

(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2021-108104
(P2021-108104A)

(43) 公開日 令和3年7月29日(2021.7.29)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 15/80 (2006.01)	G06F 15/80	5B056
G06N 3/063 (2006.01)	G06N 3/063	
G06F 17/16 (2006.01)	G06F 17/16	M

審査請求 未請求 請求項の数 20 O L 外国語出願 (全 25 頁)

(21) 出願番号 特願2020-157941 (P2020-157941)
 (22) 出願日 令和2年9月18日(2020.9.18)
 (31) 優先権主張番号 16/729, 381
 (32) 優先日 令和1年12月28日(2019.12.28)
 (33) 優先権主張国・地域又は機関
 米国 (US)

(71) 出願人 591003943
 インテル・コーポレーション
 アメリカ合衆国 95054 カリフォル
 ニア州・サンタクララ・ミッション カレ
 ッジ ブレーバード・2200
 (74) 代理人 110000877
 龍華国際特許業務法人
 (72) 発明者 カムレシュ アール. ピレイ
 アメリカ合衆国 95054 カリフォル
 ニア州・サンタクララ・ミッション カレ
 ッジ ブレーバード・2200 インテル
 ・コーポレーション内

最終頁に続く

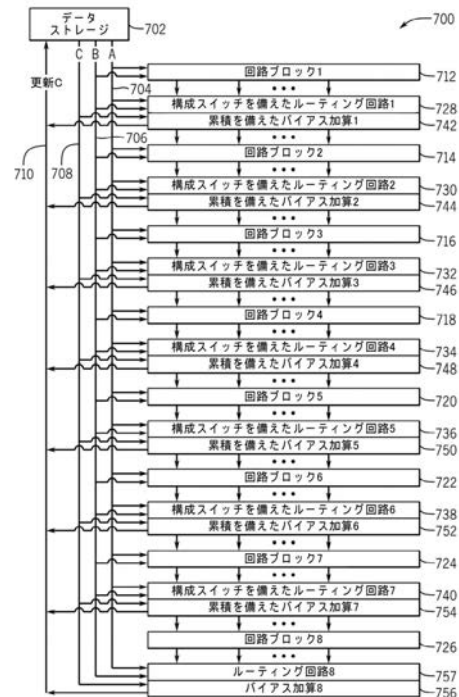
(54) 【発明の名称】 部分的読み取り／書き込みが可能な再構成可能なシストリックアレイのシステム及び方法

(57) 【要約】 (修正有)

【課題】再構成可能なシストリックアレイ回路を含むデータ処理システムを提供する。

【解決手段】再構成可能なシストリックアレイ回路700は、1又は複数のグループの処理要素を含む第1回路ブロック712と、1又は複数のグループの処理要素を含む第2回路ブロック714とを含む。再構成可能なシストリックアレイ回路は、さらに、累積値、乗算積、又はそれらの組み合わせに行列バイアスを加算する累積回路を備えた第1バイアス加算742を含む。再構成可能なシストリックアレイ回路は、さらに、第1回路ブロックから第2回路ブロックへ、第1回路ブロックから累積回路を備えた第1バイアス加算へ又はそれらの組み合わせへ微分をルーティングする第1ルーティング回路728をさらに含む。

【選択図】図7



【特許請求の範囲】

【請求項 1】

データを格納するように構成されたデータストレージと、
再構成可能なシストリックアレイ回路と、を備え、前記再構成可能なシストリックアレイ回路は、

前記データを処理するように構成された 1 又は複数のグループの処理要素を有する第 1 回路ブロックと、

前記データを処理するように構成された 1 又は複数のグループの処理要素を有する第 2 回路ブロックと、

行列バイアスを累積値に、乗算積に、又はそれらの組み合わせに加算するように構成された累積回路を備えた第 1 バイアス加算と、

微分を前記第 1 回路ブロックから前記第 2 回路ブロックに、前記第 1 回路ブロックから累積回路を備えた前記第 1 バイアス加算に、又はそれらの組み合わせにルーティングするように構成された第 1 ルーティング回路と、

を含む、システム。

【請求項 2】

前記第 1 ルーティング回路は、構成スイッチ信号を受信することに基づいて、前記微分を前記第 1 回路ブロックから前記第 2 回路ブロックに、前記第 1 回路ブロックから累積回路を備えた前記第 1 バイアス加算に、又はそれらの組み合わせにルーティングするように構成されたデマルチプレクサ及びマルチプレクサ回路を含む、請求項 1 に記載のシステム

【請求項 3】

累積回路を備えた前記第 1 バイアス加算は、クロック信号に基づいて前記乗算積を前記累積値として累積するように構成された記憶回路、及び前記行列バイアスを前記累積値に、前記乗算積に、又はそれらの組み合わせに加算するように構成された少なくとも 1 つの加算器を含む、請求項 1 又は 2 に記載のシステム。

【請求項 4】

累積回路を備えた前記第 1 バイアス加算は、N の加算器のレイテンシを含み、前記記憶回路は、N 個の記憶コンポーネントを含む、請求項 3 に記載のシステム。

【請求項 5】

前記 N 個の記憶コンポーネントは、それぞれ、フリップフロップを含む、請求項 4 に記載のシステム。

【請求項 6】

前記記憶回路は、前記 N 個の記憶コンポーネントをマルチプレクサに結合する N 個のラインを含み、前記記憶回路は、前記加算器のレイテンシが演算中に N を超えた場合に、累積値を前記 N 個の記憶コンポーネントから前記マルチプレクサに前記 N 個のラインを介して送信するように構成される、請求項 4 又は 5 に記載のシステム。

【請求項 7】

累積回路を備えた前記第 1 バイアス加算は、累積回路を備えた前記第 1 バイアス加算に入る新しい値を前記累積値に加算し、前記加算の結果を前記 N 個の記憶コンポーネントに格納するように構成される、請求項 6 に記載のシステム。

【請求項 8】

1 又は複数のグループの処理要素を有する第 3 回路ブロックと、

第 2 行列バイアスを第 2 累積値に、前記乗算積に、又はそれらの組み合わせに加算するように構成された累積回路を備えた第 2 バイアス加算と、

微分を前記第 2 回路ブロックから前記第 3 回路ブロックに、前記第 2 回路ブロックから累積回路を備えた前記第 2 バイアス加算に、又はそれらの組み合わせにルーティングするように構成された第 2 ルーティング回路と、

を備える、請求項 1 から 7 のいずれか一項に記載のシステム。

【請求項 9】

10

20

30

40

50

前記第3回路ブロックの下流に配置され、第3行列バイアスを前記第3回路ブロックから出力に加算するように構成されたバイアス加算回路を備える、請求項8に記載のシステム。

【請求項10】

前記再構成可能なシストリックアレイ回路を使用するよう、又は前記再構成可能なシストリックアレイ回路を含むように構成されたホストプロセッサ(CPU)を備え、

前記CPUは、『「M」累積によるタイル部分「N」ドット積』命令であり、前記Nは、ともにマージされている異なる行列の数であり、Mは、前記再構成可能なシストリックアレイ回路への入力として使用される不完全な行列の数である、前記命令、前記再構成可能なシストリックアレイ回路への入力として使用されるとともにマージされている前記異なる行列のサイズを指定する即値を有する「ドット積に対するタイルサイズ」命令、累積回路を備えた前記第1バイアス加算を制御する「タイル累積ドット積」命令、又はそれらの組み合わせを実行するように構成される、
請求項1から9のいずれか一項に記載のシステム。

10

【請求項11】

行列A及び行列Bに基づいてデータの1又は複数のタイルのそれぞれに対してタイルサイズを決定する段階と、

タイルサイズに基づいて、完全タイル、不完全タイル、又はそれらの組み合わせを導出する段階と、

行列Cの結果を導出するために、前記完全タイル、前記不完全タイル、又はそれらの組み合わせを、再構成可能なシストリックアレイ回路を介して処理する段階であり、前記完全タイル、前記不完全タイル、又はそれらの組み合わせを処理する段階とは、前記行列Cの結果を提供するために、前記再構成可能なシストリックアレイ回路に含まれるルーティング回路及び前記再構成可能なシストリックアレイ回路に含まれる累積回路を備えたバイアス加算又はそれらの組み合わせを適用することを含む、段階と、
を備える方法。

20

【請求項12】

前記再構成可能なシストリックアレイ回路は、N行M列のアレイサイズを含み、前記完全タイルは、N行以下及びM列以下を有する完全サイズを含み、前記不完全タイルは、N行より多く、M列より多く、又はそれらの組み合わせを有する不完全サイズを含む、請求項11に記載の方法。

30

【請求項13】

前記ルーティング回路を適用する段階は、微分を1又は複数のグループの処理要素を含む第1回路ブロックから1又は複数のグループの処理要素を含む第2回路ブロックにルーティングする段階、及び微分を前記第1回路ブロックから累積回路を備えた前記バイアス加算に、又はそれらの組み合わせにルーティングする段階を含み、

微分を前記第1回路ブロックから累積回路を備えた前記バイアス加算にルーティングする段階は、前記微分を累積回路を備えた前記バイアス加算にて受信し、前記微分を行列Cのバイアスに加算するために累積値に累積する段階を含む、請求項11又は12に記載の方法。

40

【請求項14】

前記完全タイル、前記不完全タイル、又はそれらの組み合わせを前記再構成可能なシストリックアレイ回路を介して処理する段階は、

行列Cのアドレス衝突を検出し、累積回路を備えた前記バイアス加算に通信される累積イネーブル信号を自動的にオンにするように構成されたマイクロアーキテクチャモードを適用する段階と、

前記再構成可能なシストリックアレイ回路への入力として使用されるとともにマージされている異なる行列のサイズを指定する即値を有する「ドット積に対するタイルサイズ」命令、累積回路を備えた前記バイアス加算を制御する「タイル累積ドット積」命令、又はそれらの組み合わせを実行することによりアーキテクチャモードを適用する段階と、

50

を含む、請求項 11 から 13 のいずれか一項に記載の方法。

【請求項 15】

データを格納するように構成されたデータストレージと、
再構成可能なシストリックアレイ回路と、

前記再構成可能なシストリックアレイ回路に接続されるコアのデコーダであり、単一の命令をデコードされた 1 又は複数の命令にデコードする、デコーダと、を備え、

前記 1 又は複数の命令は、

行列 A 及び行列 B を表す前記データを、前記データストレージから、前記データを処理し、前記データに基づく微分を提供するように構成された 1 又は複数のグループの処理要素を含む第 1 回路ブロックに通信し、

再構成可能なルーティング回路のスイッチングオン又はオフに基づいて、前記微分を前記第 1 回路ブロックから第 2 回路ブロックに、累積回路を備えたバイアス加算に、又はそれらの組み合わせにルーティングするように構成され、

累積回路を備えた前記バイアス加算は、行列バイアスを、累積値に、行列 A 及び行列 B の乗算積に、又はそれらの組み合わせに加算するように構成され、

前記第 1 回路ブロック、前記第 2 回路ブロック、前記再構成可能なルーティング回路、累積回路を備えた前記バイアス加算、又はそれらの組み合わせは、前記再構成可能なシストリックアレイ回路に含まれる、装置。

【請求項 16】

前記単一の命令は、デコードされると、前記再構成可能なシストリックアレイ回路への入力として使用されるとともにマージされている異なる行列のサイズを指定する即値を有する「ドット積に対するタイルサイズ」命令、累積回路を備えた前記バイアス加算を制御する「タイル累積ドット積」命令、又はそれらの組み合わせを介してアーキテクチャモードを使用する、請求項 15 に記載の装置。

【請求項 17】

前記単一の命令は、『「M」累積によるタイル部分「N」ドット積』命令を含み、N は、ともにマージされている異なる行列の数であり、M は、前記再構成可能なシストリックアレイ回路への入力として使用される不完全な行列の数である、請求項 16 に記載の装置。

【請求項 18】

前記単一の命令は、デコードされると、前記再構成可能なシストリックアレイ回路に、前記データを使用することによって $C = + A * B$ を解かせ、前記データは、前記行列 A 及び前記行列 B を表す、請求項 15 から 17 のいずれか一項に記載の装置。

【請求項 19】

前記再構成可能なシストリックアレイ回路を有する回路を備え、

前記回路は、マイクロプロセッサ、ハードウェアアクセラレータ、フィールドプログラマブルゲートアレイ (FPGA)、特定用途向け集積回路 (ASIC)、カスタムマイクロチップ、又はそれらの組み合わせを含む、請求項 15 から 18 のいずれか一項に記載の装置。

【請求項 20】

行列 A 及び行列 B に基づいてデータの 1 又は複数のタイルのそれぞれに対してタイルサイズを決定する手段と、

タイルサイズに基づいて、完全タイル、不完全タイル、又はそれらの組み合わせを導出する手段と、

行列 C の結果を導出するために、前記完全タイル、前記不完全タイル、又はそれらの組み合わせを、再構成可能なシストリックアレイ回路を介して処理する手段であり、前記完全タイル、前記不完全タイル、又はそれらの組み合わせを処理することは、前記行列 C の結果を提供するために、前記再構成可能なシストリックアレイ回路に含まれるルーティング回路及び前記再構成可能なシストリックアレイ回路に含まれる累積回路を備えたバイアス加算又はそれらの組み合わせを適用することを含む、手段と、

10

20

30

40

50

を備えるシステム。

【発明の詳細な説明】

【背景技術】

【0001】

本開示は、概して、シストリックアレイベースのアクセラレータに関し、より詳細には、部分的読み取り/書き込みを伴うシストリックアレイベースのアクセラレータに関する。

【0002】

本セクションは、読者に本開示の様々な態様に関連し得る技術分野の様々な態様を紹介することが意図されており、当該態様は、以下に説明される及び/又は特許請求の範囲に記載される。この検討は、本開示の様々な態様のより良好な理解を容易にするための背景情報を読者に提供することに役立つと考えられる。従って、これらの記述は、先行技術を承認するものとしてではなく、この観点で読むべきものであることを理解されたい。

【0003】

シストリックアレイベースのアクセラレータの使用は、ディープニューラルネットワーク(DNN)ベースのアプリケーションで有用なそれらのようなより効率的な計算を提供してよい。シストリックアレイベースのDNNアクセラレータは、アプリケーションの計算エンジンを提供するために、数百の算術演算ユニット、例えば処理要素(PE)を使用してよい。DNNアクセラレータは、通常且つ固定サイズの密行列乗算に対してより最適化されてよい。例えば、算術演算ユニットのシストリックアレイ実装が使用されて、性能を向上し、表面積を減らし、電力の利益を得てよい。従って、特定のDNNアクセラレータは、非常に規則的なデータフロー用に最適化された高密度2次元(2D)アレイを使用してよい。多くのDNNアクセラレータは、比較的低速又は非効率であり得る。

【図面の簡単な説明】

【0004】

【図1】本開示の実施形態による、再構成可能なシストリックアレイベースのアクセラレータの回路を有する1又は複数のプロセッサを含むデータ処理システムのブロック図である。

【0005】

【図2】本開示の実施形態によるシストリックアレイシステムの例のブロック図である。

【0006】

【図3】本開示の実施形態による部分的なバイアス累積サポートを含む再構成可能なシストリックアレイシステムを実行するために使用され得るスケジューラの実施形態のブロック図である。

【0007】

【図4】本開示の実施形態による、図3の再構成可能なシストリックアレイシステムのさらに詳細を示すブロック図である。

【0008】

【図5】本開示の実施形態による、再構成可能なルーティング回路及び累積回路を備えたバイアス加算の実施形態を示す概略図である。

【0009】

【図6】本開示の実施形態による、さらに詳細を示す累積回路を備えたバイアス加算の実施形態の概略図である。

【0010】

【図7】本開示の実施形態による、複数の再構成可能なルーティング回路及び累積回路を備えたバイアス加算を有する再構成可能なシストリックアレイシステムを示すブロック図である。

【0011】

【図8】本開示の実施形態による再構成可能なシストリックアレイシステムの回路を実行するのに好適な処理を示すフローチャートである。

10

20

30

40

50

【発明を実施するための形態】

【0012】

1又は複数の具体的な実施形態が以下に説明される。これらの実施形態の簡潔な記載を提供するために、実際の実装のすべての特徴が本明細書において説明されるわけではない。任意のエンジニアリング又は設計プロジェクトにおけるように、任意のそのような実際の実装の開発においては、システム関連及びビジネス関連の制約の順守などの開発者の特定の目的を達成すべく、多数の実装固有の決定がなされなければならない、当該制約は実装によって変動され得ることを理解されたい。さらに、そのような開発への取り組みは複雑且つ時間を消費するものであり得るが、それにもかかわらず、本開示の利益を有する当業者には設計、製造、及び製造のルーチン作業であることを理解されたい。

10

【0013】

本明細書で説明される技術は、ディープニューラルネットワーク(DNN)で使用されるそれらの計算のような特定の計算を向上するのに有用な特定のストリックアレイ技術を含む。ストリックアレイは、処理ユニットがセル又はノードと称されてよい密に接続された処理ユニットの同種のネットワークを含んでよい。各ノードには、様々な計算を提供するために使用されてよい融合型積和ユニット(FMA)のような処理要素(PE)を含んでよい。データはストリックアレイに入り、アレイのFMAを通じて、例えば隣接するFMA間を流れてよく、データフローの結果は、特定のアプリケーション、例えばDNNアプリケーションの計算として提供されてよい。DNNストリックアレイアクセラレータは、通常且つ固定サイズの密行列乗算に対してより最適化されてよい。例えば、DNNストリックアレイアクセラレータは、非常に規則的なデータフロー用により最適化された高密度2次元アレイを使用してよい。DNNストリックアレイアクセラレータを介して解決される非常に大きい又は小さいのいずれかの問題及び/又は提供された規則的データフローにうまくマッピングされない問題は、部分的な結果の複数の読み取り/書き込み、及び/又はストリックアレイ内のPEの大きく低下した使用率を引き起こし得る。

20

【0014】

ディープラーニングアプリケーションは、密DNNと疎DNNを含むように分類されてよい。密及び疎DNNの両方に対して、実行のいくつかの部分は、所与のストリックアレイに対する規則的なデータフロー上に完全にマッピングされてよいが、すべてではない。例えば、密DNNの場合、問題のサイズは非常に大きくなってよく、行列A、B、及びCの配列計算が $C + = A * B$ のような方程式を含む場合、各行列は複数のタイル(例えば、2Dデータ構造)に分割されて、行列をストリックアレイに「合致」させてよい。例えば、X次元にxPE及びY次元にyPEを有するストリックアレイにおいてCの単一のタイルを計算するために、Aに対するX次元及びBに対するY次元におけるすべての対応するx、yタイルに沿った計算が使用されてよい。X及びY次元の計算は、各それぞれのタイル乗算から生成される部分的な結果が書き出され、次いで、累積の単一の「チェーン」内のすべてのタイルの完了まで、さらなる処理(例えば、他の部分的な結果との累積)のために読み戻すことを必要としてよい。行列はタイルより小さくてよく(例えば、すべてのタイルよりも少ないスペースを使用する)、行列はタイルと同一のサイズであってよく、又は行列は複数のタイルを使用してよい(例えば、行列はいずれか1つのタイルのサイズより大きい)ことに留意されたい。従って、本明細書で称されるタイル又はタイルデータは、N列M行、いくつかの場合に $N = M$ を有するデータのアレイを含んでよい。行及び/又は列は、本明細書で「グループ」と称されてよい。

30

40

【0015】

疎DNNの場合、「ブロックパーシティ」処理が存在してよく、この場合に行列は任意のサイズの密なブロックにより表される。そのような密なブロック表現は、ゼロが表される必要がなくでよいので、行列内の多くの又はほとんどのゼロに亘って「スキップ」することを可能にしてよい。しかし、ブロックパーシティの副作用は、一般的な行列乗算(GEMM)の微分のような特定の微分を計算する場合に、入力行列に小さい及び/又は

50

不規則なサイズのブロックが見つけれられる可能性がある。すべてのディープラーニングアプリケーション（疎並びに密）に対して、部分的読み取り／書き込みを縮小し、シストリックアレイ上で不規則な幅で複数の行列乗算を実行し、それによりPEの利用が高くなることは有益であろう。さらに、現在、シストリックアレイのサイズに完全に合致するであろう密行列乗算に対する性能、面積、及び／又は電力不利益の支払いを最小に抑えつつ、PEの利用を向上することは有益であろう。

【0016】

本明細書で説明される技術は、部分的な累積サポートを備えた再構成可能なシストリックアレイを含む。累積サポートは、既存のタイルストレージから分離した、例えばスケジューラを介して複数の行列乗算を処理するのに好適な累積ストレージを含んでよい。スケジューラは、行列が実行のために送信される順序をスケジューリングしてよく、新しい命令（複数可）（例えば、マクロ命令）が使用されて、再構成可能なシストリックアレイを通じてデータフローを実行してよい。2つ（又はより多く）の命令が、さらに下で説明されるように計算間において宛先を使用（又は上書き）しないで同一の宛先を有する場合に、複数の行列乗算命令に亘ってシストリックアレイの宛先をチェックし、ソフトウェアの介入なしにハードウェアベースの計算を有効にするために使用されるマイクロアーキテクチャ機能が提供されてよい。再構成可能なシストリックアレイは、スケジューラによってスケジューリングされているタイルに基づいて有効にされてよい累積ロジックシステムを含む。累積ロジックシステムは、解かれる問題の終わりまで部分値を累積し、最終出力をストレージ（例えば、メモリ、バッファ、レジスタなど）に書き込んでよい。再構成可能なシストリックアレイを提供することにより、ハードウェアベースの計算がより柔軟になってよく、それとともに、ハードウェア及びストレージ（例えば、タイルレジスタファイル）間のデータ転送がさらに削減されるため、DNNベースのアプリケーションのような特定のアプリケーションに対する利用を向上し、データ転送を減少してよい。

10

20

【0017】

上記を念頭に置いて、図1は、本開示の実施形態による、1又は複数のプロセッサ（複数可）102を含むデータ処理システム100のブロック図である。データ処理システム100は、示されているものより多い又はより少ないコンポーネント（例えば、電子ディスプレイ、ユーザインタフェース構造、特定用途向け集積回路（ASIC））を含んでよい。データ処理システム100は、データ処理システム100へのデータ処理要求（例えば、DNN計算、機械学習、ビデオ処理、音声認識、画像認識、データ圧縮、データベース検索ランキング、バイオインフォマティクス、ネットワークセキュリティパターンの認識、空間ナビゲーションなどを実行する）を管理してよいINTEL（登録商標）第10世代プロセッサ（例えば、アイスレークプロセッサ）のような1又は複数のプロセッサ102を介して特定のコード又はコンピュータ命令を実行してよい。本明細書における命令という用語は、マクロ命令、例えば実行のためにプロセッサ102に提供される命令、又はマイクロ命令、例えばマクロ命令をデコードするプロセッサ102のデコーダから得られる命令を指してよいことに留意されたい。デコーダは、プロセッサ102のコアに含まれてよい。

30

【0018】

プロセッサ（複数可）102は、メモリ及び／又は記憶回路104と通信し得る。メモリ及び／又は記憶回路104は、ランダムアクセスメモリ（RAM）、リードオンリメモリ（ROM）、1又は複数のハードドライブ、フラッシュメモリ、又は任意の他の好適な光、磁気、若しくはソリッドステート記憶媒体のような有形で非一時的な機械可読媒体であってよい。メモリ及び／又は記憶回路104は、プロセッサ実行可能制御ソフトウェア、構成ソフトウェア、システムパラメータ、構成データ等のようなデータ処理システム100により処理されるデータを保持してよい。

40

【0019】

また、データ処理システム100は、データ処理システム100が他の電子デバイスと通信することを許可するネットワークインタフェース106を含んでよい。いくつかの実

50

施形態では、データ処理システム 100 は、様々な異なる要求を処理するデータセンタの一部であってよい。例えば、データ処理システム 100 は、ネットワークインタフェース 106 を介してデータ処理要求を受信して、DNN 計算、機械学習、ビデオ処理、音声認識、画像認識、データ圧縮、データベース検索ランキング、バイオインフォマティクス、ネットワークセキュリティパターンの認識、空間ナビゲーション、又はいくつかの他の専門的なタスクを実行してよい。また、データ処理システム 100 は、ディスプレイデバイス（例えば、コンピュータモニタ）、キーボード、マウス、スピーカ、音声入力デバイスなどのような、情報の入力及び / 又は表示に有用な 1 又は複数の入力 / 出力システム 108 を含んでよい。

【0020】

示される実施形態では、プロセッサ 102 は、動作可能及び / 又は通信可能に再構成可能なシストリックアレイシステム 110 に接続されてよい。再構成可能なシストリックアレイシステム 110 は、複数の処理要素（PE）、及び再構成可能なシストリックアレイシステム 110 の PE のいくつか（又はすべて）を通じてデータ（例えば、データフロー）を再構成可能に移動するために使用されてよい再構成可能なルーティングシステム 112 を含む、データをルーティングするのに好適な特定の回路を含んでよい。従って、DNN アプリケーションに使用されるデータのようなデータは、例えば、プロセッサ 102 を介して再構成可能なシストリックアレイシステム 110 に提供されてよく、次いで、再構成可能なシストリックアレイシステム 110 は、例えば、再構成可能なルーティングシステム 112 を介して、以下でさらに説明するように改善されたデータフローをより柔軟に導出してよい。再構成可能なシストリックアレイシステム 110 は、さらに、特定のバイアスデータを累積及び加算するのに好適な、累積システム 114 を備えるバイアス加算を含んでよい。例えば、累積システム 114 を備えたバイアス加算は、解かれる問題の終わりまで部分的な計算値（例えば、行列バイアス値）を累積し、最終出力をストレージに書き込んでよい。

【0021】

シストリックアレイシステムを記載することは有益であってよい。ここで図 2 を参照すると、図は、シストリックアレイシステム 200 の処理要素（PE）を通じたデータフローを介して、DNN ベースの問題のような特定の問題を解くために使用されてよいシストリックアレイシステム又は回路 200 を示すブロック図である。例えば、シストリックアレイシステム 200 は、 $C + = A * B$ （例えば、更新

【数 1】

$$c_{ij} = c_{ij} + \sum_{l=0}^{K-1} a_{il} * b_{lj}$$

ただし、K は行列の行高さである）のような様々な計算を計算するのに使用されてよい。

【0022】

示される実施形態では、データストレージ（例えば、複数のレジスタ、キャッシュ、バッファ等を有するレジスタファイル）202 は、タイルデータのような行列 A、B、C のデータを格納するのに使用されてよい。データストレージは、ライン 204、206、208、及び 210 を使用して、行列 A のタイルデータ、行列 B のタイルデータ、行列 C のタイルデータ、更新された行列 C のタイルデータをそれぞれ通信してよい。ライン 204、206、208、及び 210 のそれぞれは、複数の導管を含んでよいことに留意されたい。すなわち、ライン 204、206、208、及び 210 は、それぞれポートであってよく、各ポートは複数の導管又はラインを有してよい。ルーティング回路 212 は、行列 A の行 0 及び列 0 に対応する値 $A[0][0]$ を受信してよく、次いで、ルーティング回路 212 は、第 1 値 $A[0][0]$ を、処理要素 214、216、218 などのようなシストリックアレイシステム 200 の第 1 行内の処理要素にブロードキャストしてよい。ルーティング回路 212 は、さらに、B 内の第 1 行の値を表す値 $B[0][0]$ 、 $B[0][1]$ 、 $B[0][2]$ 、...、 $B[0][K]$ を受信して、値を処理要素 214、216

10

20

30

40

50

、218などにブロードキャストしてよい。例えば、処理要素214は、値 $B[0][0]$ を受信してよく、処理要素216は、値 $B[0][1]$ を受信してよく、処理要素218は、値 $B[0][K]$ を受信してよい。次いで、所与の行に対する処理要素のいくつか又はすべては、受信した入力に基づいて、乗算及び加算演算のような特定の演算の結果を出力してよい。次いで、例えば、処理要素214は、乗算 $A[0][0] * B[0][0]$ の結果を出力してよく、処理要素216は、乗算 $A[0][0] * B[0][1]$ の結果を出力してよく、処理要素218は、乗算 $A[0][0] * B[0][K]$ の結果を出力してよい。次いで、処理要素214、216、218の出力は、ルーティング回路220に送信されてよい。

【0023】

ルーティング回路220は、行列Aの行0及び列1に対応する値 $A[0][1]$ を受信してよく、次いで、ルーティング回路220は、値 $A[0][1]$ を、処理要素222、224、226などのようなシストリックアレイシステム200の第2行内の処理要素にブロードキャストしてよい。同様に、ルーティング回路220は、さらに、B内の第2行の値を表す値 $B[1][0]$ 、 $B[1][1]$ 、 $B[1][2]$ 、...、 $B[1][K]$ を受信して、値を処理要素222、224、226などにブロードキャストしてよい。例えば、処理要素222は、値 $B[1][0]$ を受信してよく、処理要素224は、値 $B[1][1]$ を受信してよく、処理要素226は、値 $B[1][K]$ を受信してよい。次いで、所与の行に対する処理要素のいくつか又はすべては、受信した入力に基づいて、乗算演算のような特定の演算の結果を出力してよい。次いで、例えば、処理要素222は、処理要素214の出力に加算された乗算 $A[0][1] * B[1][0]$ の結果を出力して、出力 $A[0][1] * B[1][0] + A[0][0] * B[0][0]$ に到達してよい。次いで、同様に、処理要素224は、処理要素216の出力に加算された乗算 $A[0][1] * B[1][1]$ の結果を出力して、出力 $A[0][1] * B[1][1] + A[0][0] * B[0][1]$ に到達してよい。次いで、同様に、処理要素226は、処理要素218の出力に加算された乗算 $A[0][1] * B[1][K]$ の結果を出力して、出力 $A[0][1] * B[1][K] + A[0][0] * B[0][K]$ に到達してよい。そのような積和演算は、融合型積和と称されてよく、各処理要素に含まれる融合型積和ユニット(FMA)を使用してよい。次いで、処理要素222、224、226の出力は、ルーティング回路228に送信されてよい。

【0024】

同様の態様で、ルーティング回路228及び230は、行列Aのデータ $A[0][2]$ 及び $A[0][3]$ をそれぞれ受信し、データをそれらのそれぞれの行の処理要素、例えば、ルーティング回路228に対する処理要素232、234、236及びルーティング回路230に対する処理要素238、240、242にブロードキャストしてよい。同様に、ルーティング回路228及び230は、行列の第3及び第4行に対する行列Bのデータを受信して、第3行のデータを処理要素232、234、236に、第4行のデータを処理要素238、240、242にそれぞれパスしてよい。また、処理要素232、234、236、238、240、及び242は、FMA機能、従って、シストリックアレイシステム200における行列Aの入力、行列Bの入力、及び前の処理要素の出力を含む受信した入力に基づく上述のような乗算及び加算を提供してよい。実際、示されるすべての処理要素は、融合型積和ユニットを含んでよい。

【0025】

次いで、バイアス加算回路244は、例えば、前に行列A、Bに対して実行された演算、例えば、 $C += A * B$ (例えば、行列Cからのバイアスを処理要素238、240、242からのそれぞれの結果に加算する)を用いて行列Cを加算及び/又は更新するために使用されてよい。例えば、ライン(複数可)208を介して受信された行列Cの値は、処理要素238、240、242などの出力に加算されて、ライン(複数可)210を介して更新行列Cとして格納されてよい。され得る。シストリックアレイシステム200の実施形態は、4行の処理要素を有するものとして示されているが、他の実施形態は、より多

10

20

30

40

50

い行又はより少ない行を含んでよいことを理解されたい。特定の実施形態では、シストリックアレイシステム 200 は、行ごとに 32 個の処理要素を使用してよい。例えば、128 の行列幅を有する疎 DNN ワークロードを処理する場合、処理される行列は、タイルごとに 32 列を有する 4 つのタイルに分割されてよい。すべてのタイルは、例えばデータストア 202 に書き込まれた部分的な結果を有し、次いで、次のタイルの結果に加算するために読み戻されてよい。従って、4 つの書き込み及び 4 つの読み取りをリーチタイルに使用して、疎 DNN ワークロードを完了してよい。データストア 202 の容量が増加するにつれて、使用される電力及びレイテンシが増大し得る。

【0026】

4、16、及び / 又は 36 の行列サイズを有するワークロードのような疎 DNN ワークロードの微分中、行列は、ブロックスパースティ、疎列 / 行 (CSC / CSR) の圧縮、直接インデキシング / ステップインデキシングなどのような技術を介してスパースティ圧縮を受けて、サイズ 32 の行列を得てよい。シストリックアレイシステム 200 がゼロを使用して「パディング」される場合、シストリックアレイシステム 200 は、第 1 パス上で (36 要素から) 32 要素幅を有するファイル幅タイルを処理してよく、残りの 4 つの要素を有するタイルが続ぎ、16 要素のタイルが続ぎ、次いで 4 要素のタイルが続く。従って、シストリックアレイシステム 200 の全体的な効率性は、 43.75% であってよく、これは、 $32 / 32 = 100\%$ 、 $4 / 32 = 12.5\%$ 、 $16 / 32 = 50\%$ 、及び $4 / 32 = 12.5\%$ の平均を見つけることによって計算されてよい。例えば、部分的な累積サポートを有する再構成可能なシストリックアレイを使用することにより、疎 DNN 並びに疎 DNN ワークロードの両方の処理を向上することが有益であってよい。

【0027】

ここで図 3 を参照すると、図は、スケジューラ 302 を介して複数の行列乗算を処理するのに好適な部分的なバイアス累積サポート (例えば、既存のタイルストレージとは別個のバイアス累積ストレージ) を含む再構成可能なシストリックアレイ回路又はシステム 300 の実施形態のブロック図である。例えば、スケジューラは、ホストプロセッサ (CPU)、例えばプロセッサ 102 内のソフトウェアとして、ハードウェア回路として、又は再構成可能なシストリックアレイシステム 300 に動作可能に接続されたそれらの組み合わせとして実装されてよい。示される実施形態では、スケジューラ 302 は、行列、例えば、タイプ A 304、B 306、及び / 又は C 308 の行列が、再構成可能なシストリックアレイシステム 300 への処理のために送信される順序をスケジューリングしてよい。

【0028】

スケジューラ 302 は、再構成可能なシストリックアレイシステム 300 を介して実行するためにタイルを送信する前に、行列 A 304、B 306 の特定のタイルをリオーダしてよい。また、スケジューラ 302 は、タイルをサブタイルにサイズ変更又は「ブレイク」して、再構成可能なシストリックアレイシステム 300 に含まれるバイアス累積ストレージ及びロジックを利用してよい。サブタイルに分割されていないタイルは「完全」タイルと称されてよく、完全タイルの処理にはバイアス累積を使用しなくてよい。一例では、結果行列 (例えば、行列 C 308) を通信するための x 個の読み取り / 書き込みポートがある場合、スケジューラは、任意の所与の時間に x 個より多くない完全タイルをスケジューリングしてよい。サブタイルに分割されているタイルは、「不完全」タイルと称されてよい。不完全タイルは、例えば、最後のサブタイルがスケジューリングされて、最終結果がストレージに書き出されるまで、バイアス累積ストレージ内に累積されてよい。本明細書で説明されるシステム及び方法は、以下でさらに説明されるように、完全及び不完全タイルの両方を処理し、どのタイルが完全又は不完全であることを示し、タイルの次元を示す新しいマクロ命令を含んでよい。

【0029】

また、本明細書に説明されるシステム及び方法は、スケジューラ 302 の出力に基づいてより小さい行列サイズを有する場合に行列データの再レイアウトをサポートしてよく、それにより、例えば、複数の A 行列 304 をサイドバイサイドに格納及び / 又は処理して

10

20

30

40

50

いる間に単一の A タイルがフェッチされてよい。示される実施形態では、A 1 及び A 2 は同一の A 行列に属してよく、A 1' 及び A 2' は別の A 行列に属してよい。アプリケーションに応じて、同一の B 行列を複製又はコピーするか、複数の B 行列を「繋ぎ合わせ」することのいずれかにより、B タイルが形成されてよい。図示された例では、B 1 及び B 1' は異なる B 行列からのものである。しかし、B 1 は、レプリケートされて、それにより特定のアプリケーションに対して $B 1 = B 1'$ となつてよい。

【0030】

特定の実施形態では、タイプ C 308 の行列は、入力バッファから読み取られてよく、入力バッファの帯域幅は、サイクルあたり x 読み取りに限定されてよい。従つて、スケジューラ 302 は、再構成可能なシストリックアレイ 300 のすべてのパスで実行するために最大 N 個の完全タイルをスケジューリングしてよく、従つて、C タイプの行列 308 の帯域幅の利用を向上する。従来の行列乗算では、 $C 1 + = A 1 * B 1 + A 1' * B 1' + \dots$ である。しかし、A と A' が「接着」又はともにマージされた異なる完全タイルの行列である場合、異なる演算が使用されてよい。代わりに、ハードウェア（又はソフトウェア）は、出力要素当たりより少ない演算を実行してよく、例えば $C 1 + = A 1 * B 1$ 及び $C 1' + = A 2 * B 1'$ である。しかし、通常の行列乗算より多くの出力要素があつてよい。これらの追加の出力要素は、バイアス累積回路内のストレージ又はレジスタに格納されてよいが、複数の独立した宛先が使用されて、例えばスケジューラ 302 から来る「完全」及び「不完全」タイルビットに基づいてストレージに書き込んでよい。

10

20

【0031】

前に言及されたように、本明細書で説明されるシステム及び方法は、バイアス加算累積を伴う再構成可能な行列乗算に好適な 1 又は複数のマクロ命令を提供してよい。新しい命令セットは、TPNDPMAC、『「M」累積によるタイル部分「N」ドット積』（tile partial 'N' dot product with 'M' accumulate）を含んでよい。ここで、N はともにマージされていてよい異なる行列の数であり、M は不完全な行列の数である（例えば、バイアス累積回路を使用してよい行列）。例えば、再構成可能なシストリックアレイシステム 300 に入力するための A としてマージされる 2 つの行列と、バイアス累積ロジックを使用するであろう 1 つと、1 つの B タイルにマージされる 2 つの行列との場合に、使用する命令は TP2DP1AC であろう。

30

【0032】

一実施形態では、命令に対するフォーマットは、TPNDPMAC tsrcdest、tsrc1、tsrc2 である。N = 1 の場合に、tsrcdest により指し示される単一の行列 C のソース / 宛先があつてよい。N > 1 の場合、複数の C タイルが連続してよく、tsrcdest（例えば、tsrcdest が tmm0 で N = 2 の場合は tmm0 及び tmm1）で開始し、tsrc1 が続き、次いで tsrc2 で複数のレジスタのグループを選択する。TPNDPMAC 命令は、図 4 に関連して説明されるように、再構成可能なシストリックアレイシステム 300 を使用して実装されてよい。

40

【0033】

図 4 は、特定のルーティング再構成及びバイアス累積に好適な再構成可能なシストリックアレイ回路又はシステム 300 の実施形態を示すブロック図である。示される実施形態では、再構成可能なシストリックアレイシステム 300 の特定のコンポーネントは、シストリックアレイシステム 200 に見つけられるそれらと同様に振る舞つてよい。例えば、データストレージ（例えば、複数のレジスタを有するレジスタファイル）402 は、タイルデータのような行列タイプ A 304、B 306、C 308 のデータを格納するのに使用されてよい。データストレージ 402 は、ライン 404、406、408、及び 410 を使用して、行列 A のタイルデータ、行列 B のタイルデータ、行列 C のタイルデータ、更新された行列 C のタイルデータをそれぞれ通信してよい。ライン 404、406、408、及び 410 のそれぞれは、複数の導管を含んでよいことに留意されたい。すなわち、ライン 404、406、408、及び 410 は、それぞれポートであつてよく、各ポートは複数の導管又はラインを有してよい。ルーティング回路 412 は、行列 A の行 0 及び列 0 に

40

50

対応する値 $A[0][0]$ を受信してよく、次いで、ルーティング回路 412 は、第 1 値 $A[0][0]$ を、処理要素 414、416、418 などのようなシストリックアレイシステム 200 の第 1 行内の処理要素にブロードキャストしてよい。ルーティング回路 412 は、さらに、B 内の第 1 行の値を表す値 $B[0][0]$ 、 $B[0][1]$ 、 $B[0][2]$ 、...、 $B[0][K]$ を受信して、値を処理要素 414、416、418 などにブロードキャストしてよい。例えば、処理要素 414 は、値 $B[0][0]$ を受信してよく、処理要素 416 は、値 $B[0][1]$ を受信してよく、処理要素 418 は、値 $B[0][K]$ を受信してよい。次いで、所与の行に対する処理要素のいくつか又はすべては、受信した入力に基づいて、乗算演算のような特定の演算の結果を出力してよい。次いで、例えば、処理要素 414 は、乗算 $A[0][0] * B[0][0]$ の結果を出力してよく、処理要素 416 は、乗算 $A[0][0] * B[0][1]$ の結果を出力してよく、処理要素 418 は、乗算 $A[0][0] * B[0][K]$ の結果を出力してよい。次いで、処理要素 414、416、418 の出力は、ルーティング回路 420 に送信されてよい。

10

【0034】

ルーティング回路 420 は、データを処理要素 422、424、426 にルーティングしてよく、データは順番に処理要素 414、416、及び 418 からカスケード「ダウン」するため、FMA 技術を適用してデータを乗算及び加算してよい。同様に、ルーティング回路 428 は、データを処理要素 430、432、434 にルーティングしてよく、データは順番に処理要素 422、424、及び 426 からカスケード「ダウン」するため、FMA 技術を適用してデータを乗算及び加算してよく、ルーティング回路 436 は、データを処理要素 438、440、442 にルーティングしてよく、データは順番に処理要素 430、432、及び 434 からカスケード「ダウン」するため、FMA 技術を適用してデータを乗算及び加算してよい。

20

【0035】

示される実施形態は、再構成可能なルーティング回路 444 (例えば、構成スイッチを備えたルーティング回路) を含む。ルーティング回路 412、420、428、436 と異なり、再構成可能なルーティング回路 444 は、少なくとも 2 つの演算モードに基づいて異なる方法でデータをルーティングしてよい。例えば、第 1 の演算モードでは、再構成可能なルーティング回路 444 に含まれる構成スイッチは、オンされてよく、導出されるドット積 (例えば、 $A * B$) のチェーンの「ブレーク」が生じ、新しいチェーンを開始してよい。構成スイッチがオフにされる場合、再構成可能なシストリックアレイ 300 は、1 出力を有する単一のパイプラインとして振る舞ってよい。従って、値が、処理のためにパイプラインのトップ (例えば、再構成可能なシストリックアレイシステム 300 の第 1 行) にて挿入される場合、結果は「カスケード」し、オンされる構成スイッチを有する再構成可能なルーティング回路 444 に結果が遭遇するまで下方に流れてよい。このステージにて、パイプラインは、結果値を「ブレーク」して、累積回路 446 を備えた第 1 バイアス加算に書き込まれてよい。結果値に対応する行列 C の要素に加算した後、更新値は、書き出されてよく、パイプ内の次のステージは、前の処理要素の出力値がゼロであったかのように、ロードされる。従って、オンされる構成スイッチを有する再構成可能なルーティング回路 444 とのカスケディング値の遭遇は、新しいパイプラインの開始と見なされてよい。複数の再構成可能なルーティング回路 444 が使用されてよく、例えば、再構成可能なルーティング回路 444 は、第 4 行ごとなど、8 行の再構成可能なシストリックアレイシステム 300 内に配置されてよく、従って、複数の再構成可能なルーティング回路 444 が使用されてよいことを理解されたい。

30

40

【0036】

一実施形態では、第 1 の演算モードにある場合、データストア 402 の第 1 の複数のレジスタに格納された値は、単一入力の 2 次元行列 A を表してよく、データストア 402 の第 2 の複数のレジスタに格納された値は、単一入力の 2 次元行列 B を表してよく、一方、データストア 402 の第 3 の複数のレジスタに格納された値は、単一入力の 2 次元行列 C を表してよい。第 2 の演算モードにある場合、データストア 402 の第 1 の複数のレジス

50

タに格納された値は、複数入力の2次元行列A及びA'を表してよく、データストア402の第2の複数のレジスタに格納された値は、複数入力の2次元行列B及びB'を表してよく、一方、データストア402の第3の複数のレジスタに格納された値は、複数入力の2次元行列C及びC'を表してよい。

【0037】

特定の実施形態では、第1の演算モードの実行中、再構成可能なシストリックアレイシステム300は、タイルA及びタイルBからそれぞれのルーティング回路に値を送信してよい。例えば、演算は、第1の演算モードの場合に、タイルAからの行列AにタイルBからの行列Bを乗算し、次いで、タイルCからの行列Cの対応する値にそれぞれの結果を加算し、第2の演算モードの場合に、タイルAからの行列AにタイルBからの行列Bを乗算し、次いで、タイルCからの行列Cの対応する値にそれぞれの結果を加算するとともに、タイルAからの行列A'にタイルBからの行列B'を乗算し、次いで、タイルCからの行列C'の対応する値にそれぞれの結果を加算してよい。第1の演算モードでは、処理要素438、440、442の出力は、累積回路446を備えた第1パイアス加算をバイパスしてよく、処理要素448、450、452に直接提供されてよい。次いで、処理要素448、450、452は、上述のような乗算及び加算を適用し、次いで、それぞれの出力を累積回路454を備えた第2パイアス加算に提供してよい。次いで、累積回路454を備えた第2パイアス加算は、処理要素448、450、452から提供された出力を使用して、行列Cを更新してよい。

10

【0038】

第2の演算モードでは、処理要素438、440、442の出力は、累積回路446を備えた第1パイアス加算により使用されて、例えば、特定の値を加算し格納してよい。前に言及したように、再構成可能なルーティング回路444がオンされた構成スイッチを有する場合、再構成可能なルーティング回路444は、入力として提供される値を乗算及び加算し、結果を更新行列Cに送信してよいが、結果（例えば、乗算及び加算の結果）を後の微分で使用するために累積してもよい。第2の演算モードでは、処理要素448、450、452は、処理要素438、440、442の出力に代えてゼロを受信してよく、従って、処理要素448、450、452にて開始する演算は新しいパイプラインとして進んでよい。累積回路454を備えた第2パイアス加算は、第1の演算モード（例えば、値の累積をバイパスする）を提供するためにスイッチオフされる、又は第2の演算モードに対してスイッチオンされる累積スイッチを有してよい。

20

30

【0039】

図5は、再構成可能なルーティング回路444及び累積回路501を備えたパイアス加算（例えば、回路446又は454と同等）の実施形態を示す概略図である。示される実施形態では、再構成可能なシストリックアレイシステム300の処理要素の行（例えば、行3）からのデータ500は、再構成可能なシストリックアレイシステム300の下流行502（例えば、行4）に提供されてよい。処理要素の下流行502は、行列Bのデータ504及び行列Aのデータ506を受信してもよい。次いで、行502内の処理要素は、例えばライン508を介して再構成可能なルーティング回路444への出力を提供してよい。

40

【0040】

再構成可能なルーティング回路444は、デマルチプレクサ510及びマルチプレクサ512を含んでよく、それにより、デマルチプレクサ510及びマルチプレクサ512の両方がスイッチとして使用される。すなわち、デマルチプレクサ510及びマルチプレクサ512は、同一の信号（例えば、構成オン又はオフ信号）を受信し、ともにデータルーティングのスイッチとして機能してよい。再構成可能なルーティング回路444がデマルチプレクサ510及びマルチプレクサ512へのセクタを介してオンされる場合、デマルチプレクサ510は行502の処理要素を介して導出される出力をライン514を介して累積回路501を備えたパイアス加算に書き込んでよい。順番に、マルチプレクサ512は、ゼロを下流行516（例えば、行5）の処理要素に、例えばライン518を介して

50

送信してよい。従って、行 5 1 6 の処理要素は、行 5 0 2 からのデータを使用しなくてよく、代わりに、行列 B のデータ 5 0 4 及び行列 A のデータ 5 0 6 を使用して出力 5 2 0 を導出してよく、次いで、その出力は次の下流行（例えば、行 6）送信されてよい。

【0041】

累積イネーブル信号 5 2 2 がオンされる場合、累積回路 5 0 1 を備えたバイアス加算は、バイアス 5 2 4 を C タイル 5 2 6 に加算するとともに結果を格納そうでなければ累積してよい。累積イネーブル信号 5 2 2 は、アドレスチェック信号 5 3 2 によって受信される累積イネーブル信号 5 3 0（例えば、マクロ命令の実行に基づく信号）の Boolean OR を導出する OR ゲート 5 2 8 を使用することによりオンされてよい。アドレスチェック信号 5 3 2 は、行列 C のタイルアドレス 5 3 4 の表示であってよい。より具体的には、マイクロアーキテクチャのサポートが提供されてよく、それにより、ハードウェア内で C タイルアドレス 5 3 4 がチェックされて、宛先の衝突が発生するか、例えば 2 つの行列演算は同一の行列 C の宛先アドレスを共有するか、について判断する。次いで、アドレスが同一の場合、累積ロジックは自動的にオンされて、例えば、宛先の上書きを防止する。最後のサブタイトルビット 5 3 6 が受信されると（例えば、スケジューラ 3 0 2 から到着すると）、累積回路 5 0 1 を備えたバイアス加算は、例えば、すべてのレジスタに亘って、すべての累積値を加算してよい。すなわち、最後のサブタイトルビット 5 3 6 は、すべてのサブタイトルがここで送信されていて、従って、任意の累積値がここで累積回路 5 0 1 を備えたバイアス加算を介して加算及び格納されてよいことを示してよい。

10

【0042】

再構成可能なルーティング回路 4 4 4 が（例えば、デマルチプレクサ 5 1 0 及びマルチプレクサ 5 1 2 へのセレクタを介して）オフにされる場合、デマルチプレクサ 5 1 0 は、行 5 0 2 の処理要素によりライン 5 3 5 を介してマルチプレクサ 5 1 2 に導出される出力を送信してよい。次いで、マルチプレクサ 5 1 2 は、行 5 0 2 の処理要素により導出される出力をライン 5 1 8 を介して下流行 5 1 6 に送信してもよい。従って、再構成可能なルーティング回路 4 4 4 をオフすると、再構成可能なルーティング回路 4 4 4 が、行 5 0 2 の処理要素及び行 5 1 6 の処理要素の間のパススルースイッチとして機能するようになってよい。再構成可能なルーティング回路 4 4 4 を提供することにより、本明細書で説明される技術は、再構成可能なシストリックアレイシステム 3 0 0 を通じてデータのより効率的なルーティングを可能にしてよい。

20

30

【0043】

図 6 は、さらに詳細に示す累積回路 5 0 1 を備えたバイアス加算の実施形態の概略図である。累積回路 5 0 1 を備えたバイアス加算は、メモリストレージの一致するステージを使用することにより、特定のレイテンシ（例えば、加算器のレイテンシ）を説明するように設計されてよい。例えば、3 つのステージが使用されて 3 のレイテンシを一致させてよく、4 つのステージが使用されて 4 のレイテンシ（例えば、 $2^3 + 1$ ）などを一致させてよい。示される実施形態では、カウンタ 6 0 0 は、レイテンシに基づいてカウントするために使用されてよい。従って、3 ビットのカウンタは 3 のレイテンシに使用されてよく、4 ビットのカウンタは 7 より大きいレイテンシなどに使用されてよい。従って、適切な値を多重選択するためのサイズが、増加してもよい。累積回路 5 0 1 を備えた加算への入力

40

【0044】

演算中、加算器 6 0 8 は、ドット積 6 0 2 をマルチプレクサ 6 1 0 からの出力と加算してよい。マルチプレクサ 6 1 0 の出力は、AND ゲート 6 1 2 からの信号を介して選択されてよい。AND ゲート 6 1 2 は、カウントリセット信号 6 1 4 及びクロックフリップフロップ 6 1 6 からの出力の間の Boolean AND を実行してよい。クロックフリッ

50

ブフロップ 616 は、ANDゲート 618 から出力されるデータを格納してよい。例えば、ANDゲート 618 は、カウンタ 600 の出力及び最後のサブタイトル信号 536 の間の Boolean AND 演算を実行してよい。累積イネーブル信号 522 がオンの場合、カウンタ 600 は、例えば、デマルチプレクサ 625 を選択することによりクロック信号 606 が送信されると、記憶回路（例えば、フリップフロップのような記憶コンポーネント）620、622、624 へのドット積 602 のストレージを指示してよい。

【0045】

次いで、最後のサブタイトル信号 536 は、ストレージ 620、622、624 が、ANDゲート 626、628、630 を通じて累積されたデータ値をパスし、加算器 632、634 を介して加算されるようにしてよい。次いで、累積イネーブル信号 522 を使用して、累積されたデータ値がマルチプレクサ 636 の出力として選択されてよい。加算器 638 は、使用されて、累積されたデータ値をマルチプレクサ 642 の加算バイアスセレクト信号を介して行列 C のバイアス 604 に加算してよい。行列 C のバイアス 604 は、ストレージ 643 から到着してよい。次いで、加算の結果は、例えばライン 410（図 4 に示される）を介して、更新された行列 C に提供されてよい。前に言及したように、累積回路 501 を備えたバイアス加算は、特定のレイテンシを考慮して設計されてよい。図示された例では、3つのストレージ 620、622、624 は、3又はより小さいレイテンシを処理してよい。しかし、演算中にレイテンシが増加し得る場合がある。レイテンシが増加した場合、ライン 644 が使用されて、例えば加算器 608 を介して、ストレージ 620、622、624 に格納されているより古い値で新しいドット積 602 を加算することにより、ループ内の値を連続的に累積してよい。累積イネーブル信号 522 がオフにされる場合、ドット積 602 は、デマルチプレクサ 646 をトラバースし、次いでマルチプレクサ 636 をトラバースし、その後、加算器 638 により行列 C のバイアス 604 に加算されてよい。累積回路 501 を備えたバイアス加算に提供することにより、本明細書で説明された技術は、密及び疎 DNN の両方をより効率的に処理するとともに、より柔軟なシストリックアレイベースの計算を提供してよい。

【0046】

前に言及したように、複数の再構成可能なルーティング回路 444 が使用されてよい。同様に、累積回路を備えた複数のバイアス加算、例えば、累積回路 501 を備えたバイアス加算が提供されてよい。ここで図 7 を参照すると、図は、複数のルーティング回路及び累積回路を備えた複数のバイアス加算（例えば、部分的なバイアス累積サポート）を含む再構成可能なシストリックアレイシステム 700 を示すブロック図である。示される実施形態では、シストリックアレイシステム 700 は、データストレージ 702（例えば、複数のレジスタを有するレジスタファイル）を含む。データストレージ 702 は、ライン 704、706、708、及び 710 を使用して、行列 A のタイルデータ、行列 B のタイルデータ、行列 C のタイルデータ、更新された行列 C のタイルデータをそれぞれ通信してよい。ライン 704、706、708、及び 710 のそれぞれは、複数の導管を含んでよいことに留意されたい。すなわち、ライン 704、706、708、及び 710 は、それぞれポートであってよく、各ポートは複数の導管又はラインを有してよい。

【0047】

示される実施形態はまた、8つの回路ブロック 712、714、716、718、720、722、724、726 を含む。回路ブロック 712、714、716、718、720、722、724、726 のそれぞれは、1又は複数行の処理要素を含んでよく、ここで、処理要素は融合型乗算加算ユニット（FMA）を含んでよい。一実施形態では、再構成可能なシストリックアレイシステム 700 が 32 行の処理要素を有する場合のように、回路ブロック 712、714、716、718、720、722、724、726 のそれぞれは 4 行の処理要素を含んでよい。データが第 1 回路ブロック 712 に入ると、データは、カスケード方式で処理されてよく、その後、回路ブロック 714、716、718、720、722、724、及び 726 を通じてカスケード順序で流れ、例えば $C + = A * B$ を計算してよい。

10

20

30

40

50

【0048】

図示のとおり、再構成可能なルーティング回路728、730、732、734、736、738、740は、回路ブロック712、714、716、718、720、722、724の下流に配置されてよい。再構成可能なルーティング回路728、730、732、734、736、738、740のそれぞれは、例えばスイッチングを介して、累積回路742、744、746、748、750、752、754を備えた下流のバイアス加算へのデータのフローを可能にしてよい。再構成可能なルーティング回路728、730、732、734、736、738、740は、さらに、例えば前述したようにスイッチオンされる場合に「新」しいパイプラインの作成を可能にしてよい。累積回路742、744、746、748、750および752を備えたバイアス加算のそれぞれは、上の図6に説明したように、ドット積へのバイアスを加算及び累積値をバイアスに加算するのに好適であってよい。ルーティング回路757は、スイッチング機能を含まなくてよく、従って、回路ブロック726からの出力値を、累積しないでバイアス加算のために直接パスすることによって、データをバイアス加算回路756に送信してよい。従って、再構成可能なシストリックアレイシステム700は、より効率的に且つ柔軟に、 $C + = A * B$ を含む様々な計算を導出してよい。

10

【0049】

本明細書で説明された技術をプログラムで使用するために、特定の命令（例えば、マクロ命令）が提供される。例えば、TPNDPMACは、TP2DP1ACのような命令又は「1が累積されたタイル部分の2ドット積」（`tile partial 2 dot product with 1 accumulate`）のプログラムによる使用をもたらしてよい。TP2DP1AC命令は、処理要素のアレイの中央にある再構成可能なルーティング回路（例えば、再構成可能なルーティング回路734）をオンすることにより、また累積回路を備えた対応するバイアス加算（例えば、累積回路748を備えたバイアス加算）をスイッチオンすることにより、ともにマージされた2つの均一サイズの行列を処理してよい。

20

【0050】

異なるサイズの行列がともにマージされる場合、TSZDP「ドット積のタイルサイズ」（`tile sizes for dot products`）マクロ命令が使用されてよい。一実施形態では、TSZDPマクロ命令は、A、B、及びCレジスタオペランドに加えて、ともにマージされる行列のサイズを指定する即値を取ってよい。別の実施形態では、サイズはエンコードされてよい。例えば、4の倍数（例えば、最大32）を有する行列のようなマージング行列がサポートされる場合、様々な行列サイズを次のようにエンコードしてよい。

30

【0051】

【表 1】

即時エンコーディング	第 1 K サイズ (K1)	第 2 K サイズ (K2)	第 3 K サイズ (K3)	第 4 K サイズ (K4)	第 5 K サイズ (K5)	第 6 K サイズ (K6)	第 7 K サイズ (K7)	第 8 K サイズ (K8)	デコーダ 出力	構成スイッ チ
0000000	32	0	0	0	0	0	0	0	0000000	すべてのス イッチがダ ウン
0000001	4	28	0	0	0	0	0	0	0000001	スイッチ1が オン
0000010	8	24	0	0	0	0	0	0	0000010	スイッチ2 がオン
0000011	4	4	24	0	0	0	0	0	0000011	スイッチ1及 び2がオン
0000100	12	20	0	0	0	0	0	0	0000100	スイッチ3 がオン
.
1111110	8	4	4	4	4	4	4	0	1111110	スイッチ1を 除くすべて のスイッチ がオン
1111111	4	4	4	4	4	4	4	4	1111111	すべてのス イッチがオ ン

10

20

【 0 0 5 2 】

表 1 は、スイッチの使用に言及し、スイッチは、順番に、図 7 の同等の再構成可能なルーティング回路の使用に言及する。例えば、スイッチ 1 は、再構成可能なルーティング回路 7 2 8 を指してよく、スイッチ 2 は、再構成可能なルーティング回路 7 3 0、再構成可能なルーティング回路 7 3 2 など指してよい。即時エンコーディング値がゼロの場合、これはタイルサイズが 3 2 であり、A 及び B 入力の両方の入力として単一の行列が使用されることを意味する。即時エンコーディング値 1 1 1 1 1 1 0 は、第 1 のスイッチを除くすべての構成スイッチを有効にしてよく、それにより、完全なシストリックアレイ 7 0 0 は 7 つの独立した小さなアレイと見なしてよく、第 1 のアレイは 8 の行列サイズを処理することが可能であり、他のすべてはそれぞれ 4 の行列サイズを処理してよい。同様に、1 1 1 1 1 1 1 の即時エンコーディング値は、すべての構成スイッチを有効にしてよく、それにより、完全なシストリックアレイ回路 7 0 0 は、各回路が 4 の行列サイズを処理可能である小さなアレイの 8 つの独立した回路と見なされてよい。

30

【 0 0 5 3 】

本明細書では T A C D P 「タイル累積ドット積」(tile accumulate dot product) と称される、構成スイッチを制御する前述の命令に基づいて累積ロジックを有効及び無効にする命令が使用されてもよい。この命令は、適切な構成スイッチ値でのみ有効であってよいことに留意されたい(すなわち、構成スイッチが有効でない場合、パイプの最後にあるルーティング回路 7 5 7 の終わりの構成スイッチを除いて、累積ロジックは有効でなくてよく、構成スイッチを含まなくてよい)。アキュムレータは、T A C D P imm__ac# 形式の即値をパスすることにより又は命令を通じてパスされる即値を介して有効にされてよい(例えば、T P 2 D P tsrcdest、tsrc1、tsrc2、imm__sz#、imm__ac#)。この T A C D P 命令は、T S Z D P 命令(T S Z D P imm__sz#、imm__ac#) とマージされてもよい。T A C D P 即時エンコーディングは次のとおりでよい。

40

【表 2】

即時エンコーディング	蓄積スイッチ 1	蓄積スイッチ 2	蓄積スイッチ 3	蓄積スイッチ 4	蓄積スイッチ 5	蓄積スイッチ 6	蓄積スイッチ 7	蓄積スイッチ 8
0000000	0	0	0	0	0	0	0	0
0000001	1	0	0	0	0	0	0	0
0000010	0	1	0	0	0	0	0	0
0000011	1	1	0	0	0	0	0	0
0000100	0	0	1	0	0	0	0	0
.
.
1111110	0	1	1	1	1	1	1	1
1111111	1	1	1	1	1	1	1	1

【 0 0 5 4 】

累積回路（例えば、回路 7 4 2、7 4 4、7 4 6、7 4 8、7 5 0、7 5 2）を備えたバイアス加算は、累積スイッチ（例えば、累積イネーブル信号 5 2 2）がオンされる場合にのみ有効になってよく、そうでない場合、特定のデータが再構成可能なシストリックアレイ 7 0 0 のこのセクション（例えば、累積回路を備えたバイアス加算）に入らなくてよい。しかし、累積イネーブル信号を有さなくてよいバイアス加算回路 7 5 6 が使用される。すなわち、ルーティング回路 7 5 7 は、データをバイアス加算回路 7 5 6 のみに直接ルーティングしてよく、スイッチング機能を提供しなくてよい。いくつかの実施形態では、累積回路 7 5 2 を備えた最後のバイアス加算は、累積が常にスイッチオンされると仮定して動作されてよい。

【 0 0 5 5 】

累積ロジックは、マイクロアーキテクチャモード及びアーキテクチャモードの 2 つのモードにより有効にされてよい。マイクロアーキテクチャモードでは、例えば図 5 に示されるアドレスチェック 5 3 2 を介して、前の宛先（例えば、タイルレジスタ tmm0）及び現在の宛先アドレスが同一であると識別された場合、再構成可能なシストリックアレイ 7 0 0 及び関連ハードウェアは、累積を可能にしてよい。アーキテクチャモードでは、累積は、ロジックが T S Z D P 命令又は T A C D P のいずれかにより制御される命令により有効にされてよい。前に言及したように、最後のグループの処理要素行（例えば、図示の例における回路ブロック 7 2 6）に関連する構成スイッチがないため、構成スイッチなしで最後のバイアス加算回路 7 5 6 のみがオンされてよい。

【 0 0 5 6 】

図 8 は、本明細書で説明される技術を実装するために使用されてよい処理 8 0 0 の実施形態を示す。処理 8 0 0 は、再構成可能なシストリックアレイ 3 0 0、7 0 0 及びマクロ命令 T P N D P M A C、T S Z D P 及び / 又は T A C D P を介してのようなハードウェア及び / 又はソフトウェアとして実装されてよい。示される実施形態では、解かれる問題のタイルサイズを決定してよい（ブロック 8 0 2）。問題は、密 D N N、疎 D N N、又はそれらの組み合わせとともに、機械学習、ビデオ処理、音声認識、画像認識、データ圧縮、データベース検索ランキング、バイオインフォマティクス、ネットワークセキュリティパターンの認識、空間ナビゲーションなどにおける問題を含んでよい。例えば、使用されるシストリックアレイの行及び列の数（例えば、再構成可能なシストリックアレイシステム 3 0 0、7 0 0）に基づいて、タイルサイズを選択して、例えば加算されたゼロを最小化することによりアレイをより快適に合致させてよい。タイルサイズが選択されると（ブロック 8 0 2）、いくつかの完全及び / 又は不完全タイルが導出されてよい（ブロック 8 0 4）。完全タイル 8 0 6 は、それら全体において使用されるシストリックアレイに合致してよいが、不完全タイル 8 0 8 はサブタイルに細分化されてよい。

【 0 0 5 7 】

次いで、完全タイル 8 0 6 及び不完全タイル 8 0 8 が処理されてよい（ブロック 8 1 0）。例えば、例えば図 5 に示されるアドレスチェック 5 3 2 を介して、前の宛先（例えば

10

20

30

40

50

、タイルレジスタ $tmm0$) 及び現在の宛先アドレスが同一であると識別された場合、マイクロアーキテクチャモードが使用されて、使用されるシストリックアレイを実行し、自動的に宛先の衝突を検出し、累積ロジックをスイッチオンしてよい。アーキテクチャモードでは、累積は、TSZDP命令及び/又はTACDPにより有効にされてよい。次いで、計算の結果が提供されてよい(ブロック812)。例えば、 $C += A * B$ に基づく最終Cは、計算された行列Cのそれぞれに提供されてよい。本明細書で説明される回路(例えば、再構成可能なシストリックアレイシステム300、700)は、ハードウェアアクセラレータの一部として、フィールドプログラマブルゲートアレイ(FPGA)として、特定用途向け集積回路(ASIC)として、カスタムマイクロチップとして、又はそれらの組み合わせとしてマイクロプロセッサに実装されてよいことを理解されたい。

10

【0058】

本開示に記載される実施形態が様々な修正及び代替的形態に影響されやすい場合がある一方で、具体的な実施形態が、図面における例により示されており、本明細書において詳細に説明されている。しかしながら、本開示は開示された特定の形態に限定されることが意図されているものではないことが理解され得る。本開示は、以下に添付される特許請求の範囲により定義されるように、本開示の趣旨及び範囲内に含まれるすべての修正、均等物、および代替物をカバーするものである。

【0059】

本明細書によれば、以下の各項目に記載の構成もまた開示される。

[項目1]

データを格納するように構成されたデータストレージと、再構成可能なシストリックアレイ回路と、を備え、前記再構成可能なシストリックアレイ回路は、

20

前記データを処理するように構成された1又は複数のグループの処理要素を有する第1回路ブロックと、

前記データを処理するように構成された1又は複数のグループの処理要素を有する第2回路ブロックと、

行列バイアスを累積値に、乗算積に、又はそれらの組み合わせに加算するように構成された累積回路を備えた第1バイアス加算と、

微分を前記第1回路ブロックから前記第2回路ブロックに、前記第1回路ブロックから累積回路を備えた前記第1バイアス加算に、又はそれらの組み合わせにルーティングするように構成された第1ルーティング回路と、を含む、システム。

30

[項目2]

前記第1ルーティング回路は、構成スイッチ信号を受信することに基づいて、前記微分を前記第1回路ブロックから前記第2回路ブロックに、前記第1回路ブロックから累積回路を備えた前記第1バイアス加算に、又はそれらの組み合わせにルーティングするように構成されたデマルチプレクサ及びマルチプレクサ回路を含む、項目1に記載のシステム。

[項目3]

累積回路を備えた前記第1バイアス加算は、クロック信号に基づいて前記乗算積を前記累積値として累積するように構成された記憶回路、及び前記行列バイアスを前記累積値に、前記乗算積に、又はそれらの組み合わせに加算するように構成された少なくとも1つの加算器を含む、項目1に記載のシステム。

40

[項目4]

累積回路を備えた前記第1バイアス加算は、Nの加算器のレイテンシを含み、前記記憶回路は、N個の記憶コンポーネントを含む、項目3に記載のシステム。

[項目5]

前記N個の記憶コンポーネントは、それぞれ、フリップフロップを含む、項目4に記載のシステム。

[項目6]

前記記憶回路は、前記N個の記憶コンポーネントをマルチプレクサに結合するN個のラ

50

インを含み、前記記憶回路は、前記加算器のレイテンシが演算中にNを超えた場合に、累積値を前記N個の記憶コンポーネントから前記マルチプレクサに前記N個のラインを介して転送するように構成される、項目4に記載のシステム。

[項目7]

累積回路を備えた前記第1バイアス加算は、累積回路を備えた前記第1バイアス加算に入る新しい値を前記累積値に加算し、前記加算の結果を前記N個の記憶コンポーネントに格納するように構成される、項目6に記載のシステム。

[項目8]

1又は複数のグループの処理要素を有する第3回路ブロックと、
第2行列バイアスを第2累積値に、前記乗算積に、又はそれらの組み合わせに加算するように構成された累積回路を備えた第2バイアス加算と、

微分を前記第2回路ブロックから前記第3回路ブロックに、前記第2回路ブロックから累積回路を備えた前記第2バイアス加算に、又はそれらの組み合わせにルーティングするように構成された第2ルーティング回路と、を備える、項目1に記載のシステム。

[項目9]

前記第3回路ブロックの下流に配置され、第3行列バイアスを前記第3回路ブロックから出力に加算するように構成されたバイアス加算回路を備える、項目8に記載のシステム。

[項目10]

前記再構成可能なシストリックアレイ回路を使用するよう、又は前記再構成可能なシストリックアレイ回路を含むように構成されたホストプロセッサ(CPU)を備え、

前記CPUは、『「M」累積によるタイル部分「N」ドット積』命令であり、前記Nは、ともにマージされている異なる行列の数であり、Mは、前記再構成可能なシストリックアレイ回路への入力として使用される不完全な行列の数である、前記命令、前記再構成可能なシストリックアレイ回路への入力として使用されるとともにマージされている前記異なる行列のサイズを指定する即値を有する「ドット積に対するタイルサイズ」命令、累積回路を備えた前記第1バイアス加算を制御する「タイル累積ドット積」命令、又はそれらの組み合わせを実行するように構成される、項目1に記載のシステム。

[項目11]

行列A及び行列Bに基づいてデータの1又は複数のタイルのそれぞれに対してタイルサイズを決定する段階と、

タイルサイズに基づいて、完全タイル、不完全タイル、又はそれらの組み合わせを導出する段階と、

行列Cの結果を導出するために、前記完全タイル、前記不完全タイル、又はそれらの組み合わせを、再構成可能なシストリックアレイ回路を介して処理する段階であり、前記完全タイル、前記不完全タイル、又はそれらの組み合わせを処理することは、前記行列Cの結果を提供するために、前記再構成可能なシストリックアレイ回路に含まれるルーティング回路及び前記再構成可能なシストリックアレイ回路に含まれる累積回路を備えたバイアス加算又はそれらの組み合わせを適用することを含む、段階と、を備える方法。

[項目12]

前記再構成可能なシストリックアレイ回路は、N行M列のアレイサイズを含み、前記完全タイルは、N行以下及びM列以下を有する完全サイズを含み、前記不完全タイルは、N行より多く、M列より多く、又はそれらの組み合わせを有する不完全サイズを含む、項目11に記載の方法。

[項目13]

前記ルーティング回路を適用する段階は、微分を1又は複数のグループの処理要素を含む第1回路ブロックから1又は複数のグループの処理要素を含む第2回路ブロックにルーティングする段階、及び微分を前記第1回路ブロックから累積回路を備えた前記バイアス加算に、又はそれらの組み合わせにルーティングする段階を含む、項目11に記載の方法。

。

10

20

30

40

50

[項目 1 4] 微分を前記第 1 回路ブロックから累積回路を備えた前記バイアス加算にルーティングする段階は、前記微分を累積回路を備えた前記バイアス加算にて受信し、前記微分を行列 C のバイアスに加算するために累積値に累積する段階を含む、項目 1 3 に記載の方法。

[項目 1 5]

前記完全タイル、前記不完全タイル、又はそれらの組み合わせを前記再構成可能なシストリックアレイ回路を介して処理する段階は、

行列 C のアドレス衝突を検出し、累積回路を備えた前記バイアス加算に通信される累積イネーブル信号を自動的にオンにするように構成されたマイクロアーキテクチャモードを適用する段階と、

前記再構成可能なシストリックアレイ回路への入力として使用されるとともにマージされている前記異なる行列のサイズを指定する即値を有する「ドット積に対するタイルサイズ」命令、累積回路を備えた前記バイアス加算を制御する「タイル累積ドット積」命令、又はそれらの組み合わせを実行することによりアーキテクチャモードを適用する段階と、を含む、項目 1 1 に記載の方法。

[項目 1 6]

データを格納するように構成されたデータストレージと、

再構成可能なシストリックアレイ回路と、

前記再構成可能なシストリックアレイ回路に接続されるコアのデコーダであり、単一の命令をデコードされた 1 又は複数の命令にデコードする、デコーダと、を備え、

前記 1 又は複数の命令は、

行列 A 及び行列 B を表す前記データを、前記データストレージから、前記データを処理し、前記データに基づく微分を提供するように構成された 1 又は複数のグループの処理要素を含む第 1 回路ブロックに通信し、

再構成可能なルーティング回路のスイッチングオン又はオフに基づいて、前記微分を前記第 1 回路ブロックから第 2 回路ブロックに、累積回路を備えたバイアス加算に、又はそれらの組み合わせにルーティングするように構成され、

累積回路を備えた前記バイアス加算は、行列バイアスを、累積値に、行列 A 及び行列 B の乗算積に、又はそれらの組み合わせに加算するように構成され、

前記第 1 回路ブロック、前記第 2 回路ブロック、前記再構成可能なルーティング回路、累積回路を備えた前記バイアス加算、又はそれらの組み合わせは、前記再構成可能なシストリックアレイ回路に含まれる、装置。

[項目 1 7]

前記単一の命令は、デコードされると、前記再構成可能なシストリックアレイ回路への入力として使用されるとともにマージされている異なる行列のサイズを指定する即値を有する「ドット積に対するタイルサイズ」命令、累積回路を備えた前記バイアス加算を制御する「タイル累積ドット積」命令、又はそれらの組み合わせを介してアーキテクチャモードを使用する、項目 1 6 に記載の装置。

[項目 1 8]

前記単一の命令は、『「M」累積によるタイル部分「N」ドット積』命令を含み、N は、ともにマージされている異なる行列の数であり、M は、前記再構成可能なシストリックアレイ回路への入力として使用される不完全な行列の数である、項目 1 7 に記載の装置。

[項目 1 9]

前記単一の命令は、デコードされると、前記再構成可能なシストリックアレイ回路に、前記データを使用することによって $C = + A * B$ を解かせ、前記データは、前記行列 A 及び前記行列 B を表す、項目 1 6 に記載の装置。

[項目 2 0] 前記再構成可能なシストリックアレイ回路を有する回路を備え、

前記回路は、マイクロプロセッサ、ハードウェアアクセラレータ、フィールドプログラマブルゲートアレイ (F P G A)、特定用途向け集積回路 (A S I C)、カスタムマイクロチップ、又はそれらの組み合わせを含む、項目 1 6 に記載の装置。

10

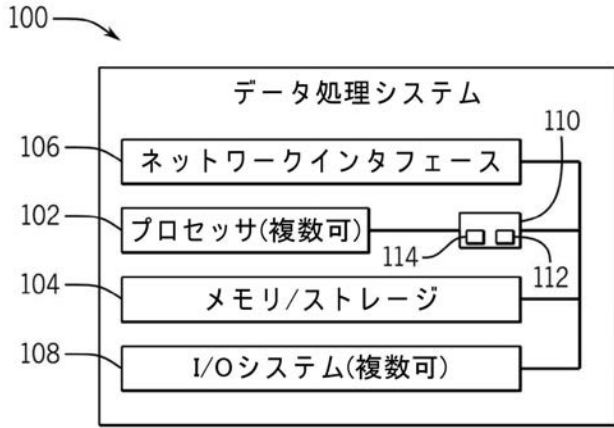
20

30

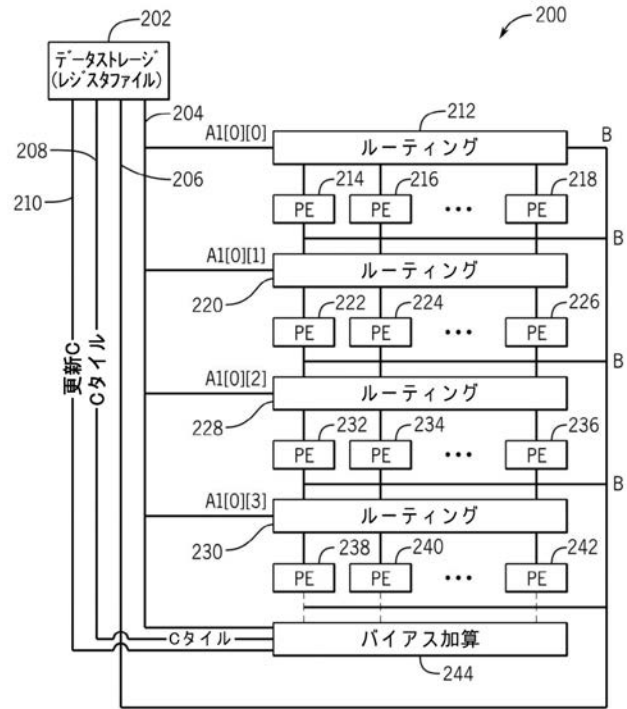
40

50

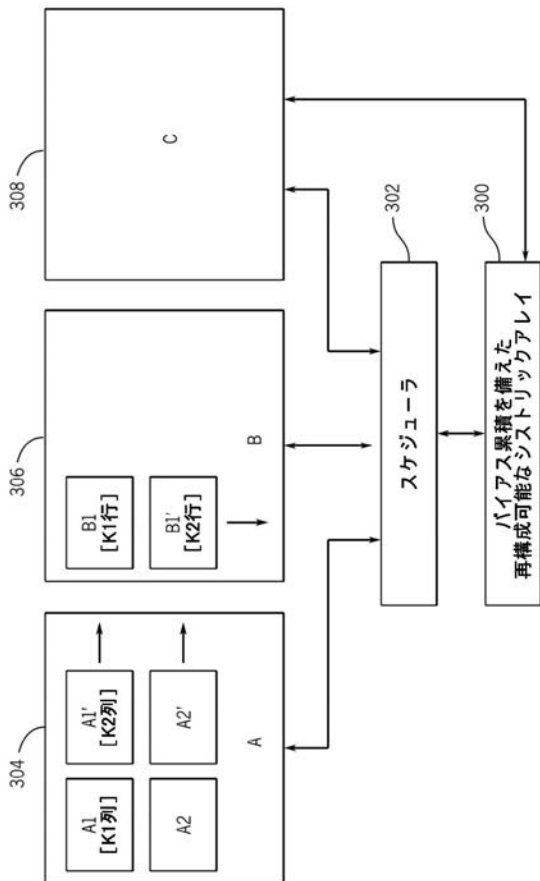
【 図 1 】



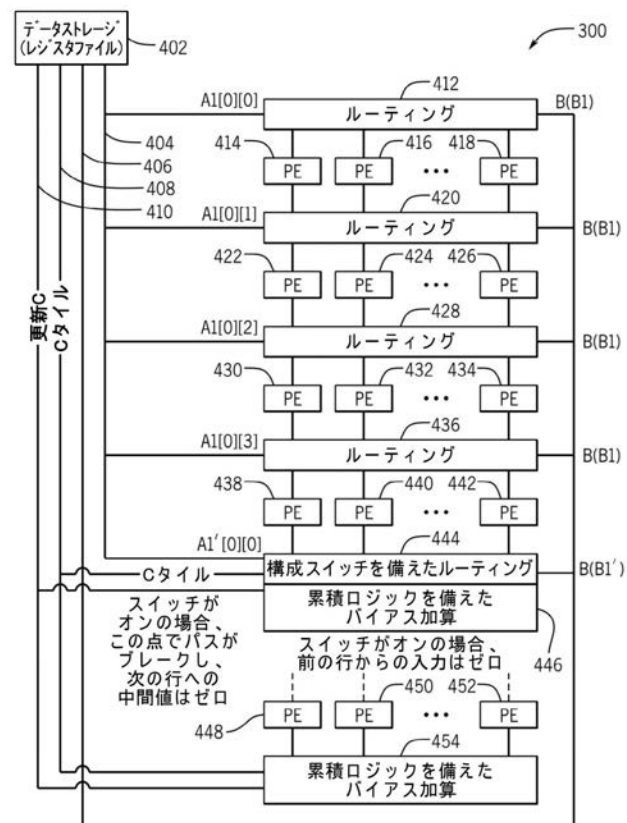
【 図 2 】



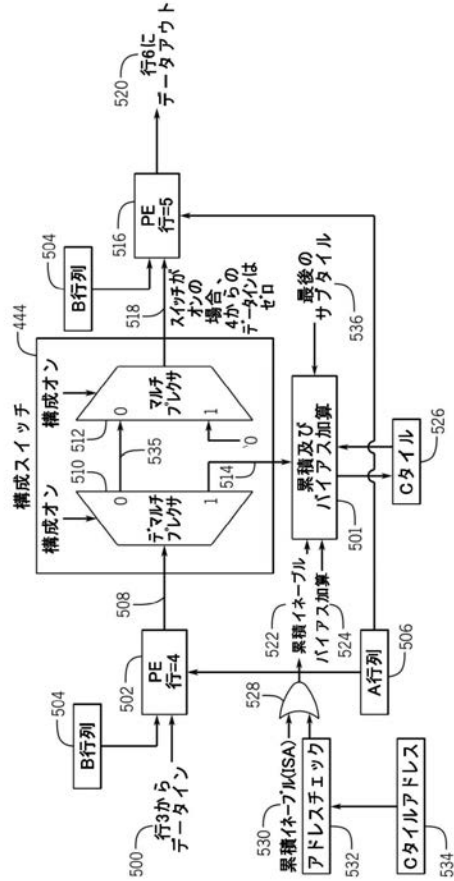
【 図 3 】



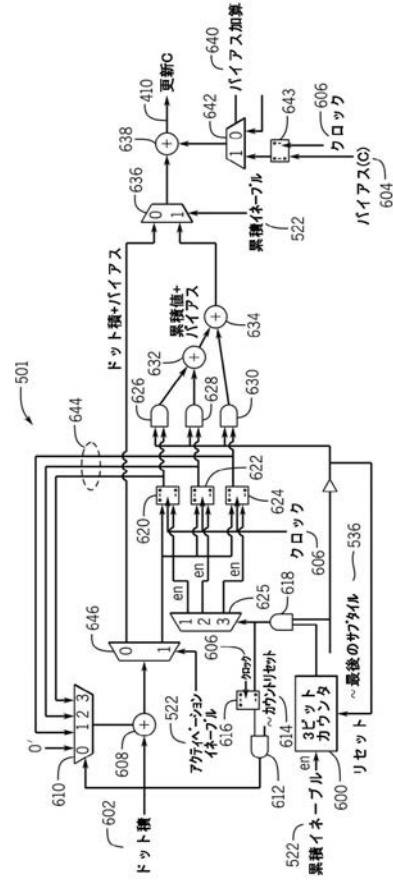
【 図 4 】



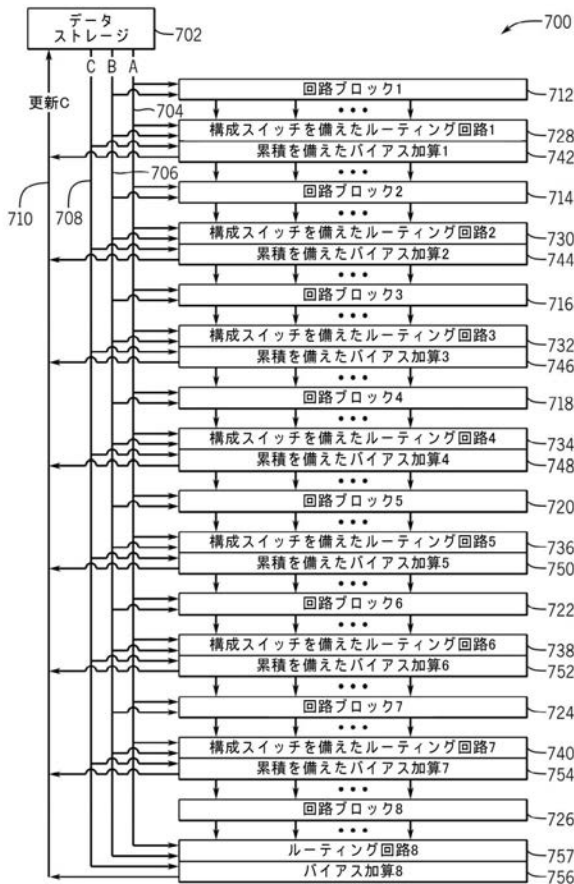
【 図 5 】



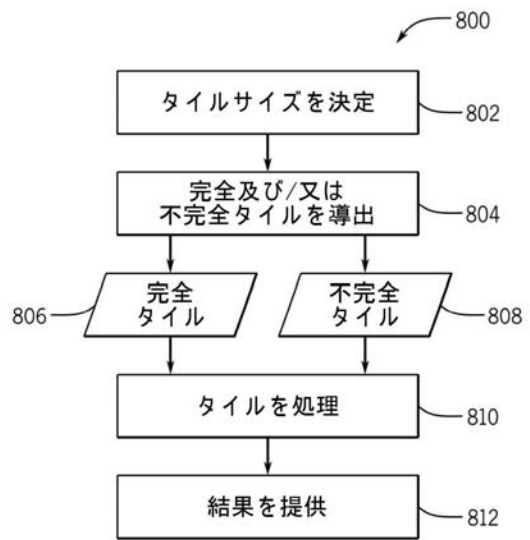
【 図 6 】



【 図 7 】



【 図 8 】



フロントページの続き

(72)発明者 ガーブリート シン カルシ
アメリカ合衆国 95054 カリフォルニア州・サンタクララ・ミッション カレッジ ブーレ
バード・2200 インテル・コーポレーション内

(72)発明者 クリストファー ジャスティン ヒューズ
アメリカ合衆国 95054 カリフォルニア州・サンタクララ・ミッション カレッジ ブーレ
バード・2200 インテル・コーポレーション内

Fターム(参考) 5B056 BB42 EE03 FF01 FF02 FF10 FF16

【外国語明細書】

2021108104000001.pdf