ（循環数）は、その処理ノードの後段に接続される全処理ノードのコンポリューションカーネルのサイズに基づいて決定することができる。

30

【 0 0 5 0 】

例えば図 7 において、メモリ領域 7 0 3 a の場合、コンポリューションカーネル 7 0 9 a と 7 0 9 d の単位演算対象画像領域の高さのうち大きい方の値をメモリ領域 7 0 3 a のリングバッファの最低限必要な高さとする。このように定めると、メモリ領域 7 0 3 a によって形成されたリングバッファに納められた画素データを用いて、コンポリューションカーネル 7 0 9 a と 7 0 9 d のどちらのコンポリューション演算も可能となる。同様にメモリ領域 7 0 3 b、7 0 3 c に最低限必要な高さは、夫々コンポリューションカーネル 7 0 9 b / 7 0 9 e、コンポリューションカーネル 7 0 9 c / 7 0 9 f のカーネルサイズから決定できる。同様に、メモリ領域 7 0 7 a、7 0 7 b は、夫々コンポリューションカーネル 7 1 2 a、7 1 2 b の単位演算対象画像領域の高さから決定できる。

40

【 0 0 5 1 】

なお、ここで決定しているのは、あくまで次階層が演算する上で最低限必要な高さであるので、他の必要性があればさらに大きい高さのリングバッファを用いてももちろんよい。例えば、図 7 では、メモリ領域 7 0 3 a、7 0 3 b、7 0 7 a、7 0 7 b はカーネルサイズから規定される次階層演算のために最低限必要な高さをリングバッファ高さとしている。しかし、メモリ領域 7 0 3 c は、次階層演算に使われるのみならず、CPU 6 8 が判

50

定処理に使用する特徴検出データでもあるため、入力画像データと同じサイズのページバッファを割り当てている。すなわちCPU68は、最終階層の特徴データ713だけでなく、メモリ領域703cに格納された第1階層の特徴データも参照して検出対象画像の存在を判定することができる。

【0052】

この様に、本実施形態のCNN処理部を用いると、中間層の特徴面をネットワークの結合状態及び目的に応じて最適なサイズのバッファ（メモリ領域703a、703b、703c、707a、707b）にアサインする事が可能となる。

【0053】

ところで、図7に示した本実施形態のCNN演算処理では、中間処理ノードの演算結果を保持するメモリ領域703a、703b、707a、707bをバンドバッファとしている。このように、中間処理ノードの出力バッファ（中間バッファ）をバンドバッファとすることにより、図4で説明した従来のCNN演算処理よりも、メモリ使用容量を減らしている。しかしながら、ネットワーク構成と演算シーケンスによっては、必ずしもバンドバッファによる中間バッファがトータルのメモリ使用量を最小にするとは限らない。以下、そのようなケースについて説明する。

【0054】

図23は、N個の階層からなるCNNに対し、中間バッファとしてバンドバッファを割り当てた場合の例を示す図である。図23において、各円は論理的な処理ノードを表し、各円の右に付随する四角がそれぞれの処理ノードに割り当てられたバッファを示している。ここでは入力画像である第0階層と出力層である第N階層の処理ノードには、1ページ分のフレームバッファ（ページバッファ）を割当て、他の中間層のノードにはバンドバッファを割り当てている。なお、本明細書では、中間層の出力結果保持バッファを特に中間バッファと呼んでいる。また、一部の円の左側に付随する四角は、その処理ノードのコンボリューションカーネルを例示するものであるが、図が複雑になるので一部のみしか記述していない。

【0055】

第0階層のノードは入力層であり、上述した通り便宜上処理ノードを割り当てているが、実際に演算処理を行うわけではなく、単に入力画像データが入力バッファメモリに格納されている状態を示す。

【0056】

また、最終層（第N階層）に割り当てられているバッファは出力バッファであって、この例では2つの処理ノードに対し入力画像サイズと同等のサイズのバッファがそれぞれ割り当てられている。もちろんこれら出力バッファの大きさは、このCNNによる演算結果を使用する、例えばCPU68による後処理部の都合によって定めればよいもので、コンボリューション演算の都合にはよらない。例えば、最終出力結果からある範囲の重心を取ってその座標に検出対象が存在するとする後処理を行うならば、出力バッファの大きさは重心を取る範囲の演算結果のみを保持できる大きさで十分である。

【0057】

図23のCNNでは、中間バッファ（第1階層～第N-1階層の出力バッファ）を、図7で説明したのと同様のバンドバッファとして割り当てている。それぞれのバンドバッファの横幅は入力画像幅と同じであり、縦高さは先に説明した通り、接続される次階層の処理ノード（隣接上層処理ノード）のコンボリューションカーネルのサイズによって定まる最低限必要な高さとしている。尚、ここでは説明を簡単にするため、中間層の結果は後処理には用いないものとする。

【0058】

このように中間バッファを全てバンドバッファとして割り当てるとき、CNN全体の処理シーケンスは、ライン単位の演算処理を、処理ノードを切り替えて順に行っていくことになる。例えば第1階層の第1処理ノードのバンドバッファの高さが5だったとすると、5ライン分の演算処理が終わりその結果が格納されると、この処理ノードの演算結果を参

10

20

30

40

50

照する第2階層の処理ノードが、1ライン分の単位演算処理を行うことができる。逆に第2階層の処理ノードの処理が終わらない限り、バンドバッファ内の5ラインの内の1ラインも破棄することはできない。そのため、新たな演算結果の格納場所がないため、第1階層の第1処理ノードは次の単位演算を行うことができない。こうして、第1階層の第1処理ノードは、その演算結果を必要とする全ての第2階層の処理ノードが1ライン分の単位演算処理を完了したら、第1階層の第1処理ノードは次の1ラインの単位演算処理を行うことができるということになる。

【0059】

従って、CNN全体で処理が1ライン単位ずつ進んでいくことになるので、全ての中間バッファは、基本的に同時に存在している必要がある。ここで、

全階層数(入力層除く)： N

階層番号変数(中間層)： $l = 1, 2, \dots, N - 1$

各階層特徴数： F_l

特徴番号(注目階層の注目特徴)： $f = f_l = 1, \dots, F_l$

前階層の特徴番号： $f' = f_{l-1}$

入力画像サイズ(水平方向, 垂直方向)： I_x, I_y

カーネルサイズ(水平方向, 垂直方向)： $W_x(l, f, f'), W_y(l, f, f')$

バンドバッファ高さ： $B_y(l, f)$

バンド中間バッファ必要容量： S_B

とすると、全ての中間バッファを最低限必要なサイズのバンドバッファとして割り当てるときに必要となるトータルのサイズ S_B は、

【数2】

$$B_y(l, f) = \max (W_y(l+1, f_{l+1}, f_l) | 1 \leq l \leq N-1, 1 \leq f_{l+1} \leq F_{l+1})$$

$$S_B = I_x \times \sum_{l=1}^{N-1} \sum_{f=1}^{F_l} B_y(l, f)$$

... (2)

として求めることができる。

【0060】

上記計算では、

- ・ネットワーク演算を構成する全ての処理ノードの各々に必要な中間バッファのサイズを当該処理ノードの後段に接続される処理ノードが必要とするデータ量に設定し、
- ・それら中間バッファのサイズを合計することにより、ネットワーク演算に必要なメモリ量を算出している。

【0061】

一方、図24は、図23で示したCNNと同じネットワーク構成のCNNに対し、中間バッファとしてページバッファを用いた場合の例を示している。入力層および出力層のバッファは図23と同様、ネットワーク構成や演算シーケンスによらず定まるもので、ここではやはりページバッファとしてある。

【0062】

図24では、第1階層から第 $N - 1$ 階層までの夫々の処理ノードに入力画像サイズと同等のサイズのページバッファが割り当てられている。しかしながら、これらのページバッファは、全て同時に存在している必要はない。

【0063】

まず第1階層の各ノードに注目すると、それぞれのノードにおける演算処理は、第0階層すなわち入力画像のみを参照して行われる。本実施形態において入力画像データは全てメモリ上に格納されているので、各ノードは必要なときに必要な範囲のデータを参照可能である。また、各ノードの出力格納バッファとして、1ページ分のバッファが用意されているので、バンドバッファの場合と異なりラインの上書きを気にすることなくいつでも演算結果を格納可能である。

【0064】

10

20

30

40

50

次に第2階層の各処理ノードに注目すると、これらのノードは第1階層の処理結果を参照する。従って、第1階層の各処理ノードが、1ページ分の処理を終えた後であれば、必要な参照データは入力画像と同様ページで保持された状態となるので、所望のときに所望の範囲を参照することができる。以下の層も同様で、基本的に前段の階層の処理が完了しており参照データがページ分揃っていれば処理可能である。

【0065】

逆に、各階層の処理ノードは演算処理を行うに当たって、前々階層の処理結果は必要としない。つまり、ある処理ノードの演算処理を開始するに当たって、前階層の全処理ノードの演算が1ページ分完了しているならば、更にその前の階層の処理ノードの中間バッファは開放してしまっても構わないことになる。

10

【0066】

従って、各処理ノードでライン単位の演算処理を順次行うに当たって、まずは第1階層に属する処理ノードでのみ順次単位演算を行う。そして、第1階層の全ての処理ノードで1ページ分の演算処理が完了した後に、第2階層の処理ノードでの単位演算処理を開始する。これを順に次の階層にたいしても繰り返してゆく、というように、単位演算処理シーケンスを組むと、中間バッファは1時期に連続する2階層のみで存在していれば、最終層までの演算が可能となる。より一般化して言うと、

- ・ 1つの階層に属する全ての処理ノードの出力が生成された後に次階層に属する処理ノードの演算処理を開始するように制御し、
- ・ N番目の階層に属する全ての処理ノードの出力が生成された後、N - 1番目の階層に属する全ての処理ノードが使用していた中間バッファ領域を開放してN + 1番目以降の階層に属する処理ノードの中間バッファ領域に割り当て可能にする。

こうして、1階層の演算処理が完了する毎に、その前の階層のページバッファを、次階層の中間バッファとして再利用するようにすることにより、トータルの中間バッファ必要サイズを減らすことが可能となる。以下、この方式をページバッファ方式という。

20

【0067】

ここで、ページ中間バッファ必要容量： S_p
とすると、

【数3】

$$S_p = (I_x \times I_y) \times \max(F_l, F_{l+1}) \quad | \leq l \leq N-1$$

30

... (3)

として、ページで中間バッファを構成した際（ページバッファ方式を採用した際）のトータルの必要サイズを求めることができる。この計算は、ネットワーク構成において連続する2つの階層の組に属する全ての処理ノードが生成する演算結果データのサイズの合計を全ての組について計算し、その内の最大となるサイズを必要なメモリ量とするものである。

【0068】

(2)、(3)式から分かる通り、バンド中間バッファ必要容量 S_B とページ中間バッファ必要容量 S_p のどちらが小さくなるかは、CNNネットワークの各論理的処理ノードの接続構造と、各ノードのコンボリューションカーネルサイズに依存する。一般的に、ネットワークの階層数が少なくカーネルサイズが小さければ、バンドバッファ方式の割当てが有利となるが、階層数が多くかつ各階層に属する処理ノード数が比較的少なければページバッファ方式の割当てが有利となる。

40

【0069】

通常、システムとして、中間バッファに用いることのできるメモリの上限サイズは一定である。この条件サイズをMとする。特に本実施形態のような専用のパターン検出装置においては、中間バッファに割り当てるメモリサイズは、小さければ小さいほど良いという訳ではなく、割り当て可能なサイズM以下であれば全く問題ない。

【0070】

また、上述のようにバンドバッファ方式の場合と、ページバッファ方式の場合では、単

50

位演算の処理シーケンスが異なる。特にページバッファ方式では、一つの処理ノードでの単位演算を1ページ分連続して行うことも可能である。このようなシーケンスにすると、実装によっては、処理ノード切り替えのオーバーヘッドを省略でき、バンドバッファ割当て方式よりトータル演算時間を若干短くできる可能性がある。従って、 S_B と S_P のどちらもM以下になる場合は、より小さい方を選ぶのではなく、ページバッファ方式を優先する方がよい。

【0071】

ところで、図11で示すように、ページバッファ割当て方式の場合、最終出力層の出力バッファも他の中間バッファと領域を兼用することが可能である。そこで、(3)式の変わりに、以下の(4)式を用いる。

【0072】

【数4】

$$S_p' = (I_x \times I_y) \times \max(F_l, F_{l+1}) \quad | \quad l \leq N$$

... (4)

【0073】

また、中間バッファだけでなく最終出力層も含めて割当て可能なサイズをM' とすると、以下の(5)式のように表される。

【0074】

【数5】

$$M' = M + (I_x \times I_y) \times F_N$$

... (5)

【0075】

すなわち、ページバッファ方式での割当てが可能かどうかを判断するのに、(3)式の代わりに(4)式を用いて、M' と比較するようにしてもよい。つまりページバッファ方式の場合は、実質割当て可能なサイズを増やすことができる。このようにすると、さらにメモリ使用効率を上げられる。

【0076】

以下、図8および図26、27に示すフローチャートを用いて、本実施形態の階層型ネットワークによる演算処理動作を詳細に説明する。図8はCPU68での一連の検出処理動作を示すフローチャートである。CPU68におけるソフトウェアによる処理は、所定の設定処理等を行った後、画像入力部61やCNN処理部63等のハードウェア処理回路部をドライブする。

【0077】

まず、ステップS801において、CPU68は、検出処理の開始に先立ち、変数やレジスタ等の各種初期化処理を実行する。各処理部はCPU68の初期化指示に従って内部のレジスタ、メモリ等を初期化する。またここでは、入力画像サイズ全域を処理領域とする初期設定も行う。本処理を含め、以降CPU68は、ブリッジ65及び画像バス64を介してCNN処理部63や他のハードウェア回路部にアクセスし、CPUバスアクセス制御部113を介して所定のレジスタやメモリにデータを設定することが可能となる。また、本初期化処理においては、演算部101内の重み係数記憶部1205(図12)へ、本パターン検出装置で対応する全ての検出対象に対する全ての処理ノードのコンボリューションカーネルの係数データをロードする。本パターン検出装置は、コンボリューションカーネルの重み係数、シーケンス情報、ネットワーク構成情報等を入れ替える事で、同一のハードウェアで様々な検出対象に対応することができるが、そのためのカーネルをここで一通りロードしておく。

【0078】

次に、ステップS802において、CPU68は、ユーザからの検出対象の選択を受け付ける。ここでは、本パターン検出装置が対応する全ての検出対象の中から、画像中より

10

20

30

40

50

検出したい所望の検出対象がユーザにより選択される。

【 0 0 7 9 】

検出対象が決定したら、ステップ S 8 0 3 でネットワーク構成管理部 1 0 8 のネットワーク構成情報設定部 1 1 0 にネットワーク構成情報を設定する。ネットワーク構成情報はネットワークの結合関係を指定するテーブルであり、レジスタファイル或いは R A M 等により構成される。ネットワーク構成情報も重み係数と共に検出対象に応じて異なる値となる。

【 0 0 8 0 】

図 9 はネットワーク構成情報の一例を示す図であり、図 1 4 に示した C N N の構成情報を表現している。図 9 において、「対象処理ノード」は図 1 4 に示すネットワークの論理的な第 0 から第 8 処理ノードに対応する。なお、論理的な処理ノードとは、演算部 1 0 1 を時分割利用する事により実現する論理的な演算処理の単位である。ここで第 0 処理ノードは入力画像データ面に対応する処理ノードであって便宜上第 0 処理ノードとしているが、実際には演算は実行しない。

【 0 0 8 1 】

「隣接下層処理ノード数」とは、処理ノードが演算実行時に必要とする下位層の接続数を示す。例えば第 4 処理ノードの場合、3 つの下位層に接続される。ネットワーク構成管理部 1 0 8 では、当該隣接下層処理ノード数に応じてメモリアクセスと演算を制御する。

【 0 0 8 2 】

「隣接下層処理ノード」は処理ノードの演算時に必要とする下位層の処理ノードを指定する情報である。例えば第 4 処理ノードは、第 1 処理ノード、第 2 処理ノード、第 3 処理ノードに接続される。つまり、第 4 処理ノードの演算時は、第 1 ~ 3 処理ノードの演算結果が参照データとして使用される。

【 0 0 8 3 】

「演算種別」は実行する演算の種別を示す情報であり、C N N 演算の場合、演算種別に応じて重み係数を選択することになる。すなわち、演算部 1 0 1 が図 1 2 に示す構成の場合、ここでの演算種別番号が係数を選択するための「カーネル選択信号」に相当する。C N N 処理の場合、各処理ノードは、「カーネル選択信号」に応じて選択した、それぞれ異なる重み係数を用いて、コンボリューション演算を実行する。

【 0 0 8 4 】

「参照データ幅」は「演算種別」に対応するコンボリューションカーネルの幅に相当し、「参照データ高さ」はコンボリューションカーネルの高さに相当する。

【 0 0 8 5 】

また、「処理開始ライン」は、該処理ノードにおいて有効な演算出力が可能な画像位置を表す。先に説明したように、コンボリューション演算はカーネルサイズの範囲の周辺画素を掃き寄せる演算であるので、カーネル参照範囲が参照画像データの有効範囲外にはみ出す場合には、有効な演算結果を得ることができない領域となる。横方向には図 2 2 の領域 2 2 0 6 が有効範囲外の領域（無効領域）に相当するが、同様に縦方向にも無効領域が存在する。これは処理ノードのカーネルサイズに依存し、階層を経ることによって蓄積される。本実施形態のパターン検出装置の C N N では、ライン単位に処理を行うので、無効領域のラインは計算をスキップした方が全体の処理を高速化できる。つまり、「処理開始ライン」までは、その処理ノードにおける演算の無効領域となるので、処理をスキップすることができる。ちなみに、本実施形態では入力画像の最初のラインの番号を 0 として開始している。また、同様に開始ラインのみならず終了ラインも存在するが、これは入力画像のライン数 I_y から「処理開始ライン」を引いてさらに 1 を減じたライン番号となる。

【 0 0 8 6 】

尚、本実施形態では、入力画像に対して各処理ノードは演算可能な最大範囲を演算するように設定しているが、最終的に演算結果を利用する処理ノードの利用したい演算範囲から、前階層に向かって順に演算範囲を逆算してももちろんかまわない。この場合は終了ラインも合わせて情報として持つようにすると、より演算を無駄無く行える。

10

20

30

40

50

【 0 0 8 7 】

次にステップ S 8 0 4 において、CPU 6 8 は、中間バッファ割当て方式を選択する。この処理の詳細を示すのが、図 2 6 のフローチャートである。中間バッファ割当て方式の選択処理では、まずステップ S 2 6 0 1 において、CPU 6 8 は、設定されている入力処理画像サイズ I_x 、 I_y を取得する。そしてステップ S 2 6 0 2 において、CPU 6 8 は、中間バッファとして利用可能なサイズ M を計算する。このサイズ M は、ワークメモリ 1 0 2 で利用可能な容量から、入力画像バッファサイズや出力バッファサイズその他の処理に必要なサイズを差し引いた値となる。

【 0 0 8 8 】

そして、ステップ S 2 6 0 3 において、CPU 6 8 は、まずページバッファ方式で中間バッファを割り当てた場合に必要な容量 S_p を計算する。この容量 S_p は、上述した式 (3) で求めることができる。ステップ S 2 6 0 4 において、この S_p が M 以下であるかどうかを確認する。 M 以下であれば、ステップ S 2 6 0 5 において、CPU 6 8 は、ページバッファ方式を今回の中間バッファ割当て方式として選択し、本処理終了する。尚、ここで使用する値として、上述したように、 S_p と M の代わりに S_p' と M' を用いてももちろん良く、このようにするとよりページバッファ方式の選択される率が高まる。

【 0 0 8 9 】

一方、 M より S_p が大きい場合には、処理はステップ S 2 6 0 4 からステップ S 2 6 0 6 へ進む。ステップ S 2 6 0 6 において、CPU 6 8 は、バンドバッファ方式で割り当てた場合の中間バッファ容量 S_B を計算する。これは上述の式 (2) により求めることができる。次に、ステップ S 2 6 0 7 において、 S_B が M 以下であるかどうかを確認する。 M 以下であれば、ステップ S 2 6 0 8 において、CPU 6 8 は、バンドバッファ方式を今回の中間バッファ割当て方式として選択し、本処理を終了する。このように、本実施形態では、複数種類のバッファ割り当て方法のうち、必要となるバッファの合計サイズが中間バッファとして利用可能なメモリ量以下 (M 以下) になるバッファ割り当て方法が選択される。

【 0 0 9 0 】

S_p 、 S_B が共に M より大きい場合には、複数種類のバッファ割り当て方法のうち、必要となるバッファの合計サイズが中間バッファとして利用可能なメモリ量以下 (M 以下) になるバッファ割り当て方法がない。従って、そのままでは本パターン検出装置で演算処理を行うことができない。この場合、ステップ S 2 6 0 9 において、CPU 6 8 は、元の入力画像サイズを 2 分割したサイズを新たな処理画像サイズとして再設定する。この分割処理は、図 2 5 に示すように、カーネルのオーバーヘッドを考慮して行われる。図 2 5 において (a) は元の入力画像を示し、太線枠は最終層における累積カーネル領域 2 5 0 1 を示す。なお、ここで累積カーネル領域とは、最終層処理ノードの参照範囲が有効になるように、順に隣接下層処理ノードを遡って入力画像まで行き着いたときの参照範囲のことである。

【 0 0 9 1 】

これを図 2 5 の (b) の点線領域に示すように、左右に 2 分割することを考える。左右それぞれの点線領域を有効とするには、累積カーネル領域 2 5 0 1 を考慮すると、網掛けで示した領域まで入力する必要があることが分かる。

【 0 0 9 2 】

従って、オーバーヘッドを考慮した 2 分割とは、図 2 5 の (c) に示すものとなる。すなわち、入力画像サイズ (a) の単純な 2 分割の領域に、累積カーネル領域 2 5 0 1 の $1/2$ (奇数時切捨て) の幅を加えたものとなる。このように分割して演算した結果の両点線枠内の領域を合成すれば、図 2 5 の (d) に示すように、元の入力画像に CNN による演算を行ったのと同様の演算結果を得ることができる。なお、このように入力データを分割して処理する方式をタイル処理と呼ぶこともある。

【 0 0 9 3 】

本実施形態においては、図 2 5 に示したような横方向の分割処理を優先的に行う。なぜ

10

20

30

40

50

ならば、縦方向の分割処理では、ページ方式の中間バッファ使用容量は減るが、バンド方式の中間バッファ方式は変わらないからである。横方向の分割では、どちらの方式であっても一時的なトータル使用容量も削減される。分割可能なサイズは、累積カーネル領域の大きさによって決まる（領域以上のサイズが必要）ので、横方向がこれ以上分割できなくなったら、縦方向の分割を行う。尚、本処理では、分割をするたびに初期値 1 の分割カウンタをインクリメントするとともに、分割方向も合わせて記憶している。

【 0 0 9 4 】

ステップ S 2 6 1 0 において、分割後のサイズが累積カーネルサイズより大きいことを確認し、再度ステップ S 2 6 0 1 からの処理を繰り返す。もし、上述のように縦方向分割を駆使しても、分割後サイズが累積カーネルサイズよりも小さくなってしまった場合、その検出処理は本 C N N 演算装置で実行することができない。よって、処理はステップ S 2 6 1 1 に進み、C P U 6 8 は、処理不可能の判断を行う。これは全処理ノードのカーネルサイズと割当て可能な中間バッファ容量最大値によって決まるので、通常はこのようにことにならない構成を取る。

【 0 0 9 5 】

図 8 に戻り、ステップ S 8 0 5 において、C P U 6 8 は、ステップ S 8 0 4 の選択処理において処理不可能と判別されていないことを確認する。処理不可能と判別されていた場合、処理はステップ S 8 0 5 からステップ S 8 1 9 へ進み、終了判定処理（後述）が実行される。但し、通常はそのようなことのない構成をとる。

【 0 0 9 6 】

ステップ S 8 0 6 において、C P U 6 8 は、ステップ S 8 0 4 の選択処理で選択された方式がページバッファ方式であったかどうかを判定する。ページバッファ方式の場合、処理はステップ S 8 0 7 のページバッファ方式シーケンス設定およびステップ S 8 0 8 のページバッファ方式中間バッファ設定へ進む。一方、バンドバッファ方式が選択されていた場合、処理はステップ S 8 0 9 のバンドバッファ方式シーケンス設定およびステップ S 8 1 0 のバンドバッファ方式中間バッファ設定へと進む。

【 0 0 9 7 】

ステップ S 8 0 7 のページバッファ方式シーケンス設定処理と、ステップ S 8 0 9 のバンドバッファ方式シーケンス設定処理において、C P U 6 8 は、シーケンス制御部 1 0 9 のシーケンス情報設定部 1 1 1 にシーケンス情報を設定する。シーケンス情報とは、時分割処理する演算処理単位（本実施形態ではライン単位）のシーケンス動作を規定するテーブル情報であり、R A M 等に保持される。なお、上述したように、シーケンス情報は R A M 6 9 に格納され、C P U 6 8 から C N N 処理部 6 3 ないに書き込まれる。

【 0 0 9 8 】

図 1 0 はシーケンス情報テーブルの例であって、図 1 4 に示した C N N ネットワークで処理を行う場合の、(a) がバンドバッファ方式、(b) がページバッファ方式の場合のシーケンスを示している。ここで、「シーケンス番号」はライン単位での処理の順序を示す番号で、テーブルの配列 Index に相当するので実際に数字を保持する必要はない。「処理ノード番号」は図 1 4 の処理ノード番号であって、シーケンス番号に対応する論理的な実行処理ノードを示す。即ち、シーケンス番号 1 のとき、第 1 処理ノードに対して演算処理を実行し、シーケンス番号 2 においては第 2 処理ノードのライン単位演算処理を実行する。処理ラインは、該シーケンス番号の処理のとき該処理ノードが出力するラインの位置を示している。各処理ノードの処理ラインは、図 9 で説明した「処理開始ライン」より順にスタートし、1 ライン処理が完了する毎にインクリメントするものである。従って、本実施形態のように、テーブル上にシーケンス番号に対応する番号として明記しなくとも、各処理ノード毎に現在の処理ラインをレジスタ等で記憶しておくようにしてもよい。

【 0 0 9 9 】

図 1 0 の (a) は、バンドバッファ方式の場合のシーケンス情報テーブルの例であり、1 ライン単位処理するごとに、処理ノードを切り替えている。例えば、シーケンス番号 “ 2 1 ” までは、処理ノード番号 “ 1 ” ~ “ 3 ” までの演算を行っており、シーケンス番号

“ 2 2 ”において処理ノード番号 4 の処理を開始している。これは、シーケンス番号 “ 2 1 ” までで、処理ノード番号 4 の 1 ライン分の演算処理に必要な参照画像が演算されたためである。このようにバンドバッファ方式では、演算可能になった処理ノードから速やかにライン単位演算を行っていくことにより、参照される側のバッファをライン単位で順次開放可能にしている。これにより、中間バッファを必要最小限の循環数のリングバッファとして構成することが可能となる。

【 0 1 0 0 】

図 1 0 の (b) は、ページバッファ方式の場合のシーケンス情報テーブルの例であり、特定処理ノードの単位演算を連続して行い、該特定処理ノードの演算が有効領域全域について完了した後に、次の処理ノードの演算を開始している。このようにすると、図 9 および図 1 4 の CNN ネットワークにおいて、例えば処理ノード番号 “ 4 ” と “ 5 ” の全有効領域演算が完了すると、処理ノード番号 “ 1 ” ~ “ 3 ” の出力に割り当てられていた中間バッファ領域は不用となる。従って、最終階層まで演算の完了していないこの時点で、このバッファ領域を開放して同領域を処理ノード番号 7、8、9 の出力用として割り当てることができるようになる。

10

【 0 1 0 1 】

ステップ S 8 0 8 のページバッファ方式による中間バッファ設定およびステップ S 8 1 0 のバンドバッファ方式による中間バッファの設定では、各処理ノードに必要な中間バッファ領域が割り当てられる。当該処理は、論理的な処理ノードの数に対応する数のレジスタセット (リングバッファ設定部 1 0 4 -1 ~ 1 0 4 -n) に値を設定することにより行われる。図 1 4 および図 9 に示したネットワーク構造の場合、8 個のリングバッファ設定部 (1 0 4 -1 ~ 1 0 4 -8) に所定の値を設定する必要がある。

20

【 0 1 0 2 】

リングサイズ設定部 1 0 6 には対応する論理的な処理ノードのリングバッファの高さ (リングカウンタの循環数) が設定される。バンドバッファ方式の場合、この循環数は当該処理ノードの出力を参照する次階層の処理ノード (隣接上層処理ノード) のコンポリューションカーネルのうちの最大の高さに相当する。一方、ページバッファ方式の場合は、リングサイズ設定部 1 0 6 に設定されるリングバッファの高さは、入力画像のライン数 (もしくは各処理ノードの有効領域ライン数) と同じになり、当該リングバッファは実質循環させずに使用されることになる。オフセットアドレス設定部 1 0 7 には、対応する処理ノードのリングバッファ先頭アドレスが設定される。

30

【 0 1 0 3 】

図 1 1 は図 1 4 に示すネットワークを実現する場合の処理ノードとオフセットアドレス及びリングバッファの高さの関係の一例を示すメモリマップである。図 1 1 において、(a) がバンドバッファ方式の場合のメモリマップを、(b) がページバッファ方式の場合のメモリマップを示す。ADR x (x : 0 ~ 8) はオフセットアドレス、BH x (x : 0 ~ 3) がリングバッファの高さ (循環数) に相当する。I_x は入力画像データの幅を示す。図中、第 0 処理ノードのための領域とは、入力画像データを保持するページバッファ領域である。つまり BH0 は入力画像データの高さ I_y に等しい。また本実施形態において BH1 は、次階層 (第 2 階層) の第 4 処理ノード及び第 5 処理ノードのカーネルサイズのうちの大きい方、すなわち図 9 の構成情報テーブルから 9 ライン分のサイズが設定される。同様に BH2 には 1 3、BH3 には 1 7 ライン分のサイズが設定される。一方、ページバッファ方式の場合は、全ての領域の高さは BH0 となり、またいくつかの領域は複数の処理ノードによって利用される。本実施形態ではこのように、メモリ 1 0 2 を所定の領域に分割して各領域をサイズの異なるリングバッファ或いはフレームバッファとして利用する。

40

【 0 1 0 4 】

以上の各種設定処理により、後述のステップ S 8 1 4 で実行される CNN 演算処理におけるネットワーク演算処理の実行手順が決定されることになる。すなわち、選択された中間バッファ割り当て方式 (上記例ではバンドバッファ方式かページバッファ方式のいずれか) に応じてネットワーク演算処理における処理ノードの実行順が決定される。以上の各

50

種設定を終了すると、ステップS 8 1 1から、画像入力部6 1、前処理部6 2、CNN処理部6 3に対して各処理の開始を指示する。

【0105】

まずステップS 8 1 1において、CPU 6 8より処理開始の指示を受けた画像入力部6 1は、1フレーム分の画像データを取得し、図示しない内部バッファに格納する。画像入力部6 1は、画像データの格納が終了するとCPU 6 8に対して画像取得終了割り込みを発生する。CPU 6 8はこの割り込みを検知すると、DMAC 6 6を起動して取得した画像データを前処理部6 2の内部メモリ（図示しない）に転送する。前処理部6 2は画像データの転送が終了すると、前処理を開始する。前処理部6 2は、例えば、予め指定するコントラスト補正情報に従って画像データのコントラストを補正する。前処理部6 2は補正処理を終了するとCPU 6 8に対して割り込みを発生する。CPU 6 8はこの割り込みを検知すると、再びDMAC 6 6を起動し、前処理部6 2によって補正された画像データをCNN処理部6 3内のメモリ102の入力画像バッファ（図11の第0処理ノード領域に相当）に転送する。

10

【0106】

次に、ステップS 8 1 2において、CPU 6 8は、分割処理カウンタを設定する。これは上述のステップS 8 0 4において、中間バッファ割当て方式の選択処理を行った際に、入力画像（処理画像）が分割指定された場合に1よりも大きい値となる。分割無しの場合は1が設定される。そしてステップS 8 1 3において、CPU 6 8は、分割無しの場合は画像全域を、分割有りの場合は初回の処理領域（例えば図25の左右どちらか）を、処理領域として設定する。

20

【0107】

次に、ステップS 8 1 4において、CPU 6 8がCNN処理部6 3に対し演算開始トリガを送ることにより、検出処理が開始される。CNN処理部6 3におけるハードウェア処理について、図27のフローチャートを参照して、以下に説明する。

【0108】

まず、ステップS 2 7 0 1において、シーケンス制御部109は、処理ノードを決定する。シーケンス制御部109は、上述したシーケンス情報設定部111に記載されたシーケンス情報テーブルを、ライン単位演算毎に上から辿り、毎回の処理ノードを決定する。図10の(a)に示す例の場合、初回のシーケンスでは処理ノード番号1を選択する。シーケンス制御部109はシーケンス回数をカウントするシーケンスカウンタを有し、シーケンス単位（この場合ライン単位の処理毎）でカウントアップする。このカウンタは初期値1であって、図10のシーケンス情報テーブルのIndexとして使用できる。つまり、シーケンス制御部109は、シーケンスカウンタをアドレスとしてシーケンス情報テーブル（図10）を参照する事により処理対象処理ノードを決定する。

30

【0109】

ステップS 2 7 0 2において、処理ノードの演算に必要な参照データをメモリ102から読み出す。より具体的には、まず、ネットワーク構成管理部108がシーケンス制御部109の出力するシーケンス指示情報に従って参照データに対応するリングバッファ設定部104を選択する。すなわち、リングバッファ設定部104-1~104-nの何れかを選択する。例えば、ステップS 2 7 0 1で第1処理ノードが選択された場合、図9に示すネットワーク構成情報テーブルの内容に従って、「接続ノード数が1」「接続先ノードが第0処理ノード」「演算種別1」が決定される。ネットワーク構成管理部108はこのネットワーク構成情報テーブルの内容に従ってノード選択信号をセレクタ1121, 1122に出力し、参照すべきリングバッファに対応したリングバッファ設定部の出力を選択する。例えば初回のシーケンスでは対象処理ノードが第1処理ノードであるので、その隣接下層ノードである第0処理ノードに対応する選択信号を出力する。選択されたリングバッファ設定部の情報（この場合、第0処理ノードに対応したリングカウンタ値、オフセットアドレス値）に従ってメモリアクセス制御部103は読み出すメモリの先頭アドレスを生成する。

40

50

【 0 1 1 0 】

図 1 6 はメモリアクセス制御部 1 0 3 の内部を説明する図である。また、図 1 5 はメモリアクセス制御部 1 0 3 の参照データの読み出し動作を説明する図である。

【 0 1 1 1 】

図 1 5 において、1 5 0 1 はバンドバッファ方式におけるリングバッファ、1 5 0 2 は演算するコンボリューションカーネルの参照ウインドウに相当する大きさ、1 5 0 3 はコンボリューションカーネルの重み係数列の様子を説明する図である。ここではコンボリューションカーネルサイズが 6×6 の場合について例示している。重み係数列 1 5 0 3 において、W00 ~ W05 は 1 行目のデータ列に対する重み係数列、W10 ~ W15 は 2 行目のデータ列に対する重み係数列であり、以下同様に各データ列に対する重み係数列を示している。コンボリューション演算時は当該係数値と対応する位置の参照データの積和演算処理が実行される。

10

【 0 1 1 2 】

WIDTH は特徴面の幅（即ち本実施形態の場合、入力画像データの幅 I_x に相当）、L3 ~ L8 は特徴面の 3 行目から 8 行目のラインデータであることを示す。また A1 ~ A6 は夫々対応するラインの先頭メモリアドレスである。

【 0 1 1 3 】

メモリアクセス制御部 1 0 3 において、制御部 1 6 0 1 はネットワーク構成管理部 1 0 8 の出力する動作制御信号に従って各処理部及びメモリに対するコマンド信号（Read/Write 制御信号）を生成する。1 6 0 2 は列カウンタであり、今回の演算で使用するカーネルサイズに等しい参照ウインドウ 1 5 0 2（図 1 5 の点線枠で示された範囲）のデータバッファ上の位置（横方向の位置）を示す。アドレス変換部 1 6 0 5 が生成する行先頭アドレス A1 ~ A6 と、列カウンタ 1 6 0 2 及び後述のウインドウカウンタ 1 6 0 7 のカウンタ値を加算器 1 6 0 3 で加算する事で、リングバッファ内の各行内のデータをアクセスするためのメモリアドレスが生成される。1 6 0 7 はウインドウカウンタであり、参照範囲の各行に対して横方向に連続する参照画素を取り出すための参照ウインドウ幅（カーネル幅に相当）カウンタである。ウインドウカウンタ 1 6 0 7 は、参照ウインドウ幅分の画素をカウントすると 0 にリセットされる。尚、列カウンタ 1 6 0 2 やウインドウカウンタ 1 6 0 7 はネットワーク構成管理部 1 0 8 が保持するネットワーク構成情報（図 9）の内容に従って演算種別の変更毎に設定される。

20

30

【 0 1 1 4 】

すなわち、メモリアドレスは、アドレス変換部 1 6 0 5 の生成する行先頭アドレス、参照ウインドウの列位置を指定する列カウンタ 1 6 0 2 のカウンタ値、参照ウインドウ内の画素位置を指定するウインドウカウンタ 1 6 0 7 の各出力を加算した値である。メモリアクセス制御部 1 0 3 は、このようにして生成されたメモリアドレスをメモリ 1 0 2 に対して逐一出力する。

【 0 1 1 5 】

例えば、アドレス変換部 1 6 0 5 には、オフセットアドレス設定部 1 0 7 からの各処理ノードに割り当てられたオフセットアドレスと、リングカウンタ 1 0 5 からの最終行カウンタ値が入力される。そして最終行カウンタ値の示す行を縦方向の最後の位置として、カーネル高さ分のバッファ内の各行の先頭アドレス A1 ~ A6 を順次出力する。ここで、最終行カウンタ値は、図 1 5 に示すバンドバッファ方式の場合、リングバッファの値と一致し、リングバッファ内で最新の行が入っている位置を示すものとなる。最終行カウンタ値から逆算することにより、図 1 5 の例では L3 の入っている行の先頭アドレスすなわち A3 から、順に A4, A5, A6, A1, A2 と出力される。尚、特に図示はしないがページバッファ方式の場合は、各シーケンスに対し図 1 0 に示した処理ライン番号に、カーネル高さ（参照データ高さ）の $1/2$ （切捨て）を加えた値を最終行カウンタ値とする。なお、1 ライン分のデータ投入によりリングカウンタがインクリメントされ、アドレス変換部 1 6 0 5 には、リングカウンタからの最終行カウンタ値が入力される。したがって、アドレス変換部 1 6 0 5 は、カーネル幅のカウントを行う毎に次の行の先頭アドレスを発生することになる。

40

50

【 0 1 1 6 】

双方向制御部 1 6 0 4 はメモリデータバスの双方向制御を司るバッファであり、制御部 1 6 0 1 の出力する制御信号に従ってデータバスの方向制御をする。1 6 0 6 はコンボリユーシオン演算に必要な参照データを一時的に保持するキャッシュメモリである（以下、参照データキャッシュという）。上記のアドレス変換結果に基づいて得られた参照ウィンドウ内の参照データは、参照データキャッシュ 1 6 0 6 に格納される。制御部 1 6 0 1 はウィンドウカウンタ 1 6 0 7 を更新しながら参照データキャッシュ 1 6 0 6 を制御する事で、列方向に連続する参照データをキャッシュする。アドレス変換部 1 6 0 5 による先頭アドレスの出力順に従って、参照データキャッシュ 1 6 0 6 には元の正しいライン順でデータが格納される。

10

【 0 1 1 7 】

以上のようにして、メモリアクセス制御部 1 0 3 によってメモリ 1 0 2 からキャッシュへの参照データ群の読み出しを終了すると、処理はステップ S 2 7 0 3 に進む。ステップ S 2 7 0 3 において、演算部 1 0 1 は、コンボリユーシオン演算処理を開始する。ネットワーク構成管理部 1 0 8 は、構成情報テーブルに記録された「演算種別」情報（カーネル選択信号に対応する）に従って演算部 1 0 1 の重み係数を指定し、演算部 1 0 1 を駆動する。演算部 1 0 1 の乗算器 1 2 0 1 は、メモリアクセス制御部 1 0 3 の参照データキャッシュ 1 6 0 6 に格納された参照データを読み出し、演算種別情報で指定された重み係数を用いてコンボリユーシオン演算処理を実行する。そして、ステップ S 2 7 0 4 において、演算部 1 0 1 の累積加算器 1 2 0 2 は、コンボリユーシオン演算処理の演算結果を累積加算する。

20

【 0 1 1 8 】

次に、ステップ S 2 7 0 5 において、全ての接続先処理ノードの参照データについてコンボリユーシオン演算処理を実施したか否かを判定し、コンボリユーシオン演算が未実施の参照データがあれば処理をステップ S 2 7 0 4 に戻す。例えば、図 9 に示す例において、処理対象ノードが第 4 処理ノードの場合、接続先ノード数は 3 である。この場合、構成情報テーブルの内容に従って、第 1 処理ノードの結果、第 2 処理ノードの結果、第 3 処理ノードの結果に対するコンボリユーシオン演算処理が順次実行され、累積加算器 1 2 0 2 に累積、保持される。各処理ノードに対する参照データの読み出しと演算のシーケンスは前述した方法と同じである。即ち、メモリアクセス制御部 1 0 3 は処理ノード毎に異なるリングカウンタ値、オフセットアドレス値等の情報に従って、必要な参照データ群をメモリ 1 0 2 から参照データキャッシュ 1 6 0 6 に読み出す。そして、演算部 1 0 1 は当該キャッシュデータに対してコンボリユーシオン演算を実行する。

30

【 0 1 1 9 】

全ての接続先ノードに対する演算を終了すると処理はステップ S 2 7 0 5 からステップ S 2 7 0 6 に進む。ステップ S 2 7 0 6 において、非線形処理部 1 2 0 3 は、累積加算器 1 2 0 2 の出力を非線形変換する。

【 0 1 2 0 】

次に、ステップ S 2 7 0 7 において、CNN 処理部 6 3 は、演算部 1 0 1 で得られた非線形変換結果をメモリ 1 0 2 に格納する。より具体的には、まずネットワーク構成管理部 1 0 8 は自身の処理ノードに関するリングバッファ設定部 1 0 4 を選択する。例えば、第 1 処理ノードを演算している場合、第 1 処理ノードに対応するリングバッファ設定部 1 0 4 を選択する。メモリアクセス制御部 1 0 3 はここで指定されたリングバッファ設定部 1 0 4 のリングカウンタ 1 0 5 が示す行の次の行を先頭アドレスとしてメモリアドレスを生成する。ページバッファ方式であっても、リングカウンタには常に最新の格納済み行番号が示されている。なお、ライト動作時は制御部 1 6 0 1 によって、ウィンドウカウンタ 1 6 0 7 は 0 に初期化されている。メモリアクセス制御部 1 0 3 は生成した先頭アドレスに演算結果を書き込む。書き込みを終了すると列カウンタ 1 6 0 2 の値を 1 インクリメントする。列カウンタ 1 6 0 2 は 1 つの演算結果書き込み毎にインクリメントする。従って、次の処理時は 1 列分（1 画素分）ずれた領域の参照データ群が読み出される。

40

50

【 0 1 2 1 】

図 1 7 はここで説明した演算の様子をネットワーク構成管理部 1 0 8、メモリアクセス制御部 1 0 3、演算部 1 0 1 別に模式的にタイムチャート化した図である。上段がネットワーク構成管理部 1 0 8 の動作を示し、メモリアクセス制御部 1 0 3 及び演算部 1 0 1 はネットワーク構成管理部 1 0 8 の指示に従って各処理を実行する。

【 0 1 2 2 】

上述したように、ステップ S 2 7 0 1 において、ネットワーク構成管理部 1 0 8 は、シーケンス制御部 1 0 9 からのシーケンス制御指示情報に従って処理ノードを選択する (1 7 0 1)。そして、ネットワーク構成情報テーブルを参照して、接続ノード数を設定する (1 7 0 2)。続いて、ネットワーク構成管理部 1 0 8 は、選択された参照ノードに関する情報 (リングカウンタ値、オフセットアドレス値等) をメモリアクセス制御部 1 0 3 に通知し、参照データの読み出しを指示する (1 7 0 3)。メモリアクセス制御部 1 0 3 は、通知されたリングカウンタ値、オフセットアドレス値を用いてメモリ 1 0 2 から参照データを読み出し、参照データキャッシュ 1 6 0 6 にキャッシュする (1 7 0 4 , 1 7 0 5)。メモリアクセス制御部 1 0 3 による参照データの読み出しが完了すると、ネットワーク構成管理部 1 0 8 は、演算部 1 0 1 に対して、演算の開始を指示する。演算部 1 0 1 は、参照データキャッシュ 1 6 0 6 にキャッシュされた参照データを読み出してコンボリューション演算処理を実行する (1 7 0 6 , 1 7 0 7)。演算部 1 0 1 におけるコンボリューション演算処理が完了すると、ネットワーク構成管理部 1 0 8 は、次の参照ノードについて同様の処理 (1 7 0 9 ~ 1 7 1 3) を繰り返す。全ての参照ノードについてコンボリューション演算を完了すると、ネットワーク構成管理部 1 0 8 は、演算部 1 0 1 に非線形変換処理を実行させ (1 7 1 4)、特徴面における 1 画素の演算結果を得る。この演算結果をメモリ 1 0 2 に格納するために、ネットワーク構成管理部 1 0 8 は、上記処理ノードに関する情報 (リングカウンタ値、オフセットアドレス値等) をメモリアクセス制御部 1 0 3 に通知し、演算結果の書込みを指示する。メモリアクセス制御部 1 0 3 は、通知されたリングカウンタ値、オフセットアドレス値を用いてメモリ 1 0 2 に 1 行分の演算結果を書き込む (1 7 1 5 , 1 7 1 6 , 1 7 1 7)。そして、列カウンタを 1 6 0 2 をインクリメントする (1 7 1 8)。

【 0 1 2 3 】

以上の処理を 1 ライン分繰り返し (S 2 7 0 8)、処理を終了すると、処理はステップ S 2 7 0 8 からステップ S 2 7 0 9 へ進む。ステップ S 2 7 0 9 において、ネットワーク構成管理部 1 0 8 は処理中の演算ノードに対応するリングバッファ設定部 1 0 4 のリングカウンタ 1 0 5 をインクリメントする。リングカウンタ 1 0 5 の更新は 1 ラインの処理終了毎に行われる。リングカウンタ 1 0 5 はカウント値がリングサイズ設定部 1 0 6 の値に等しくなった場合 0 に初期化される。つまり、リングカウンタ 1 0 5 はリングサイズを基準にして循環する。但し、ページバッファ方式の場合はリングサイズ設定部 1 0 6 の値が入力画像サイズ高さと同様になるため、実質 0 に戻ることはない。この様に、メモリ 1 0 2 に対するアクセスを論理的な処理ノード毎にリングカウンタ 1 0 5 の動作に伴って処理する事でメモリ 1 0 2 上の所定の領域をサイズ (循環数) の異なる複数のリングバッファとして独立に使用する事ができる。即ち図 1 1 で示すメモリマップ上の領域を夫々リングバッファとして利用する事になる。

【 0 1 2 4 】

次にステップ S 2 7 1 0 において、CNN 処理部 6 3 は、全ての処理ノードが演算を終了したか否かを判定する。ここではシーケンス情報テーブル (図 1 0) に記された最後のシーケンスまでを終了したか否かが判定される。なお、シーケンス制御部 1 0 9 は図示しないシーケンスカウンタを予め設定されたシーケンス数と比較する事で終了判定を行っても良い。或いは、シーケンス制御部 1 0 9 は、テーブルの最後に付加された、予め定められた Termination データを検出することによっても終了判定を行うようにしてもよい。演算が終了していない場合、処理はステップ S 2 7 1 0 からステップ S 2 7 0 1 に戻る。そして、CNN 処理部 6 3 は、シーケンスカウンタを更新し、カウンタ値に対応するテーブ

10

20

30

40

50

ルを参照する事で次に処理する処理ノード番号を取得する。処理ノードを決定すると、シーケンス指示情報に従ってネットワーク構成管理部108は次の処理ノードに対する処理を開始する。異なる処理ノードを処理する場合も、リングバッファ及び演算に関する各種パラメータが異なるだけであり、前述した処理と同様の動作が繰り返される。

【0125】

尚、図10に示した様に、演算処理は下位層から順次リングバッファに特徴データを格納しながら処理を進めるが、バンドバッファ方式の場合とページバッファ方式の場合では処理の順序が異なる。

【0126】

以上、ステップS2701～S2710の処理を繰り返す事で、所定のCNNネットワークに基づく各特徴面の演算が、ライン単位で時分割処理しながら実行される。そして、CNN処理部63は全てのシーケンスを終了すると、ステップS2711において、CPU68に対して割り込みを発生する。

10

【0127】

図8に戻り、CPU68は割り込みを検知すると、ステップS814のCNN演算処理が完了したと見なし、ステップS815において、出力画像の取得処理を行う。この処理において、CPU68は、DMAC66を起動してCNN処理部63から必要な演算結果をRAM70に転送する。本実施形態では、最終層の第8処理ノードの出力結果を吸い上げている。尚、ステップS812で設定した分割処理カウンタが1よりも大きい場合（入力画像が分割されて処理された場合）は、RAM70に用意する出力データ格納領域の中で、ステップS804で定めた領域に対応する位置へとデータを転送する。

20

【0128】

ステップS816において、CPU68は、分割カウンタをデクリメントする。そして、ステップS817において、分割カウンタが0になるまでステップS813からの一連の処理を繰り返す。この結果、最終的に入力画像一面分に対応するCNN演算結果がRAM70上に格納される。

【0129】

演算結果がRAM70上に格納されると、ステップS818において、CPU68は、判定処理を実行する。この判定処理では、RAM70上に吸い上げた最終層処理ノードの出力である特徴データを利用して対象物の検出状況を判定する。例えば、所定のしきい値で特徴データを2値化しその重心を取得する等の方法で対象物の有無を判定するという判定処理が行われる。

30

【0130】

以上で入力画像に対する1検出対象の検出処理が完了する。次の入力画像や検出対象を変更しての処理を行わないならば、本処理を終了する（ステップS819）。一方、検出対象を変更する場合、処理はステップS820からステップS802に戻り、各種パラメータが再設定される。そして、上述の処理を繰り返し、検出対象に応じた重み係数/ネットワーク構成情報/シーケンス情報をそれぞれ更新する。更にリングバッファ設定部104のリングカウンタも新たな重み係数及びネットワーク構成情報に応じて再設定する。これにより、論理的な処理ノードは、検出対象に応じて、メモリ102を異なるサイズのリングバッファとしてマッピングし処理を行う。

40

【0131】

一方、ステップS820において検出対象を変更しない場合は、ステップS821に処理を進める。ステップS821とS822において、CPU68は、リングバッファ設定部104-1～104-nのリングカウンタ105及びシーケンス制御部109の内部カウンタ等を初期化する。そして、処理をステップS811に戻し、画像データの取得から再開する。即ち次のフレーム画像に対して同じ検出処理を実行する。

【0132】

以上、第1実施形態によれば、論理的な処理ノードに毎にリングバッファを制御するリングバッファ設定部104が設けられ、ネットワーク構成と目的に応じてリングバッファ

50

のサイズ（循環数）が設定される。そして、この構成において、複数の中間バッファ割当て方式（ページバッファ方式、バンドバッファ方式）から最適なものをネットワーク構成に基づいて選択可能としている。この構成により、同一のハードウェアでより多くの種類のコンポリュショナルニューラルネットワーク等の階層的な演算処理を処理する事が可能になる。更に、当該複数の中間バッファ割当て方式に優先順位を設けたことにより、同一条件化でより高速な演算の可能な方式を選択することができる。また、どの中間バッファ割当て方式でも所定の範囲内の容量に収まらない場合に、入力画像を分割してタイル処理することにより、より多くの種類の階層演算に対応可能としている。

【 0 1 3 3 】

< 第 2 実施形態 >

第 1 実施形態では全ての論理的な処理ノード毎にリングバッファのサイズを設定可能な構成について説明したが、本発明はこれに限るわけではない。例えば、階層毎にリングバッファのサイズを設定する構成とすることも可能である。第 2 実施形態では、そのような構成について説明する。

【 0 1 3 4 】

図 19 は階層毎にリングバッファのサイズを規定する場合の CNN 処理部 63 の構成を示す。図 18 は図 19 に示す CNN 処理部 63 で実現される CNN ネットワーク構成の一例を示す図である。図 18 では、階層毎にのみリングバッファのサイズが異なっている様子が示されている。即ち、第 1 階層の演算結果を格納するためのメモリ領域 1803a ~ 1803c と、第 2 階層の演算結果を格納するためのメモリ領域 1807a ~ 1807b を夫々同じサイズのリングバッファで構成する。図 18 は図 7 と比較して、メモリ領域 1803c のバッファサイズが異なっていることが分かる。

【 0 1 3 5 】

以下、第 1 実施形態との違いについて説明する。第 2 実施形態の CNN 処理部 63 は、論理的な処理ノード毎にリングバッファ設定部 194-1 ~ 194-n を有する。以下、リングバッファ設定部 194-1 ~ 194-n の任意の 1 つを指す場合は、リングバッファ設定部 194 と記載する。各リングバッファ設定部 194 は、第 1 実施形態のリングカウンタ 105 とオフセットアドレス設定部 107 に対応するリングカウンタ 1951 とオフセットアドレス設定部 1971 を有する。但し、第 2 実施形態のリングバッファ設定部 194 は、第 1 実施形態のリングバッファ設定部 104 が有していたリングサイズ設定部 106 を有していない。その代わりに、第 2 実施形態による CNN 処理部 63 は、階層型ネットワーク演算の論理的な階層毎に、リングサイズ設定部を有する。図 19 の例では、階層数が 3 までに対応するべく、2 つのリングサイズ設定部 1961a, b が設けられている。

【 0 1 3 6 】

リングサイズ設定部 1961a, b は夫々複数のリングバッファ設定部 194-1 ~ 194-n に接続されている。本例では、リングサイズ設定部 1961a は第 1 階層 1806 の処理ノードに対応した複数のリングバッファ設定部 194 に接続され、リングサイズ設定部 1961b は第 2 階層 1810 の処理ノードに対応した複数のリングバッファ設定部 194 に接続される。即ち、リングバッファ設定部 194-1 ~ 194-n がリングサイズ設定部 1961a, b によってグルーピングされている。

【 0 1 3 7 】

図 20 にリングサイズ設定部 1961a, b とリングバッファ設定部 194 の関係を示す。第 1 階層 1806 のためのリングバッファの制御に利用するリングバッファ設定部 194 として、リングサイズ設定部 1961a が接続されたリングバッファ設定部が選択される。一方、第 2 階層のリングバッファの制御に利用するリングバッファ設定部 194 には、リングサイズ設定部 1961b が接続されたリングバッファ設定部が選択される。演算時、処理ノードに対応するリングバッファ設定部 194 の選択は、ネットワーク構成管理部 108 が保持する管理テーブル情報に従って行われる。

【 0 1 3 8 】

以下、処理フローに関して、図 8 を用いて第 1 実施形態との違いを説明する。第 2 実施

10

20

30

40

50

形態では、ステップ S 8 0 8 およびステップ S 8 1 0 において、リングサイズ設定部 1 9 6 1 a、b への設定が階層毎に行われる。また、階層毎にグルーピングされたリングバッファ設定部 1 9 4 の中から構成するネットワークに対応するリングバッファ設定部 1 9 4 を選択し、オフセットアドレスを設定する。リングサイズ設定部 1 9 6 1 a、b には、図 1 8 のネットワーク構成では、第 1 階層、第 2 階層のリングバッファ高さに相当する値が設定される。ページバッファ方式の場合、この値は入力画像のサイズと同じである（ステップ S 8 0 8）。バンドバッファ方式の場合、この値は次階層の処理ノードのコンポリューションカーネルの内の最大のものの高さとなる。

【 0 1 3 9 】

他の処理は第 1 実施形態と同じであり説明を省略する。図 1 8 に示すネットワークを処理する場合、以上の設定でネットワーク構成管理部 1 0 8 が所定の論理処理ノードに対応するリングバッファ設定部 1 9 4 を選択しながら処理を進める事で第 1 実施形態と同様にライン単位で処理が実行される。

【 0 1 4 0 】

以上のように、第 2 実施形態によれば、リングサイズ設定部を各リングバッファ設定部に設けず、階層毎に設けるようにしたので、リングサイズ設定部を構成するレジスタの数を削減する事が可能になる。

【 0 1 4 1 】

< 第 3 実施形態 >

上記第 1、第 2 実施形態では CNN 処理部 6 3 をハードウェアで実現する場合について説明したが、本発明はソフトウェアにより実現する場合にも適用することが可能である。図 2 1 にソフトウェアで実現する場合の構成例を示す。図 2 1 に示す構成は図 6 に示す構成から CNN 処理部 6 3 を取り除き、ランダムアクセス可能な高速メモリである RAM 2 1 0 1 を追加したものであるためその違いについて説明する。

【 0 1 4 2 】

CPU 6 8 は前処理部 6 2 の終了割り込みを受け付けると DMA C 6 6 を起動して前処理部 6 2 内のメモリに格納された補正後の画像データを RAM 2 1 0 1 に転送する。CPU 6 8 は RAM 2 1 0 1 に格納した画像データに対して、図 2 7 に示したステップ S 2 7 0 1 ~ ステップ S 2 7 1 0 の処理をソフトウェアにより実行する。その場合、CNN 演算処理の動作に必要なワークメモリとして RAM 2 1 0 1 を使用する。即ち、CPU 6 8 は、RAM 2 1 0 1 上に図 1 1 で示すメモリマップを構成し、処理ノードに対応する各メモリ領域をリングバッファとして使用する。もちろんリングバッファはハード構成の場合と同様、バンドバッファ方式とページバッファ方式のいずれかの割当て方式を選択して使用できる。

【 0 1 4 3 】

尚、第 1 実施形態の CNN 処理部 6 3 に存在するリングバッファ設定部 1 0 4 等はソフトウェア上の変数として構成され、具体的には RAM 7 0 上にアサインされる。

【 0 1 4 4 】

以上の第 3 実施形態によれば、CNN 処理部 6 3 をハードウェアにより構成する場合と同様に、処理に必要なバッファメモリを削減する事が可能になる。図 2 1 に示す構成の場合、RAM 2 1 0 1 を少ないメモリで実現する事が出来る。また、RAM 2 1 0 1 を用意せずに RAM 7 0 をワークメモリとして利用する場合であっても同様である。

【 0 1 4 5 】

< 他の実施形態 >

上記実施形態では、リングカウンタを使用して、メモリ 1 0 2 の所定の連続領域をライン単位で循環しながら使用方法について説明したが、本発明はこのようなメモリの使用方法に限るわけではない。例えば、リングカウンタに対応するメモリアドレステーブルを有し、当該テーブルを参照する事で、不連続な領域を所定の処理単位に割り当てながら処理する等の方法でも良い。即ち、本発明で規定するリングバッファとは狭義のリングバッファ或いは循環バッファに限定するものではない。

10

20

30

40

50

【0146】

上記実施形態では、シーケンス情報テーブルの内容に従って論理的な処理ノードをライン単位で時分割処理する場合について説明したが、他の方法を適用しても良い。例えば、読み取り側バッファと書き込みバッファの利用状況から適応的にスケジューリングする等の方法でも良い。特にページ割当て方式を選択した場合、処理単位自体をライン単位からページ単位に切り替えるような実装としてももちろんかまわない。

【0147】

また実施形態ではページバッファ方式をバンドバッファ方式より優先的に選択する例を紹介したが、もちろん優先順はこれに限るものではない。例えば、演算結果データ群の最初のラインが出力されるまでのレイテンシを重視する場合には、バンドバッファ方式を優先することも可能である。このようにすると、例えば、全演算結果が出力される前に途中で演算を打ち切るような場合に有利にできることは明らかである。なお、バンドバッファ方式を優先する場合は、図26において、S2603～S2605の処理と、S2606～S2608の処理を入れ替えれば良い。

10

【0148】

また、中間バッファの割り当て方法としては、上述のバッファ割当て方式に限らず、他のバッファ割当て方式や、複数の方式を混在させるようにしても良いことは明らかである。すなわち、

- ・各処理ノードの演算結果データを保持するための中間バッファを各処理ノードへの割り当てるための複数種類のバッファ割り当て方法を用意しておく、
- ・これらのバッファ割り当て方法のそれぞれについて、ネットワーク演算に必要なメモリ量を当該ネットワーク演算の構成に基づいて算出し、
- ・算出されたメモリ量に基づいて、複数種類のバッファ割り当て方法のうちの1つを選択するように構成することが可能である。

20

【0149】

ここで、算出された必要なメモリ量が、メモリ102の中間バッファの割り当てに利用可能なメモリ容量以下となるバッファ割り当て方法が複数存在した場合は、予め定められた優先順位に従って使用すべきバッファ割り当て方法が選択される。また、算出された必要なメモリ量が、上記割り当てに利用可能なメモリ容量以下となるバッファ割り当て方法が存在しない場合には、上述したように入力データを分割して処理する。

30

【0150】

また、上記各実施形態では、特徴抽出結果を入力層と同じ解像度で保持するが、特徴面を入力面に対してサブサンプリングする場合に対しても同様に適用可能である。

【0151】

また、上記各実施形態では、最も効率の良い処理単位としてライン単位でシーケンス制御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。1ライン以下の単位やブロック単位でシーケンス制御する場合にも適用可能であることは明らかであり、その構成も当業者には明らかである。

【0152】

また、上記各実施形態では、本発明をコンポリューションニューラルネットワーク演算に適用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。所定の参照領域を必要とする様々な階層的な演算処理に対して適用することが可能である。更に、2次元演算に限るわけでもない。

40

【0153】

また、上記実施形態ではリングバッファ設定部104-1～104-nをレジスタとして提供する場合について説明したが、メモリとして他のパラメータメモリやワークメモリ102と共有する構成にしても良い。その場合、回路リソースをより有効に利用できる。即ち、より柔軟なネットワーク構成を実現する事が可能になる。

【0154】

また、上記実施形態では、参照データキャッシュ1606に2次元の参照データを全て

50

取り込んだ後に演算部 101 を起動する場合について説明したが、これに限られるものではない。例えば、コンボリューション演算等の場合、参照ウィンドウ内の 1 行単位で演算部 101 を駆動する様に制御する事も可能である。この場合、ウィンドウカウンタ 1607 のデクリメントに伴う連続する列方向参照データの読み出しを終了すると、次のラインの参照データ読み出し処理開始前に演算部 101 を駆動し、読み出したデータに対するコンボリューション演算を実行する。演算が終了すると次のラインの参照データ読み出しを開始する。以上の処理を繰り返す。この場合、参照データキャッシュのサイズが参照データ幅分のみで良いため、少ないキャッシュ容量で実現する事が出来る。

【0155】

また、第 2 実施形態では、リングサイズ設定部 1961a, b を階層毎に有する場合について説明したが、リングバッファ設定部 104 (194) を階層毎に有する構成とすることも可能である。その場合は階層処理単位でリングカウンタを更新する様に制御する。論理的な処理ノードの動作シーケンスに制限が生じる(必ず階層単位でシーケンスする必要が生じる)が、より回路規模を削減する事が出来る。

【0156】

また、上記実施形態では、リングサイズ設定部が任意に設定可能な場合(レジスタや RAM により構成する場合)について説明したが、全て或いは一部が固定的な値として指定される構成とすることもできる。

【0157】

以上、実施形態を詳述したが、本発明は、例えば、システム、装置、方法、プログラムもしくは記憶媒体等としての実施態様をとることが可能である。具体的には、複数の機器から構成されるシステムに適用しても良いし、また、一つの機器からなる装置に適用しても良い。

【0158】

尚、本発明は、ソフトウェアのプログラムをシステム或いは装置に直接或いは遠隔から供給し、そのシステム或いは装置のコンピュータが該供給されたプログラムコードを読み出して実行することによって前述した実施形態の機能が達成される場合を含む。この場合、供給されるプログラムは実施形態で図に示したフローチャートに対応したコンピュータプログラムである。

【0159】

従って、本発明の機能処理をコンピュータで実現するために、該コンピュータにインストールされるプログラムコード自体も本発明を実現するものである。つまり、本発明は、本発明の機能処理を実現するためのコンピュータプログラム自体も含まれる。

【0160】

その場合、プログラムの機能を有していれば、オブジェクトコード、インタプリタにより実行されるプログラム、OS に供給するスクリプトデータ等の形態であっても良い。

【0161】

コンピュータプログラムを供給するためのコンピュータ読み取り可能な記憶媒体としては以下が挙げられる。例えば、フロッピー(登録商標)ディスク、ハードディスク、光ディスク、光磁気ディスク、MO、CD-ROM、CD-R、CD-RW、磁気テープ、不揮発性のメモリカード、ROM、DVD(DVD-ROM, DVD-R)などである。

【0162】

その他、プログラムの供給方法としては、クライアントコンピュータのブラウザを用いてインターネットのホームページに接続し、該ホームページから本発明のコンピュータプログラムをハードディスク等の記録媒体にダウンロードすることが挙げられる。この場合、ダウンロードされるプログラムは、圧縮され自動インストール機能を含むファイルであってもよい。また、本発明のプログラムを構成するプログラムコードを複数のファイルに分割し、それぞれのファイルを異なるホームページからダウンロードすることによっても実現可能である。つまり、本発明の機能処理をコンピュータで実現するためのプログラムファイルを複数のユーザに対してダウンロードさせる WWW サーバも、本発明に含まれる

10

20

30

40

50

ものである。

【 0 1 6 3 】

また、本発明のプログラムを暗号化してCD-ROM等の記憶媒体に格納してユーザに配布するという形態をとることもできる。この場合、所定の条件をクリアしたユーザに、インターネットを介してホームページから暗号を解く鍵情報をダウンロードさせ、その鍵情報を使用して暗号化されたプログラムを実行し、プログラムをコンピュータにインストールさせるようにもできる。

【 0 1 6 4 】

また、コンピュータが、読み出したプログラムを実行することによって、前述した実施形態の機能が実現される他、そのプログラムの指示に基づき、コンピュータ上で稼動しているOSなどとの協働で実施形態の機能が実現されてもよい。この場合、OSなどが、実際の処理の一部または全部を行ない、その処理によって前述した実施形態の機能が実現される。

【 0 1 6 5 】

さらに、記録媒体から読み出されたプログラムが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書き込まれて前述の実施形態の機能の一部或いは全てが実現されてもよい。この場合、機能拡張ボードや機能拡張ユニットにプログラムが書き込まれた後、そのプログラムの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPUなどが実際の処理の一部または全部を行なう。

【 図面の簡単な説明 】

【 0 1 6 6 】

【 図 1 】 第 1 実施形態の階層的演算処理装置の構成を説明するブロック図である。

【 図 2 】 階層結合型ニューラルネットワークの構成例を説明する図である。

【 図 3 】 ニューロンの構成を示す図である。

【 図 4 】 CNNのネットワーク構成例を説明する図である。

【 図 5 】 CNNの特徴抽出の例を説明する図である。

【 図 6 】 実施形態による、階層的演算処理を利用した画像処理装置の構成を説明するブロック図である。

【 図 7 】 実施形態に関するCNNのネットワークの構成を説明する図である。

【 図 8 】 第 1 実施形態の画像処理装置の動作を説明するフローチャートである。

【 図 9 】 ネットワーク構成情報テーブルのデータ構成例を示す図である。

【 図 1 0 】 シーケンス情報テーブルのデータ構成例を示す図である。

【 図 1 1 】 メモリの割り当てに関する例を示す図である。

【 図 1 2 】 演算処理部 1 0 1 の構成例を説明するブロック図である。

【 図 1 3 】 リングバッファの例を説明する図である。

【 図 1 4 】 処理ノードの論理的な接続構成を説明する図である。

【 図 1 5 】 参照データ群の読み出しを説明する図である。

【 図 1 6 】 メモリアクセス制御部 1 0 3 の構成を説明する図である。

【 図 1 7 】 CNN演算単位の動作タイミングを説明する図である。

【 図 1 8 】 第 2 実施形態のCNNネットワークの構成を説明する図である。

【 図 1 9 】 第 2 実施形態の演算処理装置の構成を説明するブロック図である。

【 図 2 0 】 第 2 実施形態のリングバッファ設定部とリングサイズ設定部の関係を説明する図である。

【 図 2 1 】 第 2 実施形態の演算処理装置の構成を説明するブロック図である。

【 図 2 2 】 ライン単位で演算部 1 0 1 が処理を実行する場合の様子を模式的に説明する図である。

【 図 2 3 】 バンドバッファ方式による中間バッファ割当ての様子の一例を示す図である。

【 図 2 4 】 ページバッファ方式による中間バッファ割当ての様子の一例を示す図である。

【 図 2 5 】 実施形態による分割処理の一例を示す図である。

10

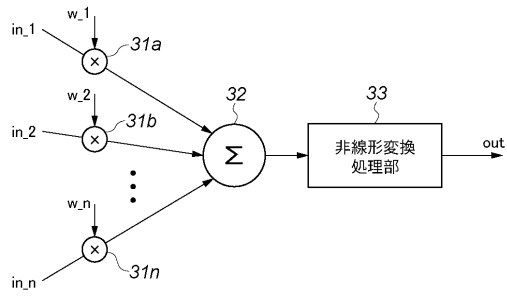
20

30

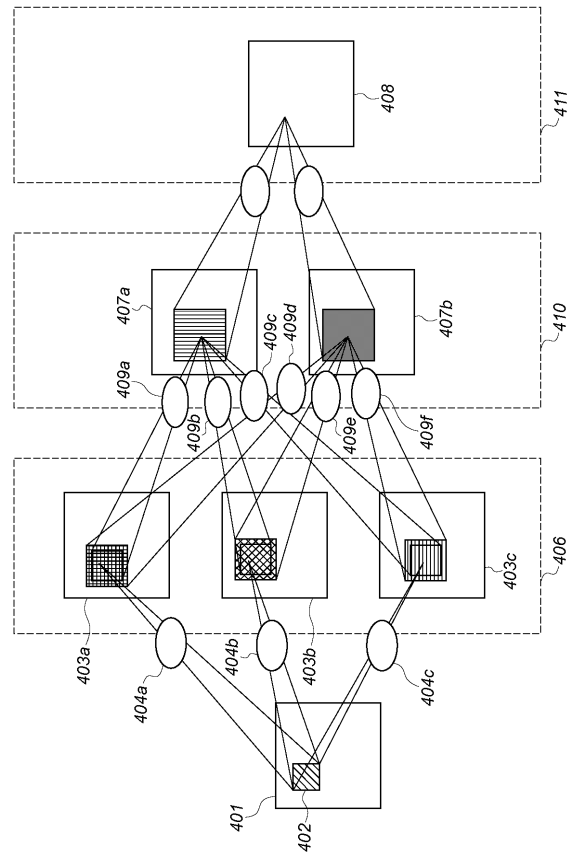
40

50

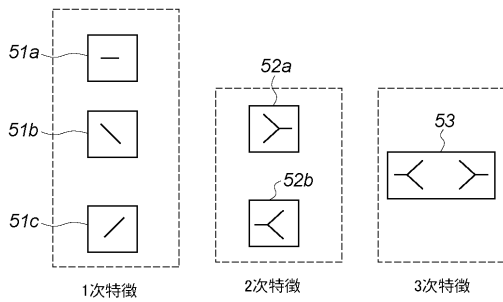
【図3】



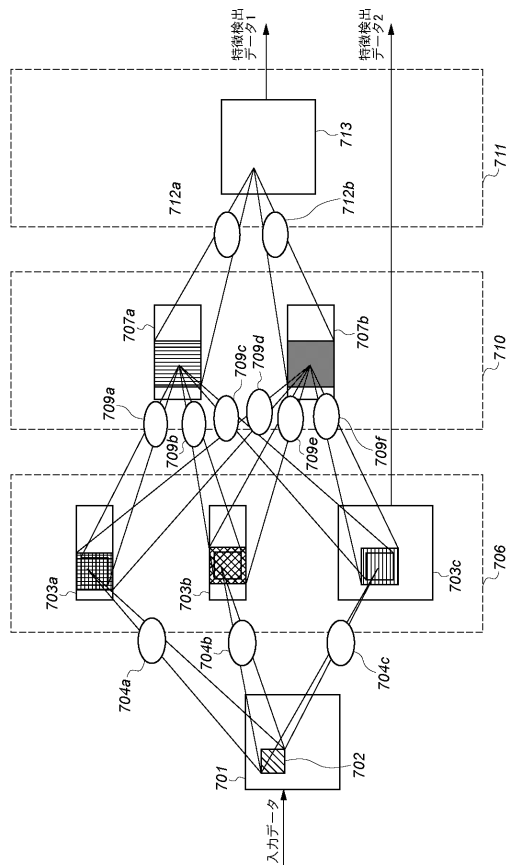
【図4】



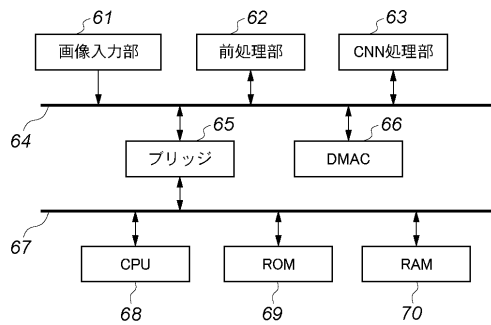
【図5】



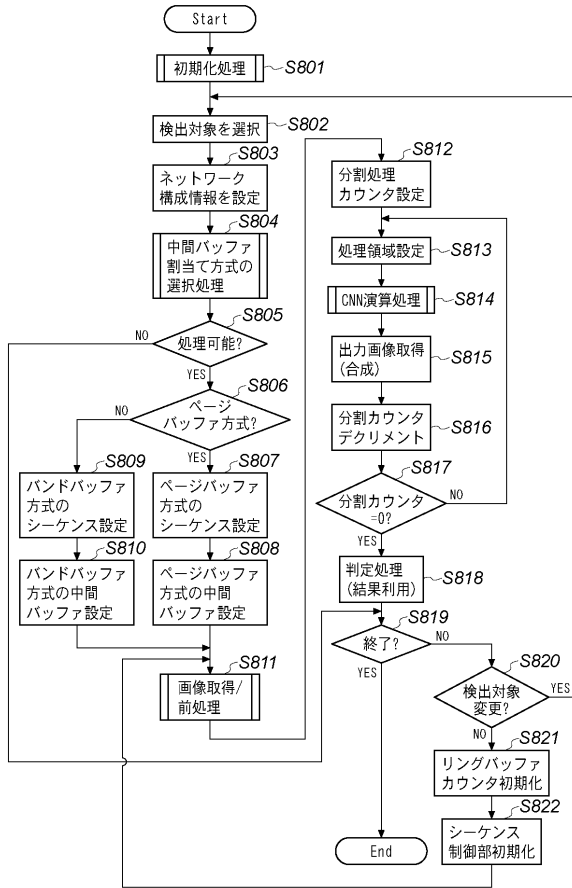
【図7】



【図6】



【図 8】



【図 9】

対象処理ノード	隣接下層処理ノード数	隣接下層処理ノード	演算種別	参照データ幅	参照データ高さ	処理開始ライン (有効領域オフセット)
第0処理ノード	1	第0処理ノード	演算1	1	1	0
第1処理ノード	1	第0処理ノード	演算2	3	3	1
第2処理ノード	1	第0処理ノード	演算3	3	3	1
第3処理ノード	3	第1処理ノード	演算4	7	7	4
第4処理ノード	3	第1処理ノード	演算5	7	7	4
第5処理ノード	3	第2処理ノード	演算6	7	7	4
第6処理ノード	3	第2処理ノード	演算7	9	9	5
第7処理ノード	3	第2処理ノード	演算8	9	9	5
第8処理ノード	2	第3処理ノード	演算9	9	9	5
第9処理ノード	2	第3処理ノード	演算10	11	11	10
第10処理ノード	2	第3処理ノード	演算11	11	11	10
第11処理ノード	2	第4処理ノード	演算12	13	13	11
第12処理ノード	2	第4処理ノード	演算13	13	13	11
第13処理ノード	2	第5処理ノード	演算14	17	17	17
第14処理ノード	2	第5処理ノード	演算15	17	17	17

【図 10】

シーケンス番号	処理ノード番号	処理ライン
1	1	1
2	2	1
3	3	1
4	1	2
5	2	2
6	3	2
...
19	1	7
20	2	7
21	3	7
22	4	4
...
27	1	9
28	2	9
29	3	9
30	4	6
31	5	5
...
4	16	4
5	15	5
6	10	6
...
3	21	3
4	18	4
5	17	5
6	12	6
7	11	7
...
3	37	3
4	34	4
5	33	5
6	28	6
7	27	7
8	19	8
...
1	N-2	1
2	N-2	2
3	N-2	3
4	N-5	4
5	N-6	5
6	N-11	6
7	N-12	7
8	N-20	8

(a)

シーケンス番号	処理ノード番号	処理ライン
1	1	1
2	1	2
3	1	3
4	1	4
...
N-4	1	N-4
N-3	1	N-3
N-2	1	N-2
N-1	2	1
N	2	2
N+1	2	3
...
2N-3	2	N-3
2N-4	2	N-2
...
3	1	3
...
3	2	3
...
4	4	4
...
4	5	4
...
4	6	4
...
4	N-4	4
...
4	N-5	4
5	5	5
...
5	6	5
...
5	N-6	5
...
6	N-11	6
7	N-11	7
...
7	N-11	7
...
7	N-12	7
...
8	N-19	8
8	N-20	8

(b)

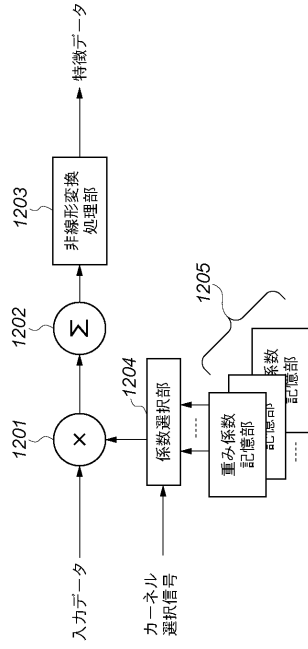
【図 11】



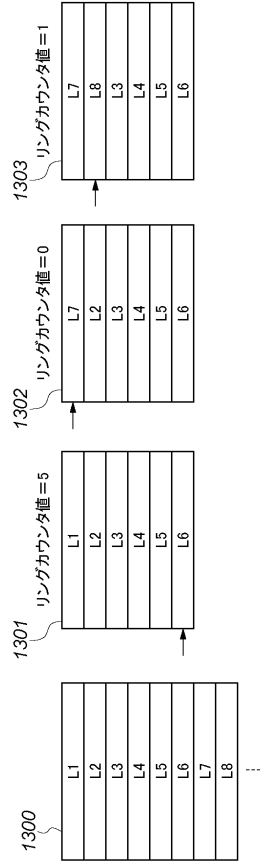
(a)

(b)

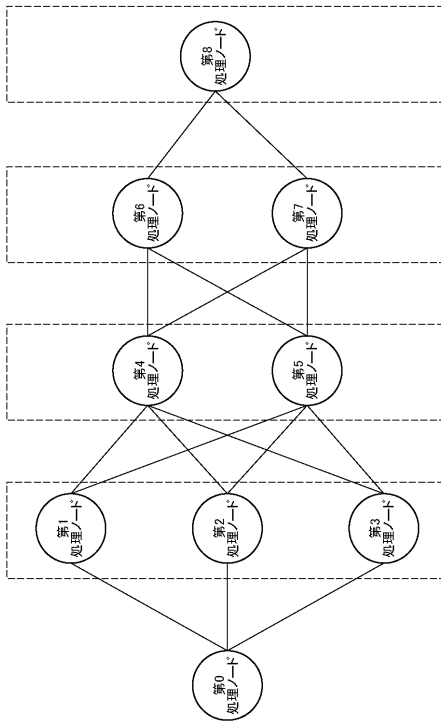
【図 1 2】



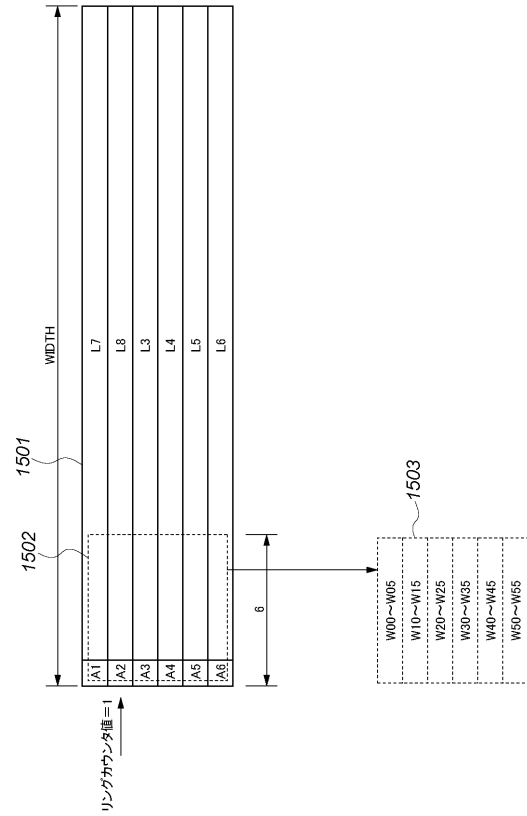
【図 1 3】



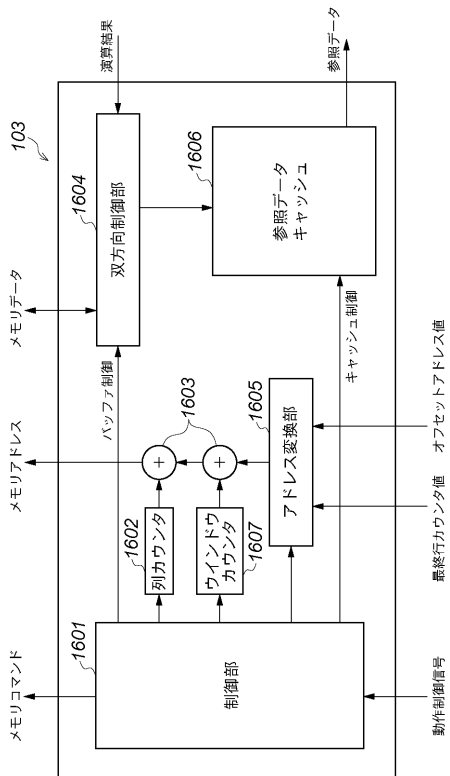
【図 1 4】



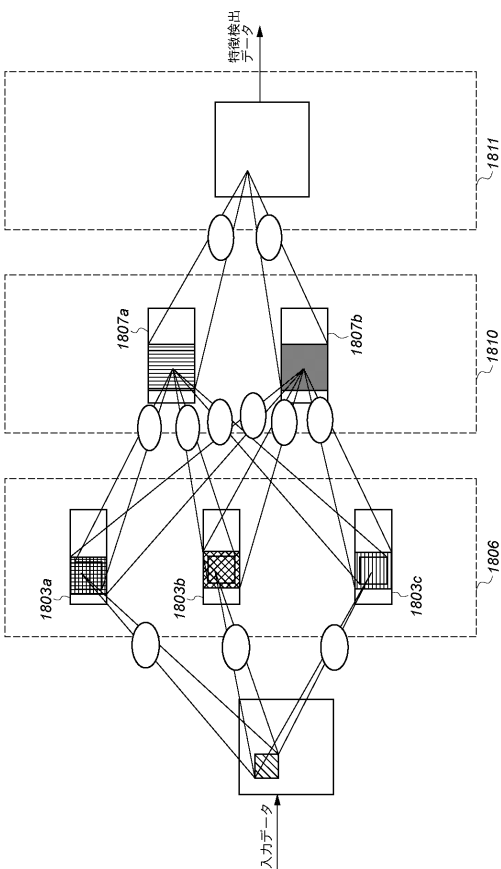
【図 1 5】



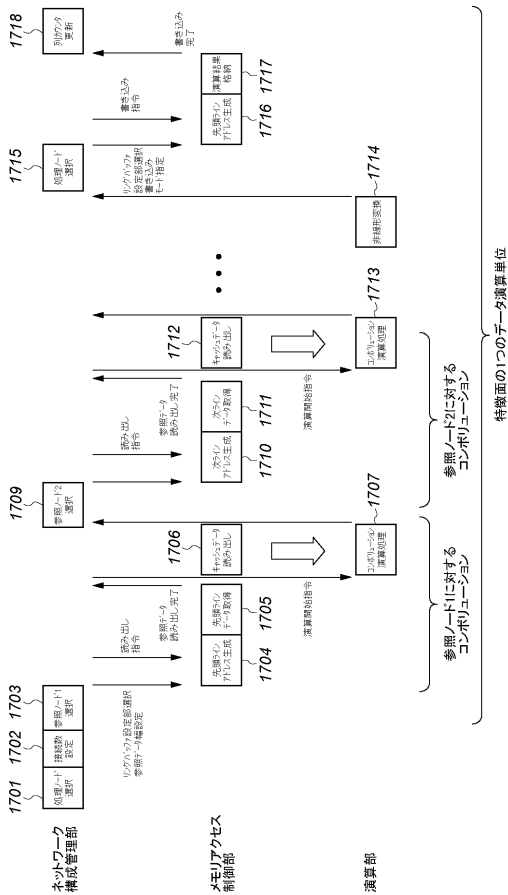
【図 16】



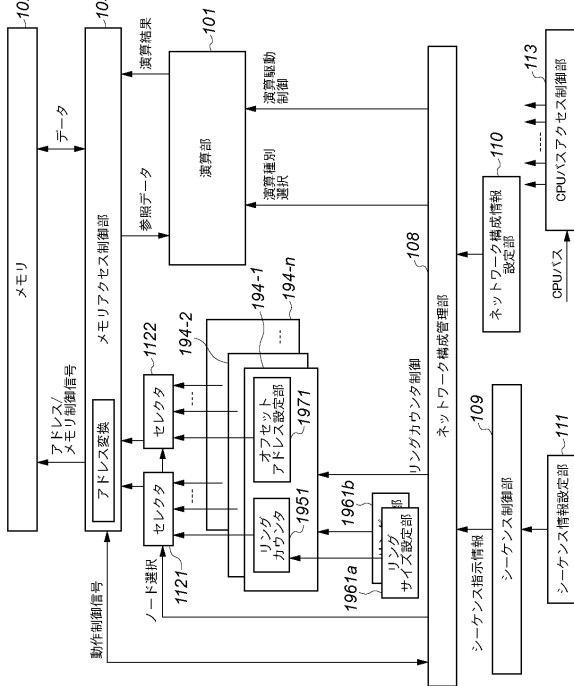
【図 18】



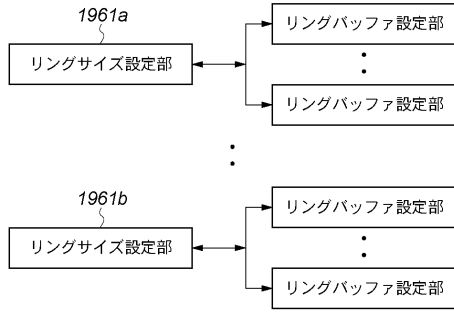
【図 17】



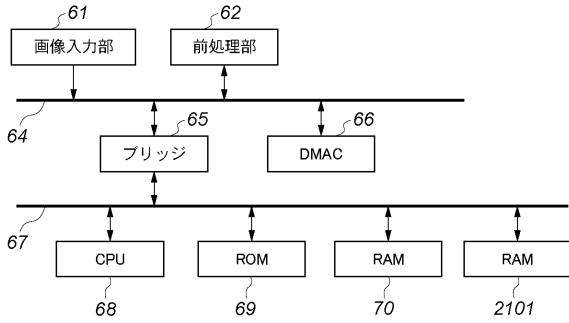
【図 19】



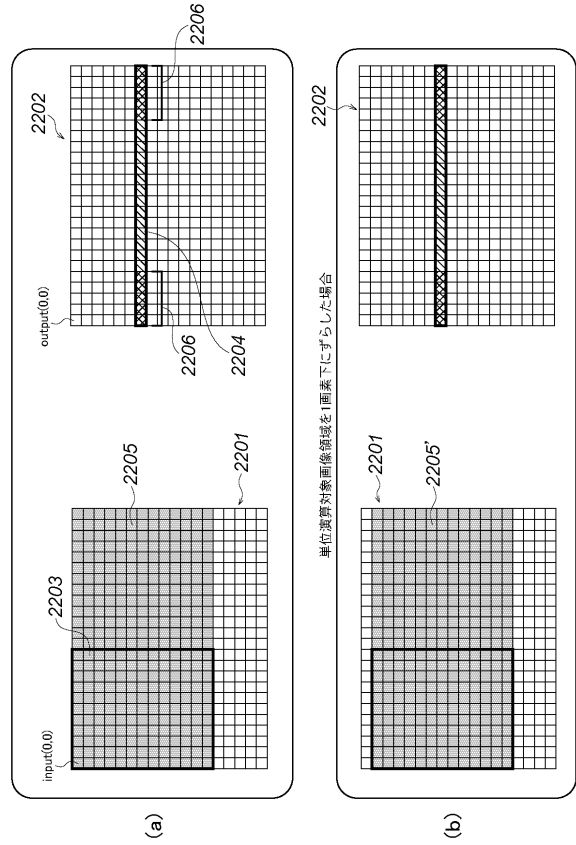
【図20】



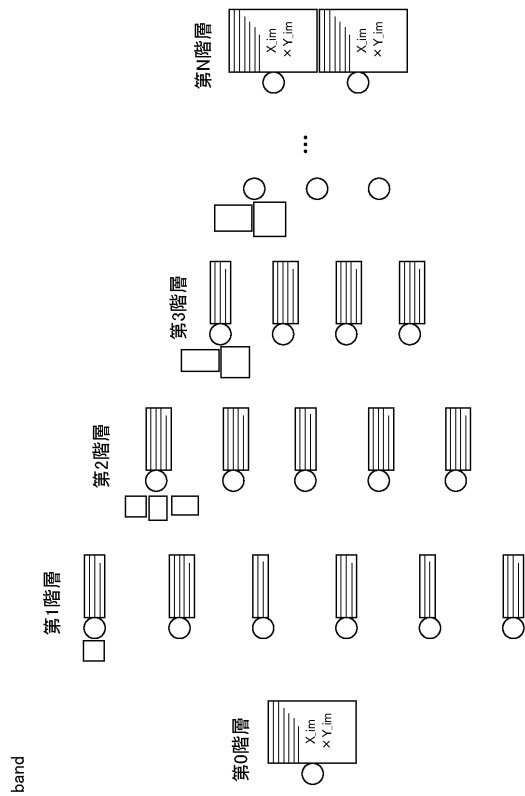
【図21】



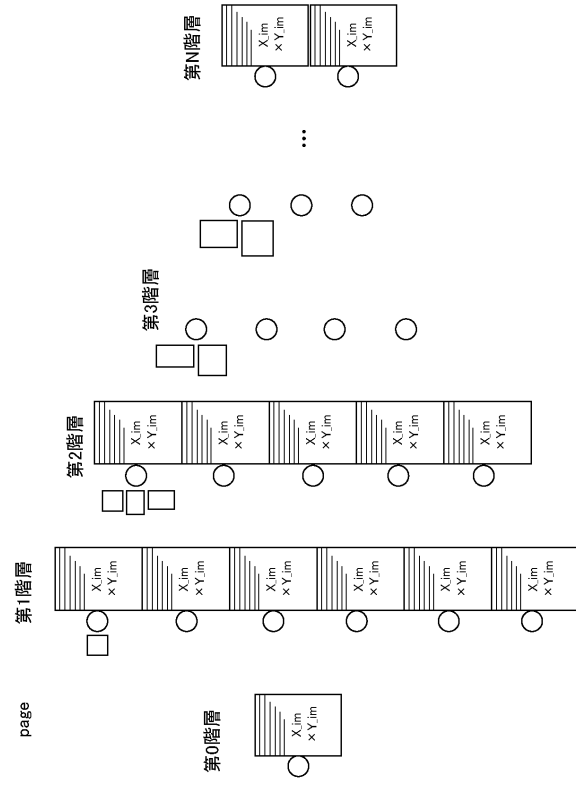
【図22】



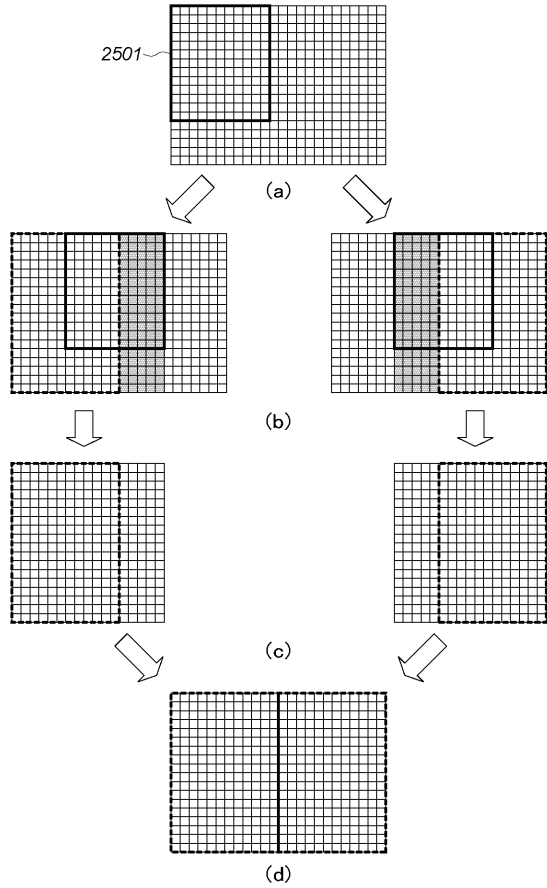
【図23】



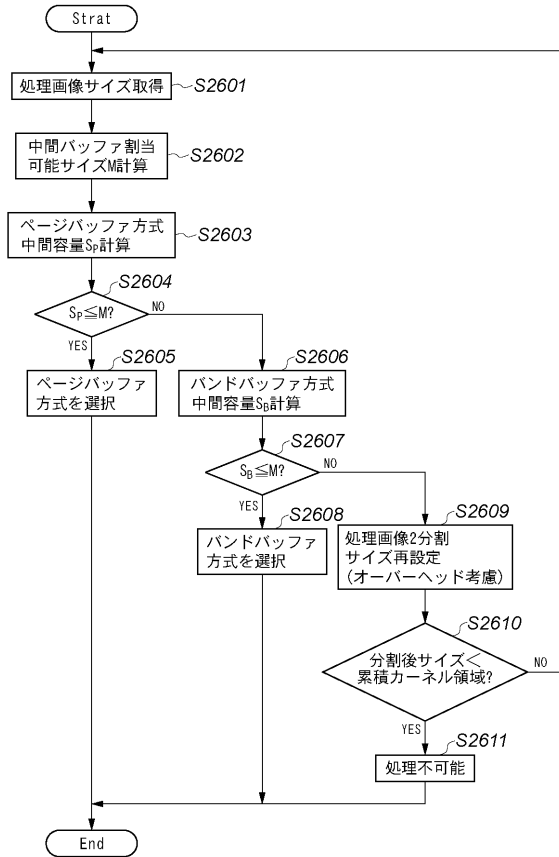
【図24】



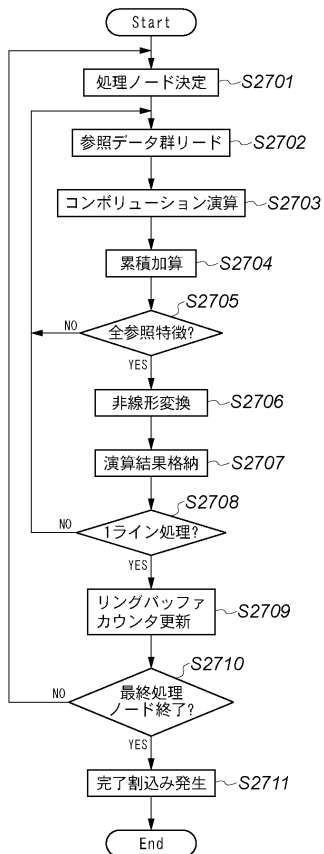
【図 25】



【図 26】



【図 27】



フロントページの続き

- (72)発明者 伊藤 嘉則
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 加藤 政美
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 山本 貴久
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

審査官 塚田 肇

- (56)参考文献 特開平05-108593(JP,A)
特許第2729987(JP,B2)

- (58)調査した分野(Int.Cl., DB名)
- | | |
|------|-------|
| G06N | 3/063 |
| G06N | 3/00 |
| G06T | 7/00 |