

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2019-3414

(P2019-3414A)

(43) 公開日 平成31年1月10日(2019.1.10)

| (51) Int.Cl. | | F I | テーマコード (参考) |
|--------------------|------------------|-------------|-------------|
| G06F 15/80 | (2006.01) | G06F 15/80 | 5B045 |
| G06F 17/10 | (2006.01) | G06F 17/10 | A 5B056 |
| G06N 3/063 | (2006.01) | G06N 3/063 | |
| G06F 15/173 | (2006.01) | G06F 15/173 | 683D |

審査請求 未請求 請求項の数 10 O L (全 19 頁)

| | | | |
|-----------|------------------------------|----------|--|
| (21) 出願番号 | 特願2017-117686 (P2017-117686) | (71) 出願人 | 506301140 公立大学法人会津大学 福島県会津若松市一箕町大字鶴賀字上居合 90番地 |
| (22) 出願日 | 平成29年6月15日(2017.6.15) | (74) 代理人 | 100094525 弁理士 土井 健二 |
| | | (74) 代理人 | 100094514 弁理士 林 恒徳 |
| | | (72) 発明者 | 富岡 洋一 福島県会津若松市一箕町大字鶴賀字上居合 90番地 公立大学法人会津大学内 |
| | | (72) 発明者 | スタニスラフ セドゥーキン 福島県会津若松市一箕町大字鶴賀字上居合 90番地 公立大学法人会津大学内 |
| | | Fターム(参考) | 5B045 GG13 5B056 AA05 BB26 FF01 FF02 FF05 |

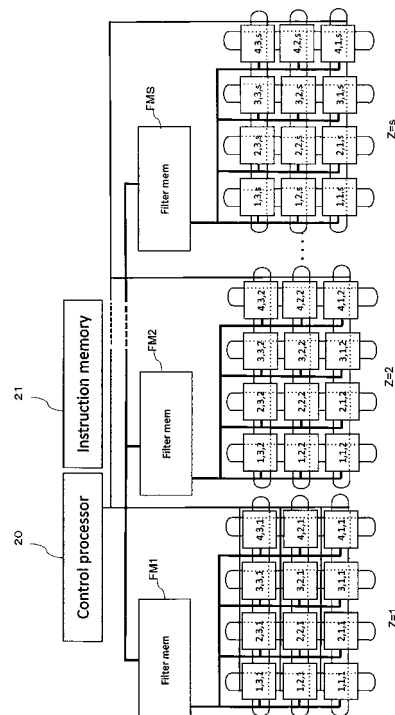
(54) 【発明の名称】 データ処理装置、及びこれにおけるデータ処理方法

(57) 【要約】

【課題】演算素子プロセッシングエレメントを三次元的に多数配置し、並列性を保ったまま省力で高速の計算を行えるデータ処理方法及びデータ処理装置を提供する。

【構成】乗算及び加算機能を有する複数のプロセッシングエレメントを3次元方向に有するデータ処理装置であって、それぞれ複数の前記プロセッシングエレメントが2次元方向に配置され、Z軸方向に積層された複数の2次元面を有し、前記複数の2次元面のそれぞれに対応して特徴重みが配置されるフィルタメモリを有し、入力データがZ軸方向の上位面の2次元面から配置され、一の面に配置されたプロセッシングエレメントで前記乗算機能により前記入力データと前記特徴重みの積を順次演算して2次元畳み込みデータを演算し、更に下面から転送されるデータと自身のデータを加算する演算を行い、当該演算結果を隣接する上面のプロセッシングエレメントに転送することを特徴とする。

【選択図】 図4



【特許請求の範囲】**【請求項 1】**

乗算及び加算機能を有する複数のプロセッシングエレメントを 3 次元方向に有するデータ処理装置であって、

それぞれ複数の前記プロセッシングエレメントが 2 次元方向に配置され、Z 軸方向に積層された複数の 2 次元面を有し、

前記複数の 2 次元面のそれぞれに対応して特徴重みが配置されるフィルタメモリを有し

、
入力データが Z 軸方向の上位面の 2 次元面から配置され、

一の面に配置されたプロセッシングエレメントで前記乗算機能により前記入力データと前記特徴重みの積を順次演算して 2 次元畳み込みデータを演算し、更に下面から転送されるデータと自身のデータを加算する演算を行い、当該演算結果を隣接する上面のプロセッシングエレメントに転送する、

ことを特徴とするデータ処理装置。

10

【請求項 2】

請求項 1 において、

前記 2 次元方向に配置されたプロセッシングエレメントはトーラスネットワークに接続され、Z 軸方向には、上下面に隣接するプロセッシングエレメントが、ネットワークで双方向に接続される、

ことを特徴とするデータ処理装置。

20

【請求項 3】

請求項 1 において、

前記 2 次元畳み込みデータは、隣接するプロセッシングエレメントからの転送データと自身のデータを加算演算し、更にシフト方向に隣接するプロセッシングエレメントに前記加算演算結果を転送する、

ことを特徴とするデータ処理装置。

【請求項 4】

請求項 1 において、

前記 Z 軸方向の最上位面にあるプロセッシングエレメントは、下面の複数のプロセッシングエレメントから転送される 2 次元畳み込みデータと自身のデータを加算して 2 . 5 次元畳み込みデータを演算する、

ことを特徴とするデータ処理装置。

30

【請求項 5】

請求項 4 において、

前記 2 . 5 次元畳み込みデータは、順次下面のプロセッシングエレメントにシフトされる、

ことを特徴とするデータ処理装置。

【請求項 6】

請求項 1 乃至 4 の何れか 1 項において、

前記特徴重みは、前記入力データの配置された 2 次元面の数で分割され、前記入力データの配置された面毎に対応するフィルタメモリに配置し、前記畳み込みの演算の際、前記フィルタメモリに配置された特徴重みを、対応する面の全てのプロセッシングエレメントにブロードキャストする、

ことを特徴とするデータ処理装置。

40

【請求項 7】

乗算及び加算機能を有する複数のプロセッシングエレメントを 3 次元方向に有するデータ処理装置におけるデータ処理方法であって、

前記データ処理装置は、それぞれ複数の前記プロセッシングエレメントが 2 次元方向に配置され、Z 軸方向に積層された複数の 2 次元面を有し、前記複数の 2 次元面のそれぞれに対応して特徴重みが配置されるフィルタメモリを有し、

50

入力データを前記 Z 軸方向の上位面の 2 次元面から配置する工程と、
 一の面に配置されたプロセッシングエレメントで前記乗算機能により前記入力データと前記特徴重みの積を順次演算して 2 次元畳み込みデータを演算する工程と、
 更に下面から転送されるデータと自身のデータを加算する演算を行い、当該演算結果を隣接する上面のプロセッシングエレメントに転送する工程を、
 有することを特徴とするデータ処理方法。

【請求項 8】

請求項 7 において、
 さらに、前記 2 次元畳み込みデータは、隣接するプロセッシングエレメントからの転送データと自身のデータを加算演算し、更にシフト方向に隣接するプロセッシングエレメントに前記加算演算結果を転送する工程を、
 有することを特徴とするデータ処理方法。

10

【請求項 9】

請求項 7 において、
 前記 Z 軸方向の最上位面にあるプロセッシングエレメントは、下面の複数のプロセッシングエレメントから転送される 2 次元畳み込みデータと自身のデータを加算して 2 . 5 次元畳み込みデータを演算する工程を、
 有することを特徴とするデータ処理方法。

【請求項 10】

請求項 9 において、
 前記 2 . 5 次元畳み込みデータを、順次下面のプロセッシングエレメントにシフトする工程を、
 有することを特徴とするデータ処理方法。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ処理装置、及びこれにおけるデータ処理方法に関し、特に、畳み込みニューラルネットワークにおける畳み込み演算に適したデータ処理装置、及びこれにおけるデータ処理方法に関する。

【背景技術】

30

【0002】

ニューラルネットワークに畳み込み(圧縮処理: Convolution)を追加した畳み込みニューラルネットワーク(Convolutional Neural Network: 以下適宜 CNN と表記)が、特に画像認識に有効な機械学習として広く認識されている。

【0003】

図 1 は、CNN のシステム構成の概略を示す図である。入力データに対して、複数の層(レイヤー L1 - L5)構造で処理を行う。

【0004】

図 1 では、レイヤー L1、L2 のそれぞれは、畳み込み層(Convolutional Layer)、プーリング層(Pooling Layer)を含みこれを繰り返す。

40

【0005】

畳み込み層は、入力データに対して フィルタ(kernel)特徴を乗算する(特徴量を畳み込む)層である。入力データが画像データである場合、入力データ(画像)に対して、それぞれ異なるフィルタ特徴を乗算してフィルタの数に対応する画像を得ている。複数のフィルタを使うことにより入力画像のさまざまな特徴が捉えられ、特徴量の畳み込みによって画像内のパターンが検出出来る。

【0006】

プーリング層は、畳み込み層の直後に置かれ、レイヤーを縮小して扱いやすくし、抽出された特徴の位置感度を低下させる。

【0007】

50

CNNは、次いで、レイヤーL3 - L5により全結合した多層パーセプトロンを配置して入力データ（画像）を認識する。

【0008】

ここで、序盤のレイヤーで行う畳み込み演算には膨大な計算回数が必要である。このためかかる部分の省電力化が非常に重要な課題となっている。

【0009】

しかし、複数の演算素子（PE：プロセッシングエレメント）を、アレイ状に配置するアレイ型の並列演算処理素子とすると、周辺機能ブロックとの類似度計算をするために多くの配線資源や転送時間が必要となる。

【0010】

かかる点に鑑みて、本発明者等は、先にPE間の通信でのデータ衝突を回避し、かつPEを特定の方向に偏ることなく増加させることが可能な拡張性の高いデータ処理装置を提案している(特許文献1)。

【0011】

かかる先の発明技術では、n次元のネットワークを構成するn次元の方向に配置された全てのPEが、転送クロックに同期してデータを入出力する。そして、データを入出力する方向であるシフト方向に隣接する第1の隣接PEから第1のデータを受け取るとともに、反対側に隣接する第2の隣接PEに第2のデータを出力し、隣接するPE間のデータ転送レートがシフト方向によらず等しいという特徴を有する演算装置である。

【先行技術文献】

【特許文献】

【0012】

【特許文献1】特許第5939572号公報

【発明の概要】

【発明が解決しようとする課題】

【0013】

これまでのCNN演算のための技術は、上記特許文献1に提案の発明に従う場合であっても、並列演算において個々のPEの処理量が大きくなるもの、即ち、技術的に最速であるがエネルギーに乏しいPEあるいはコアの数が、メモリに蓄積されるデータの数よりはるかに小さい。

【0014】

換言すれば、各演算ステップにおけるアクティブなプロセッサ、メモリ動作の数が基本的にCNNアルゴリズムにおける可能性より小さいものであった。

【0015】

結果として、CNNの解決のための時間が、最小値よりはるかに大きく、解決すべきエネルギーが高くなる。

【0016】

かかる点に鑑みて、本発明の目的は、演算素子PEを三次元的に多数配置し、並列性を保ったまま省電力で高速の計算を行えるデータ処理装置、及びこれにおけるデータ処理方法を提供することにある。

【課題を解決するための手段】

【0017】

上記目的を達成する本発明に従う第1の側面は、乗算及び加算機能を有する複数のプロセッシングエレメントを3次元方向に有するデータ処理装置であって、それぞれ複数の前記プロセッシングエレメントが2次元方向に配置され、Z軸方向に積層された複数の2次元面を有し、前記複数の2次元面のそれぞれに対応して特徴重みが配置されるフィルタメモリを有し、入力データがZ軸方向の上位面の2次元面から配置され、一の面に配置されたプロセッシングエレメントで前記乗算機能により前記入力データと前記特徴重みの積を順次演算して2次元畳み込みデータを演算し、更に下面から転送されるデータと自身のデータを加算する演算を行い、当該演算結果を隣接する上面のプロセッシングエレメントに

10

20

30

40

50

転送することを特徴とする。

【0018】

上記目的を達成する本発明に従う第1の側面において、第1の態様として、前記2次元方向に配置されたプロセッシングエレメントはトラスネットワークに接続され、Z軸方向には、上下面に隣接するPEが、ネットワークで双方向に接続されることを特徴とする。

【0019】

上記目的を達成する本発明に従う第1の側面において、第2の態様として、前記2次元畳み込みデータは、隣接するプロセッシングエレメントからの転送データと自身のデータを加算演算し、更にシフト方向に隣接するプロセッシングエレメントに前記加算演算結果を転送することを特徴とする。

10

【0020】

上記目的を達成する本発明に従う第1の側面において、第3の態様として、前記Z軸方向の最上位面にあるプロセッシングエレメントは、下面の複数のプロセッシングエレメントから転送される2次元畳み込みデータと自身のデータを加算して2.5次元畳み込みデータを演算することを特徴とする。

【0021】

上記目的を達成する本発明に従う第1の側面における第3の態様において、第4の態様として、前記2.5次元畳み込みデータは、順次下面のプロセッシングエレメントのシフトされることを特徴とする。

20

【0022】

上記目的を達成する本発明に従う第1の側面における上記何れかの態様において、前記特徴重みは、前記入力データの配置された2次元面の数で分割され、前記入力データの配置された面毎に対応するフィルタメモリに配置し、前記畳み込み演算の際、前記フィルタメモリに配置された特徴重みを、対応する面の全てのプロセッシングエレメントにブロードキャストすることを特徴とする。

【0023】

上記目的を達成する本発明に従う第2の側面は、乗算及び加算機能を有する複数のプロセッシングエレメントを3次元方向に有するデータ処理装置におけるデータ処理方法であって、前記データ処理装置は、それぞれ複数の前記プロセッシングエレメントが2次元方向に配置され、Z軸方向に積層された複数の2次元面を有し、前記複数の2次元面のそれぞれに対応して特徴重みが配置されるフィルタメモリを有し、入力データを前記Z軸方向の上位面の2次元面から配置する工程と、一の面に配置されたプロセッシングエレメントで前記乗算機能により前記入力データと前記特徴重みの積を順次演算して2次元畳み込みデータを演算する工程と、更に下面から転送されるデータと自身のデータを加算する演算を行い、当該演算結果を隣接する上面のプロセッシングエレメントに転送する工程を有することを特徴とする。

30

【0024】

上記目的を達成する本発明に従う第2の側面において、第1の態様として、前記2次元畳み込みデータは、隣接するプロセッシングエレメントからの転送データと自身のデータを加算演算し、更にシフト方向に隣接するプロセッシングエレメントに前記加算演算結果を転送する工程を有することを特徴とする。

40

【0025】

上記目的を達成する本発明に従う第2の側面において、第2の態様として、前記Z軸方向の最上位面にあるプロセッシングエレメントは、下面の複数のプロセッシングエレメントから転送される2次元畳み込みデータと自身のデータを加算して2.5次元畳み込みデータを演算する工程を有することを特徴とする。

【0026】

上記目的を達成する本発明に従う第2の側面において、第3の態様として、前記2.5次元畳み込みデータを、順次下面のプロセッシングエレメントにシフトする工程を有する

50

ことを特徴とする。

【発明の効果】

【0027】

上記本発明に従う特徴構成により、処理されるデータと同数のプロセッサで並列演算を行うことで、最小実行ステップ数でCNNの各層の計算を実行できる。このため、リアルタイム処理に求められる実行時間制約を達成できる最小の動作クロック周波数で実行可能であり、リアルタイムかつ低消費電力の計算を行える。

【図面の簡単な説明】

【0028】

【図1】CNNのシステム構成の概略を示す図である。

10

【図2】TAP(Tensor Array Processor)における3次元アレイ状のPME(Processor in Memory)を示す図である。

【図3】各PMEの機能構成例ブロック図である。

【図4】TAPの3次元アレイをそれぞれのXY面に展開して示す図である。

【図5】図4における一つの面に属するPMEを拡大して示す図である。

【図6】畳み込み演算による演算結果の変化を示す図である。

【図7A】本発明のデータ処理方法におけるある層の初期状態を表す図である。

【図7B】1番目のフィルタについて2次元畳み込みを行っている状態を表す図である。

【図7C】下の面から1番目の2次元畳み込み演算データが転送され、自身の畳み込み演算結果と加算することで2.5次元畳み込み演算の結果を求める処理を示す図である。

20

【図7D】1番目の2.5次元畳み込みの結果(白丸)を下の面にデータシフトした状態を表す図である。

【図7E】2番目のフィルタについて2次元畳み込みを行っている状態を表す図である。

【図7F】下の面から転送される2番目のフィルタに対する2次元畳み込み演算データと自身の畳み込み演算結果との加算により2.5次元畳み込み演算の結果を求める処理を示す図である。

【図7G】1番目と2番目の2.5次元畳み込みの結果(白丸)を下の面にデータシフトした状態を表す図である。

【図7H】3番目のフィルタについて2次元畳み込みを行っている状態を表す図である。

【図7I】下の面から転送される3番目のフィルタに対する2次元畳み込み演算データと自身の畳み込み演算結果と加算により2.5次元畳み込み演算の結果を求める処理を示す図である。

30

【図7J】1番目から3番目の2.5次元畳み込みの結果(白丸)を下の面にデータシフトした状態を表す図である。

【図7K】4番目のフィルタについて2次元畳み込みを行っている状態を表す図である。

【図7L】下の面から転送される4番目のフィルタに対する2次元畳み込み演算データと自身の畳み込み演算結果と加算により2.5次元畳み込み演算の結果を求める処理を示す図である。

【図7M】4枚のフィルタを用いた畳み込み層の計算結果を表す図である

【図8】畳み込み演算の例を具体的数値で説明する図である。

40

【図9】特徴重みの縦横を異なるものとした時のデータ上のシフト方向を考察する図である。

【図10】本発明に従う2.5次元畳み込み演算の様子を示すタイムチャート図である。

【発明を実施するための形態】

【0029】

以下に、本発明の実施例を添付の図面に従い説明する。これらの実施例は本発明の理解を容易とするためのものであり、本発明の適用は、これら実施例に限定されるものではない。また、本発明の保護の範囲は、特許請求の範囲と同一又は類似の範囲にも及ぶ。

【0030】

本発明に従うデータ処理装置は、3次元アレイ状に配置されたそれぞれメモリ機能を有

50

する演算素子であるPME(Processor in Memory)とネットワークから構成されるシステムであり、以降TAP(Tensor Array Processor)と称する。

【 0 0 3 1 】

図 2 は、かかるTAPにおける 3 次元アレイ状のPMEを示す図であり、それぞれ計算モジュールを有する複数のPMEが 3 次元 (X,Y,Z) 方向に積層配列されている。

【 0 0 3 2 】

かかる構造は、半導体技術により、3次元プロセッサとして作成可能である。すなわち、複数のPMEがX,Y方向に配列された2次元半導体面をZ方向に積み重ねて1チップで3次元構造とすることが可能である。

【 0 0 3 3 】

図 2 において、複数のPMEが、 $W_s \times H_s \times C_s$ の3次元アレイ状 (X軸方向に W_s 個、Y軸方向に H_s 個、Z軸方向に C_s 個) に配列されている。各XY面では、各PMEがトラスネットワーク (L_x 、 L_y) に接続され、各PMEの有するデータをX軸正方向、X軸負方向、Y軸正方向、Y軸負方向の4方向にデータシフトする機能を有する。

【 0 0 3 4 】

また、Z軸方向では、上下面に隣接するPMEが、ネットワーク L_z で双方向に接続され、Z軸正方向 (上面方向) とZ軸負方向 (下面方向) のデータ転送を同時に行うことが可能に構成されている。かかるデータシフトの方向制御及び、そのための共通シフトクロックは、後にデータ処理装置を展開図で示す制御プロセッサにより供給される。

【 0 0 3 5 】

図 3 は、各PMEの機能構成例ブロック図であり、畳み込み層の計算に必要な乗算及び加算機能ブロック 1 0 とプーリング層の計算に必要なMAX (最大値) 演算機能ブロック 1 1 を有している。さらに、必要に応じて、追加の機能ブロックを添えることは可能である。

【 0 0 3 6 】

図 4 は、上記TAPの3次元アレイをそれぞれのXY面に展開して示す図である。かかるTAPは、システムとして共通の制御プロセッサ 2 0 と、指示メモリ 2 1 を有している。

【 0 0 3 7 】

3次元アレイのそれぞれのXY面 ($Z=1, Z=2, \dots, Z=S$) に対応してフィルタメモリFM1-FMSを有し、フィルタメモリFM1-FMSのそれぞれは、対応する同一面に存在する全てのPMEと接続されている。計算の際に、フィルタメモリFM1-FMSからフィルタの特徴重みを対応する面の全てのPMEにブロードキャストすることが可能である。

【 0 0 3 8 】

ここで、本発明に従うデータ処理方法をCNNの畳み込み演算処理に用いる場合を想定する。

【 0 0 3 9 】

指示メモリ 2 1 には、事前の学習により得られた各層のフィルタサイズ及びフィルタ数が畳み込みニューラルネットワーク構造として入力される。これに基づき、制御プロセッサ 2 0 により、各面のフィルタメモリFM1-FMSに対応する重み、及び共通のクロック信号の供給等が行われる。

【 0 0 4 0 】

PMEは、畳み込み層において、自身の有するデータと対応するフィルタメモリからブロードキャストされる重みとを乗算し、更に隣接するPMEから転送されるデータとを加算して、その結果を反対側に隣接するPMEに転送する。プーリング層においては、自身のデータと隣接するPMEから転送されるデータの最大値を求め、隣接するPMEに順次転送する。

【 0 0 4 1 】

先に、説明した様に指示メモリ 2 1 からの予め学習によって得られた指示データに基づき、上記のPMEによる演算と転送の方向及び共通シフトのためのタイミングクロックが、制御プロセッサ 2 0 から全てのPMEに送られる。

【 0 0 4 2 】

図 5 は、かかる図 4 における一つの面に属するPMEを拡大して示す図である。この例で

10

20

30

40

50

は、 $W_s=4$, $H_s=3$ の場合で、 $Z=k$ の面を示している。ファイルメモリFM_kから共通に、 $Z=k$ の面にある全てのPMEにフィルタ（特徴重み）が供給される。

【0043】

かかる構成のTAPに、列数 W_0 × 行数 H_0 × チャンネル数 C_0 の3次元テンソルデータがCNNの第1層の入力となる。

【0044】

この入力データに繰り返しCNNの畳み込み層の計算を適用することで、各層のデータサイズが変化する。

【0045】

ここで、 k 層目のデータサイズを $W_k \times H_k \times C_k$ とする。本発明のシステムでは、

【0046】

【数1】

$$W_s \geq \max_k W_k, \quad H_s \geq \max_k H_k, \quad C_s \geq \max_k C_k$$

【0047】

であると想定する。

【0048】

入力のチャンネル数がこの値より小さい場合は、TAPの上位のXY面から順に入力データを配置する。入力データが配置されたPMEは活性面（active）、そうでなければ非活性面（inactive）となる。ただし、PMEが非活性面であってもデータの転送は行われる。

【0049】

例えば、入力データとして一枚の画像データを考えた時、次のように想定することが出来る。一枚の画像データを同じ大きさの領域ごとに区切り複数の領域データ（チャンネル）として切り出し、各領域データをTAPの最上位の面から順に該当の面にあるPMEに配置していく。

【0050】

この時、フィルタ（特徴重み）は、次のように処理される。一つの特徴重みを前記画像データの配置される面の数に対応して分割し、それぞれの分割特徴重みを対応する面のファイルメモリFMに格納する。そして、計算時に対応するフィルタメモリFMに配置されている特徴重みが当該面に属する全てのPMEにブロードキャストされる。

【0051】

それぞれのPMEは、ブロードキャストされた特徴重みと自身のデータとの積を演算する。さらに、PMEは、一方向の隣接するメモリ要素から転送されるデータを前記の積の演算結果に加え、反対方向に隣接するPMEに転送する。かかる処理を繰り返し、2次元畳み込みを行う。なお、かかる場合の転送制御は、先に述べた特許文献1の発明に従い実行される。

【0052】

さらに、本発明では、特徴として、入力データに対し2次元畳み込みを行ったデータが配置されたTAPの最下面から最上面まで、それぞれ2次元畳み込みデータを上方向に転送する。この時、各面のPMEは自身のデータと一つ下の面からの畳み込み演算結果を足し合わせ、その結果を一つ上の面のPMEに転送する。

【0053】

最終的に最上位面の2次元アレイプロセッサで、全ての2次元アレイの2次元畳み込み結果を足し合わせた2.5次元畳み込み演算結果を得ることが出来る。

【0054】

さらに、後に詳述するように、最上位面で得られた2.5次元畳み込み演算結果は、順次下面にシフトされる。

【0055】

かかる畳み込み演算による演算結果の変化を図6に示す。図6(1)に示すように、N

10

20

30

40

50

$\times M \times C_{in}$ 個のPMEに配置された入力データが、畳み込み演算の結果 $N \times M \times C_{out}$ 個のPMEに畳み込み演算結果が得られる。このとき C_{out} の大きさは、特徴重み(kernel)の数に依存する。

【0056】

この結果がCNNの一つの層の畳み込み演算処理結果のデータであり、次いで、図6(2)に示すようにプーリング層の演算処理を行ってレイヤーを縮小して扱いやすくする。同時に、この演算結果は次の層の入力になる。

【0057】

ここで、本発明のデータ処理装置において実行されるデータ処理に従う畳み込み演算処理の特徴を理解容易のために、更に図7A~7Mにおいて各面における変化を模式的に示す。

10

【0058】

図7A~7Mにおいて、活性面は入力データが配置された面である。図7Aは、初期状態を表す。実戦の直方体が活性面のPMEを表し、破線の直方体が付加性面のPMEを表す。以下、図7B~7Mにおいて同様である。さらに、灰色の丸で占められる表示は入力データ D_{in} を表している。図7Bは1番目のフィルタに対して各面において2次元畳み込みを行っている状態を示す。図の矢印は、各面にあるPMEに対するデータの転送方向を表し、黒丸は計算結果を示す。以下、図7C~7Mにおいて同様である。

【0059】

図7Cは、下の面からその2次元畳み込み演算データが転送され、自身の2次元畳み込み演算結果と加算することで2.5次元畳み込み演算の結果を求める処理を示している。

20

【0060】

図7Dは、2番目のフィルタに対する畳み込み計算を行う準備として、この2.5次元畳み込み演算結果を下の面にシフトする状態を示している。このシフトは、不活性面を含めて行われる。

【0061】

図7Eは2番目のフィルタに対して各面において2次元畳み込みを行っている状態を示す。図7Fは下の面から2番目のフィルタに対する2次元畳み込み演算データが転送され、自身の畳み込み演算結果と加算することで、2.5次元畳み込み演算の結果を求める処理を示している。

30

【0062】

同様に図7G~図7Lは3番目、4番目のフィルタに対する畳み込み計算の様子を示している。すなわち、図7Gは、1番目と2番目の2.5次元畳み込みの結果(白丸)を矢印のように上の面から下の面にデータシフトした状態を示している。

【0063】

図7Hは、3番目のフィルタについて2次元畳み込みを行っている状態を示している。図7Iは、下から3番目のフィルタに対する2次元畳み込み演算データが転送され、自身の畳み込み演算結果と加算することで2.5次元畳み込み演算の結果を求める処理を示している。

【0064】

図7Jは、1番目から3番目の2.5次元畳み込み結果(白丸)を上から下の面にデータシフトした状態を示している。この際、不活性面にもデータがシフトされている。

40

【0065】

図7Kは、4番目のフィルタに対する2次元畳み込み演算を行っている状態を示している。図7Lは、図7Iの処理と同様であるが、下から4番目のフィルタに対する2次元畳み込み演算データが転送され、自身の畳み込み演算結果と加算することで2.5次元畳み込み演算お結果を求める処理を示している。

【0066】

図7Mは、最終的に4枚のフィルタを用いた畳み込み層の計算結果を表し、これが次層の入力となる。

50

【 0 0 6 7 】

ここで、畳み込み演算を式で表すと下記(1)式のようになる。 b_0 はバイアス定数項である。バイアス b_0 は、畳み込み演算の結果を一定値増加、減少するために使用される。このバイアス b_0 とフィルタの重みはともに、CNNの学習時に自動的に決定される。

【 0 0 6 8 】

【 数 2 】

$$y_o^{(l)}(i, j) = b_o + \sum_{c=1}^{C^l} \sum_{\Delta i = -\lfloor \frac{\omega}{2} \rfloor}^{\lfloor \frac{\omega}{2} \rfloor} \sum_{\Delta j = -\lfloor \frac{\omega}{2} \rfloor}^{\lfloor \frac{\omega}{2} \rfloor} w_{o,c}^{(l)}(\Delta i, \Delta j) \cdot x_c^{(l)}(is - \Delta i, js - \Delta j), \quad \dots (1)$$

10

【 0 0 6 9 】

ただし、 s は自然数であり、畳み込み計算を行うときのストライドを表す。さらに、簡単化のため、本発明の説明ではストライドが1のときのみを説明しているが、ストライドが2以上であっても本発明の適用可能は、否定されない。

【 0 0 7 0 】

【 数 3 】

$$w_{o,c}^{(l)}(\Delta i, \Delta j)$$

【 0 0 7 1 】

は、第1層の o 番目のフィルタの重み、

【 0 0 7 2 】

【 数 4 】

$$x_c^{(l)}(is - \Delta i, js - \Delta j)$$

【 0 0 7 3 】

は、第1層の入力データである。

【 0 0 7 4 】

それぞれの面にあるPMEは、(1)式の後半部分

【 0 0 7 5 】

【 数 5 】

$$\sum_{\Delta i = -\lfloor \frac{\omega}{2} \rfloor}^{\lfloor \frac{\omega}{2} \rfloor} \sum_{\Delta j = -\lfloor \frac{\omega}{2} \rfloor}^{\lfloor \frac{\omega}{2} \rfloor} w_{o,c}^{(l)}(\Delta i, \Delta j) \cdot x_c^{(l)}(is - \Delta i, js - \Delta j),$$

30

【 0 0 7 6 】

の計算を行う。このとき、 $C=C^l$ の2次元畳み込み演算をTAPの一番上の面のPMEが計算しており、同様に $C=C^l - 1$ の2次元畳み込み演算をその一つ下の面のPMEが計算している。各面で計算した上記の後半部分の計算結果を足し合わせることで(1)式全体の計算をしている。

40

【 0 0 7 7 】

図8は、上記の畳み込み演算の例を具体的数値で説明する図であり、2つの上下面の場合を例にしている。

【 0 0 7 8 】

一の面(Ch1)で入力(Input) x_{60} と重み(kernel) w_{61} を矢印方向に移動しながら乗算し、同時に下の面(Ch2)で入力(Input) x_{62} と重み(kernel) w_{63} を矢印方向に移動しながら乗算する。これにより、それぞれ2次元畳み込み演算結果 64 が得られる。

50

【 0 0 7 9 】

一の面 (Ch1) の初期時点での入力 6 0 と重み 6 1 との乗算結果は、次のようであり、
 $(-3 * 1) + (-2 * 2) + (1 * 2) + (3 * 2) = 1$

次いで、1桁分矢印方向にシフトした時の入力 6 0 と重み 6 1 との乗算結果は、次のようである。

【 0 0 8 0 】

$(1 * 1) + (-3 * 2) + (-2 * 3) + (2 * 2) + (1 * 2) + (3 * 1) = -2$

これらは、2次元畳み込み演算結果 6 4 に示される通りである。

【 0 0 8 1 】

一方、下の面 (Ch2) の初期時点での入力 6 2 と重み 6 3 との乗算結果は、次のようであり。

【 0 0 8 2 】

$(-2 * 2) + (3 * 3) + (-3 * 1) + (1 * 3) = 5$

次いで、1桁分矢印方向にシフトした時の入力 6 0 と重み 6 1 との乗算結果は、次のようである。

【 0 0 8 3 】

$(2 * 2) + (-2 * 3) + (3 * 3) + (2 * 1) + (3 * 3) + (1 * 2) = 1$ である。

【 0 0 8 4 】

これらは、2次元畳み込み演算結果 6 5 に示される通りである。

【 0 0 8 5 】

ついで、前記一の面 (Ch1) では、自身の2次元畳み込み演算結果 6 4 を得て、更に下の面 (Ch2) から転送される2次元畳み込み演算結果 6 5 が転送される。したがって、それら2次元畳み込み演算結果 6 4 及び 6 5 とパイアス $b_0 = 1$ とを加算して2 . 5次元畳み込み演算結果 6 6 に示すように求める。

【 0 0 8 6 】

上記の様に、入力データ上で重みを順次所定桁数分ずつシフトして乗算及び加算を繰り返すことにより2次元畳み込み演算結果が得られる。

【 0 0 8 7 】

この際、指示メモリ 2 1 に格納されている指示に基づき、PMEからのデータ転送の方向がデータを一筆書きに転送し、無駄な転送をなくし、同じPMEに複数のデータ転送が行われないように制御され、データの衝突を回避することが出来る。

【 0 0 8 8 】

ここで、上記図 8 に示す例では特徴重みを縦横 3×3 、即ち縦横の長さが同じ $w \times w$ としているが、縦横の長さが異なる様に一般化することが出来、これを $w_1 \times w_2$ として表す。

【 0 0 8 9 】

図 9 は、特徴重みの縦横を異なるものとした時のデータ上のシフト方向を考察する図である。図 9 (1) は、縦横長さが同じ奇数で、中心に向かう様に一筆書きでシフトすることが出来る。図 9 (2) は、縦の長さ、横の長さのいずれか一方が偶数であり、図 9 (1) と同様に、全ての点をちょうど 1 回ずつ通ハミルトンパスが存在する。これに対し、図 9 (3) は、縦の長さも横の長さも奇数の場合で有り、1 個のPMEは 2 回通過することになるので、無駄な転送が発生する。

【 0 0 9 0 】

ここで、上記説明したように2 . 5次元の畳み込み演算の結果が得られるが、このデータはTAPの一番上のPMEが保有している。本発明に従うアルゴリズムでは、TAPの最下面を $Z=1$ 、最上面を $Z=C_s$ とすると、1層目で0番目のフィルタを用いたときの2 . 5次元畳み込み演算の結果

【 0 0 9 1 】

10

20

30

40

50

【数 6】

$$y_0^{(l)}(i, j)$$

を

【数 7】

$$Z = C_s - C^{l+1} + 0$$

【0092】

の面状のPMEに配置し、次の層の畳み込みの計算の準備に整える。このため、最上面で計算結果が得られる度に、各PMEが保有している2．5次元畳み込み演算結果を下面方向に1回シフトする(図7B参照)。

10

【0093】

図10は、更に本発明に従う2．5次元畳み込み演算の様子を示すタイムチャート図である。このタイムチャートでは、 $PME(i, j, 1)$, $PME(i, j, 2)$, ... $PME(i, j, C_s)$ の動作を表している。また、この図では C_{in} 個の面がactiveである。

【0094】

タイムチャートにおいて、各面の墨塗り部分Aで2次元畳み込み演算を行っている。この計算結果が終了した次のステップでその計算結果とbiasBを足し合わせて上の面のPMEにデータを転送する(上方向矢印)。

【0095】

次の面のPMEは、下の面のPMEから転送されたデータと自身の2次元畳み込み演算結果を足し合わせて、その結果Cを更に一つ上の面のPMEに転送することを繰り返す。最終的に一番上の面のPMEでの演算結果が、2．5次元畳み込みの演算結果Dとなる。

20

【0096】

次いで、この、2．5次元畳み込みの演算結果Dが、一つの重みについて2．5次元畳み込み演算が終わる都度、下向矢印の方向に下の面にシフトされる。この際、上の面からシフトされるデータはPMEでは、それを保存するだけで、その他の処理は行われない。2．5次元畳み込み演算を求めるために、一度だけシフトを行う。

【0097】

上記の動作を C_{out} 回繰り返し、一番上の面から C_{out} までに2．5次元畳み込み演算結果が保持される。

30

【0098】

ここで、CNNの各層において、上記したように2．5次元畳み込み処理が行われた後、プーリング(pooling)演算を行なって、次の層の入力データとされる。

【0099】

プーリング演算は、2次元畳み込み演算の時と同じ方法で周辺のPMEが持つデータを受け取り次の式(2)で最大値を計算する。

【0100】

【数 8】

$$x_0^{(l+1)}(i, j) = \max_{\substack{-\lfloor \frac{\omega'}{2} \rfloor \leq \Delta i \leq \lfloor \frac{\omega'}{2} \rfloor \\ -\lfloor \frac{\omega'}{2} \rfloor \leq \Delta j \leq \lfloor \frac{\omega'}{2} \rfloor}} \left\{ \text{act} \left(y_0^{(l)}(is' + \Delta i, js' + \Delta j) \right) \right\} \quad \dots (2)$$

40

【0101】

ただし、 s' は2以上の自然数であり、畳み込み計算を行うときのストライドを表す。さらに、actは活性化関数であり、例えば

【0102】

【数 9】

$$act(\cdot) = \max(0, \cdot)$$

【0103】

が用いられる。

【0104】

また、プーリングでは、2次元畳み込みと同様に上面側にデータ転送を行うが、各PMEは、自身の有するデータ x と隣接した下面のPMEから受け取った y_{in} を用いて

【0105】

【数10】

$$y_{out} = \max(y_{in}, x) \dots (3)$$

【0106】

を計算する。

【0107】

この計算を行いながら先に、図9で説明した様にデータを転送することにより周辺のPMEの持つデータの最大値を求める。

【0108】

以上説明したように、本発明に従うデータ処理装置は、CNNにおけるデータ処理装置として使用される場合は、CNNの構造（学習によって得られた各層のフィルタサイズ、フィルタ数）が入力として与えられる。各フィルタの重みがTAPの各面のフィルタメモリ上に与えられる。さらに、CNNの主入力データがTAPのPMEに配置される。

【0109】

各面のPMEは、自身のデータと重みを乗算して2次元畳み込みデータを、周辺から転送されるデータとを加算して上位の面上のPMEに送る。したがって、最上位の面にあるPMEで、全ての下面の2次元畳み込み演算結果を加算することにより並列性を保ったまま省力で高速の2.5次元畳み込み演算結果を得ることが出来る。

【符号の説明】

【0110】

PME メモリ要素

10 乗算及び加算機能ブロック

11 MAX（最大値）演算機能ブロック

20 制御プロセッサ

21 指示メモリ

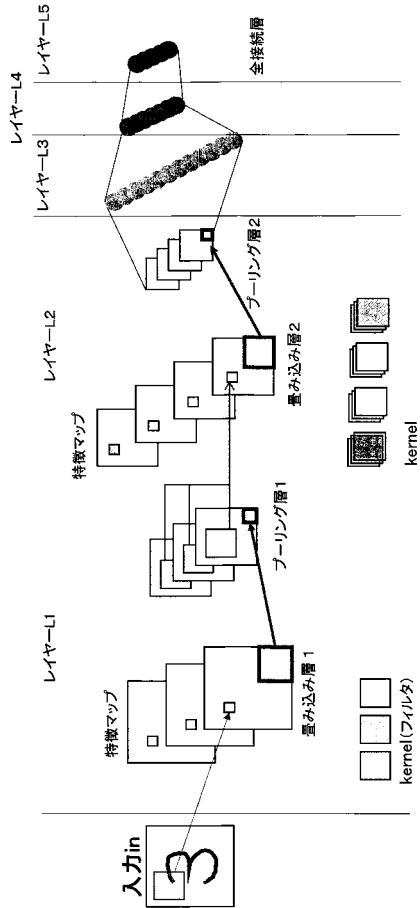
FM1-FMS フィルタメモリ

10

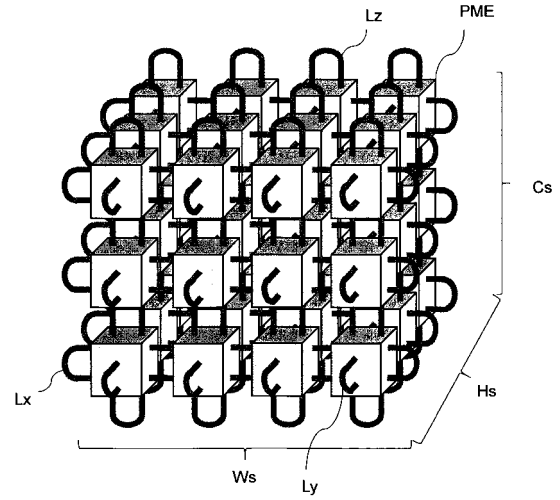
20

30

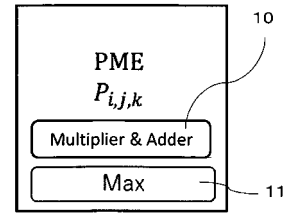
【図1】



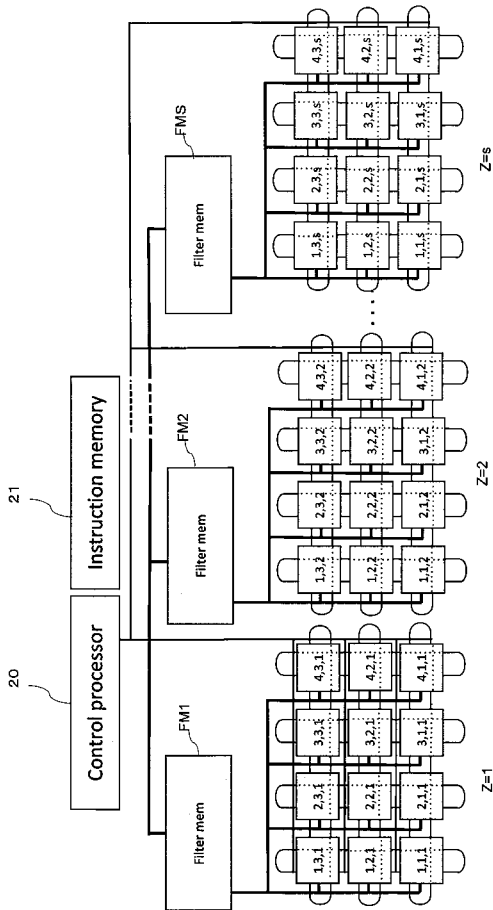
【図2】



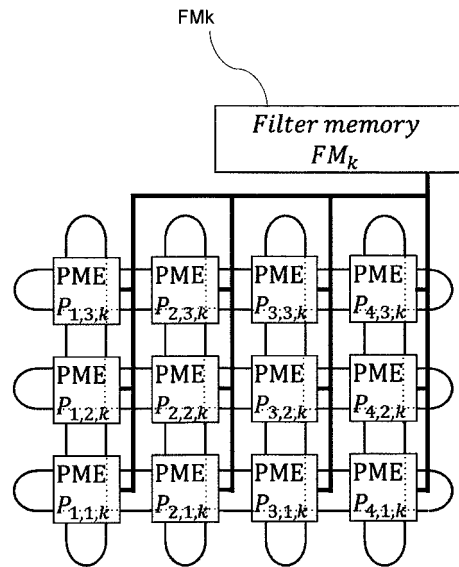
【図3】



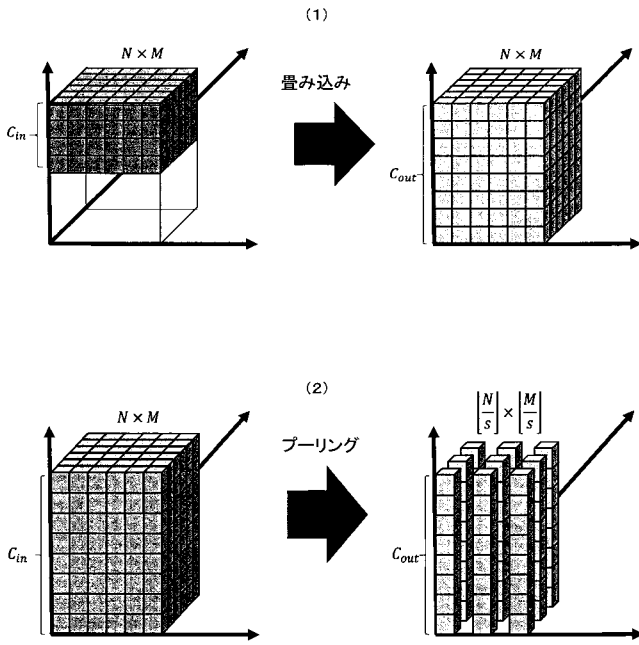
【図4】



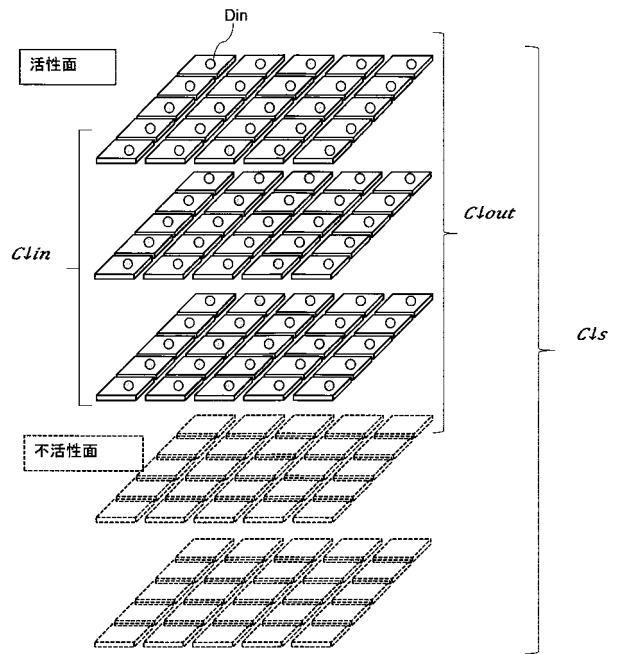
【図5】



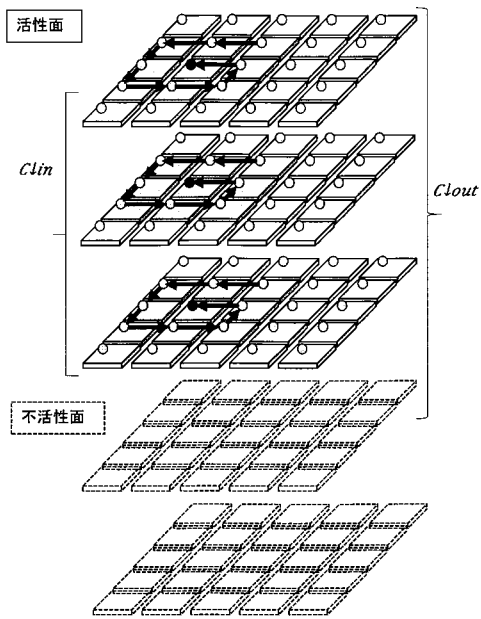
【 図 6 】



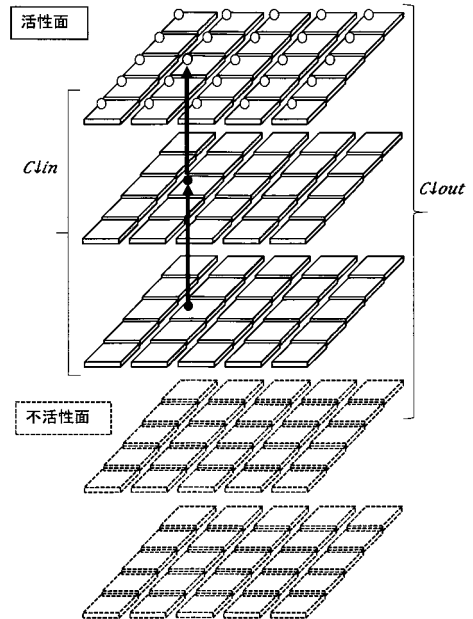
【 図 7 A 】



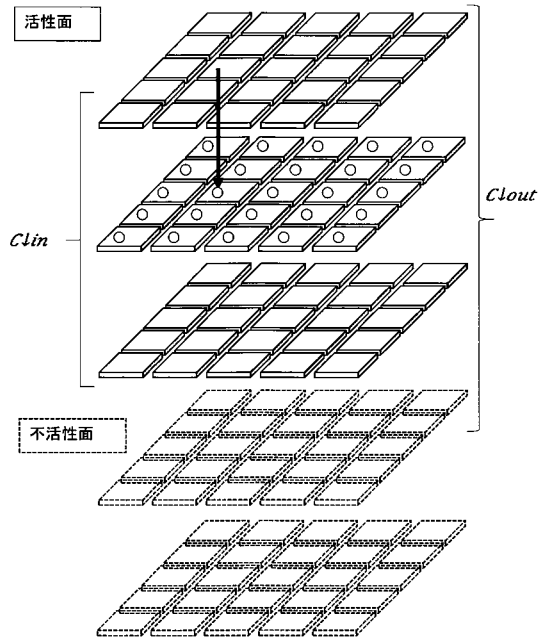
【 図 7 B 】



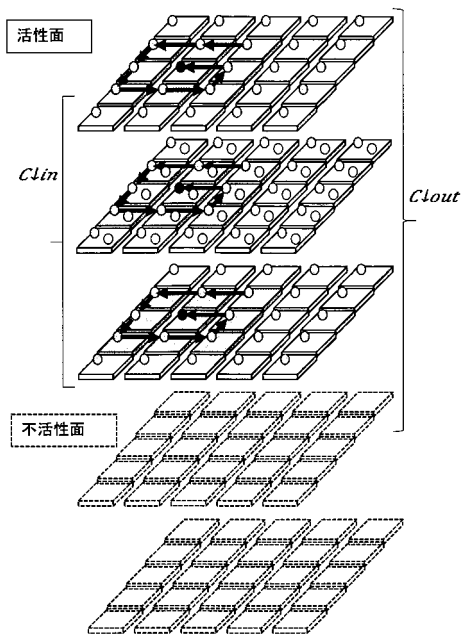
【 図 7 C 】



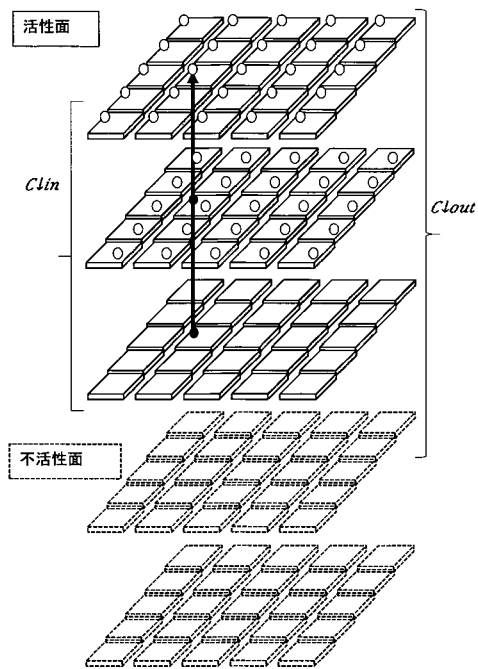
【 図 7 D 】



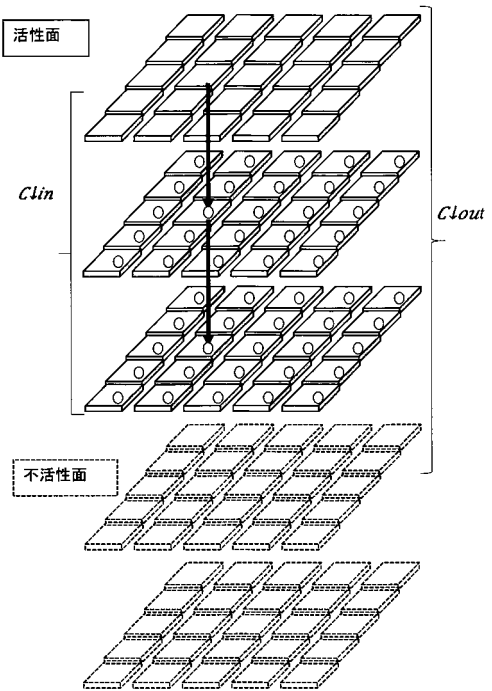
【 図 7 E 】



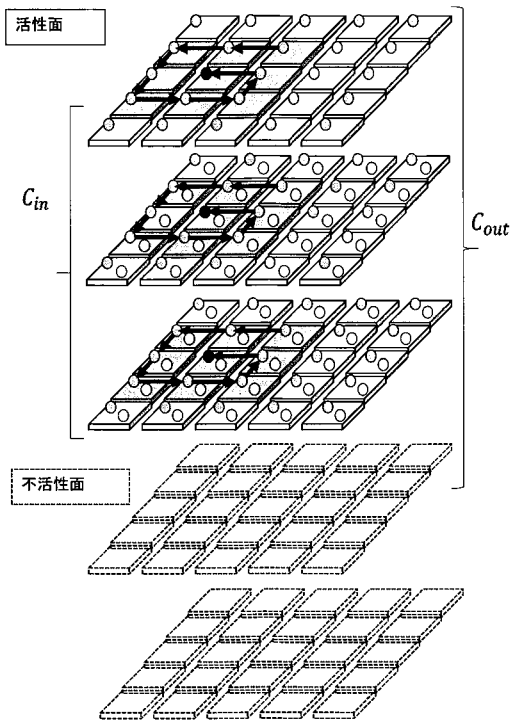
【 図 7 F 】



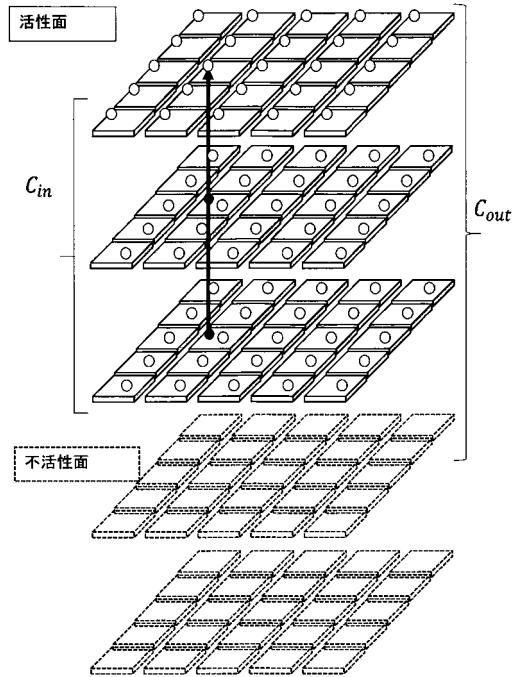
【 図 7 G 】



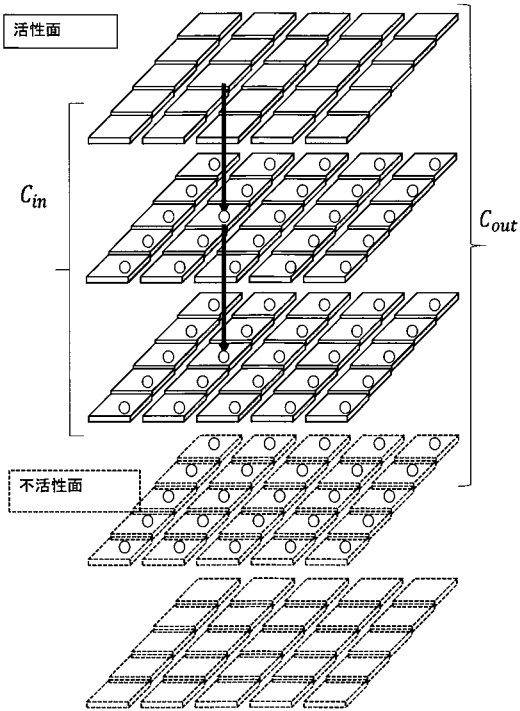
【 図 7 H 】



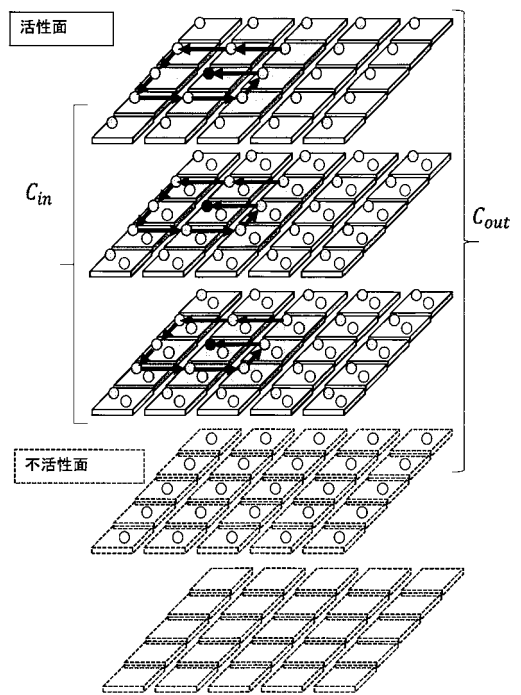
【 図 7 I 】



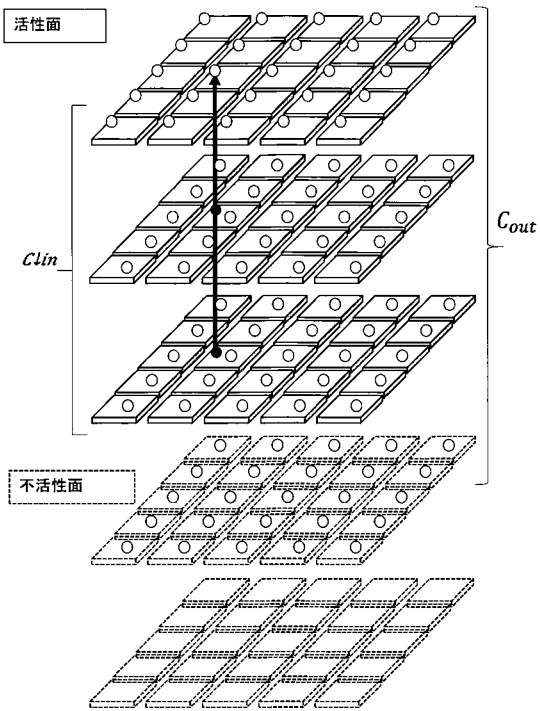
【 図 7 J 】



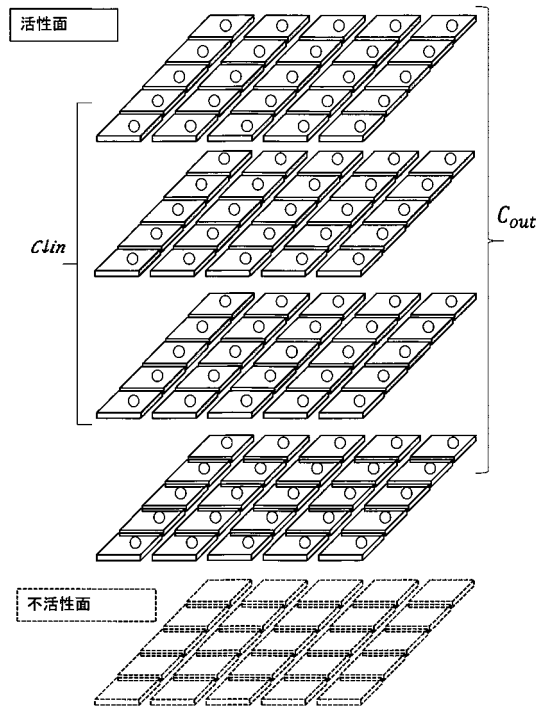
【 図 7 K 】



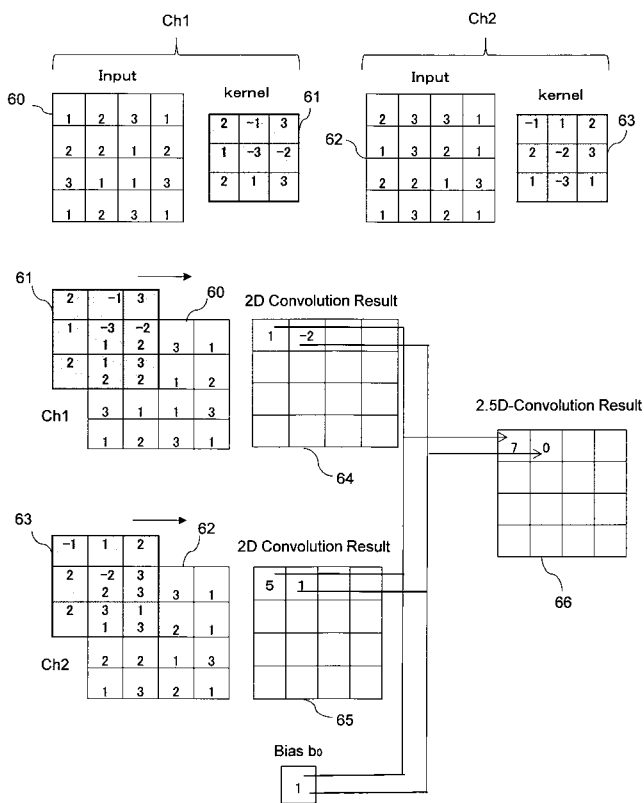
【 図 7 L 】



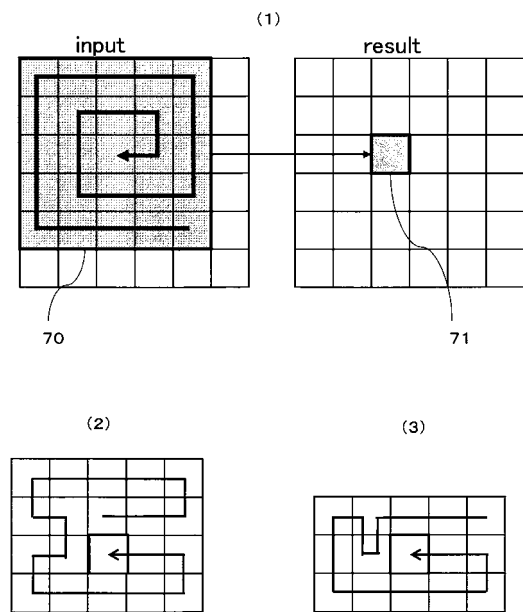
【 図 7 M 】



【 図 8 】



【 図 9 】



【 図 10 】

