

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7374236号

(P7374236)

(45)発行日 令和5年11月6日(2023.11.6)

(24)登録日 令和5年10月26日(2023.10.26)

(51)国際特許分類

F I

G 0 6 F 17/16 (2006.01)

G 0 6 F 17/16

M

G 0 6 N 3/063(2023.01)

G 0 6 N 3/063

請求項の数 20 外国語出願 (全17頁)

(21)出願番号	特願2022-2202(P2022-2202)	(73)特許権者	510192916
(22)出願日	令和4年1月11日(2022.1.11)		テスラ, インコーポレイテッド
(62)分割の表示	特願2020-503780(P2020-503780)		アメリカ合衆国 テキサス州 7 8 7 2 5
)の分割		, オースティン, テスラ ロード 1
原出願日	平成30年7月19日(2018.7.19)	(74)代理人	110000659
(65)公開番号	特開2022-64892(P2022-64892A)		弁理士法人広江アソシエイト特許事務所
(43)公開日	令和4年4月26日(2022.4.26)	(72)発明者	パノン, ピーター ジョセフ
審査請求日	令和4年2月8日(2022.2.8)		アメリカ合衆国 カリフォルニア州 9 4
(31)優先権主張番号	62/536,399		0 6 2, ウッドサイド, グレンウッド
(32)優先日	平成29年7月24日(2017.7.24)		アベニュー 1 2 8
(33)優先権主張国・地域又は機関	米国(US)	(72)発明者	ハード, ケヴィン アルテア
(31)優先権主張番号	15/710,433		アメリカ合衆国 カリフォルニア州 9 4
(32)優先日	平成29年9月20日(2017.9.20)		0 6 1, レッドウッド シティ, セント
(33)優先権主張国・地域又は機関			フランシス ストリート 1 1 2 5
	最終頁に続く	(72)発明者	タルベス, エミル
			最終頁に続く

(54)【発明の名称】 加速数学エンジン

(57)【特許請求の範囲】

【請求項 1】

センサデータを受信するように構成された第1入力回路と、
複数のフィルタのうちの1つ以上のフィルタを受信するように構成された第2入力回路と、

前記センサデータおよび前記フィルタを受信するように構成された複数のサブ回路であって、各サブ回路は、算術論理ユニットを備え、前記サブ回路は、前記センサデータおよび前記フィルタを畳み込みするように構成された、複数のサブ回路と、を備え、

前記センサデータおよび前記フィルタを畳み込むために、前記サブ回路は、前記サブ回路を介して、前記1つ以上のフィルタで前記センサデータの個々のサブセットを順に畳み込みするように構成され、残りのフィルタのうちの1つ以上が、続いて、畳み込みのために受信され、前記サブ回路の各々が、前記畳み込みに関連する各出力ピクセルを保存する
行列プロセッサ。

【請求項 2】

前記センサデータは、画像データ、ライダデータ、超音波データまたはレーダデータを含む、請求項1に記載の行列プロセッサ。

【請求項 3】

前記受信されたセンサデータは、線形化されたセンサデータを表す再フォーマットされたオペランドを含む、請求項1に記載の行列プロセッサ。

【請求項 4】

10

20

前記サブ回路のうちの 1 つ以上は、エンコーダをさらに備える、請求項 1 に記載の行列プロセッサ。

【請求項 5】

前記サブ回路の少なくとも一部は、特定のエンコーダを共有し、前記特定のエンコーダはブースエンコーダである、請求項 4 に記載の行列プロセッサ。

【請求項 6】

前記行列プロセッサは、冗長データを識別するように構成された状態機械を実装する、請求項 1 に記載の行列プロセッサ。

【請求項 7】

冗長データの識別は、前記複数のフィルタのうちの個々のフィルタおよび / または 1 つ以上のストライドのうちの個々のストライドに関連付けられた各サイズを含む入力に基づいている、請求項 6 に記載の行列プロセッサ。

10

【請求項 8】

第 1 のサブセットの畳み込みは、前記フィルタの 1 つ以上での前記第 1 のサブセットの畳み込みを決定することを含み、1 つ以上の残りのサブセットは、残りのフィルタの 1 つ以上で順に畳み込みされる、請求項 1 に記載の行列プロセッサ。

【請求項 9】

前記行列プロセッサは、タイルのアレイを備え、前記タイルは、前記サブ回路の各サブセットを備える、請求項 1 に記載の行列プロセッサ。

【請求項 10】

20

出力アレイの各列は、前記複数のフィルタのうちの各フィルタで畳み込みされた前記センサデータの個々のサブセットを含む、請求項 1 に記載の行列プロセッサ。

【請求項 11】

センサデータをフォーマットし、前記フォーマットされたセンサデータを行列プロセッサに提供するように構成された第 1 論理回路と、

複数のフィルタのうちの 1 つ以上のフィルタを前記行列プロセッサに提供するように構成された第 2 論理回路と、

複数のサブ回路を備える行列プロセッサであって、前記サブ回路は、前記センサデータの個々のサブセットを、前記 1 つ以上のフィルタで順に畳み込みするように構成され、残りのフィルタの 1 つ以上が、続いて畳み込みのために受信され、前記サブ回路の各々が、前記畳み込みに関連する各出力ピクセルを保存する、行列プロセッサと、
を備える、システム。

30

【請求項 12】

前記センサデータの個々のサブセットを順に畳み込みすることは、前記 1 つ以上のフィルタでの前記センサデータの第 1 のサブセットの畳み込みを決定することを含み、前記複数のフィルタのうちの残りのフィルタの 1 つ以上は、前記第 1 のサブセットでの畳み込みのために受信される、請求項 11 に記載のシステム。

【請求項 13】

前記システムは出力アレイを備え、前記出力アレイの各列は、前記複数のフィルタのうちの各フィルタで畳み込みされた前記センサデータの個々のサブセットを含む、請求項 11 に記載のシステム。

40

【請求項 14】

前記第 1 論理回路は、前記センサデータの一部を格納する複数のデータレジスタを備え、前記複数のデータレジスタは、センサデータから取得された入力領域のサイズに対応する第 1 幅を有する、請求項 11 に記載のシステム。

【請求項 15】

前記入力領域は前記センサデータの個々のサブセットに対応する、請求項 14 に記載のシステム。

【請求項 16】

前記システムは、センサデータを複数のベクトルに線形化するように構成され、各ベク

50

トルは前記センサデータの各サブセットを表す、請求項 11 に記載のシステム。

【請求項 17】

行列プロセッサによって実行される方法であって、

第 1 論理回路から、複数のサブセットを含むセンサデータを受信するステップと、

第 2 論理回路から、複数のフィルタのうちの 1 つ以上のフィルタを受信するステップと、

前記行列プロセッサの複数のサブ回路を使用して、前記 1 つ以上のフィルタで、個々のサブセットを順に畳み込みするステップであって、残りのフィルタのうちの 1 つ以上が、前記個々のサブセットの順の畳み込みのために、続いて受信され、前記サブ回路の各々が、前記畳み込みに関連する各出力ピクセルを保存する、ステップと、

を含む、方法。

10

【請求項 18】

前記 1 つ以上のフィルタで第 1 のサブセットを畳み込みすることに続いて、第 2 のサブセットが前記 1 つ以上のフィルタで畳み込みされる、請求項 17 に記載の方法。

【請求項 19】

前記行列プロセッサは、前記複数のフィルタのうちの 1 つ以上のフィルタの各々を順に受信する、請求項 17 に記載の方法。

【請求項 20】

出力アレイの各列に、前記列に関連付けられたフィルタで畳み込みされた前記センサデータの個々のサブセットを出力するステップをさらに含む、請求項 17 に記載の方法。

【発明の詳細な説明】

20

【技術分野】

【0001】

[関連出願の相互参照]

本出願は、2017 年 7 月 24 日に提出された「Accelerated Mathematical Engine」と題された米国仮出願第 62 / 536399 (20150 - 2154P (P0822 - 1PUS)) の優先権を主張し、発明者として Peter Joseph Bannon、Kevin Altair Hurd、および Emil Talpes を掲載する。前述の特許文書は、その全体が参照により本明細書に組み込まれる。

【0002】

30

また、本出願は、2017 年 9 月 20 日に提出された「Accelerated Mathematical Engine」と題され、発明者として Peter Joseph Bannon、Kevin Altair Hurd、および Emil Talpes を掲載した共通所有米国特許出願第 15 / 710433 (20150 - 2154 (N0822 - 1NUS)) に優先権を主張する。前述の特許文書のそれぞれは、その全体が参照により本明細書に組み込まれる。

【0003】

本開示は、大量のデータで動作する加速数学エンジンに関し、より具体的には、行列乗算演算に基づいて複雑な畳み込み演算を実行する加速数学エンジンに関する。

【背景技術】

40

【0004】

当業者は、時間に敏感で複雑な数学的演算を実施するために使用される一般的なプロセッサおよびシステムに対する速度および性能のますます高まる要求を認識するであろう。これらの一般的なシステムは、大量のデータを処理し、複雑な数学的演算を実行するために使用されるため、計算リソースと計算速度は、これらの計算を実行する既存の一般的なハードウェア設計の機能によって制限される。たとえば、行列演算を実行する汎用コンピュータ装置およびプロセッサは、特定の状況下でこれらの演算をタイムリーに実行できない場合がある。デジタル信号処理操作を実行する多くの従来の乗算器は、一連のソフトウェアおよびハードウェア行列操作手順（アドレス生成、転置、ビットごとの加算およびシフトなど）に依存しており、時間に敏感なシステム内においてボトルネックになり得る。

50

多くの場合、これらの操作ステップでは、プロセッサの算術関数を使用して中間結果を生成する必要があるが、さまざまな場所から中間結果を保存およびフェッチして操作を完了するため、計算時間が無駄になる。

【 0 0 0 5 】

図 1 は、従来の乗算器システムの一例を示している。乗算器システム 1 0 0 は、計算ユニット 1 0 2、レジスタ 1 0 4、キャッシュ 1 0 6、およびメモリ 1 0 8 を含むスカラマシンである。動作中、計算ユニット 1 0 2 は、レジスタ 1 0 4 およびキャッシュ 1 0 6 を使用して、メモリ 1 0 8 に格納されたデータを検索する。典型的には、計算ユニット 1 0 2 は、例えば乗算を加算に変換し、結果を何らかの内部レジスタに出力することにより、結果の行列を得るために入力行列で行列乗算を含む様々な計算手順を実行できる CPU または GPU などのマイクロプロセッサである。

10

【 0 0 0 6 】

たとえば、画像の出力ピクセルを表す内積は、通常、2 つの行列の個々の行列要素をドット乗算して部分的な結果を取得し、最終結果を加算して最終的な内積を取得することによって生成される。個々の行列要素の乗算、つまりスカラ乗算は、通常、ドット乗算を一連の個々のサブ操作に分割することにより、個々のデータ要素に対して実行される。その結果、単一の算術演算を完了するために、レジスタ 1 0 4、キャッシュ 1 0 6、およびメモリ 1 0 8 の 1 つまたは複数に部分積を格納およびフェッチする必要がある。

【 0 0 0 7 】

畳み込みなどの計算が要求されるアプリケーションでは、多くの場合、計算ユニット 1 0 2 にソフトウェア関数を組み込み、畳み込み演算を代替の行列乗算演算に変換するために使用する必要がある。これは、データを 2 つの行列に並べ替えて再フォーマットし、それをそのまま行列乗算できるようにすることで実現される。しかし、各スカラ演算を実行するために必要なデータをレジスタから何度も再格納および再フェッチする必要があるため、スカラマシン 1 0 0 でデータを効率的に共有または再利用するメカニズムは存在しない。これらの操作の複雑さと管理オーバーヘッドは、畳み込み演算の対象となる画像データの量が増加するにつれて著しく大きくなる。

20

【 0 0 0 8 】

算術演算を完了するためにレジスタ 1 0 4、キャッシュ 1 0 6、およびメモリ 1 0 8 から中間結果を格納およびフェッチする追加された非効率的なステップと相まって、スカラマシン 1 0 0 内のデータの多くを再利用できないことは、乗算器システム 1 0 0 などの既存のシステムの欠点の一部にすぎない。

30

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 9 】

したがって、必要なのは、行列数学演算を迅速かつ効率的に実行できる高計算スループットのシステムと方法である。

【 図面の簡単な説明 】

【 0 0 1 0 】

本発明の実施形態を参照し、その例を添付の図面に示すことができる。これらの図は、限定ではなく例示を目的としている。本発明をこれらの実施形態の文脈で一般的に説明するが、本発明の範囲をこれらの特定の実施形態に限定することを意図するものではないことを理解されたい。図の項目は縮尺どおりではない場合がある。

40

【 0 0 1 1 】

【 図 1 】 従来の乗算器システムの例を示す。

【 0 0 1 2 】

【 図 2 】 本開示の様々な実施形態に従って算術演算を実行するための例示的な行列プロセッサアーキテクチャを示す。

【 0 0 1 3 】

【 図 3 】 図 2 に示される行列プロセッサアーキテクチャの例示的な構成の詳細を示す。

50

【 0 0 1 4 】

【図 4】図 3 に示される論理回路の例示的な乗加算回路の実装を示す。

【 0 0 1 5 】

【図 5】本開示の様々な実施形態による例示的な畳み込み演算を示す。

【 0 0 1 6 】

【図 6】本開示の様々な実施形態による例示的な畳み込み演算の詳細を示す。

【図 7】本開示の様々な実施形態による例示的な畳み込み演算の詳細を示す。

【図 8】本開示の様々な実施形態による例示的な畳み込み演算の詳細を示す。

【 0 0 1 7 】

【図 9】本開示の様々な実施形態による例示的な逆畳み込み演算を示す。

10

【 0 0 1 8 】

【図 10】本開示の様々な実施形態による、畳み込みニューラルネットワークをより高速にするために算術演算を実行するプロセスを示す図である。

【発明を実施するための形態】

【 0 0 1 9 】

以下の説明では、説明の目的で、本発明の理解を提供するために特定の詳細が述べられている。ただし、これらの詳細なしで本発明を実施できることは当業者には明らかであろう。さらに、当業者は、以下で説明する本発明の実施形態が、プロセス、装置、システム、デバイス、または有形のコンピュータ可読媒体上の方法などの様々な方法で実装され得ることを認識するであろう。

20

【 0 0 2 0 】

図に示される構成またはモジュールは、本発明の例示的な実施形態の例示であり、本発明を曖昧にすることを避けることを意図している。また、この説明全体を通して、構成はサブユニットを含み得る別個の機能ユニットとして説明されるが、当業者は、さまざまな構成またはその一部が別個の構成に分割されてもよく、または単一のシステムまたは構成内に統合されることを含め、一緒に統合されてもよいことがわかるであろう。本明細書で説明する機能または動作は、構成として実装できることに留意されたい。構成は、ソフトウェア、ハードウェア、またはそれらの組み合わせで実装できる。多くの構成は、多くのサブ構成の相互接続を通じて形成される。これらの論理的に異なるサブ構成を集約して他のサブ構成と組み合わせることができ、ここで説明したものと同様のまたは同一の機能を集約構成レベルで提供できる場合は、ここで示されているものと動作が論理的に異なるサブ構成を選択できる（たとえば、アクティブな *h i g h* 信号をアクティブな *l o w* 信号に、AND ゲートを反転入力 NOR ゲートなどに置き換えることができる、など）。

30

【 0 0 2 1 】

さらに、図内の構成またはシステム間の接続は、直接接続に限定されるものではない。むしろ、これらの構成間のデータは、中間構成によって変更、再フォーマット、または変更され得る。また、追加の接続またはより少ない接続を使用してもよい。また、「結合」、「接続」、または「通信結合」という用語には、直接接続、1 つまたは複数の中間デバイスを介した間接接続、および無線接続が含まれることを理解されたい。

【 0 0 2 2 】

40

本明細書における「一実施形態」、「好ましい実施形態」、「実施形態」、または「いくつかの実施形態」への言及は、実施形態に関連して説明される特定の特徴、構造、特性、または機能が本発明の少なくとも 1 つの実施形態に含まれ、2 つ以上の実施形態に含まれてもよいことを意味する。また、本明細書の様々な場所での上記のフレーズの出現は、必ずしもすべてが同じ実施形態または複数の実施形態を指しているわけではない。

【 0 0 2 3 】

本明細書のさまざまな場所での特定の用語の使用は、説明のためであり、限定するものとして解釈されるべきではない。サービス、機能、またはリソースは、単一のサービス、機能、またはリソースに限定されず、これらの用語の使用は、分散または集約され得る関連するサービス、機能、またはリソースのグループ化を指す。

50

【 0 0 2 4 】

「含む」、「含んでいる」、「備え」、「備えている」という用語は、限定されていない (open) 用語であると理解され、以下の如何なるリストも例であり、リストされた項目に限定されることを意図するものではなく、追加項目とともに、部分集合または上位集合の項目を含み得る。ここで使用されているあらゆる見出しは、組織的な目的のみに使用され、説明または如何なる請求項の範囲をも制限するために使用されることはない。この特許文献に記載されている各文献は、その全体が参照により本明細書に組み込まれる。

【 0 0 2 5 】

さらに、当業者は、(1) 特定のステップをオプションで実行できること、(2) ステップは、ここに記載されている特定の順序に限定されない場合があること、(3) 特定のステップが異なる順序で実行される場合があること、(4) 特定のステップを同時に実行できること、を認識しなければならない。

【 0 0 2 6 】

本明細書の実施形態は主に畳み込みの文脈で説明されているが、当業者は、逆畳み込みおよび他の行列演算も行列 - 行列型乗算演算として構成できること、したがって本発明の原理も同様に逆畳み込みに適用できることを理解するであろう。さらに、本開示の様々な実施形態に従って、他のタイプの数学的演算を実施してもよい。

【 0 0 2 7 】

図 2 は、本開示の様々な実施形態に従って算術演算を実行するための例示的な行列プロセッサアーキテクチャを示している。システム 200 は、論理回路 232、234、キャッシュ/バッファ 224、データフォーマッタ 210、重みフォーマッタ 212、データ入力行列 206、重み入力行列 208、行列プロセッサ 240、出力アレイ 226、後処理ユニット 228、および制御ロジック 250 を含む。行列プロセッサ 240 は、算術論理ユニット (Arithmetic Logic Units: ALU)、レジスタ、および、いくつかの実施形態ではエンコーダ (ブースエンコーダなど) を含む複数のサブ回路 242 を含む。論理回路 232 は、N 個の入力演算子およびデータレジスタを表す回路であり得る。論理回路 234 は、M 個の重みオペランドを行列プロセッサ 240 に入力する回路であり得る。論理回路 232 は、画像データオペランドを行列プロセッサ 240 に入力する回路であり得る。重み入力行列 208 およびデータ入力行列 206 は、SRAM デバイスを含む様々なタイプのメモリに格納されてもよい。当業者は、様々なタイプのオペランドが行列プロセッサ 240 に入力され得ることを認識するであろう。

【 0 0 2 8 】

特定の実施形態による動作において、システム 200 は、システム内の冗長動作を削減し、ハードウェア固有のロジックを実装してデータおよび重みの大きなセットにわたって特定の数学的演算を実行することにより畳み込み演算を加速する。この加速は、行列プロセッサ 240 内の数学的演算のタイミングを大規模に計るだけでなく、画像データと重みとを取得して行列プロセッサ 240 に入力する方法 (および対応するハードウェア構成) の直接的な結果である。

【 0 0 2 9 】

実施形態において、図 2 のフォーマッタ 210、212 の例はインラインフォーマッタとして実装される。特定の実施形態では、フォーマッタ 210、212 は別個の構成であり、他の実施形態では、フォーマッタ 210、212 は一緒におよび/または 1 つ以上の他の構成と統合される。それぞれがハードウェアで実装され、行列をオペランド上のベクトルに変換して、行列プロセッサ 240 内で動作する。他の実施形態では、フォーマッタ 210、212 はソフトウェアで実装されるが、これは通常、速度の低下をもたらす。データフォーマッタ 210 は、データ入力行列 206 を含む 2 次元または三次元 (例えば、 $3 \times 3 \times 3$ 立方) のデータを単一のベクトルまたは行もしくは列で表されるストリングに変換し、それによりデータ入力行列 206 を線形化またはベクトル化する。詳細には、フォーマッタ 210 はデータ入力行列 206 を受け取り、行列プロセッサ 240 によって処理される入力データを準備する。実施形態において、これは、行列プロセッサ 240 が出

10

20

30

40

50

カピクセルを生成するときに畳み込み計算の一部として行列乗算を効率的に実行できるように、行列プロセッサ 240 のハードウェア要件に従ってデータ入力行列 206 のパラメータを適切なフォーマットにマッピングすることによって達成される。

【0030】

一例として、行列プロセッサ 240 が 96 行 96 列を含むと仮定すると、 96×96 フォーマットにマッピングされたデータは、行列プロセッサ 240 をその全計算能力まで利用させ、したがって、好ましい効率を提供する。その場合、フォーマッタ 210 は 96 列幅の出力を生成する必要がある。同様に、フォーマッタ 212 は、重み入力行列 208 に基づいて 96 行幅の出力を生成する必要がある。

【0031】

実施形態では、フォーマッタ 210 は、いくつかのマルチプレクサまたはスイッチを使用して、データ入力行列 206 の一部またはすべてをフェッチし、そこから異なる要素を選択して、行列プロセッサ 240 の列に従って並べられるデータを生成する。実施形態では、選択により、データ入力行列 206 からの適切なデータが、規定のクロックサイクルで各列に渡されることが保証される。実施形態において、重みが静的である場合、それらはオフラインで事前にフォーマットされ、メモリに格納され、一度だけフェッチされ、フォーマッタ 212 を使用せずに修正されたベクトル化されたフォーマットで行列プロセッサ 240 に直接供給されてもよい。他の実施形態では、様々なフォーマットおよびフェッチ操作に従って、重みを動的に調整し、行列プロセッサ 240 に供給してもよい。実施形態では、行列プロセッサ 240 は、さまざまなサイズの列および行の入力を可能にする。すなわち、行列プロセッサ 240 は、任意のサイズの $N \times M$ 計算を計算するように設計されている。

【0032】

他の実施形態では、データ入力行列 206 の列の数（例えば、 X ）が行列プロセッサ 240 の列の数よりも大きくなるように行列プロセッサ 240 の列の数が制限（例えば、 N 列に）された場合（すなわち、 $X > N$ ）の場合、制御ロジック 250 は、データ入力行列 206 を複数のサブ行列に分割し、各サブ行列が行列プロセッサ 240 によって計算されるようにすることができる。そのような場合、各行列プロセッサ 240 は異なるスレッドで実行されていてもよい。たとえば、データ入力行列 206 が 192×96 のデータポイントで構成され、行列プロセッサが 96 列 96 行を有している場合（すなわち、1 クロックサイクルで 96×96 の計算が発生する場合）、制御ロジック 250 はデータ入力行列 206 を 2 つのサブ行列（データ入力行列 206 の左半分とデータ入力行列 206 の右半分など）に分割してもよい。各サブ行列は、 96×96 のデータポイントで構成される。個別にスレッド化された各行列プロセッサ 240 は、送信されるサブ行列の出力チャネルを計算し、全チャネルからの値（つまり 192 の値）を保持するのに十分大きいはずである最終出力アレイ 260 に結果を配置する。より一般的には、データ入力行列 206 は、任意の数のサブ行列に分割され、それぞれが別個のスレッドで実行される異なる行列プロセッサ 240 に送信されてもよい。出力アレイ 260 と同様に、データ入力行列 206、データフォーマッタ 210、キャッシュ/バッファ 224、論理回路 232、および後処理ユニット 228 は、同様により大きなデータを収容できなければならない。

【0033】

代替実施形態では、制御ロジック 250 が内積に沿って計算を分割することにより、複数の行列プロセッサ 240 間で CNN を計算することができる。内積のセグメントがそれぞれ異なる行列プロセッサ 240 で計算され、次いで、入力された積が一緒に加算されて出力ベクトルが計算され、出力ベクトルは出力アレイ 260 に格納される。

【0034】

高速行列乗算に適した代替形式にデータを再配置することにより畳み込み演算を行列乗算に変換するための、CPU または GPU によって実行されるフォーマット機能の一般的なソフトウェア実装とは異なり、本開示のさまざまなハードウェア実装は、データを臨機応変にフォーマットし、実行できるようにする。たとえば、サイクルごとに 96 個のデー

10

20

30

40

50

タが有効になり、行列の非常に多数の要素を並列処理できるため、データを行列演算に効率的にマッピングすることができる。実施形態では、 $2N$ 個のフェッチされた入力データについて、 $2N^2$ の計算データが単一のクロックサイクルで取得され得る。このアーキテクチャは、複数のデータ入力にわたって多数の数学的演算を実行する並列で効率的な同期プロセスを提供するだけでなく、一般的なプロセッサアーキテクチャで使用される読み取り操作またはフェッチ操作の数を効果的に削減することにより、処理速度を大幅に向上させる。

【0035】

実施形態では、任意の数の列および行を有する行列プロセッサ240の効率を高めるために、フォーマッタ212、214は、入力行列データの異なる形状を行列プロセッサ240に適した列および行に再フォーマットすることができる。実施形態では、異なる入力サイズを有する行列の処理に対応するために、フォーマットが動的に実行される。実施形態において、入力チャネルを含む再フォーマットされた行列は、キャッシュ/バッファ224に供給される。

10

【0036】

キャッシュ/バッファ224は、様々なデータが再利用されるため、 k を畳み込みカーネル幅として、データ入力行列206から $1/k$ 回だけデータをフェッチすることができる。たとえば、特定のサイクルで行がフェッチされると、特定の列はその行のすべてのデータにアクセスできる。いくつかの実施形態において、キャッシュ/バッファ224は、SRAMに再アクセスしてSRAMからデータ読み取る必要なく、畳み込みによって再利用され得るデータの部分的なコピーを格納する部分的なバッファであってもよい。

20

【0037】

行列プロセッサ240が計算を完了すると、結果のセットは、例えば、行列プロセッサ240の最下行のアキュムレータから、例えば、内積を受信するシフトレジスタを効果的に形成する出力フリップフロップ（図示せず）にシフトされ得る。実施形態では、出力チャネルに対応する行から、例えばクロックサイクルごとに1つ、出力アレイ226に結果をプルまたはシフトすることは、状態機械（図示せず）によって達成され得る。状態機械は、例えば、データをSRAMおよび/または後処理ユニット228に送信する前に、出力チャネル上で追加の動作を実行してもよい。行列プロセッサ240の内部動作は、以下により詳細に説明される。

30

【0038】

実施形態では、行列プロセッサ240は、行列プロセッサ240を介して出力アレイ226に渡される結果のコピーを保存することにより並列処理を可能にするシャドーレジスタを含む。実施形態では、演算結果を出力レジスタからシャドーレジスタに移動することは、次の値のセットをALUにロードすることを伴う。

【0039】

累積が完了すると、畳み込みが開始され、前の畳み込みのすべてのデータが出力アレイ226に出力される前に蓄積がやり直され得る。その結果、各クロックサイクルにおいて、行列プロセッサ240のデータが1行下に移動し、各サイクルにおいて最後の行が出力アレイ226に出力されるようになる。実際、この動作モードは、SRAMにデータを保存するなどの追加の処理操作とは無関係に、中断することなく連続した各サイクルで新しい計算が行われることを保証する。

40

【0040】

後処理ユニット228は、ハードウェア加速ブリーディングユニット、メモリからデータを取得してSRAMなどにデータ（重みや結果など）を保存するダイレクトメモリアクセス（direct memory access：「DMA」）の一部であり得るDRAMなどの多数のデバイス（図示せず）を含むか、またはそれらと相互作用し得る。デバイスは、システム200内のフォーマッタ210、212および他の構成を管理する制御ロジック250によって部分的または全体的に制御されてもよい。

【0041】

50

データを読み取るためのアドレスを生成し、結果を書き込み、畳み込みの後続のステップで使用されるデータを取得する場所と実行する方法とを計算するために、システム 200 が畳み込みのどこにあるかを追跡するシーケンサなど、管理機能を実行する補助デバイスは図 2 には示されていない。

【0042】

特定の実施形態では、重み入力行列 208 は物理的に分割され、行列プロセッサ 240 の 2 つの異なる側から重みを駆動し、2 次元アレイが、重み入力行列 208 のデータの一部をそれぞれが受け取る 2 つの領域（例えば、左側と右側）に分割されるようにする。このような実装では、重みが既知であるという事実を活用することにより、データの待ち時間が短縮される。実施形態では、ピーク電力消費を低減するために、重みとデータとの乗算が特定のサイクル数にわたって分散されるように動作タイミングが選択されてもよい。この効率的な動作タイミングにより、行列プロセッサによって実行される読み取り操作の数の減少や、行列内（サブ回路間など）でのデータ移動の効率の改善など、エネルギー消費ステップが削減される。

10

【0043】

実施形態では、冗長データを識別するように構成された状態機械（図示せず）を使用することができる。識別された冗長データは列全体で再利用できるため、データを再フェッチする必要はない。状態機械は、例えば、画像サイズ、フィルタサイズ、ストライド、チャネル数、および同様のパラメータに関連する入力に基づいて、実行されるデータをどのように、およびどこにシフトするかを決定するように構成されてもよい。

20

【0044】

実施形態では、ブースエンコーダは、行列プロセッサ 240 の乗算アーキテクチャのいくつかの要素にわたって共有される。ブースエンコーダは、当技術分野で知られている如何なるブースエンコーダであってもよく、2 つの数値を乗算し、2 つの数値の 1 つを、たとえば、8 ビット値から 12 ビットまたは乗算器ロジックでの乗算演算を容易にするその他の値にエンコードするために使用されてもよく、したがって、より高速である。実施形態では、ブースエンコーダは、同一のコード化された代替の重み値をすべての列にわたって共有するように、行全体にわたって並列に適用されてもよい。すべての列にオペランドをロードすることにより、行全体にわたって 1 クロックサイクルで乗算を実行できる。したがって、N 個の計算要素で同じデータ（たとえば、重み）を共有するために再エンコード

30

【0045】

図 3 は、図 2 に示される行列プロセッサアーキテクチャの例示的な構成の詳細を示す。実施形態では、行列プロセッサ 300 は、各軸線上において所定のベクトル長に対応し得る。図 3 に示されるように、行列プロセッサ 300 は、行列形式で配置された 6×6 タイル 302 のアレイを備えてもよい。各タイル 302 は、サブ回路 350 を順に含む行列 320 を含む得る。図 4 を参照して以下で詳細に説明するように、各サブ回路 350 は、算術演算を実行できるセルであってもよい。実施形態では、サブ回路 350 は、乗算、累積、およびシフト演算を同時に実行する。

40

【0046】

実施形態では、行列プロセッサ 300 の複数の行および列を利用して $N \times N$ タイル出力を生成することにより、算術演算が並列化される。たとえば、96 の行サイズと 96 の対応する列サイズとは、 $2 * 9216$ の数学計算の出力を容易にする。他の実施形態では、行の数と列の数とは異なっていてもよい。つまり、N 行 M 列であり得、 $N \times M$ タイル出力が生成され得る。たとえば、行サイズが 96 で対応する列サイズが 192 の場合、1 クロックサイクルで $2 * 18,432$ の計算結果が生成される。

【0047】

図 4 は、図 3 に示されるサブ回路の例示的な乗加算回路の実装を示す。図 4 に示すよう

50

に、乗加算回路 4 0 0 は、乗算器 4 3 0、加算器 4 3 2、論理 4 3 4、4 3 6、4 3 8、アキュムレータ 4 2 4、シャドーレジスタ 4 2 8、および出力レジスタ 4 4 0 を含む。実施形態では、アキュムレータ 4 2 4 は累算レジスタとして実装されてもよい。

【 0 0 4 8 】

実施形態では、アキュムレータ 4 2 4 は、レジスタを含む A L U のセットと、A L U の出力を受信するように構成され得るシャドーレジスタ 4 2 8 とを含み得る。

【 0 0 4 9 】

動作中、乗算器 4 3 0 は、重み 4 0 2 およびデータ 4 0 4 を受信して乗算し、それらから積を生成する。各積は、乗算器 4 3 0 から積を受け取ることに応答して、積をアキュムレータ 4 2 4 の現在値に加算する加算器 4 3 2 に提供されてもよい。

10

【 0 0 5 0 】

実施形態において、アキュムレータ 4 2 4 は、例えば出力レジスタ 4 4 0 に格納される累算値を生成する。累積値は畳み込みの結果であり、図 2 を参照して述べたように、2 つのフォーマットされた行列の内積に対応し得る。

【 0 0 5 1 】

実施形態では、結果 4 5 0 を出力するシャドーレジスタ 4 2 8 に出力レジスタ 4 4 0 の結果のコピーを提供し、アキュムレータ 4 2 4 に再度アクセスして新しい計算を開始することができる。実施形態では、図 4 の乗加算回路 4 0 0 は、乗算、加算演算、およびシフト演算を同時に、すなわち単一のサイクル内で実行することができ、それにより各サイクルで発生する演算の総数を倍にすることができる。

20

【 0 0 5 2 】

実施形態において、C l e a r A c c 信号 4 0 8 は、例えば、乗算器 4 3 0 が乗算演算を実行したときに、累積演算がやり直せるように、アキュムレータ 4 2 4 の内容を消去する。実施形態において、R e s u l t E n a b l e 信号 4 1 2 は、データ 4 0 4 が有効であるという判定に応じて起動される。アキュムレータ 4 2 4 は、データを累積して保存し、データを累積して消去し、または単にデータを消去し得ることが理解される。

【 0 0 5 3 】

実施形態では、結果は、単一のクロックサイクルで、すなわち、中間の実行動作および保存動作を必要とせずに、出力レジスタ 4 4 0 からシャドーレジスタ 4 2 8 に移動される。

【 0 0 5 4 】

30

図 5 は、本開示の様々な実施形態による例示的な畳み込み演算を示している。畳み込み 5 0 0 は、入力画像 5 0 2 の入力チャンネル I C、重み 5 3 2、内積 5 1 4、出力チャンネル O C、およびアキュムレータ 5 4 0 を含む。

【 0 0 5 5 】

実施形態では、畳み込み演算 5 0 0 は、例えば入力画像 5 0 2 内の小さな特徴を検出するために、個々のフィルタ（すなわち、重み）5 3 2 を入力画像 5 0 2 に適用する。異なる順序で一連の異なる特徴を分析することにより、入力画像 5 0 2 でマクロ特徴を識別することができる。他の実施形態では、入力 5 0 2 は非画像データである。例えば、入力 5 0 2 は、超音波、レーダ、ライダー（L I D A R）、または他のセンサデータなどの非画像センサデータであってもよい。入力 5 0 2 は、一般的な数学的計算または当業者に知られている他のタイプのデータでもよい。

40

【 0 0 5 6 】

畳み込み 5 0 0 は、各入力チャンネル I C が異なる情報のセットを含み得るため、各入力チャンネル I C に対して異なる重みのセット 5 3 2 を使用してもよく、各重み行列 5 3 2 は、異なる特徴の識別を助けるように設計され得る。実施形態では、畳み込み 5 0 0 は、長方形の入力行列 5 0 4 に長方形の重み行列 5 3 2 を乗算して、部分的な内積を取得する。次に、出力画像内の出力ピクセル 5 1 4 を表す累積された内積 5 1 4（すなわち整数）を生成するために、部分内積が加算器 5 4 6 によって合計され得る。

【 0 0 5 7 】

実施形態では、出力チャンネル O C の各ピクセルは、乗算器 5 4 2 および加算器 5 4 4 に

50

よって生成される。実施形態において、部分的な内積の値は、入力画像 5 0 2 の領域 5 0 4 への重み行列 5 3 2 の全体の適用に対応する。換言すれば、各重み 5 3 2 は、乗算器 5 4 2 により領域 5 0 4 とドット乗算されて部分的な内積を生成し、次いで、部分的な内積がアキュムレータ 5 4 0 に累積されて、畳み込みを表す累積出力を生成する。

【0058】

1 つまたは複数の入力チャンネル IC、たとえば各色（たとえば RGB）に 1 つを使用することができる。例えば、各畳み込みは、各色に 1 つずつ、3 つの異なる行列を表す重み 5 3 2 を使用してもよい。各出力チャンネル OC 5 1 2 は、入力データ 5 0 2 の異なる特徴を表す異なるフィルタまたは重み 5 3 2 を使用して生成され得る。出力チャンネルの数は、機能の数に依存し得る。畳み込みの数は、出力チャンネル OC の数に入力チャンネル IC の数

10

【0059】

図 5 に示すように、入力行列 5 0 4 は、3 つの入力チャンネル、すなわち $3 \times 3 \times IC$ にわたって 3×3 の重み行列 5 3 2 と組み合わせることができる $K \times K \times Ky$ （すなわち、 3×3 ）行列であり、深さは一致し、出力プレーンにおいて単一の要素、内積 5 1 4 を生成する。出力チャンネル 5 1 2 の各内積 5 1 4 は、ドット乗算の結果である。

【0060】

図 6 ~ 図 8 は、本開示の様々な実施形態による例示的な畳み込み演算の詳細を示している。畳み込み 6 0 0 は、入力データ行列 6 0 2、重みデータ行列 6 0 4、アレイ 6 0 6、および内積 6 3 0 を含む。実施形態では、アレイ 6 0 6 は、図 2 および図 3 に示されるような行列プロセッサアーキテクチャである。

20

【0061】

図 6 の入力データ行列 6 0 2 は、実施形態において、図 5 に示される矩形入力行列 5 0 4 などの入力行列を線形化することにより得られ、入力行列のベクトル化形態を得ることができる列 6 1 0 を含む。同様に、重みデータ行列 6 0 4 は、図 5 の矩形重み行列 5 3 2 などの重み行列のベクトル化された形式であり得る行 6 2 0 を含む。例として、 3×3 入力行列と 3 つの入力チャンネルとは、入力データ行列 6 0 2 で使用するために 27 要素列 6 1 0 を生成できる $3 \times 3 \times 3 = 27$ 要素を含むベクトルに再フォーマットできる。逆に、

30

【0062】

実施形態において、図 5 において長方形として描かれた入力チャンネルおよび入力重みは、例えば、図 2 を参照して論じられたフォーマットによって、行列乗算器 / プロセッサ（図 2 の要素 2 4 0 として示されている）に提供されるベクトルフォーマット（例えば、96 要素を有するベクトル）に再フォーマットされ、 96×96 要素の内積演算を並行して実行できる。詳細には、各入力チャンネルの長方形として図 5 に示されている入力データ 5 0 4 および入力重み 5 3 2 は、ベクトル形式に再フォーマットされる。

40

【0063】

実施形態では、図 6 に入力データ 6 0 2 および入力重み 6 0 4（例えば、それぞれ 96 個の要素を含む）として示される結果のベクトル形式は、 96×96 要素の内積演算を並列に実行する行列プロセッサまたは行列乗算器 2 4 0 に提供される。実施形態では、出力チャンネルの計算において、同じ入力データのセットを使用するが重みの異なるセット（すなわち、フィルタ）を使用して同じ出力ピクセルが生成され、入力データを読み取ること

【0064】

さらに、入力データ行列 6 0 2、重みデータ行列 6 0 4、およびアレイ 6 0 6 は、図 6

50

に示されているものとは異なる数の列および行を有し得ることが理解される。特に、入力データ行列 6 0 2 および重みデータ行列 6 0 4 の形状は、アレイ 6 0 6 の任意の調停構成の列および行に対応するようにフォーマットされてもよい。加えて、重みデータ行列 6 0 4 が既知の状況では、フォーマッタを使用せずに行 6 2 0 を生成し、ベクトル化された形式で保存することができる。

【 0 0 6 5 】

実施形態では、図 6 の内積 6 3 0 は、列 6 1 0 に対応するベクトルに行 6 2 0 に対応するベクトルをドット乗算することにより生成される。実施形態では、図 7 に示されるように、次の内積 6 3 2 は、列 6 1 2 に対応するベクトルと行 6 2 0 に対応するベクトルをドット乗算することにより取得され得る。当業者が認識するように、アレイ 6 0 6 の第 1 行のすべての内積が満たされると、アレイ 6 0 6 の第 2 行の内積は、入力データ行列 6 0 2 の第 1 列 6 1 0 の要素を重みデータ行列 6 0 4 の 2 行目などとドット乗算することにより計算され得る。

【 0 0 6 6 】

図 6 ~ 図 8 は単に例示の目的を果たすだけであり、上述のドット乗算が同時に実行されてワンショット行列 - 行列乗算演算を生成できることに留意することが重要である。

【 0 0 6 7 】

図 9 は、本開示の様々な実施形態による例示的な逆畳み込み演算を示している。逆畳み込みシステム 9 0 0 は、入力画像 9 0 2 の入力チャンネル IC、重み 9 2 2、内積 9 0 4、9 0 6、および出力チャンネル OC を含む。当業者は、逆畳み込み演算 9 0 0 が、実際には、畳み込み演算、例えば図 5 に示される畳み込みの数学的転置（ほぼ逆）であることを認識するであろう。当業者はさらに、通常の畳み込みニューラルネットワークに使用される手順と同様の手順を適用することにより、ニューラルネットワークを使用して逆畳み込み演算 9 0 0 を学習できることを認識するであろう。簡潔にするために、図 5 のものと同様の構成の説明または機能はここでは繰り返さない。

【 0 0 6 8 】

実施形態では、図 9 の逆畳み込み演算 9 0 0 は、重み 9 2 2 を使用して内積 9 0 4 9 0 6 を分解することにより行列 9 1 2 を再構築する。畳み込み演算と同様に、逆畳み込み 9 0 0 は、各入力チャンネル IC に対して異なるセットの重み 9 2 2 を使用してもよい。実施形態では、逆畳み込み 9 0 0 を画像に有利に適用して、例えばアーチファクトに対するロバスト性を改善するために画像逆畳み込みを実行することができる。他のアプリケーションは、画像データの分析および復元などを含み得る。

【 0 0 6 9 】

図 10 は、本開示の様々な実施形態による畳み込みニューラルネットワークを加速するために算術演算を実行するプロセスを示す。

【 0 0 7 0 】

算術演算を実行するためのプロセス 1 0 0 0 は、データ行列の行を表すことができるオペランドの第 1 セットが第 1 論理回路から受信されるとき、ステップ 1 0 0 2 で始まる。このオペランドの第 1 セットは、オペランドが行列プロセッサへの入力と整列するようにベクトル化できる。特定の実施形態では、ベクトル化されたオペランドのサイズは、軸線に沿った行列プロセッサへの入力の数に直接関係する。

【 0 0 7 1 】

ステップ 1 0 0 4 で、重み行列の列を表し得るオペランドの第 2 セットが、第 2 論理回路から受信される。オペランドのこの第 2 セットは、オペランドが行列プロセッサへの対応する入力内で整列するようにベクトル化される。特定の実施形態では、ベクトル化されたオペランドのサイズは、異なる軸線に沿った行列プロセスへの入力の数に直接関係する。

【 0 0 7 2 】

ステップ 1 0 0 6 で、オペランドの第 1 セットは、オペランドの第 2 セットとドット乗算され、1 つ以上の内積を取得する。特定の実施形態では、オペランドのセットにわたるこのセット操作は、単一のクロックサイクルで実行される。

10

20

30

40

50

【 0 0 7 3 】

ステップ 1 0 0 8 で、内積を使用して、画像をフィルタで畳み込み、畳み込み結果を生成することができる。

【 0 0 7 4 】

ステップ 1 0 1 0 で、畳み込み結果がさらに処理されて、画像出力が強化される。このさらなる処理は、非線形関数、正規化操作、またはプーリング操作を使用して発生し得る。

【 0 0 7 5 】

当業者は、計算システムまたはプログラミング言語が本発明の実施に重要ではないことを認識するであろう。当業者はまた、上記の多くの要素が物理的および／または機能的にサブモジュールに分離されるか、または一緒に結合され得ることを認識するであろう。

10

【 0 0 7 6 】

以下の請求項の要素は、複数の依存関係、構成、および組み合わせを含むなど、さまざまに配置できることに留意されたい。例えば、実施形態において、様々な請求項の主題は他の請求項と組み合わせられてもよい。

【 0 0 7 7 】

前述の例および実施形態は例示であり、本発明の範囲を限定するものではないことを当業者は理解するであろう。明細書を読み、図面を検討することで当業者に明らかなすべての置換、強化、同等物、組み合わせ、および改良は、本発明の真の精神および範囲内に含まれることが意図されている。

20

30

40

50

【図面】

【図 1】

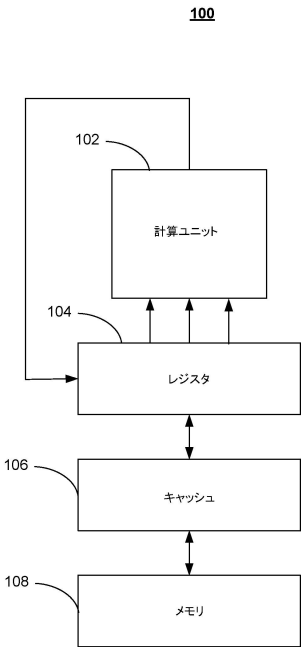


FIGURE 1
(従来技術)

【図 2】

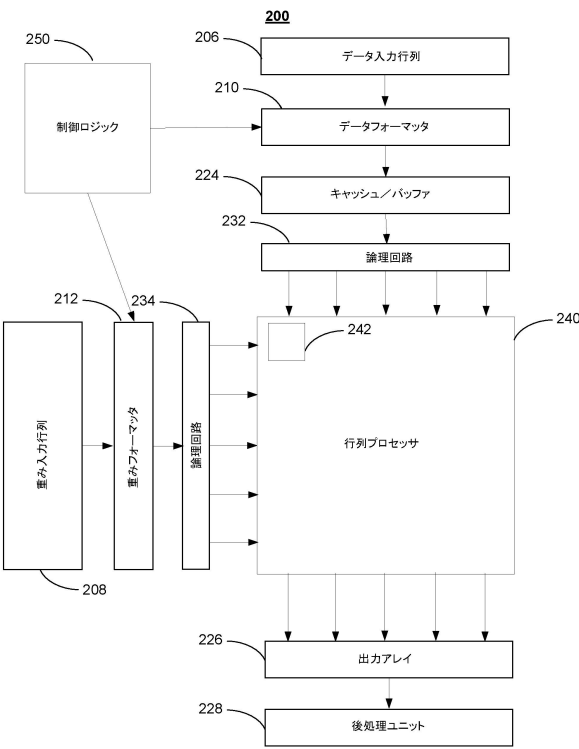


FIGURE 2

【図 3】

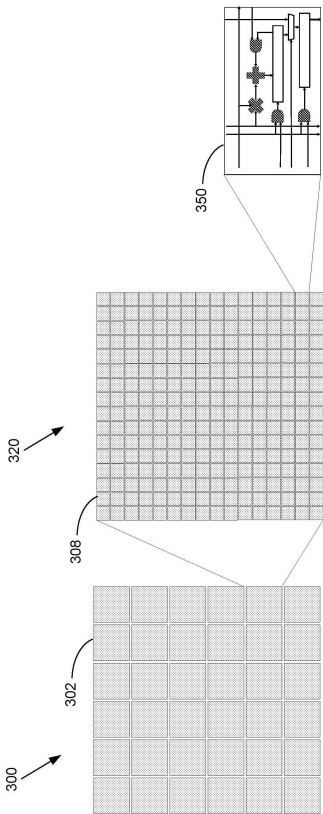


FIGURE 3

【図 4】

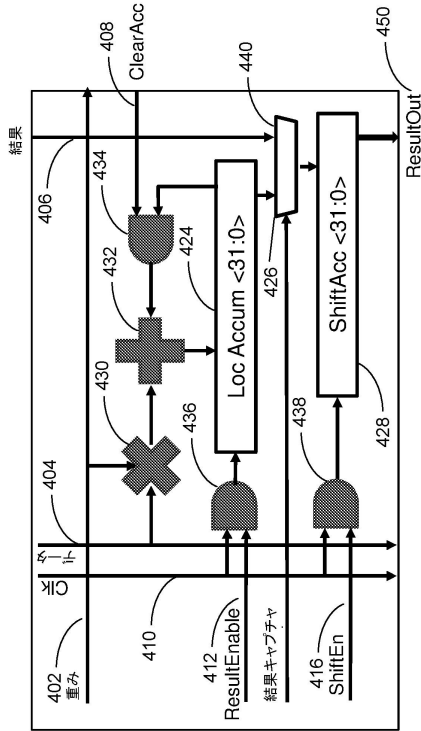


FIGURE 4

10

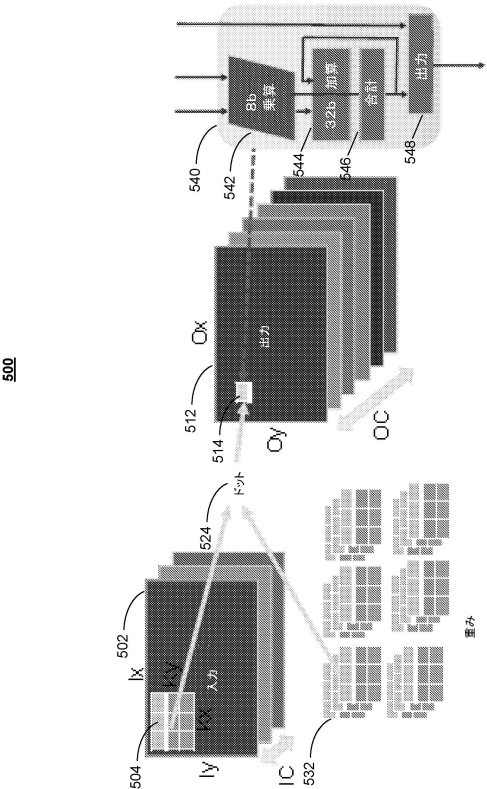
20

30

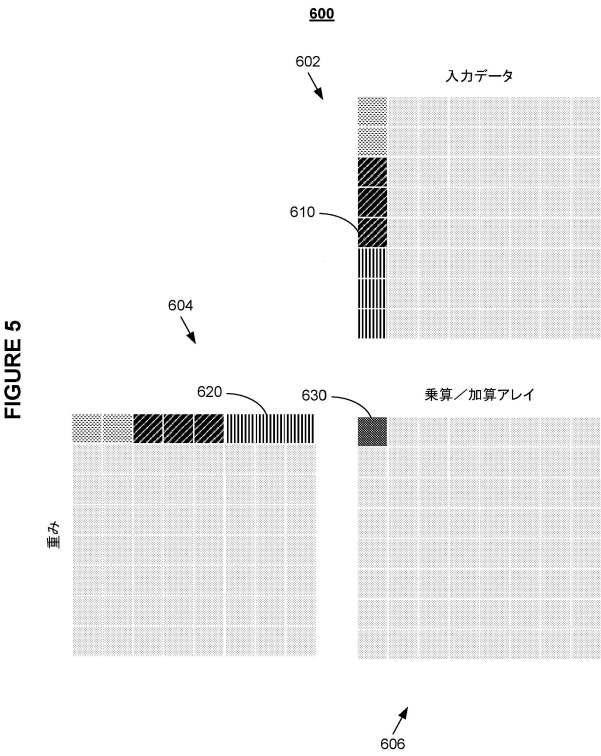
40

50

【図 5】



【図 6】



10

20

FIGURE 6

【図 7】

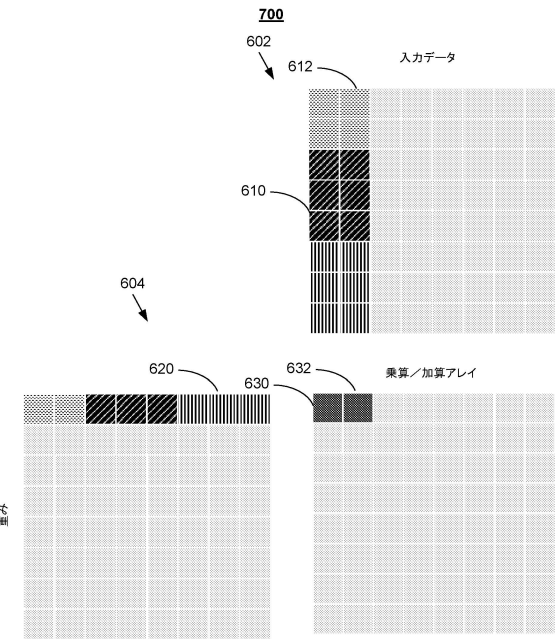


FIGURE 7

【図 8】

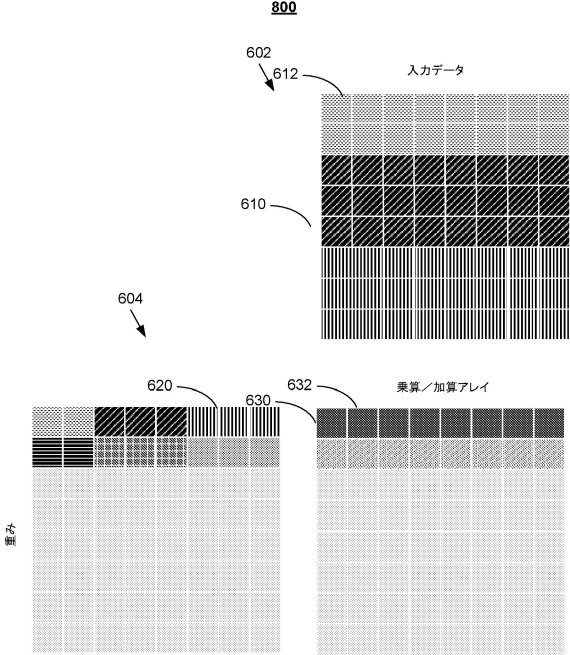


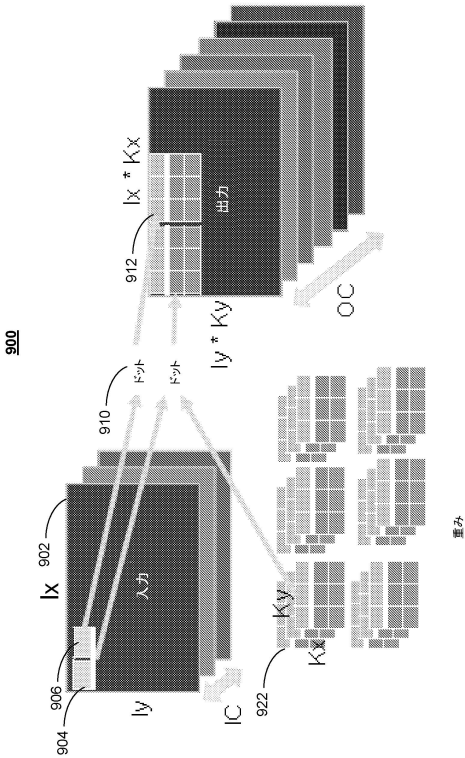
FIGURE 8

30

40

50

【図 9】



【図 10】

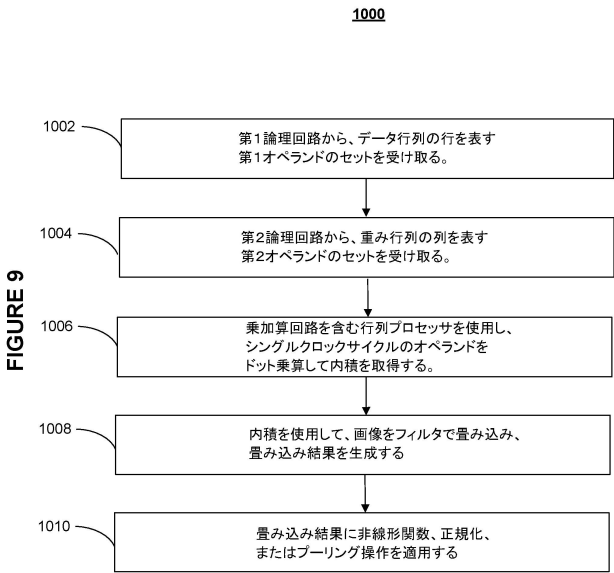


FIGURE 10

フロントページの続き

米国(US)

アメリカ合衆国 カリフォルニア州 9 4 4 0 3 , サン マテオ , パーバンク アベニュー 1 6 0

審査官 坂庭 剛史

(56)参考文献

米国特許出願公開第 2 0 1 7 / 0 0 9 7 8 8 4 (U S , A 1)

米国特許出願公開第 2 0 1 6 / 0 3 4 2 8 9 1 (U S , A 1)

特開 2 0 1 0 - 0 7 9 8 4 0 (J P , A)

特開 2 0 1 5 - 0 5 6 1 2 4 (J P , A)

米国特許出願公開第 2 0 1 6 / 0 3 7 9 1 0 9 (U S , A 1)

(58)調査した分野 (Int.Cl. , D B 名)

G 0 6 F 1 7 / 1 6

G 0 6 N 3 / 0 6 3