

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2021-508125
(P2021-508125A)

(43) 公表日 令和3年2月25日(2021.2.25)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 17/16 (2006.01)	G06F 17/16	M 5B056
G06F 15/80 (2006.01)	G06F 15/80	
G06F 7/57 (2006.01)	G06F 7/57	

審査請求 有 予備審査請求 未請求 (全 42 頁)

(21) 出願番号 特願2020-536531 (P2020-536531)
 (86) (22) 出願日 平成30年10月19日 (2018.10.19)
 (85) 翻訳文提出日 令和2年8月5日 (2020.8.5)
 (86) 国際出願番号 PCT/CN2018/111077
 (87) 国際公開番号 W02019/128404
 (87) 国際公開日 令和1年7月4日 (2019.7.4)
 (31) 優先権主張番号 201711499179.X
 (32) 優先日 平成29年12月29日 (2017.12.29)
 (33) 優先権主張国・地域又は機関 中国 (CN)

(71) 出願人 504161984
 ホアウェイ・テクノロジーズ・カンパニー・リミテッド
 中華人民共和国・518129・グアンドン・シェンツェン・ロンガン・ディストリクト・バンティアン・(番地なし)・ホアウェイ・アドミニストレーション・ビルディング
 (74) 代理人 110000877
 龍華国際特許業務法人

最終頁に続く

(54) 【発明の名称】 行列乗算器

(57) 【要約】

本発明の実施形態は行列乗算器を開示し、データコンピューティング技術の分野に関し、それにより計算のために2つの行列をブロックに分割する。行列乗算器は、第1メモリ、第2メモリ、演算回路、およびコントローラを含み、ここで、演算回路、第1メモリ、および第2メモリはバスを使用してデータ通信を実行してよく、コントローラは、予め設定されたプログラムまたは命令に従って、第1行列および第2行列を、ブロックに分割されるように制御し、演算回路を、コントローラのブロック分割結果に基づいて、第1メモリおよび第2メモリの対応するブロックで乗算演算を実行するように制御するように構成される。行列乗算器は、2つの行列で乗算演算を実行するように構成されてよい。

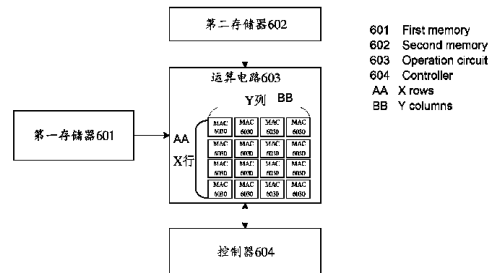


図 6

【特許請求の範囲】

【請求項 1】

M * K 行列である第 1 行列を格納するように構成された第 1 メモリと、
 K * N 行列である第 2 行列を格納するように構成された第 2 メモリと、
 前記第 1 メモリおよび前記第 2 メモリに接続される演算回路と、
 前記演算回路に接続されるコントローラと、を含む、
 行列乗算器であって、

前記演算回路は、X 行 * Y 列からなる演算ユニットを含み、各前記演算ユニットは、ベクトル乗算回路および加算回路を含み、前記行列乗算回路は、前記第 1 メモリによって送られる行ベクトルのデータおよび前記第 2 メモリによって送られる列ベクトルのデータを
 10 受信し、前記 2 つのベクトルを乗算するように構成され、前記加算回路は、前記 2 つのベクトルの乗算によって取得された結果を加算し、同一の演算ユニットの計算結果を累算し、各演算ユニットの演算結果を取得するように構成され、

前記コントローラは、以下の動作、すなわち、

前記第 1 行列を、サイズが X * L であるサブブロックを単位とするブロックに分割し、同じサイズの S * R 個のサブブロックを取得し、前記 S * R 個のサブブロックのうち第 s 行第 r 列におけるサブブロックは $A_{s,r}$ 、 $s = (1, 2, 3, \dots, \text{および } S)$ 、および $r = (1, 2, 3, \dots, \text{および } R)$ で表される、動作と、

前記第 2 行列を、サイズが L * Y であるサブブロックを単位とするブロックに分割し、同じサイズの R * T 個のサブブロックを取得し、R * T 個のサブブロックのうち第 r 行第 t 列におけるサブブロックは、 $B_{r,t}$ 、 $r = (1, 2, 3, \dots, \text{および } R)$ 、 $t = (1, 2, 3, \dots, \text{および } T)$ で表される、動作と
 20

を実行するように構成され、

前記コントローラは、さらに以下の動作、すなわち、

任意のサブブロック $A_{s,r}$ の X 個の行ベクトルにおける第 x 行と、対応するサブブロック $B_{r,t}$ の Y 個の列ベクトルにおける第 y 列とを、X 行 * Y 列からなる演算ユニットの第 x 行第 y 列において前記演算ユニットに入力し、それにより、処理を実行する動作を実行するように構成され、 $x = (1, 2, 3, \dots, \text{および } X)$ 、 $y = (1, 2, 3, \dots, \text{および } Y)$ であり、前記任意のサブブロック $A_{s,r}$ における r と、前記対応するサブブロック $B_{r,t}$ における r とは同じ値を有する、
 30

行列乗算器。

【請求項 2】

前記コントローラは、以下の動作、すなわち、

前記任意のサブブロック $A_{s,r}$ の前記 X 個の行ベクトルにおける前記第 x 行と、前記対応するサブブロック $B_{r,t}$ の前記 Y 個の列ベクトルにおける前記第 y 列とを、同じクロックサイクルにおいて並行して、X 行 * Y 列からなる前記演算ユニットの第 x 行第 y 列において前記演算ユニットに入力し、それにより前記処理を実行する、動作

を実行するように具体的に構成される、請求項 1 に記載の行列乗算器。

【請求項 3】

前記コントローラはさらに、前記任意のサブブロック $A_{s,r}$ の行ベクトルを、x 個の行番号の昇順で、X 行 * Y 列からなる前記演算ユニットに対応する第 x 行に連続的に入力するように制御するように構成され、近接する行ベクトルが同じ列で異なる行の演算ユニットに入る時点の間の差は 1 クロックサイクルであり、前記コントローラはさらに、前記対応するサブブロック $B_{r,t}$ の列ベクトルを、y 個の列番号の昇順で、X 行 * Y 列からなる前記演算ユニットに対応する第 y 行に連続的に入力するように同時に制御するように構成され、近接する列ベクトルが同じ行で異なる列の演算ユニットに入る時点の間の差は 1 クロックサイクルである、
 40

請求項 1 または 2 に記載の行列乗算器。

【請求項 4】

前記コントローラはさらに、s および r の値を変更されないままにして、t の値を、少
 50

なくとも2つの連続するサブブロック乗算計算サイクルにおいて変更されるように制御するように構成され、その結果、前記第1メモリは、前記少なくとも2つの連続するサブブロック乗算計算サイクル内で同じサブブロック $A_{s,r}$ を再使用し、前記サブブロック乗算計算サイクルは、1つのサブブロック $A_{s,r}$ および対応するサブブロック $B_{r,t}$ 上での行列乗算演算を完了させるように、 X 行* Y 列の前記演算ユニットによって使用される時間である、

請求項1から3のいずれか一項に記載の行列乗算器。

【請求項5】

前記行列乗算器はさらに、前記演算回路に接続された第3メモリを含み、

前記コントローラは、前記ベクトル乗算回路および前記加算回路の演算結果を前記第3メモリに格納するように、 X 行* Y 列の前記演算ユニットを制御するように構成される、請求項1から4のいずれか一項に記載の行列乗算器。

10

【請求項6】

前記行列乗算器はさらに、前記第1メモリおよび前記第2メモリに接続される第4メモリと、前記第3メモリに接続される第5メモリとを含み、

前記コントローラはさらに、前記第1行列および前記第2行列の乗算演算を実行する前に、

前記第4メモリから、前記第1行列および前記第2行列のデータソースを、それぞれ前記第1メモリおよび前記第2メモリに移動させ、前記第3メモリから、前記計算結果を前記第5メモリに移動させるように、制御するように構成される、

20

請求項5に記載の行列乗算器。

【請求項7】

前記ベクトル乗算回路は L 個の乗算器を含み、前記加算回路は入力数が $L+1$ である加算木を含む、

請求項1から6のいずれか一項に記載の行列乗算器。

【請求項8】

前記第1メモリ、前記第2メモリ、前記演算回路、および前記コントローラはバスインタフェースユニットを使用して接続される、

請求項1から7のいずれか一項に記載の行列乗算器。

【請求項9】

$S =$

【数22】

$$\begin{cases} M/X, M\%X = 0 \\ \left[\frac{M}{X} \right] + 1, M\%X \neq 0 \end{cases}$$

および

$R =$

【数23】

$$\begin{cases} K/L, K\%L = 0 \\ \left[\frac{K}{L} \right] + 1, K\%L \neq 0 \end{cases}$$

40

であり、

$M\%X = 0$ のとき、計算は前記第1行列の第 $(M+1)$ 行から第 $(S * X - M)$ 行まで実行されず、結果の値には0が割り当てられ、 $K\%Y = 0$ のとき、計算は前記第1行列の第 $(K+1)$ 行から第 $(R * Y - K)$ 行まで実行されず、結果の値には0が割り当てられる、

請求項1から8のいずれか一項に記載の行列乗算器。

【請求項10】

50

R =

【数 2 4】

$$\begin{cases} K/L, K\%L = 0 \\ \left[\frac{K}{L}\right] + 1, K\%L \neq 0 \end{cases}$$

および

T =

【数 2 5】

$$\begin{cases} N/Y, N\%Y = 0 \\ \left[\frac{N}{Y}\right] + 1, N\%Y \neq 0 \end{cases}$$

であり、

$K\%Y = 0$ のとき、計算は前記第 1 行列の第 $(K + 1)$ 列から第 $(R * Y - K)$ 列まで実行されず、結果の値には 0 が割り当てられ、 $N\%X = 0$ のとき、計算は前記第 1 行列の第 $(N + 1)$ 行から第 $(T * X - N)$ 行まで実行されず、結果の値には 0 が割り当てられる、

請求項 1 から 8 のいずれか一項に記載の行列乗算器。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はコンピューティング技術の分野に関し、特に行列乗算器に関する。

【背景技術】

【0002】

現在、2つの行列 A および B の積は、以下の 2つの方式のいずれにおいて計算され得る。

【0003】

方式 1：計算がベクトルプロセッサを使用して実行される。

【0004】

C = A * B であり、ベクトルプロセッサによって同時に計算できる要素の数は M であることが前提とされる。図 1 を参照すると、ベクトルプロセッサは、行列 A の第 i 行（要素 $A_{i1}, A_{i2}, \dots, A_{i(M-1)},$ および A_{iM} を含む）におけるベクトルをソースレジスタ Reg 0 にロードし、そして、行列 B の第 j 列（要素 $B_{j1}, B_{j2}, \dots, B_{j(M-1)},$ および B_{jM} を含む）におけるベクトルをレジスタ Reg 1 にロードし、その結果、Reg 0 および Reg 1 に対応する要素の間の乗算が実装できる。最終的に、加算木を使用することによって累算演算が完了し、行列 C の第 i 行第 j 列におけるデータ C_{ij} が、計算を通じて取得され、行列 C は複数回の計算を実行することによって取得され得る。

【0005】

方式 2：計算速度をさらに増加させるように、行列の乗算演算は、2次元計算アレイを使用して完了され得る。

【0006】

例えば、2次元計算アレイは $N * N$ シストリックアレイであり得る。方式 1 において、2個の $N * N$ 行列の乗算演算を完了するには、 N^3 個の乗算演算が必要である。ベクトルプロセッサは、各クロックサイクルにおいて M 個の要素の間の乗算の計算が可能なので、1つの乗算演算が完了するために必要な期間は N^3 / M 個のクロックサイクルである。方式 2 において、2個の $N * N$ 行列の乗算演算を完了するには、 N^3 個の乗算演算が必要である。シストリックアレイは N^2 個の演算ユニットを有するので、1つの行列演算を完了させるために必要な期間は、 $N^3 / N^2 = N$ 個のクロックサイクルである。

10

20

30

40

50

方式 1 および方式 2 の両者において、 $N * N$ 行列の乗算演算を完了させるためには長い時間がかかり、比較的固定されて柔軟性のないコンピューティングサイズをもたらす。

【発明の概要】

【0007】

本発明の実施形態は、行列乗算器および関連するデバイスを提供し、それにより、行列乗算の最中の、柔軟性のない計算および低い効率という問題を解決する。

【0008】

第 1 の態様によると、本発明の実施形態は行列乗算器を提供し、行列乗算器は、 $M * K$ 行列である第 1 行列を格納するように構成された第 1 メモリと、 $K * N$ 行列である第 2 行列を格納するように構成された第 2 メモリと、第 1 メモリおよび第 2 メモリに接続される演算回路と、演算回路に接続されたコントローラと、を含み、

演算回路は X 行 $* Y$ 列からなる演算ユニットを含み、各演算ユニットはベクトル乗算回路および加算回路を含み、行列乗算回路は、第 1 メモリによって送られる行ベクトルのデータおよび第 2 メモリによって送られる列ベクトルのデータを受信し、2 つのベクトルを乗算するように構成され、加算回路は、2 つのベクトルの乗算によって取得された結果を加算し、同一の演算ユニットの計算結果を累算し、各演算ユニットの演算結果を取得するように構成される、演算回路と、

コントローラは以下の動作、すなわち、

第 1 行列を、サイズが $X * L$ であるサブブロックを単位とするブロックに分割し、同じサイズの $S * R$ 個のサブブロックを取得し、 $S * R$ 個のサブブロックのうち第 s 行第 r 列におけるサブブロックは $A_{s,r}$ 、 $s = (1, 2, 3, \dots, \text{および } S)$ 、および $r = (1, 2, 3, \dots, \text{および } R)$ で表される、動作と、

第 2 行列を、サイズが $L * Y$ であるサブブロックを単位とするブロックに分割し、同じサイズの $R * T$ 個のサブブロックを取得し、 $R * T$ 個のサブブロックのうち第 r 行第 t 列におけるサブブロックは、 $B_{r,t}$ 、 $r = (1, 2, 3, \dots, \text{および } R)$ 、 $t = (1, 2, 3, \dots, \text{および } T)$ で表される、動作とを実行するように構成され、

コントローラは、さらに以下の動作、すなわち、

任意のサブブロック $A_{s,r}$ の X 個の行ベクトルにおける第 x 行と、対応するサブブロック $B_{r,t}$ の Y 個の列ベクトルにおける第 y 列とを、 X 行 $* Y$ 列からなる演算ユニットの第 x 行第 y 列において演算ユニットに入力し、それにより、処理を実行する動作を実行するように構成され、 $x = (1, 2, 3, \dots, \text{および } X)$ 、 $y = (1, 2, 3, \dots, \text{および } Y)$ であり、任意のサブブロック $A_{s,r}$ における r と、対応するサブブロック $B_{r,t}$ における r とは同じ値を有する。

【0009】

本発明の実施形態は、行列乗算器を提供し、ここで、行列乗算器は行列乗算ブロック分割方法、すなわち、 MNK フラクタルを完了して、行列乗算器 60 における内部コントローラ 604 の制御ロジックを使用することによって、乗算のために大きい行列を単位行列（具体的には、 $X * L * L * Y$ 行列）に分割するようにコントローラを使用する。コントローラ 604 の制御ロジックは、各クロックサイクルにおいて、単位行列乗算タスクを演算回路 603 に送り、その結果、データがパイプライン方式で実行され、 X 行 $* Y$ 列の演算ユニットがフルロード状態で動作する。行列乗算の効率が増大し、ニューラルネットワークアルゴリズムを大幅に改善する適用効果を実現される。本発明のこの実施形態において提供される行列乗算器は、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および FC 演算を実行し得る。

【0010】

可能な実装において、コントローラは、以下の動作を実行するように具体的に構成される。

任意のサブブロック $A_{s,r}$ の X 個の行ベクトルにおける第 x 行と、対応するサブブロック $B_{r,t}$ の Y 個の列ベクトルにおける第 y 列とを、同じクロックサイクルにおいて並行し

10

20

30

40

50

て、X行*Y列からなる演算ユニットの第x行第y列において演算ユニットに入力し、それにより演算を実行する。

【0011】

可能な実装において、コントローラはさらに、任意のサブブロック $A_{s,r}$ の行ベクトルを、x個の行番号の昇順で、X行*Y列からなる演算ユニットに対応する第x行に連続的に入力するように制御するように構成され、近接する行ベクトルが同じ列で異なる行の演算ユニットに入る時点の間の差は1クロックサイクルであり、コントローラはさらに、対応するサブブロック $B_{r,t}$ の列ベクトルを、y個の列番号の昇順で、X行*Y列からなる演算ユニットに対応する第y行に連続的に入力するように同時に制御するように構成され、近接する列ベクトルが同じ行で異なる列の演算ユニットに入る時点の間の差は1クロックサイクルである。

10

【0012】

可能な実装において、コントローラはさらに、sおよびrの値を変更されないままにして、tの値を少なくとも2つの連続するサブブロック乗算計算サイクルにおいて変更されるように制御するように構成され、その結果、第1メモリは、少なくとも2つの連続するサブブロック乗算計算サイクル内で同じサブブロック $A_{s,r}$ を再使用し、サブブロック乗算計算サイクルは、1つのサブブロック $A_{s,r}$ および対応するサブブロック $B_{r,t}$ 上での行列乗算演算を完了させるようにX行*Y列の演算ユニットによって使用された時間である。

20

【0013】

可能な実装において、行列乗算器はさらに、演算回路に接続された第3メモリを含み、コントローラは、ベクトル乗算回路および加算回路の演算結果を第3メモリに格納するように、X行*Y列の演算ユニットを制御するように構成される。

【0014】

可能な実装において、行列乗算器はさらに、第1メモリおよび第2メモリに接続される第4メモリと、第3メモリに接続される第5メモリとを含み、コントローラはさらに、第1行列および第2行列の乗算演算を実行する前に、第4メモリから、第1行列および第2行列のデータソースを、それぞれ第1メモリおよび第2メモリに移動させ、第3メモリから、計算結果を第5メモリに移動させるように、制御するように構成される。

30

【0015】

可能な実装において、ベクトル乗算回路はL個の乗算器を含み、加算回路は入力数がL+1である加算木を含む。

【0016】

可能な実装において、第1メモリ、第2メモリ、演算回路、およびコントローラはバスインタフェースユニットを使用して接続される。

【0017】

可能な実装において、 $S =$

【数1】

$$\begin{cases} M/X, M\%X = 0 \\ \lceil \frac{M}{X} \rceil + 1, M\%X \neq 0 \end{cases}$$

40

および $R =$

【数2】

$$\begin{cases} K/L, K\%L = 0 \\ \lceil \frac{K}{L} \rceil + 1, K\%L \neq 0 \end{cases}$$

であり、

50

$M \% X = 0$ のとき、計算は第 1 行列の第 $(M + 1)$ 行から第 $(S * X - M)$ 行まで実行されず、結果の値には 0 が割り当てられ、 $K \% Y = 0$ のとき、計算は第 1 行列の第 $(K + 1)$ 行から第 $(R * Y - K)$ 行まで実行されず、結果の値には 0 が割り当てられる。

【0018】

可能な実装において、 $R =$

【数 3】

$$\begin{cases} K/L, K \% L = 0 \\ \lfloor \frac{K}{L} \rfloor + 1, K \% L \neq 0 \end{cases}$$

および $T =$

【数 4】

$$\begin{cases} N/Y, N \% Y = 0 \\ \lfloor \frac{N}{Y} \rfloor + 1, N \% Y \neq 0 \end{cases}$$

であり、

$K \% Y = 0$ のとき、計算は第 1 行列の第 $(K + 1)$ 列から第 $(R * Y - K)$ 列まで実行されず、結果の値には 0 が割り当てられ、 $N \% X = 0$ のとき、計算は第 1 行列の第 $(N + 1)$ 行から第 $(T * X - N)$ 行まで実行されず、結果の値には 0 が割り当てられる。

【0019】

可能な実装において、行列乗算器はさらに直接メモリアクセスユニットを含み、直接メモリアクセスユニットは、第 1 行列を第 1 メモリに移動する前に第 1 行列に行列転置演算を実行するか、または、第 2 行列を第 2 メモリに移動する前に第 2 行列に行列転置演算を実行するように構成される。

【0020】

可能な実装において、コントローラは第 1 行列の任意のサブブロックを、行形式で第 1 メモリに格納されるように制御し、または、第 2 行列の任意のサブブロックを、行形式で第 2 メモリに格納されるように制御する。このようにして、サブブロックが素早く読み出されることができ、サブブロックはフレキシブルに、かつ素早く転置されることができ

【0021】

第 2 の態様によれば、本願は電子デバイスを提供し、電子デバイスは、

第 1 の態様の任意の実施例において提供されるセキュアな要素と、チップに結合された個別デバイスとを含み得る。

【0022】

第 3 の態様によれば、本願はシステムオンチップを提供し、システムオンチップは第 1 の態様の任意の実施例において提供されたチップを含む。システムオンチップはチップを含み得、または、チップおよび別の個別デバイスを含み得る。

【図面の簡単な説明】

【0023】

本発明の実施形態または背景における技術的解決方法をより明確に説明するように、以下に、本発明の実施形態または背景を説明するために必要とされる添付の図面を簡潔に説明する。

【0024】

【図 1】先行技術における、2 つの行列の積を計算する処理の概略図である。

【0025】

【図 2】先行技術における、畳み込みカーネルを重み行列に変換する概略図である。

【0026】

【図 3】先行技術における、入力データを入力行列に変換する概略図である。

10

20

30

40

50

【0027】

【図4】先行技術における、2つの行列の乗算演算を実行するための方法の概略図である。

【0028】

【図5】先行技術における、TPUシストリックアレイの概略図である。

【0029】

【図6】本発明の実施形態による、行列乗算アクセラレータの構造図である。

【0030】

【図7】本発明の実施形態による、演算ユニット6030の構造図である。

【0031】

【図8】本発明の実施形態による、行列をブロックに分割する概略図である。

【0032】

【図9】本発明の実施形態による、特定の演算回路603の配線の概略図である。

【0033】

【図10】本発明の実施形態による、特定の演算回路603の配線の概略図である。

【0034】

【図11】本発明の実施形態による、ベースが4である行列乗算器の入力フォーマットを示す図である。

【0035】

【図12】 $M = 2$ 、 $N = 2$ 、および $K = 2$ のときの、 $T = 0$ の時点での行列乗算器のパイプライン実行の概略図である。

【0036】

【図13】 $M = 2$ 、 $N = 2$ 、および $K = 2$ のときの、 $T = 1$ の時点での行列乗算器のパイプライン実行の概略図である。

【0037】

【図14】 $M = 2$ 、 $N = 2$ 、および $K = 2$ のときの、 $T = 7$ の時点での行列乗算器のパイプライン実行の概略図である。

【0038】

【図15】 $M = 2$ 、 $B = 2$ 、および $K = 2$ のときの、 $T = 11$ の時点での行列乗算器のパイプライン実行の概略図である。

【0039】

【図16】本発明の実施形態による、別の行列乗算器の構造図である。

【0040】

【図17】本発明の実施形態による、さらに別の行列乗算器の構造図である。

【0041】

【図18】本発明の実施形態による、命令非同期実行シーケンスの概略図である。

【発明を実施するための形態】

【0042】

以下に、本発明の実施形態の添付の図面を参照して、本発明の実施形態を説明する。

【0043】

本願の本明細書、特許請求の範囲、および添付の図面において、用語「第1」、「第2」、「第3」、「第4」および同様のものは、異なるオブジェクトとの間の区別を意図するものであり、特定の順序を意味するものではない。加えて、用語「を含む」、「を有する」およびそれらの任意の他の変形は、非限定的な含有を含めることを意図するものである。例えば、一連の段階またはユニットを含む処理、方法、システム、製品、およびデバイスは、列挙された段階またはユニットに限定されるものではなく、任意で、列挙されていない段階またはユニットをさらに含むか、任意で、当該処理、方法、製品、またはデバイスにもともと備わった別の段階またはユニットをさらに含む。

【0044】

本明細書で「実施形態」に言及することは、実施形態を参照して説明される特定の特性

10

20

30

40

50

、構造、および特徴が、本願の少なくとも1つの実施形態に含まれ得ることを意味する。本明細書の様々な位置に示される語句は、必ずしも同じ実施形態を参照しなくともよく、独立した、または別の実施形態から排除された追加の実施形態ではない。本明細書において説明される実施形態は、別の実施形態と組み合わせられてよいことが、当業者には明確および暗示的に理解される。

【0045】

本明細書で使用される「コンポーネント」、「モジュール」および「システム」などの用語は、コンピュータ関連のエンティティ、ハードウェア、ファームウェア、ハードウェアおよびソフトウェアの組み合わせ、ソフトウェア、または実行中のソフトウェアを意味するように使用される。例えば、コンポーネントは、プロセッサ上で動作する処理、プロセッサ、オブジェクト、実行可能なファイル、実行スレッド、プログラム、および/またはコンピュータであってよいが、それらに限定されるものではない。図に示されるように、コンピューティングデバイス上で動作するアプリケーションと、コンピューティングデバイスとの両者が、コンポーネントであってよい。1または複数のコンポーネントが、処理および/または実行スレッド内に存在してよく、コンポーネントは、1つのコンピュータに、および/または2つ以上のコンピュータの間で分散されて配置されてよい。加えて、これらのコンポーネントは、様々なデータ構造を格納する、様々なコンピュータ読み出し可能な媒体から実行されてよい。例えば、コンポーネントは、ローカルな、および/またはリモートな処理を使用することによって、例えば、1または複数のデータパケットを有する信号（例えば、ローカルシステムの、分散システムの、および/または、信号を使用することによって他のシステムとインタラクトするインターネットなどのネットワークにわたる、別のコンポーネントとインタラクトする2つのコンポーネントからのデータ）によって、送信してよい。

10

20

【0046】

次に、解決される必要がある技術的問題と、本願の応用的シナリオが提供される。近年、画像分類、画像認識、音声認識、および他の関連する分野において、畳み込みニューラルネットワークが良好な性能であるので、畳み込みニューラルネットワークは、学界および産業界において研究および開発のホットスポットとなっている。畳み込みニューラルネットワークは主に、畳み込み演算および全結合（fully-connected、FC）演算を含む。畳み込み演算の演算量は、通常、ネットワークの全演算量の70%より多くを占有することがある。

30

【0047】

畳み込み演算は、行列乗算演算と厳格に同等ではない。しかしながら、畳み込み演算は、適切なデータ調整によって、行列乗算演算に変換され得る。通常、畳み込みニューラルネットワークには複数の畳み込みカーネルがある。畳み込みカーネルは3次元であり、3次元のデータを含む。方向xおよびyは、データの長さおよび幅を表し、方向zは、データの深さとみなされ得る。畳み込みカーネルは実際にはフィルタ（filter）であり、主に画像から異なる特徴を取り出すように構成される。図2を参照すると、畳み込みカーネルは、実質的に、一連の重みの組み合わせである。K個の畳み込みカーネルがあると前提する。K個の畳み込みカーネルの同じ位置で方向zにN個の要素が取り出され、その結果、 $N * K$ の重み行列（weight matrix）が取得できる。畳み込みカーネルは、行列乗算器の仕様（具体的には、行列乗算器によって計算できる行列の行の数および列の数）に基づいて重み行列の形態で行列乗算器のメモリに予め格納されてよく、その結果、畳み込みカーネルは、行列乗算器が行列乗算演算を実行するとき呼び出される。本発明の実施形態において、「*」は「乗算」を表す。

40

【0048】

図3を参照すると、畳み込みカーネルのストライド（stride）（本発明のこの実施形態において、ストライドは1である）に基づいて、行列乗算器は、方向zにおいてM個の入力点のN個のデータ、すなわち、合計で $M * N$ 個のデータを取り出し得る。入力行列（input matrix）が形成され得る。行列乗算器は、入力行列と重み行列に

50

対して乗算演算を実行する必要がある。

【0049】

FC演算は、実質的にベクトルと行列との乗算演算である。FC演算の入力はベクトル9216であり、FC演算は4096の点を必要とする。この場合、FC演算によって出力される点を取得するために、ベクトル9126と9216個の重みとに対して小数点乗算演算が実行される必要がある、4096の点すべてを取得するためには、ベクトル9216と9216×4096個の重みとに対して小数点乗算演算が実行される必要がある。

【0050】

図4は、行列 $C = A * B$ の計算式を示し、ここで、Aは $M * K$ サイズの行列を表し、Bは $K * N$ サイズの行列を表す。本発明のこの実施形態において、M、N、およびKはそれぞれ正の整数である。計算によって行列Cの1個のデータを取得するためには、小数点乗算演算は、行列Aにおける1つの行ベクトルのデータと、行列Bにおける1つの列ベクトルの対応するデータとで実行される必要がある、そして累算が実行される。言い換えれば、計算によって行列Cの1個のデータを取得するためには、N個の乗算演算が実行される必要がある。この場合、計算によって行列Cを取得するためには、 $M * N * K$ 個の乗算演算が実行される必要がある。

10

【0051】

先行技術において、シストリックアレイコンピューティング方式、例えば、機械学習のためにGoogleによってカスタマイズされた専用チップ(AASIC)、Google TPuv1は、 $256 * 256$ 2-D MACアレイを使用することによって、行列乗算および畳み込み演算(図5に示されるように)に最適化されたシストリックアレイ設計を使用する。図の各セルは1つの乗算器である。乗算器が2つの行列の要素を乗算した後、計算によって取得された結果(部分和、すなわち、行列乗算における中間結果)が、図の下部の累算ユニットに伝送され、以前の関連する累算値に累算される。このようにして、データがフルロード状態で動作するとき、シストリックアレイは、各クロックサイクルにおいて1つの行列のサイズの間接値を累算する。前述の解決手段において、計算密度が低いので、行列乗算計算効率は比較的低い。加えて、畳み込み演算の最中に、シストリックアレイのコンピューティングサイズは比較的固定されているので、シストリックアレイの演算効率を増加させるように、入力および重みは多くの形態に転換される必要がある、柔軟性のない演算をもたらす。さらに、行列乗算の最中に、パイプライン実行効果を実現するように、データは大きなサイズを有する必要がある。例えば、小さい行列における $256 * 256$ 2-Dシストリックアレイの計算効率は高くない。

20

30

【0052】

加えて、関連する特許は、 $M * K * N$ 3-D MACアレイを実装する。TPuv1およびNVDLA 2-D MACアレイ解決方法と比較すると、行列乗算計算効率は大幅に増大する。本発明は新しいハードウェアアクセラレータアーキテクチャを提供し、その結果、新しいハードウェアアクセラレータアーキテクチャは $[N * N]$ 行列乗算演算を単一のクロックサイクルで完了することができる。ハードウェアアーキテクチャにおいて、処理エンジン(PE)に含まれる数の個数は $N * N * N$ であり、加算木に含まれる数の個数は $N * N$ である。加えて、大きい行列をより小さい行列に分割する計算方法も、また提供される。しかしながら、前述の解決手段において、ハードウェアによってサポートされるサイズとなるように、行列サイズが追加される必要がある。このことは、データ帯域幅を浪費し、計算効率を低減させる。行列が人為的に大きい行列および小さい行列に分割される場合、ソフトウェアプログラミングは複雑であり、また、関連するソフトウェアプログラミング量も格段に増大する。加えて、アクセラレータが単方向に周期的方式のみ行列の要素をロードすることができ、ソフトウェアは独立して行列を分割する必要があるので、計算モードは単一でありフレキシブルでない。さらに、行列Aおよび行列Bのメモリがすべてのデータを収容できなくなると、繰り返しの読み出しが発生する。したがって、バッファサイズはサービアルゴリズムに比較的強く依存しており、具体的には、アクセラレータは、密結合したオンチップメモリに強く従属している。

40

50

【0053】

したがって、本願において解決されるべき技術的問題は、効率的で、フレキシブルで、低いエネルギー方式で、ハードウェアを使用することによる畳み込みニューラルネットワークにおいて、多数のデータ演算をいかに実行するかである。

【0054】

本発明のこの実施形態において提供される行列乗算器は、機械学習、ディープラーニング、および畳み込みニューラルネットワークなどの分野に適用されてよく、または、デジタル画像処理およびデジタル信号処理などの分野に適用されてよく、または、行列乗算演算に関連する他の分野に適用されてよいことが、理解されることができ。

【0055】

前述の分析に基づいて、本願は、行列乗算アクセラレータを提供し、本願において提供される技術的問題を具体的に解析および解決する。図6は、本発明の実施形態による行列乗算アクセラレータの構造図である。図6に示されるように、行列乗算器60は第1メモリ601、第2メモリ602、演算回路603、およびコントローラ604を含む。演算回路603は、バスを使用して、第1メモリ601、第2メモリ602、およびコントローラ604とデータ通信を実行し得る。演算回路603は、第1メモリ601および第2メモリ602から行列データを取り出し、ベクトル乗算および加算演算を実行するように構成される。コントローラ604は、ベクトル演算を完了するように、予め設定されたプログラムまたは命令に従って、演算回路603を制御するように構成される。第1メモリ601は第1行列を格納するように構成される。

10

20

【0056】

第1行列は $M \times K$ 行列である。行列 a が第1行列である場合、第1行列 a の第 i 行第 j 列の要素は a_{ij} と表されてよく、ここで、 $i = (1, 2, 3, \dots, \text{および} M)$ であり、 $j = (1, 2, 3, \dots, \text{および} K)$ である。

【0057】

本発明のこの実施形態において説明される第1メモリ601、および、以下に説明される関連する行列乗算器の第2メモリ602、第3メモリ606、および内部メモリは、それぞれ、レジスタ、ランダムアクセスメモリ(random access memory、略してRAM)、静的ランダムアクセスメモリ、フラッシュメモリ、または別の読み出しおよび書き込み可能メモリであってよい。本願において、第1行列、第2行列および演算結果のデータ型はそれぞれ、int8、fp16、またはfp32などの型であってよい。

30

【0058】

第2メモリ602は第2行列を格納するように構成され、第2行列は $K \times N$ 行列である。行列 b が第2行列である場合、第2行列 b の第 j 行第 g 列の要素は B_{jg} と表されてよく、ここで、 $j = (1, 2, 3, \dots, \text{および} K)$ であり、 $g = (1, 2, 3, \dots, \text{および} N)$ である。

【0059】

本明細書では、 M 、 K 、 N 、 X 、および Y はそれぞれ、0より大きい整数である。 M 、 N 、および K のうちいずれか2つのパラメータが等しくてもよく、等しくなくてもよい。代替的に、 M 、 N 、および K は、等しくてもよく、等しくなくてもよい。 X および Y は、等しくてもよく、等しくなくてもよい。これは、本願を具体的に限定するものではない。

40

【0060】

演算回路603は、 X 行 \times Y 列の演算ユニット6030(乗算累算ユニットMACと称されてよい)を含み得る。各演算ユニットは、独立してベクトル乗算演算を実行し得る。図6において、演算回路603が 4×4 演算ユニット6031を含む例が図に使用され、すなわち、 $X = 4$ および $Y = 4$ である。演算ユニット6030は、それぞれ第1メモリ601によって送られた行ベクトルと、第2メモリ602によって送られた列ベクトルとを受信し、行ベクトルと列ベクトルとのベクトル乗算演算を実行するように使用される、2つの入力を提供される。具体的には、1つの演算回路6030はベクトル乗算回路および

50

加算回路を含み、ここで、行列乗算回路は第1メモリ601によって送られる行ベクトルのデータと、第2メモリ602によって送られる列ベクトルのデータを受信し、2つのベクトルを乗算するように構成され、加算回路は、2つのベクトルの乗算によって取得された結果を加算し、同一の演算ユニットの計算結果を累算し、演算ユニットの演算結果を取得するように構成される。

【0061】

図7は、演算ユニット6030の構造図である。可能な実装において、ベクトル乗算回路はL個（例えば、 $L=4$ ）の乗算器を含む。加算回路は、入力数が $L+1$ である加算木を含み、具体的には、加算木はL個の乗算結果を累算し、異なるクロックサイクルの演算ユニットの計算結果を累算するように構成される。任意で、行列乗算器60は第3メモリ605をさらに含み、第3メモリ605はベクトル乗算回路および加算回路の演算結果を格納し、異なるクロックサイクルの演算結果を格納するように構成される。本願における第3メモリ605は、 $X*Y$ 個のストレージユニットを含み得、各ストレージユニットは、対応する演算ユニットが演算を実行するたびに、取得された演算結果を格納するように構成されることが、理解され得る。代替的に、各演算ユニットは、第3メモリ605の指定されたストレージスペースに対応し、ストレージスペースは、演算ユニットが演算を実行するたびに取得された演算結果を格納するように使用される。

10

【0062】

コントローラ604は、第1行列および第2行列の積を計算するように、以下の動作を実行し得る。

20

【0063】

コントローラ604は、第1行列を、サイズが $X*L$ であるサブブロックを単位とするブロックに分割し、同じサイズの $S*R$ 個のサブブロックを取得し、ここで、 $S*R$ 個のサブブロックのうち第s行第r列におけるサブブロックは、 $A_{s,r}$ と表され、 $s=(1, 2, 3, \dots, \text{および} S)$ 、 $r=(1, 2, 3, \dots, \text{および} R)$ である。すなわち、本願の行列乗算器60に関して、行列乗算器60に含まれる X 行 $*Y$ 列の行列データは、生産または送後の後に固定され、対応する乗算回路における数Lの乗算器もまた固定される。したがって、行列演算の最中、第1行列および第2行列はフラクタルである必要があり、すなわち、ブロックに分割される必要がある。分割方式は、第1行列を、 $X*L$ サブブロックを単位として使用したブロックに分割することである。本発明のこの実施形態において、ブロック分割の目的は、大きい行列を、行列乗算器のサイズに準拠した多くの小さい行列に分割し、そして、特定のシーケンスの小さい行列を計算し、関連する小さい行列の値を累算し、最終的に行列乗算結果を取得することである。このようにして、フレキシブルな計算が実行でき、後の再使用およびマルチレベルのキャッシングを円滑化し、計算効率がさらに増大することができ、データ移動帯域幅およびエネルギー消費を低減することができる。

30

【0064】

第1行列が $M*K$ 行列であり、第1行列が整数個の $X*L$ サブブロックで正確に分割できない場合が存在し得ることに、留意すべきである。したがって、 M/X または K/L が整数でないとき、演算は要素0をパディングする方式で実行されてよい。代替的に、対応する位置で全く計算が実行されず、結果の値に0が割り当てられる。具体的には、

40

【数5】

$$S = \begin{cases} M/X, M\%X = 0 \\ \lceil \frac{M}{X} \rceil + 1, M\%X \neq 0 \end{cases}$$

および

【数 6】

$$R = \begin{cases} K/L, K\%L = 0 \\ \lceil \frac{K}{L} \rceil + 1, K\%L \neq 0 \end{cases}$$

であり、 $M\%X = 0$ のとき、計算は第 1 行列の第 $(M + 1)$ 行から第 $(S * X - M)$ 行まで実行されず、結果の値には 0 が割り当てられ、 $K\%Y = 0$ のとき、計算は第 1 行列の第 $(K + 1)$ 行から第 $(R * Y - K)$ 行まで実行されず、結果の値には 0 が割り当てられる。言い換えれば、演算ユニットは対応する行および列において実体的乗算計算を実行せず、処理のために、演算が実行されたが結果が 0 であるとみなす。このようにして、対応する演算ユニットの読み出しおよび演算電力消費は低減され得る。 10

【0065】

対応して、コントローラ 604 は、第 2 行列を、サイズが $L * Y$ であるサブブロックを単位とするブロックに分割し、同じサイズの $R * T$ 個のサブブロックを取得し、ここで、 $R * T$ 個のサブブロックのうち第 r 行第 t 列におけるサブブロックは、 $B_{r,t}$ と表され、 $r = (1, 2, 3, \dots, \text{および } R)$ 、 $t = (1, 2, 3, \dots, \text{および } T)$ である。コントローラ 604 が第 1 行列を、演算回路 603 の仕様に従ってブロックに分割されるように制御した後、第 2 行列もまた、第 1 行列に合致する方式で対応して分割されることを必要とし、そうでなければ、行列乗算計算は実行され得ない。 20

【0066】

第 2 行列が $K * N$ 行列であり、第 2 行列が整数個の $L * Y$ サブブロックで正確に分割できない場合が存在し得ることに、留意すべきである。したがって、 K/L または N/Y が整数でないとき、演算は要素 0 をパディングする方式で実行されてよい。代替的に、対応する位置で全く計算が実行されず、結果の値に 0 が割り当てられる。具体的には、

【数 7】

$$R = \begin{cases} K/L, K\%L = 0 \\ \lceil \frac{K}{L} \rceil + 1, K\%L \neq 0 \end{cases}$$

および 30

【数 8】

$$T = \begin{cases} N/Y, N\%Y = 0 \\ \lceil \frac{N}{Y} \rceil + 1, N\%Y \neq 0 \end{cases}$$

であり、 $K\%Y = 0$ のとき、計算は第 1 行列の第 $(K + 1)$ 列から第 $(R * Y - K)$ 列まで実行されず、結果の値には 0 が割り当てられ、 $N\%X = 0$ のとき、計算は第 1 行列の第 $(N + 1)$ 行から第 $(T * X - N)$ 行まで実行されず、結果の値には 0 が割り当てられる。言い換えれば、演算ユニットは対応する行および列において実体的乗算計算を実行せず、処理のために、演算が実行されたが結果が 0 であるとみなす。このようにして、対応する演算ユニットの読み出しおよび演算電力消費は低減され得る。 40

【0067】

固定された仕様に従って、第 1 行列および第 2 行列が別々にブロックに分割された後、2 つの行列は、サブブロックの間の行列乗算演算を実行するように、演算回路 603 に入力されてよい。具体的な計算処理において、コントローラ 604 は、任意のサブブロック $A_{s,r}$ の X 個の行ベクトルのうちの第 x 行と、対応するサブブロック $B_{r,t}$ の Y 個の列ベクトルのうちの第 y 列とが、 X 行 $* Y$ 列からなる演算ユニットにおける第 x 行第 y 列の演算ユニットに入力され、それにより演算を実行するように制御し得、ここで、 $x = (1, 2, 3, \dots, \text{および } X)$ 、 $y = (1, 2, 3, \dots, Y)$ であり、任意のサブブロック $A_{s,r}$ における r および対応するサブブロック $B_{r,t}$ における r は同じ値を有する。サブブ 50

ブロック A_{s_r} の行ベクトルおよびサブブロック B_{r_t} の列ベクトルが演算ユニットに入力される前に、第 1 行列および第 2 行列はブロックに分割されており、すなわちフラクタルである。したがって、特定のシーケンスにおいて、演算回路 603 にサブブロック A_{s_r} および対応するサブブロック B_{r_t} を入力する複数の実装があつてよい。

【0068】

可能な実装において、サブブロック A_{s_r} および対応するサブブロック B_{r_t} における s または t の値のシーケンスにおいて連続的に演算が実行されてよい。図 8 に示されるように、例えば、第 1 行列は $M * K$ 行列であり、第 2 行列は $K * N$ 行列である。 $M = 12$ 、 $K = 6$ 、 $N = 12$ 、 $X = 4$ 、 $Y = 4$ 、および $L = 3$ であると前提される。第 1 行列および第 2 行列がブロックに分割された後、 $S = 3$ 、 $R = 2$ 、および $T = 3$ であることがわかる。この場合、第 1 行列

10

【数 9】

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \\ A_{31} & A_{33} \end{bmatrix}$$

および第 2 行列

【数 10】

20

$$B = \begin{bmatrix} B_{11} & B_{12} & B_{13} \\ B_{21} & B_{22} & B_{23} \end{bmatrix}$$

がブロック分割後に取得され、ここで、 A は $X * L$ 行列、すなわち $4 * 3$ 行列を表し、 B の各要素は実際、 $L * Y$ 行列、すなわち $3 * 4$ 行列である。

【0069】

$$C = A * B =$$

【数 11】

30

$$\begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} & A_{11}B_{13} + A_{12}B_{23} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} & A_{21}B_{13} + A_{22}B_{23} \\ A_{31}B_{11} + A_{32}B_{21} & A_{31}B_{12} + A_{32}B_{22} & A_{31}B_{13} + A_{32}B_{23} \end{bmatrix}$$

【0070】

第 1 行列および第 2 行列の乗算演算において、行列乗算演算が、任意の 1 つのサブブロック A_{s_r} 、すなわち、第 1 行列の各サブブロック A_{s_r} と、第 2 行列の対応するサブブロック B_{r_t} とで実行されることが必要である。行列乗算計算がシーケンスで最初に行われる特定のシーケンスおよび特定のサブブロックを決定する、複数の実装があり得る。

40

【0071】

方式 1：行列乗算シーケンスにおいて、例えば、サブブロックはサブブロック A_{1_1} およびサブブロック B_{1_1} であつてよい。 A_{1_1} のすべての行ベクトルと、対応する B_{1_1} のすべての列ベクトルとが、サブブロック乗算計算サイクル（第 1 のラウンドとして理解されてよい）に入力され、それにより、演算を実行する。 A_{1_2} のすべての行ベクトルおよび対応する B_{2_1} のすべての列ベクトルでの演算が、第 2 サブブロック乗算計算サイクル（第 2 のラウンドとして理解されてよい）において実行される。このようにして、演算ユニットが累算を実行した後、結果行列 C の第 1 行第 1 列における結果点 C_{1_1} の値が取得され得る。同様に、結果行列 C のすべての位置における結果点を取得され得る。実際、

50

C_{11}

【数 1 2】

$$A_{11}B_{11} + A_{12}B_{21}$$

ここで、

【数 1 3】

$$A_{11} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} \end{bmatrix}$$

10

【数 1 4】

$$A_{12} = \begin{bmatrix} a_{14} & a_{15} & a_{16} \\ a_{24} & a_{25} & a_{26} \\ a_{34} & a_{35} & a_{36} \\ a_{44} & a_{45} & a_{46} \end{bmatrix}$$

20

【数 1 5】

$$B_{11} = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \\ b_{41} & b_{42} & b_{43} \end{bmatrix}$$

30

【数 1 6】

$$B_{12} = \begin{bmatrix} b_{14} & b_{15} & b_{16} \\ b_{24} & b_{25} & b_{26} \\ b_{34} & b_{35} & b_{36} \\ b_{44} & b_{45} & b_{46} \end{bmatrix}$$

40

【0072】

言い換えれば、 C_{11} は実際、 4×4 行列である。したがって、行列計算ルールによれば、最終的に取得される行列 C は $M \times N$ 結果行列であり、すなわち、 12×12 結果行列である。

【0073】

方式 2：1つのサブブロックが具体的ルールによって再使用される。本発明のこの実施形態は、第 1 行列の 1つのサブブロック A_{s_r} および第 2 行列の対応するサブブロック B_{r_t} に、サブブロックでの行列乗算演算の実行をもたらすサブブロック再使用方式を提供する。具体的には、コントローラ 604 はさらに、 s および r の値を変更されないままに

50

して、 t の値を、少なくとも2つの連続するサブブロック乗算計算サイクルにおいて変更されるように制御するように構成され、その結果、第1メモリは、少なくとも2つの連続するサブブロック乗算計算サイクル内で同じサブブロック $A_{s,r}$ を再使用し、ここで、サブブロック乗算計算サイクルは、1つのサブブロック $A_{s,r}$ および対応するサブブロック $B_{r,t}$ での行列乗算演算を完了するように、 X 行 * Y 列の演算ユニットによって使用される時間である。

【0074】

例えば、 $M = 12$ 、 $K = 6$ 、 $N = 12$ 、 $X = 4$ 、 $Y = 4$ 、および $L = 3$ を前提とする前述した実施形態において、 $A_{1,1}$ のすべての行ベクトルおよび対応するサブブロック $B_{1,1}$ のすべての列ベクトルが、サブブロック乗算計算サイクル（第1のラウンドとして理解されてよい）に入力され、それにより演算を実行する。第2サブブロック乗算計算サイクル（第2のラウンドとして理解されてよい）において、 s および r の値は変更されないままであり、しかし、 t の値は変更される必要があり、具体的には、 $A_{1,1}$ のすべての行ベクトルおよび別の対応するサブブロック $B_{1,2}$ のすべての列ベクトルで演算が実行される。任意で、第3のサブブロック乗算計算サイクル（第3ラウンドとして理解されてよい）において、 $A_{1,1}$ のすべての行ベクトルおよびさらに別の対応するサブブロック $B_{1,3}$ のすべての列ベクトルで演算が実行される。このようにして、第1メモリの $A_{1,1}$ は、複数の連続するサブブロック乗算計算サイクルで繰り返し使用されることができ、その結果、読み出しおよび書き込みオーバーヘッドが低減され、データ移動帯域幅が低減される。

【0075】

方式1および方式2において、サブブロック乗算計算サイクルにおける、第1行列のサブブロック $A_{s,r}$ および第2行列の対応するサブブロック $B_{r,t}$ に関する計算ルールは、第1行列における任意のサブブロック $A_{s,r}$ の X 個の行ベクトルにおける第 x 行と、対応するサブブロック $B_{r,t}$ の Y 個の列ベクトルのうちの第 y 列とが、 X 行 * Y 列からなる演算ユニットにおける第 x 行第 y 列の演算ユニットに入力され、それにより演算を実行するように制御し得、ここで、 $x = (1, 2, 3, \dots, \text{および } X)$ 、 $y = (1, 2, 3, \dots, Y)$ であり、任意のサブブロック $A_{s,r}$ における r および対応するサブブロック $B_{r,t}$ における r は同じ値を有する。すなわち、サブブロック $A_{s,r}$ の任意の行ベクトルと、第2行列の対応するサブブロック $B_{r,t}$ の任意の列ベクトルは、計算のために、 X 行 * Y 列からなる演算ユニットにおける指定された演算ユニットに入力される。例えば、 $A_{1,1}$ の第2行ベクトル

【数17】

$$[a_{21} \quad a_{22} \quad a_{23}]$$

および、第2行列の対応するサブブロック $B_{1,1}$ における第3列ベクトル

【数18】

$$\begin{bmatrix} b_{13} \\ b_{23} \\ b_{33} \end{bmatrix}$$

は、演算のために、 X 行 * Y 列からなる演算ユニットにおける第2行第3列に対応する演算ユニットに入力される、などである。

【0076】

図6で示された演算回路603における演算ユニットの構成方式によると、図9は、本発明の実施形態による特定の演算回路603の配線の概略図である。

【0077】

$BUFA$ は第1行列の第1メモリ601であり、 $BUFB$ は第2行列の第2メモリ602であり、 $BUFC$ は各演算ユニット6030の計算結果を格納する第3メモリ605で

10

20

30

40

50

あり、演算回路603はX行*Y列(X=4、Y=4と前提する)からなる演算ユニット、すなわち図のMAC GRP R00C00からMAC GRP R03C03を含む。加えて、各演算ユニットMAC GRPは、X*L行列の1つの行ベクトルとL*Y行列の1つの列ベクトルで乗算演算を実行し得る。

【0078】

本発明のこの実施形態において、演算回路603は、フラクタル行列乗算ユニットと称されてよく、3-D MACアレイ(MAC Cube)およびアキュムレータ(Accumulator)を含み、以下のようなフラクタル行列乗算命令を実行するように構成される。C=A*BまたはC=A*B+C、ここで、A/B/Cは2次元行列である。Aのサイズは(M*ベース)×(K*ベース)、Bのサイズは(K*ベース)×(N*ベース)、Cのサイズは(M*ベース)×(N*ベース)である。ベースは演算回路603の基本サイズであり、すなわちX*Y、例えば、8*8、16*16、および32*32である。前述のC=A*BまたはC=A*B+C計算演算は、MNK行列乗算(および累算)と称される。実際の実行処理において、コントローラは、特定のシーケンスの組み合わせ(上で説明された方式1または方式2)におけるフラクタル方式でMNK行列乗算を完了させるように、大きい行列を、ベースサイズの基本行列に分割されるように制御する。

10

【0079】

フラクタル行列乗算ユニットの具体的なアーキテクチャは、図7に示される(ベース=4と前提する)。例えば、図7において、MACグループはN*N(4*4)からなる乗算累算グループであり、N(4)個の乗算ユニット、および、入力数がN+1(5)である累算木を含む。行列乗算に関して、乗算アキュムレータが、1つの行に1つの列を乗算して累算(すなわち、結果行列の1つの要素)を実行する演算を、実行してよい。図9において、4*4乗算累算グループの全体があり、すなわち、完全な4*4*4*4の行列乗算演算が同時に計算され得る。

20

【0080】

図9の配線の概略図において、演算回路603は、1つのサブブロックA_{s,r}および対応するサブブロックB_{r,t}での行列乗算計算を、同じクロックサイクルにおいて完了することのサポートとなり得ることが、理解されることができる。サブブロックA_{s,r}のX個の行ベクトルのすべてと、対応するサブブロックB_{r,t}のY個の列ベクトルのすべてとが、図9の配線方式において対応するBUFAおよびBUFBから同時に、対応する演算ユニット6030に到達し得るので、コントローラ604は、1つのサブブロックA_{s,r}および対応するサブブロックB_{r,t}での乗算計算を1クロックサイクルで完了するように、および、次のクロックサイクルにおいて、別のサブブロックA_{s,r}および対応するサブブロックB_{r,t}での乗算計算を完了するか、または、同じサブブロックA_{s,r}および対応する別のサブブロックB_{r,t}での行列乗算計算を完了するように、演算回路603を制御し得る。

30

【0081】

図10は、本発明の実施形態による特定の演算回路603の配線の概略図である。演算回路603において、図10に対応して、シストリックアレイ構造が提供される。具体的に、コントローラ604はさらに、任意のサブブロックA_{s,r}の行ベクトルを、x個の行番号の昇順で、X行*Y列からなる演算ユニットに対応する第x行に連続的に入力するように制御するように構成され、近接する行ベクトルが同じ列で異なる行の演算ユニットに入る時点の間の差は1クロックサイクルであり、コントローラはさらに、対応するサブブロックB_{r,t}の列ベクトルを、y個の列番号の昇順で、X行*Y列からなる演算ユニットに対応する第y行に連続的に入力するように同時に制御するようにさらに構成され、近接する列ベクトルが同じ行で異なる列の演算ユニットに入る時点の間の差は1クロックサイクルである。

40

【0082】

すなわち、各演算ユニット6030(乗算アキュムレータ)をフルに利用するように、本発明のこの実施形態におけるフラクタル行列乗算ユニットは、シストリックアレイ構造

50

を有してよい。TPUV1構造からの差は、(TPUV1におけるデータの個数は1であるが)各シストリック伝送におけるデータの個数がLであることにある。したがって、データ処理の並行性の程度は、TPUV1におけるシストリックアレイのそれより高い。

【0083】

シストリックアレイアーキテクチャに基づいて、図10に対応する配線構造において、BUFA/Bが、それぞれ、第1行列/第2行列をバッファリングするように構成されるメモリである。図10において、第1行列バッファ(BUFA)は、各クロックサイクルにおいて、行列Aにおける単位行列をX個の行に分割し、同じ行のL個の要素をシストリックアレイの演算ユニットに順次送る。同様に、第2行列バッファ(BUFB)は、各クロックサイクルにおいて、第2行列における単位行列をY個の列に分割し、同じ列のL個の要素をシストリックアレイに順次送る。具体的な時間シーケンスは以下の通りである。

10

【0084】

BUFCは、「 $A * B + C$ 」計算における「C」(オフセット)行列を格納するバッファ(L0バッファまたはバッファレジスタを使用して構築されてよい)であり、また、行列乗算の中間値がBUFCに格納されてもよい。乗算アキュムレータが乗算を完了した後、累算木は、乗算後に取得されたL個の中間値と、BUFCに格納された1つのオフセットまたは中間値を累算する。

【0085】

$M = 2$ 、 $N = 2$ 、および $K = 2$ (すなわち $8 \times 8 * 8 \times 8$ 行列乗算が使用される例。行列乗算器60のコントローラ603は、図11のフォーマットで行列乗算を分割し、全体で8個の 4×4 単位行列演算を取得する。MNK行列乗算演算に関して、分割シーケンスのための多くの可能性があり、分割シーケンスのルールは、MNK行列乗算演算が、方式1および方式2のシーケンスにおいて実行されてよいことである。方式2において再使用するデータの最大数倍のポリシーを使用することによって、データを読み出すための消費電力を低減し得ることが、理解されてよい。MNKフラクタル分割が実行された後、図12から図15に示されるように、コントローラ603の制御ロジックは8クロックサイクルに8個のフラクタルをシストリックアレイに入力する。図12は、 $M = 2$ 、 $N = 2$ および $K = 2$ のときの、時点 $T = 0$ におけるフラクタル行列乗算器のパイプライン実行を示し、図13は、 $M = 2$ 、 $N = 2$ および $K = 2$ のときの、時点 $T = 1$ における行列乗算器のパイプライン実行を示し、図14は、 $M = 2$ 、 $N = 2$ および $K = 2$ のときの、時点 $T = 7$ におけるフラクタル行列乗算器のパイプライン実行を示し、図15は、 $M = 2$ 、 $N = 2$ および $K = 2$ のときの、時点 $T = 11$ におけるフラクタル行列乗算器のパイプライン実行を示す。シストリックアレイは、 $T = 6$ のとき、すなわち、第7クロックサイクルにおいてフルロード状態で動作を開始することがわかる。最後の6クロックサイクルにおいて、単位行列がシストリックアレイから出力され、行列全体の乗算演算もまた完了する。

20

30

【0086】

任意で、図16を参照すると、行列乗算器60は、命令発送ユニット606、命令フェッチユニット607、直接メモリアクセスユニット608、ベクトルユニット609、スカラユニット610、およびバスインタフェースユニット611をさらに含んでよい。さらに、本発明のこの実施形態において提供される行列乗算器60は、コプロセッサとして使用され、中央演算処理装置(Central Processing Unit, 略してCPU)80上に載置されてよく、CPUは行列乗算器60に計算タスクを割り当てる。具体的に、CPU80は第1行列、第2行列、および外部メモリ70への関連する命令を格納し得る。行列乗算器60は、第1行列、第2行列、および外部メモリ70における関連する命令を読み出すことによって、行列乗算演算を完了し得る。外部メモリ70は、具体的には、ダブルデータレートシンクロナスダイナミックランダムアクセスメモリ(Double Data Rate Synchronous Dynamic Random Access Memory, 略してDDR)、または別の読み出しおよび書き込み可能メモリであってよい。外部メモリは、行列乗算器60からプライベートなメモリであってよい。具体的に、第1メモリ601、第2メモリ602、第3メモリ605、およ

40

50

び外部メモリ70は一般的に、オンチップメモリ(On-Chip Buffer)である。

【0087】

1.ベクトルユニット609(Vector Unit)は、様々な種類のマルチパラレルコンピューティングデバイス(例えば、浮動小数点乗算、浮動小数点加算、浮動小数点値比較)を含み、ここで、コンピューティングデバイスは、SIMD(Single Instruction multiple data)命令を実行するように構成され、統一されたバッファ(Unified Buffer)およびLOCバッファのために移動する直接のデータに責任を負う。

【0088】

2.スカラーユニット610(Scalar Unit)は、様々な種類の整数基本演算デバイス(例えば、加算、乗算、比較、およびシフト)を含む。

【0089】

3.直接メモリアクセスユニット(Direct Memory Access Unit, DMA Unit)は、各ストレージユニットにデータを移動するように、例えば、L1 RAMからL0 RAMへデータを移動するように構成される。本発明のこの実施形態における直接メモリアクセスユニットが、行列乗算器の外部メモリまたは内部メモリから、乗算演算に関与する行列データを移動するとき、直接メモリアクセスユニットは、行列がブロックに分割された後に取得された結果を格納する必要がある。例えば、 2×2 行列に関して、第1行列の第1行第1列のサブブロック $A_{11} =$

【数19】

$$\begin{bmatrix} A0 & A1 \\ A2 & A3 \end{bmatrix}$$

がサブブロックのユニットに格納され、A0、A1、A2およびA3が1つの行に格納される、などである。このようにして、第1行列または第2行列が、対応する第1メモリへと移動されてよいとき、または、第2行列が、対応する第2メモリに移動されてよいときは、ストレージは、前述の方式で実行されてよい。演算ユニットが読み出しの実行を必要とするとき、演算ユニットもまた、前述のストレージシーケンスにおいて読み出しを実行してよく、それにより、計算を円滑化する。行ベクトルが列ベクトルに転置される必要があるとき、転置はフレキシブルに、そして素早く実行されてよい。

【0090】

4.命令フェッチユニット607(Instruction Fetch Unit, IFU)は、内部でPC(プログラムカウンタ)およびIM(命令メモリ)に統合され、メインメモリからバスインタフェースユニット(BIU)611を使用して命令をフェッチし、実行手順を復号および制御する。

【0091】

5.命令発送ユニット606(Dispatch Unit)は、命令フェッチユニットによって伝送された命令を構文解析し、命令に対応するタイプ命令を4つのパイプラインユニットに提示し、ここで、パイプラインユニットは図16のスカラーユニット(Scalar Unit)、ダイレクトメモリアクセス(Direct Memory Access, DMA)ユニット、ベクトルユニット(Vector Unit)、およびフラクタル行列乗算ユニットである。命令発送ユニットが、4つのパイプラインの間の順序立てた実行を制御するためのメカニズムがある。

【0092】

パイプラインユニットには2つの型、すなわち非同期実行(Posted Execution)および同期実行があることに留意すべきである。すべてのタイプ命令は順序保持方式で伝送される。違いは、非同期実行ユニットによる命令の実行は非同期的に終了し、同期実行ユニットによる命令の実行は同期的に終了する、ということにある。スカラーユ

10

20

30

40

50

ニット (Scalar Unit) は同期実行ユニットであり、フラクタル行列乗算ユニット (Fractal Mat Mult Unit)、DMAユニット、およびベクトルユニット (Vector Unit) は非同期実行ユニットである。

【0093】

可能な実装において、直接メモリアクセスユニットに関して、本発明のこの実施形態は、構成可能なオンフライン行列転置機能を提供する。例えば、第1行列のブロック行列がメモリ (例えば、行列乗算器の外部メモリ) から、別のメモリ (第1メモリなどの、行列乗算器の内部メモリ) に移動されるとき、直接メモリアクセスユニットは、当該移動の最中に行列転置演算を実行し、転置行列の順序で転置行列を格納する。行列転置は、ニューラルネットワークトレーニングプロセスの必須の演算フェーズである。移動後の転置の実行のための共通命令と比較すると、本発明のこの実施形態における構成可能なオンフライン行列転置のための移動命令は、よりフレキシブルであり、また、ソフトウェアはより容易に、およびより簡潔にされる。詳細は以下の表に示される。

10

【0094】

共通命令：構成可能なオンフライン行列転置機能のための命令。

【表1】

LOAD_L0	X2, {X1}	LOAD_L0_to_L1.Trans {X4}, {X1}
Transpose	X3, X2	
STORE_L1	{X4}, X3	

20

【0095】

共通移動命令が、構成可能なオンフライン行列転置機能のための命令と比較される。構成可能なオンフライン行列転置機能をサポートすることによって、同じ命令が、異なるパラメータへの、より多くの応用的シナリオをサポートし得る。フラクタル行列乗算プロセッサアーキテクチャに適用可能な、構成可能なオンフライン行列転置方法が設計される。

【0096】

図17を参照すると、データ再使用を円滑化し、消費電力を低減し、密結合されたオンチップメモリへの依存を低減するように、本発明の実施形態はさらに、マルチレベルバッファを使用するストレージ構造を提供する。すべての演算ユニットは、統一されたバッファ (Unified Buffer) を使用することによって、相互作用データを読み出し/書き込みしてよい。行列乗算器内には、2つのレベルの専用バッファL1およびL0がある。L1バッファおよび統一されたバッファは通常、ダイレクトメモリアクセスDMAユニットを使用して、外部格納空間とデータを交換する。外部格納空間は複数レベルのストレージユニットを含む。例えば、行列乗算器は複数レベルのバッファを含み、L0からL1へ、そしてL2バッファへと、容量が次第に増加し、帯域幅が次第に減少し、遅延が次第に増加し、消費電力オーバーヘッドが次第に増加する。L0は、最も内側のレベルのバッファであり、MNK乗算命令の3つの行列「第1行列」「第2行列」および「結果行列」をバッファリングするように構成されてよい。L0は計算に近いので、帯域幅および遅延に関する要件はもっとも高く、データ再使用の可能性は最大である。性能を改善して消費電力を低減させるように、A Dトリガ (DFE) がL0を構築するために使用されてよい。フラクタル命令のソースおよび宛先オペランドは、L1 (図17の第5メモリ612および第4メモリ613) から来る。実行の最中に、データはL0 (例えば、図17の第1メモリ601および第2メモリ602) を使用することによって再使用される。上記のフラクタル命令のようなソフトウェアは、L1を使用することによってデータを再使用し得る。マルチレベルバッファにおけるデータ再使用は、フラクタル命令を実行するシーケンスおよびフラクタル命令の上のソフトウェアを制御するシーケンスを使用することによって実装され得る。加えて、マルチレベルバッファのデータを再使用することによって、各バッファのデータのデータ移動時間もまた隠され得る。以下の表の例は、データ再使用と、バッファの複数のレベルの間の移動とを説明し得る。

30

40

【0097】

50

以下の2つの行列：A =

【数20】

$$\begin{bmatrix} A0 & A1 \\ A2 & A3 \end{bmatrix}$$

およびB =

【数21】

$$\begin{bmatrix} A0 & A1 \\ A2 & A3 \end{bmatrix}$$

10

があり、2つの行列のデータ移動ステップが以下の表に示されると前提する。

【表2】

時点	L1からの読み出し	L0への格納	計算
1	A0, B0		
2	B1	A0, B0	A0 * B0
3	A2	A0, B0, B1	A0 * B1
4	A1	A0, A2, B0, B1	A2 * B0
5	B2	A1, A2, B0, B1	A2 * B1
6	B3	A1, A2, B1, B2	A1 * B2
7	A3	A1, A2, B2, B3	A1 * B3
8		A2, A3, B2, B3	A3 * B2
9		A2, A3, B2, B3	A3 * B3

20

30

【0098】

時点1において、コントローラ60はL1バッファから行列のA0およびB0部分を読み出し、A0およびB0部分をL0に格納する。

【0099】

時点2において、A0およびB0フラクタル行列は、L0から読み出され、演算に参与することができる。同時に、ハードウェアはL1からB1フラクタルを読み出し、B1フラクタルをL0に格納し、次の動作に関する準備を行う。加えて、データの読み出し時間もまた、計算によって隠される。この場合、ハードウェアは2つのフラクタル行列の両者を読み出す必要はなく、B1行列のみを読み出す。「A0 * B1」が時点3において行列のために計算されるとき、時点1において格納されたデータA0が再使用される。前述のリストを参照すると、データが各時間単位において再使用されることが、後の計算においてわかることができる。

40

【0100】

本発明のこの実施形態は、L1およびL0の間のデータの移動に限定されるものではないことに留意すべきである。L2（例えば、外部メモリ701および外部メモリ702）からL1バッファへとデータを移動する最中に、データもまた、帯域幅を低減させてエネルギー消費を最適化するように、再使用されてよい。本発明のこの実施形態において、行

50

列分割方式および移動シーケンスは、限定されるものではない。データ再使用は、各時間単位にデータ移動を実現するように、データ移動の最中に最大化されるべきであり、フラクタル行列計算はフルロード状態で実行される。

【0101】

本発明のこの実施形態において、マルチレベルバッファ構造、行列フラクタルデータ再使用、フラクタル命令を実行するシーケンス、およびフラクタル命令の上のソフトウェアを制御するシーケンスを使用することによって、マルチレベルバッファにおけるデータ再使用が実現でき、密結合するオンチップメモリへの依存が低減され、エネルギー効率が最適化され、ソフトウェアプログラミングの複雑性が低減される。

【0102】

本発明のこの実施形態において、行列で乗算演算を実行するための命令を実行するシーケンスは、2つの方式を含む：命令同期実行および命令非同期実行である。

【0103】

本発明のこの実施形態において、フラクタル行列乗算命令が実行される前に、例えば、行列サイズの計算、行列データの読み出し、および宛先アドレスの計算といった、一連の制御準備およびデータ準備が必要とされる。プロセッサの命令実行ポリシーが同期実行である場合、具体的には、すべての命令がシーケンスにコミット (commit) される必要がある場合、関連づけられていない命令が終了するまで、命令の実行が開始しないという可能性が非常に高い。このことは、大きくそして不要である、性能の損失をもたらすことがある。以下の手順は、命令同期実行シーケンスである：アドレス計算 制御準備 行列0の読み出し 行列0の乗算 アドレス計算 制御準備 行列1の読み出し 行列1の乗算。

【0104】

前述の実行シーケンスにおいて、第2の時間の制御準備、アドレス計算、行列1のデータの読み出しは、行列0の乗算の終了に依存せず、そのような追加の時間は、不要な待ち時間をもたらすことがある。この問題を解決するように、本発明のこの実施形態において、ハードウェア命令発送ユニット606は、マルチチャネル順序保持方式で伝送を実行し、それにより、異なる型の命令が同時にかつ順次実行されることを可能にすることを保証する。前述の例において、制御準備およびアドレス計算が順序保持方式によってスカラチャネル上で実行され、行列読み出しおよび格納が順序保持方式によってデータ移動チャネル上で実行され、行列乗算計算もまた、順序保持方式によって行列演算チャネル上で実行される。チャネルは、オーバーラップしてよいが、順序保持はされず、互いに従属した命令が、待ちフラグ (Wait Flag) を設定することによって同期されてよい。命令非同期実行ポリシーを使用することによって、命令は並列で実行され得、これにより、ますます実行効率が增大する。前述の同期実行シーケンスの例において、非同期実行ポリシーが使用される場合、効果は図18に示される。命令非同期実行シーケンスにおいて、命令は順序保持されず、依存関係を有する関連する命令が、ソフトウェアによって追加された待ち命令を使用することによって同期されてよい。フラクタル行列乗算の制御準備オーバーヘッドは、この非同期実行方式を使用することによって隠され得る。フラクタル行列乗算プログラミング方式に適用可能な非同期実行方式が設計される。

【0105】

行列乗算器が提供され、ここで、行列乗算器は行列乗算ブロック分割方法、すなわち、MNKフラクタルを完了して、行列乗算器60における内部コントローラ604の制御ロジックを使用することによって、乗算のために大きい行列を単位行列 (具体的には、 $X * L * L * Y$ 行列) に分割するようにコントローラを使用する。コントローラ604の制御ロジックは、各クロックサイクルにおいて、単位行列乗算タスクを演算回路603に送り、その結果、データがパイプライン方式で実行され、 X 行 * Y 列の演算ユニットがフルロード状態で動作する。行列乗算の効率が增大し、ニューラルネットワークアルゴリズムを大幅に改善する適用効果を実現される。本発明のこの実施形態において提供される行列乗算器は、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およびFC演算を実行し

10

20

30

40

50

得る。

【0106】

前述した実施形態のすべてまたは一部は、ソフトウェア、ハードウェア、ファームウェア、またはそれらの組み合わせによって実装されてよい。実施形態を実装するようにソフトウェアプログラムが使用されるとき、実施形態は、完全に、または部分的に、コンピュータプログラム製品の形態で実装されてよい。コンピュータプログラム製品は、1または複数のコンピュータ命令を含む。コンピュータプログラム命令がコンピュータ上でロードおよび実行されるとき、本願の実施形態による手順または機能が、すべてまたは部分的に生成される。コンピュータは汎用コンピュータ、専用コンピュータ、コンピュータネットワーク、または他のプログラマブルな装置であってよい。コンピュータ命令は、コンピュータ可読記憶媒体に格納されてよく、または、コンピュータ可読記憶媒体から別のコンピュータ可読記憶媒体に伝送されてもよい。例えば、コンピュータ命令は、ウェブサイト、コンピュータ、サーバ、またはデータセンタから、別のウェブサイト、コンピュータ、サーバ、またはデータセンタに、有線（例えば、同軸ケーブル、光ファイバ、またはデジタル加入者線（Digital Subscriber Line, 略してDSL））または無線（例えば、赤外線、無線、およびマイクロ波、または同様のものなど）方式で、伝送されてよい。コンピュータ可読記憶媒体は、コンピュータによってアクセス可能な任意の可用媒体、または、1または複数の可用媒体を統合する、サーバまたはデータセンタなどのデータストレージデバイスであってよい。可用媒体は、磁気媒体（例えば、フロッピー（登録商標）ディスク、ハードディスク、または磁気テープ）、光学媒体（例えば、DVD）、半導体媒体（例えば、ソリッドステートドライブ（Solid State Disk, 略してSSD））、または同様のものであってよい。

10

20

【0107】

本願は、実施形態を参照して説明されてきたが、保護を請求する本願の実装の処理において、当業者は、添付の図面、開示された内容、および添付の特許請求の範囲を閲覧することによって、開示された実施形態の別の変形を理解および実装し得る。請求項において、「含む（comprising）」は、別のコンポーネントまたは別の工程を排除せず、「1つの（a）」または「1つ（one）」は、複数という意味を排除しない。単一のプロセッサまたは別のユニットが、請求項において列挙される複数の機能を実装してもよい。いくつかの測定値が互いに異なる従属請求項に記録されるが、これは、これらの測定値がより良い効果を作り出すように組み合わせられ得ないことを意味しない。

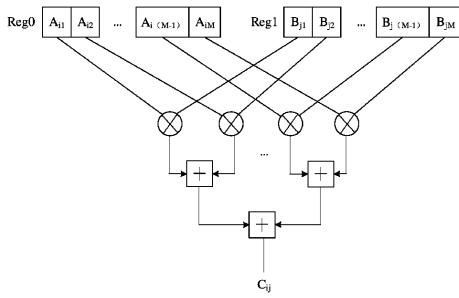
30

【0108】

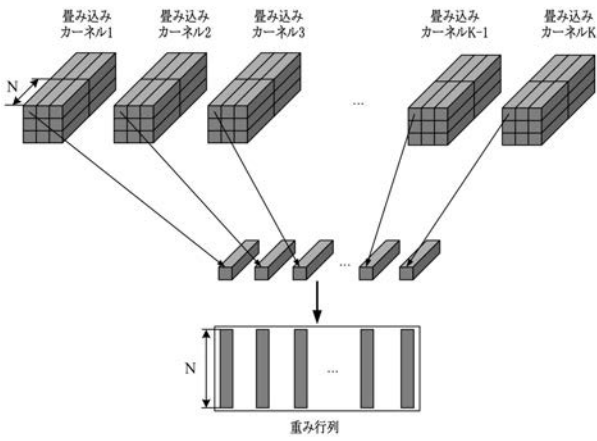
本願は具体的特徴およびそれらの実施形態を参照して説明されてきたが、明らかに、様々な修正および組み合わせが、本願の思想および範囲から逸脱することなく、行われてよい。対応して、明細書および添付の図面は単に、添付の特許請求の範囲によって画定された本願の例としての説明にすぎず、本願の範囲を含める修正、変形、組み合わせ、または均等物の、いずれかまたはすべてとみなされる。明らかに、当業者は、本願の思想および範囲から逸脱することなく、本願の様々な修正および変形を行うことができる。以下の特許請求の範囲およびそれらの均等技術によって画定される保護の範囲内に属するならば、本願は、本願へのこれらの修正および変形を含めることを意図する。

40

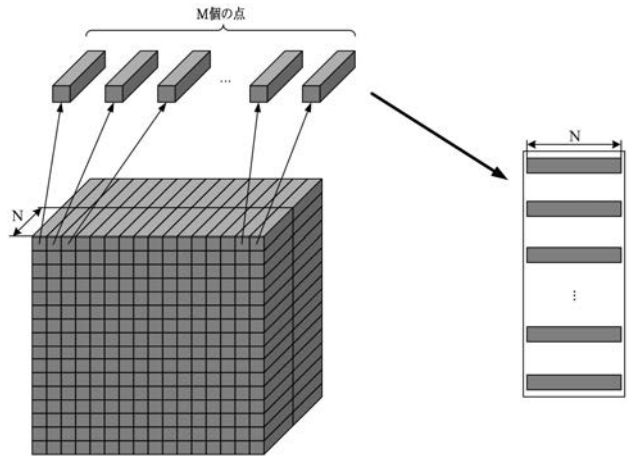
【 図 1 】



【 図 2 】



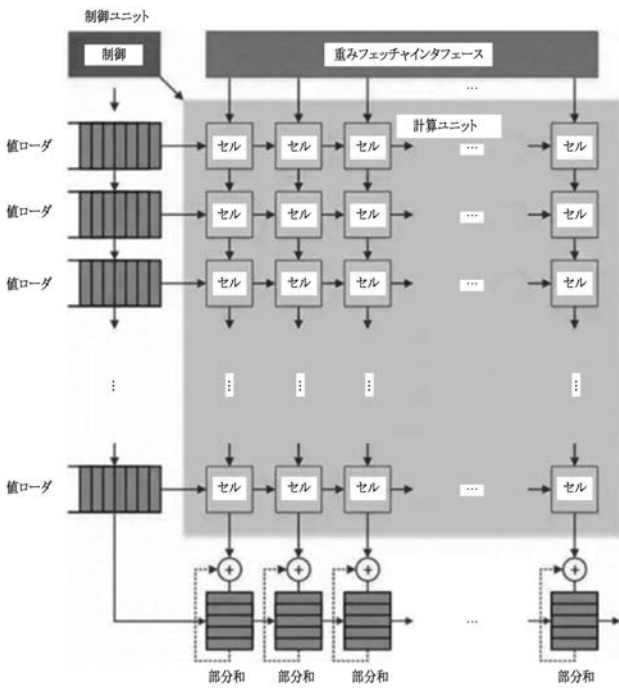
【 図 3 】



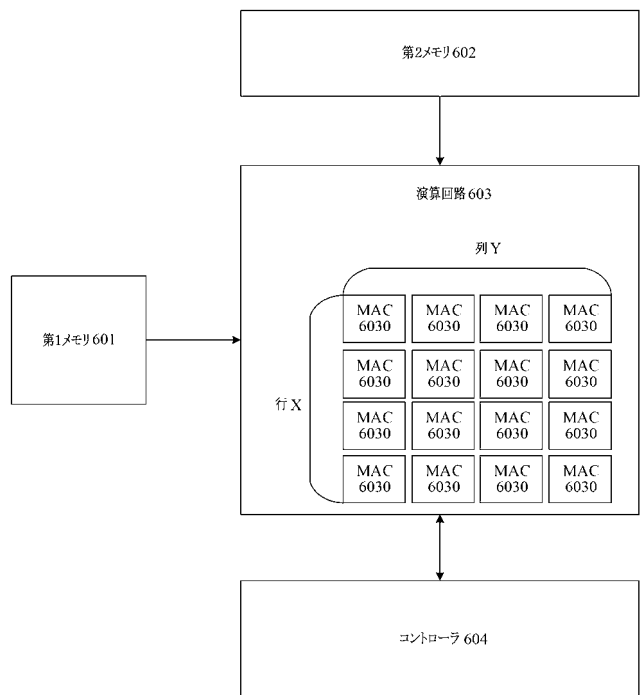
【 図 4 】

$$\begin{matrix}
 \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1K} \\ A_{21} & A_{22} & \dots & A_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ A_{M1} & A_{M2} & \dots & A_{MK} \end{bmatrix} &
 \begin{bmatrix} B_{11} & B_{12} & \dots & B_{1K} \\ B_{21} & B_{22} & \dots & B_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ B_{N1} & B_{N2} & \dots & B_{NK} \end{bmatrix} &
 = &
 \begin{bmatrix} \sum_1^N A_{1j}B_{j1} & \sum_1^N A_{1j}B_{j2} & \dots & \sum_1^N A_{1j}B_{jK} \\ \sum_1^N A_{2j}B_{j1} & \sum_1^N A_{2j}B_{j2} & \dots & \sum_1^N A_{2j}B_{jK} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_1^N A_{Mj}B_{j1} & \sum_1^N A_{Mj}B_{j2} & \dots & \sum_1^N A_{Mj}B_{jK} \end{bmatrix} \\
 \text{行列A} & \text{行列B} & & \text{行列C}
 \end{matrix}$$

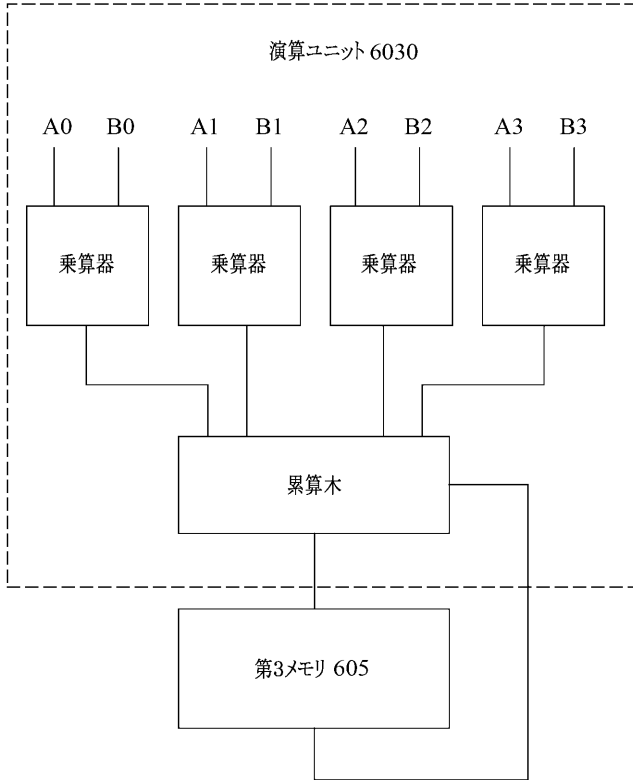
【 図 5 】



【 図 6 】



【 図 7 】



【 図 8 】

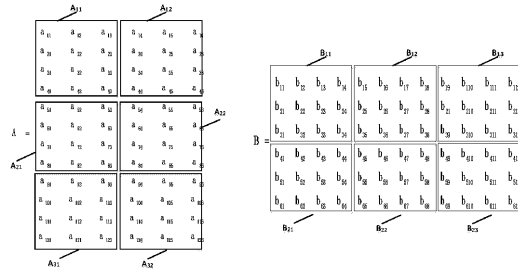


図 8

【 図 9 】

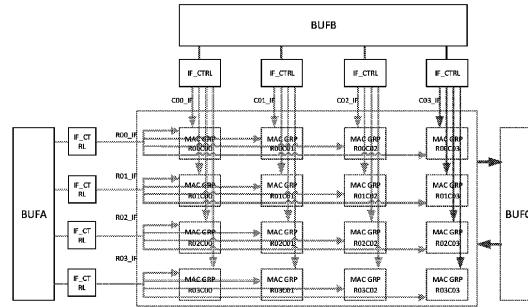


図 9

【 図 1 0 】

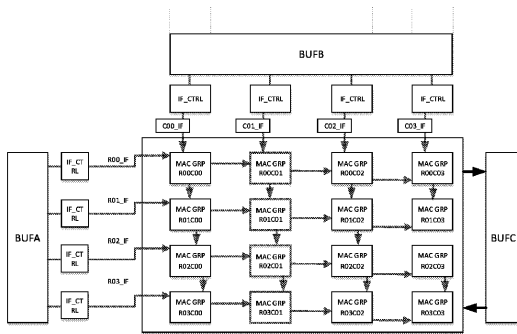


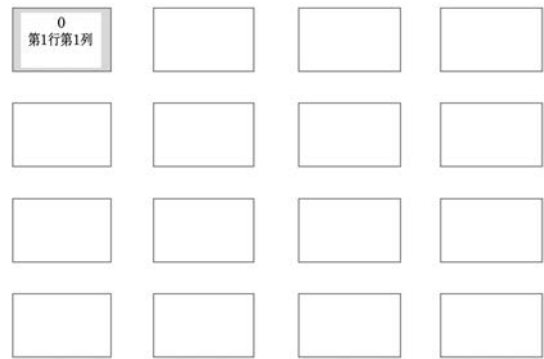
図 10

【 図 1 1 】

クロックサイクル	0	1	2	3	4	5	6
BUFA_0	[0,R0]	[1,R0]	[2,R0]	[3,R0]			
BUFA_1		[0,R1]	[1,R1]	[2,R1]	[3,R1]		
BUFA_2			[0,R2]	[1,R2]	[2,R2]	[3,R2]	
BUFA_3				[0,R3]	[1,R3]	[2,R3]	[3,R3]
クロックサイクル	0	1	2	3	4	5	6
BUFB_0	[0,C0]	[1,C0]	[2,C0]	[3,C0]			
BUFB_1		[0,C1]	[1,C1]	[2,C1]	[3,C1]		
BUFB_2			[0,C2]	[1,C2]	[2,C2]	[3,C2]	
BUFB_3				[0,C3]	[1,C3]	[2,C3]	[3,C3]

【 図 1 2 】

T = 0

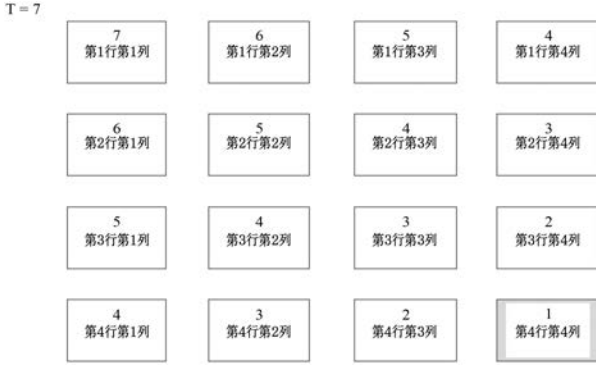


【 図 1 3 】

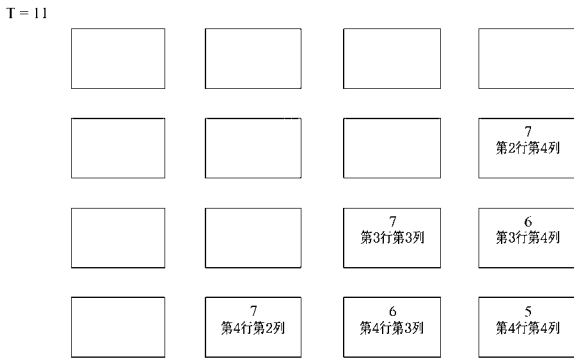
T = 1



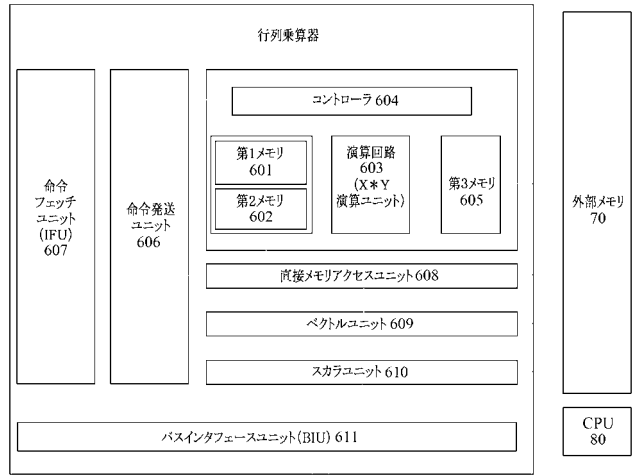
【図 14】



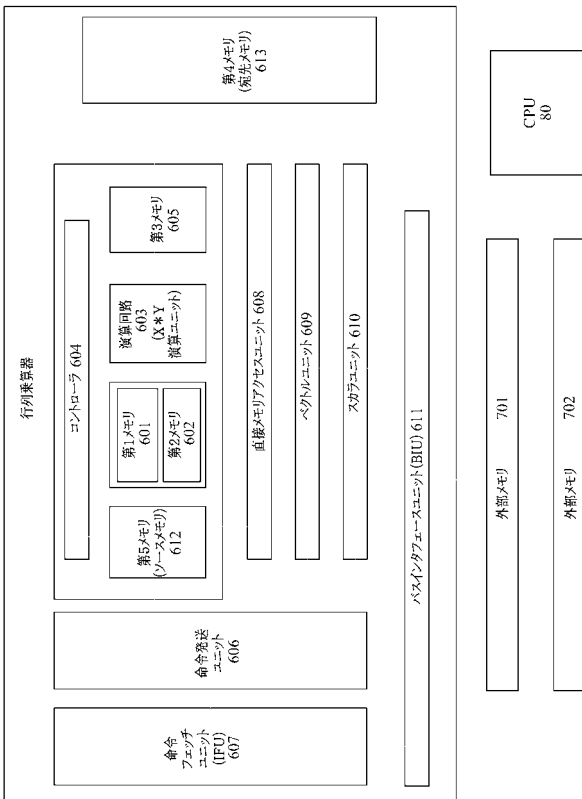
【図 15】



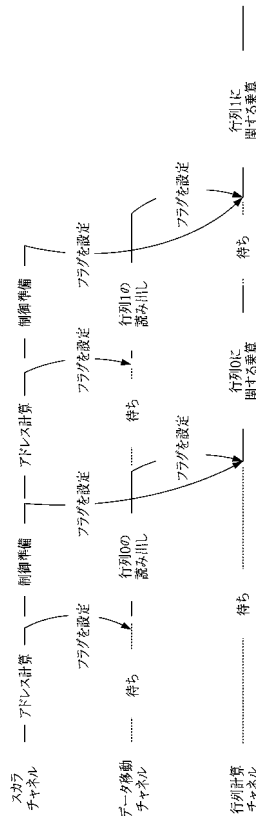
【図 16】



【図 17】



【図 18】



【手続補正書】

【提出日】令和2年8月5日(2020.8.5)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

M * K 行列である第 1 行列を格納するように構成された第 1 メモリと、
 K * N 行列である第 2 行列を格納するように構成された第 2 メモリと、
 前記第 1 メモリおよび前記第 2 メモリに接続される演算回路と、
 前記演算回路に接続されるコントローラと、を含む、
 行列乗算器であって、

前記演算回路は、X 行 * Y 列からなる演算ユニットを含み、各前記演算ユニットは、ベクトル乗算回路および加算回路を含み、前記ベクトル乗算回路は、前記第 1 メモリによって送られる行ベクトルのデータおよび前記第 2 メモリによって送られる列ベクトルのデータを受信し、前記 2 つのベクトルを乗算するように構成され、前記加算回路は、前記 2 つのベクトルの乗算によって取得された結果を加算し、同一の演算ユニットの計算結果を累算し、各演算ユニットの演算結果を取得するように構成され、

前記コントローラは、以下の動作、すなわち、

前記第 1 行列を、サイズが X * L であるサブブロックを単位とするブロックに分割し、同じサイズの S * R 個のサブブロックを取得し、前記 S * R 個のサブブロックのうち第 s 行第 r 列におけるサブブロックは $A_{s,r}$ 、 $s = (1, 2, 3, \dots, \text{および } S)$ 、および $r = (1, 2, 3, \dots, \text{および } R)$ で表される、動作と、

前記第 2 行列を、サイズが L * Y であるサブブロックを単位とするブロックに分割し、同じサイズの R * T 個のサブブロックを取得し、R * T 個のサブブロックのうち第 r 行第 t 列におけるサブブロックは、 $B_{r,t}$ 、 $r = (1, 2, 3, \dots, \text{および } R)$ 、 $t = (1, 2, 3, \dots, \text{および } T)$ で表される、動作と

を実行するように構成され、

前記コントローラは、さらに以下の動作、すなわち、

任意のサブブロック $A_{s,r}$ の X 個の行ベクトルにおける第 x 行と、対応するサブブロック $B_{r,t}$ の Y 個の列ベクトルにおける第 y 列とを、X 行 * Y 列からなる演算ユニットの第 x 行第 y 列において前記演算ユニットに入力し、それにより、処理を実行する動作を実行するように構成され、 $x = (1, 2, 3, \dots, \text{および } X)$ 、 $y = (1, 2, 3, \dots, \text{および } Y)$ であり、前記任意のサブブロック $A_{s,r}$ における r と、前記対応するサブブロック $B_{r,t}$ における r とは同じ値を有する、

行列乗算器。

【請求項2】

前記コントローラは、以下の動作、すなわち、

前記任意のサブブロック $A_{s,r}$ の前記 X 個の行ベクトルにおける前記第 x 行と、前記対応するサブブロック $B_{r,t}$ の前記 Y 個の列ベクトルにおける前記第 y 列とを、同じクロックサイクルにおいて並行して、X 行 * Y 列からなる前記演算ユニットの第 x 行第 y 列において前記演算ユニットに入力し、それにより前記処理を実行する、動作

を実行するように具体的に構成される、請求項1に記載の行列乗算器。

【請求項3】

前記コントローラはさらに、前記任意のサブブロック $A_{s,r}$ の行ベクトルを、x 個の行番号の昇順で、X 行 * Y 列からなる前記演算ユニットに対応する第 x 行に連続的に入力するように制御するように構成され、近接する行ベクトルが同じ列で異なる行の演算ユニットに入る時点の間の差は1クロックサイクルであり、前記コントローラはさらに、前記対

応するサブブロック $B_{r,t}$ の列ベクトルを、 y 個の列番号の昇順で、 X 行 * Y 列からなる前記演算ユニットに対応する第 y 列に連続的に入力するように同時に制御するように構成され、近接する列ベクトルが同じ行で異なる列の演算ユニットに入る時点の間の差は 1 クロックサイクルである、

請求項 1 または 2 に記載の行列乗算器。

【請求項 4】

前記コントローラはさらに、 s および r の値を変更されないままにして、 t の値を、少なくとも 2 つの連続するサブブロック乗算計算サイクルにおいて変更されるように制御するように構成され、その結果、前記第 1 メモリは、前記少なくとも 2 つの連続するサブブロック乗算計算サイクル内で同じサブブロック $A_{s,r}$ を再使用し、前記サブブロック乗算計算サイクルは、1 つのサブブロック $A_{s,r}$ および対応するサブブロック $B_{r,t}$ 上での行列乗算演算を完了させるように、 X 行 * Y 列の前記演算ユニットによって使用される時間である、

請求項 1 から 3 のいずれか一項に記載の行列乗算器。

【請求項 5】

前記行列乗算器はさらに、前記演算回路に接続された第 3 メモリを含み、

前記コントローラは、前記ベクトル乗算回路および前記加算回路の演算結果を前記第 3 メモリに格納するように、 X 行 * Y 列の前記演算ユニットを制御するように構成される、

請求項 1 から 4 のいずれか一項に記載の行列乗算器。

【請求項 6】

前記行列乗算器はさらに、前記第 1 メモリおよび前記第 2 メモリに接続される第 4 メモリと、前記第 3 メモリに接続される第 5 メモリとを含み、

前記コントローラはさらに、前記第 1 行列および前記第 2 行列の乗算演算を実行する前に、

前記第 4 メモリから、前記第 1 行列および前記第 2 行列のデータソースを、それぞれ前記第 1 メモリおよび前記第 2 メモリに移動させ、前記第 3 メモリから、前記計算結果を前記第 5 メモリに移動させるように、制御するように構成される、

請求項 5 に記載の行列乗算器。

【請求項 7】

前記ベクトル乗算回路は L 個の乗算器を含み、前記加算回路は入力数が $L + 1$ である加算木を含む、

請求項 1 から 6 のいずれか一項に記載の行列乗算器。

【請求項 8】

前記第 1 メモリ、前記第 2 メモリ、前記演算回路、および前記コントローラはバスインタフェースユニットを使用して接続される、

請求項 1 から 7 のいずれか一項に記載の行列乗算器。

【請求項 9】

$S =$

【数 2 2】

$$\begin{cases} M/X, M\%X = 0 \\ \lceil \frac{M}{X} \rceil + 1, M\%X \neq 0 \end{cases}$$

および

$R =$

【数 2 3】

$$\begin{cases} K/L, K\%L = 0 \\ \lceil \frac{K}{L} \rceil + 1, K\%L \neq 0 \end{cases}$$

であり、

$M \% X = 0$ のとき、計算は前記第 1 行列の第 $(M + 1)$ 行から第 $(S * X - M)$ 行まで実行されず、結果の値には 0 が割り当てられ、 $K \% Y = 0$ のとき、計算は前記第 1 行列の第 $(K + 1)$ 列から第 $(R * Y - K)$ 列まで実行されず、結果の値には 0 が割り当てられる、

請求項 1 から 8 のいずれか一項に記載の行列乗算器。

【請求項 1 0】

R =

【数 2 4】

$$\begin{cases} K/L, K\%L = 0 \\ \left\lfloor \frac{K}{L} \right\rfloor + 1, K\%L \neq 0 \end{cases}$$

および

T =

【数 2 5】

$$\begin{cases} N/Y, N\%Y = 0 \\ \left\lfloor \frac{N}{Y} \right\rfloor + 1, N\%Y \neq 0 \end{cases}$$

であり、

$K \% Y = 0$ のとき、計算は前記第 1 行列の第 $(K + 1)$ 列から第 $(R * Y - K)$ 列まで実行されず、結果の値には 0 が割り当てられ、 $N \% X = 0$ のとき、計算は前記第 1 行列の第 $(N + 1)$ 行から第 $(T * X - N)$ 行まで実行されず、結果の値には 0 が割り当てられる、

請求項 1 から 8 のいずれか一項に記載の行列乗算器。

【手続補正 2】

【補正対象書類名】明細書

【補正対象項目名】0 0 0 8

【補正方法】変更

【補正の内容】

【0 0 0 8】

第 1 の態様によると、本発明の実施形態は行列乗算器を提供し、行列乗算器は、

$M * K$ 行列である第 1 行列を格納するように構成された第 1 メモリと、

$K * N$ 行列である第 2 行列を格納するように構成された第 2 メモリと、

第 1 メモリおよび第 2 メモリに接続される演算回路と、

演算回路に接続されたコントローラと、を含み、

演算回路は X 行 * Y 列からなる演算ユニットを含み、各演算ユニットはベクトル乗算回路および加算回路を含み、ベクトル乗算回路は、第 1 メモリによって送られる行ベクトルのデータおよび第 2 メモリによって送られる列ベクトルのデータを受信し、2 つのベクトルを乗算するように構成され、加算回路は、2 つのベクトルの乗算によって取得された結果を加算し、同一の演算ユニットの計算結果を累算し、各演算ユニットの演算結果を取得するように構成される、演算回路と、

コントローラは以下の動作、すなわち、

第 1 行列を、サイズが $X * L$ であるサブブロックを単位とするブロックに分割し、同じサイズの $S * R$ 個のサブブロックを取得し、 $S * R$ 個のサブブロックのうち第 s 行第 r 列におけるサブブロックは $A_{s, r}$ 、 $s = (1, 2, 3, \dots, \text{および } S)$ 、および $r = (1, 2, 3, \dots, \text{および } R)$ で表される、動作と、

第 2 行列を、サイズが $L * Y$ であるサブブロックを単位とするブロックに分割し、同じサイズの $R * T$ 個のサブブロックを取得し、 $R * T$ 個のサブブロックのうち第 r 行第 t 列

におけるサブブロックは、 $B_{r,t}$ 、 $r = (1, 2, 3, \dots, \text{および } R)$ 、 $t = (1, 2, 3, \dots, \text{および } T)$ で表される、動作とを実行するように構成され、

コントローラは、さらに以下の動作、すなわち、

任意のサブブロック $A_{s,r}$ の X 個の行ベクトルにおける第 x 行と、対応するサブブロック $B_{r,t}$ の Y 個の列ベクトルにおける第 y 列とを、 X 行 * Y 列からなる演算ユニットの第 x 行第 y 列において演算ユニットに入力し、それにより、処理を実行する動作を実行するように構成され、 $x = (1, 2, 3, \dots, \text{および } X)$ 、 $y = (1, 2, 3, \dots, \text{および } Y)$ であり、任意のサブブロック $A_{s,r}$ における r と、対応するサブブロック $B_{r,t}$ における r とは同じ値を有する。

【手続補正 3】

【補正対象書類名】明細書

【補正対象項目名】0011

【補正方法】変更

【補正の内容】

【0011】

可能な実装において、コントローラはさらに、任意のサブブロック $A_{s,r}$ の行ベクトルを、 x 個の行番号の昇順で、 X 行 * Y 列からなる演算ユニットに対応する第 x 行に連続的に入力するように制御するように構成され、近接する行ベクトルが同じ列で異なる行の演算ユニットに入る時点の間の差は 1 クロックサイクルであり、コントローラはさらに、対応するサブブロック $B_{r,t}$ の列ベクトルを、 y 個の列番号の昇順で、 X 行 * Y 列からなる演算ユニットに対応する第 y 列に連続的に入力するように同時に制御するように構成され、近接する列ベクトルが同じ行で異なる列の演算ユニットに入る時点の間の差は 1 クロックサイクルである。

【手続補正 4】

【補正対象書類名】明細書

【補正対象項目名】0017

【補正方法】変更

【補正の内容】

【0017】

可能な実装において、 $S =$

【数 1】

$$\begin{cases} M/X, M\%X = 0 \\ \left[\frac{M}{X} \right] + 1, M\%X \neq 0 \end{cases}$$

および $R =$

【数 2】

$$\begin{cases} K/L, K\%L = 0 \\ \left[\frac{K}{L} \right] + 1, K\%L \neq 0 \end{cases}$$

であり、

$M\%X = 0$ のとき、計算は第 1 行列の第 $(M + 1)$ 行から第 $(S * X - M)$ 行まで実行されず、結果の値には 0 が割り当てられ、 $K\%Y = 0$ のとき、計算は第 1 行列の第 $(K + 1)$ 列から第 $(R * Y - K)$ 列まで実行されず、結果の値には 0 が割り当てられる。

【手続補正 5】

【補正対象書類名】明細書

【補正対象項目名】0055

【補正方法】変更

【補正の内容】

【0055】

前述の分析に基づいて、本願は、行列乗算アクセラレータを提供し、本願において提供される技術的問題を具体的に解析および解決する。図6は、本発明の実施形態による行列乗算器60の構造図である。図6に示されるように、行列乗算器60は第1メモリ601、第2メモリ602、演算回路603、およびコントローラ604を含む。演算回路603は、バスを使用して、第1メモリ601、第2メモリ602、およびコントローラ604とデータ通信を実行し得る。演算回路603は、第1メモリ601および第2メモリ602から行列データを取り出し、ベクトル乗算および加算演算を実行するように構成される。コントローラ604は、ベクトル演算を完了するように、予め設定されたプログラムまたは命令に従って、演算回路603を制御するように構成される。第1メモリ601は第1行列を格納するように構成される。

【手続補正6】

【補正対象書類名】明細書

【補正対象項目名】0057

【補正方法】変更

【補正の内容】

【0057】

本発明のこの実施形態において説明される第1メモリ601、および、以下に説明される関連する行列乗算器の第2メモリ602、第3メモリ605、および内部メモリは、それぞれ、レジスタ、ランダムアクセスメモリ(random access memory、略してRAM)、静的ランダムアクセスメモリ、フラッシュメモリ、または別の読み出しおよび書き込み可能メモリであってよい。本願において、第1行列、第2行列および演算結果のデータ型はそれぞれ、int8、fp16、またはfp32などの型であってよい。

【手続補正7】

【補正対象書類名】明細書

【補正対象項目名】0060

【補正方法】変更

【補正の内容】

【0060】

演算回路603は、 X 行 \times Y 列の演算ユニット6030(乗算累算ユニットMACと称されてよい)を含み得る。各演算ユニットは、独立してベクトル乗算演算を実行し得る。図6において、演算回路603が 4×4 演算ユニット6030を含む例が図に使用され、すなわち、 $X = 4$ および $Y = 4$ である。演算ユニット6030は、それぞれ第1メモリ601によって送られた行ベクトルと、第2メモリ602によって送られた列ベクトルとを受信し、行ベクトルと列ベクトルとのベクトル乗算演算を実行するように使用される、2つの入力を提供される。具体的には、1つの演算回路6030はベクトル乗算回路および加算回路を含み、ここで、ベクトル乗算回路は第1メモリ601によって送られる行ベクトルのデータと、第2メモリ602によって送られる列ベクトルのデータとを受信し、2つのベクトルを乗算するように構成され、加算回路は、2つのベクトルの乗算によって取得された結果を加算し、同一の演算ユニットの計算結果を累算し、演算ユニット6030の演算結果を取得するように構成される。

【手続補正8】

【補正対象書類名】明細書

【補正対象項目名】0064

【補正方法】変更

【補正の内容】

【0064】

第1行列が $M \times K$ 行列であり、第1行列が整数個の $X \times L$ サブブロックで正確に分割できない場合が存在し得ることに、留意すべきである。したがって、 M/X または K/L が

整数でないとき、演算は要素 0 をパディングする方式で実行されてよい。代替的に、対応する位置で全く計算が実行されず、結果の値に 0 が割り当てられる。具体的には、

【数 5】

$$S = \begin{cases} M/X, M\%X = 0 \\ \lceil \frac{M}{X} \rceil + 1, M\%X \neq 0 \end{cases}$$

および

【数 6】

$$R = \begin{cases} K/L, K\%L = 0 \\ \lceil \frac{K}{L} \rceil + 1, K\%L \neq 0 \end{cases}$$

であり、 $M\%X = 0$ のとき、計算は第 1 行列の第 $(M + 1)$ 行から第 $(S * X - M)$ 行まで実行されず、結果の値には 0 が割り当てられ、 $K\%L = 0$ のとき、計算は第 1 行列の第 $(K + 1)$ 列から第 $(R * Y - K)$ 列まで実行されず、結果の値には 0 が割り当てられる。言い換えれば、演算ユニットは対応する行および列において実体的乗算計算を実行せず、処理のために、演算が実行されたが結果が 0 であるとみなす。このようにして、対応する演算ユニットの読み出しおよび演算電力消費は低減され得る。

【手続補正 9】

【補正対象書類名】明細書

【補正対象項目名】0071

【補正方法】変更

【補正の内容】

【0071】

方式 1：行列乗算シーケンスにおいて、例えば、サブブロックはサブブロック A_{11} およびサブブロック B_{11} であってよい。 A_{11} のすべての行ベクトルと、対応する B_{11} のすべての列ベクトルとが、サブブロック乗算計算サイクル（第 1 のラウンドとして理解されてよい）に入力され、それにより、演算を実行する。 A_{12} のすべての行ベクトルおよび対応する B_{21} のすべての列ベクトルでの演算が、第 2 サブブロック乗算計算サイクル（第 2 のラウンドとして理解されてよい）において実行される。このようにして、演算ユニットが累算を実行した後、結果行列 C の第 1 行第 1 列における結果点 C_{11} の値が取得され得る。同様に、結果行列 C のすべての位置における結果点を取得され得る。実際、

$$C_{11} =$$

【数 12】

$$A_{11}B_{11} + A_{12}B_{21}$$

ここで、

【数 13】

$$A_{11} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} \end{bmatrix}$$

【数 1 4】

$$A_{12} = \begin{bmatrix} a_{14} & a_{15} & a_{16} \\ a_{24} & a_{25} & a_{26} \\ a_{34} & a_{35} & a_{36} \\ a_{44} & a_{45} & a_{46} \end{bmatrix}$$

【数 1 5】

$$B_{11} = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \\ b_{41} & b_{42} & b_{43} \end{bmatrix}$$

【数 1 6】

$$B_{12} = \begin{bmatrix} b_{14} & b_{15} & b_{16} \\ b_{24} & b_{25} & b_{26} \\ b_{34} & b_{35} & b_{36} \\ b_{44} & b_{45} & b_{46} \end{bmatrix}$$

【手続補正 1 0】

【補正対象書類名】明細書

【補正対象項目名】0 0 8 1

【補正方法】変更

【補正の内容】

【0 0 8 1】

図 1 0 は、本発明の実施形態による特定の演算回路 6 0 3 の配線の概略図である。演算回路 6 0 3 において、図 1 0 に対応して、シストリックアレイ構造が提供される。具体的に、コントローラ 6 0 4 はさらに、任意のサブブロック $A_{s,r}$ の行ベクトルを、 x 個の行番号の昇順で、 X 行 * Y 列からなる演算ユニットに対応する第 x 行に連続的に入力するように制御するように構成され、近接する行ベクトルが同じ列で異なる行の演算ユニットに入る時点の間の差は 1 クロックサイクルであり、コントローラ 6 0 4 はさらに、対応するサブブロック $B_{r,t}$ の列ベクトルを、 y 個の列番号の昇順で、 X 行 * Y 列からなる演算ユニットに対応する第 y 列に連続的に入力するように同時に制御するようにさらに構成され、近接する列ベクトルが同じ行で異なる列の演算ユニットに入る時点の間の差は 1 クロックサイクルである。

【手続補正 1 1】

【補正対象書類名】明細書

【補正対象項目名】0 0 8 3

【補正方法】変更

【補正の内容】

【0 0 8 3】

シストリックアレイアーキテクチャに基づいて、図 1 0 に対応する配線構造において、

B U F A / B が、それぞれ、第 1 行列 / 第 2 行列をバッファリングするように構成されるメモリである。図 10 において、第 1 行列バッファ (B U F A) は、各クロックサイクルにおいて、第 1 行列における単位行列を X 個の行に分割し、同じ行の L 個の要素をシストリックアレイの演算ユニットに順次送る。同様に、第 2 行列バッファ (B U F B) は、各クロックサイクルにおいて、第 2 行列における単位行列を Y 個の列に分割し、同じ列の L 個の要素をシストリックアレイに順次送る。具体的な時間シーケンスは以下の通りである。

【手続補正 12】

【補正対象書類名】明細書

【補正対象項目名】0085

【補正方法】変更

【補正の内容】

【0085】

M = 2、N = 2、および K = 2 (すなわち 8 × 8 * 8 × 8 行列乗算が使用される例。行列乗算器 60 のコントローラ 604 は、図 11 のフォーマットで行列乗算を分割し、全体で 8 個の 4 × 4 単位行列演算を取得する。M N K 行列乗算演算に関して、分割シーケンスのための多くの可能性があり、分割シーケンスのルールは、M N K 行列乗算演算が、方式 1 または方式 2 のシーケンスにおいて実行されてよいことである。方式 2 において再使用するデータの最大数倍のポリシーを使用することによって、データを読み出すための消費電力を低減し得ることが、理解されてよい。M N K フラクタル分割が実行された後、図 12 から図 15 に示されるように、コントローラ 604 の制御ロジックは 8 クロックサイクルに 8 個のフラクタルをシストリックアレイに入力する。図 12 は、M = 2、N = 2 および K = 2 のときの、時点 T = 0 におけるフラクタル行列乗算器のパイプライン実行を示し、図 13 は、M = 2、N = 2 および K = 2 のときの、時点 T = 1 における フラクタル行列乗算器のパイプライン実行を示し、図 14 は、M = 2、N = 2 および K = 2 のときの、時点 T = 7 におけるフラクタル行列乗算器のパイプライン実行を示し、図 15 は、M = 2、N = 2 および K = 2 のときの、時点 T = 11 におけるフラクタル行列乗算器のパイプライン実行を示す。シストリックアレイは、T = 6 のとき、すなわち、第 7 クロックサイクルにおいてフルロード状態で動作を開始することがわかる。最後の 6 クロックサイクルにおいて、単位行列がシストリックアレイから出力され、行列全体の乗算演算もまた完了する。

【手続補正 13】

【補正対象書類名】明細書

【補正対象項目名】0089

【補正方法】変更

【補正の内容】

【0089】

3. 直接メモリアクセスユニット (D i r e c t M e m o r y A c c e s s U n i t , D M A U n i t) は、各ストレージユニットにデータを移動するように、例えば、L1 RAM から L0 RAM へデータを移動するように構成される。本発明のこの実施形態における直接メモリアクセスユニットが、行列乗算器の外部メモリまたは内部メモリから、乗算演算に関与する行列データを移動するとき、直接メモリアクセスユニットは、行列がブロックに分割された後に取得された結果を格納する必要がある。例えば、2 * 2 行列に関して、第 1 行列の第 1 行第 1 列のサブブロック A₁₁ =

【数 19】

$$\begin{bmatrix} A0 & A1 \\ A2 & A3 \end{bmatrix}$$

がサブブロックのユニットに格納され、A0、A1、A2 および A3 が 1 つの行に格納

される、などである。このようにして、第 1 行列が、対応する第 1 メモリへと移動されてよいとき、または、第 2 行列が、対応する第 2 メモリに移動されてよいときは、ストレージは、前述の方式で実行されてよい。演算ユニットが読み出しの実行を必要とするとき、演算ユニットもまた、前述のストレージシーケンスにおいて読み出しを実行してよく、それにより、計算を円滑化する。行ベクトルが列ベクトルに転置される必要があるとき、転置はフレキシブルに、そして素早く実行されてよい。

【手続補正 14】

【補正対象書類名】明細書

【補正対象項目名】0098

【補正方法】変更

【補正の内容】

【0098】

時点 1 において、コントローラ 604 は L1 バッファから行列の A0 および B0 部分を読み出し、A0 および B0 部分を L0 に格納する。

【手続補正 15】

【補正対象書類名】明細書

【補正対象項目名】0108

【補正方法】変更

【補正の内容】

【0108】

本願は具体的特徴およびそれらの実施形態を参照して説明されてきたが、明らかに、様々な修正および組み合わせが、本願の範囲から逸脱することなく、行われてよい。対応して、明細書および添付の図面は単に、添付の特許請求の範囲によって画定された本願の例としての説明にすぎず、本願の範囲を含める修正、変形、組み合わせ、または均等物の、いずれかまたはすべてとみなされる。明らかに、当業者は、本願の思想および範囲から逸脱することなく、本願の様々な修正および変形を行うことができる。以下の特許請求の範囲およびそれらの均等技術によって画定される保護の範囲内に属するならば、本願は、本願へのこれらの修正および変形を含めることを意図する。

(項目 1)

M * K 行列である第 1 行列を格納するように構成された第 1 メモリと、

K * N 行列である第 2 行列を格納するように構成された第 2 メモリと、

上記第 1 メモリおよび上記第 2 メモリに接続される演算回路と、

上記演算回路に接続されるコントローラと、を含む、

行列乗算器であって、

上記演算回路は、X 行 * Y 列からなる演算ユニットを含み、各上記演算ユニットは、ベクトル乗算回路および加算回路を含み、上記行列乗算回路は、上記第 1 メモリによって送られる行ベクトルのデータおよび上記第 2 メモリによって送られる列ベクトルのデータを受信し、上記 2 つのベクトルを乗算するように構成され、上記加算回路は、上記 2 つのベクトルの乗算によって取得された結果を加算し、同一の演算ユニットの計算結果を累算し、各演算ユニットの演算結果を取得するように構成され、

上記コントローラは、以下の動作、すなわち、

上記第 1 行列を、サイズが X * L であるサブブロックを単位とするブロックに分割し、同じサイズの S * R 個のサブブロックを取得し、上記 S * R 個のサブブロックのうち第 s 行第 r 列におけるサブブロックは $A_{s,r}$ 、 $s = (1, 2, 3, \dots, \text{および } S)$ 、および $r = (1, 2, 3, \dots, \text{および } R)$ で表される、動作と、

上記第 2 行列を、サイズが L * Y であるサブブロックを単位とするブロックに分割し、同じサイズの R * T 個のサブブロックを取得し、R * T 個のサブブロックのうち第 r 行第 t 列におけるサブブロックは $B_{r,t}$ 、 $r = (1, 2, 3, \dots, \text{および } R)$ 、 $t = (1, 2, 3, \dots, \text{および } T)$ で表される、動作と

を実行するように構成され、

上記コントローラは、さらに以下の動作、すなわち、

任意のサブブロック $A_{s,r}$ の X 個の行ベクトルにおける第 x 行と、対応するサブブロック $B_{r,t}$ の Y 個の列ベクトルにおける第 y 列とを、 X 行 \times Y 列からなる演算ユニットの第 x 行第 y 列において上記演算ユニットに入力し、それにより、処理を実行する動作を実行するように構成され、 $x = (1, 2, 3, \dots, \text{および } X)$ 、 $y = (1, 2, 3, \dots, \text{および } Y)$ であり、上記任意のサブブロック $A_{s,r}$ における r と、上記対応するサブブロック $B_{r,t}$ における r とは同じ値を有する、

行列乗算器。

(項目 2)

上記コントローラは、以下の動作、すなわち、

上記任意のサブブロック $A_{s,r}$ の上記 X 個の行ベクトルにおける上記第 x 行と、上記対応するサブブロック $B_{r,t}$ の上記 Y 個の列ベクトルにおける上記第 y 列とを、同じクロックサイクルにおいて並行して、 X 行 \times Y 列からなる上記演算ユニットの第 x 行第 y 列において上記演算ユニットに入力し、それにより上記処理を実行する、動作

を実行するように具体的に構成される、項目 1 に記載の行列乗算器。

(項目 3)

上記コントローラはさらに、上記任意のサブブロック $A_{s,r}$ の行ベクトルを、 x 個の行番号の昇順で、 X 行 \times Y 列からなる上記演算ユニットに対応する第 x 行に連続的に入力するように制御するように構成され、近接する行ベクトルが同じ列で異なる行の演算ユニットに入る時点の間の差は 1 クロックサイクルであり、上記コントローラはさらに、上記対応するサブブロック $B_{r,t}$ の列ベクトルを、 y 個の列番号の昇順で、 X 行 \times Y 列からなる上記演算ユニットに対応する第 y 行に連続的に入力するように同時に制御するように構成され、近接する列ベクトルが同じ行で異なる列の演算ユニットに入る時点の間の差は 1 クロックサイクルである、

項目 1 または 2 に記載の行列乗算器。

(項目 4)

上記コントローラはさらに、 s および r の値を変更されないままにして、 t の値を、少なくとも 2 つの連続するサブブロック乗算計算サイクルにおいて変更されるように制御するように構成され、その結果、上記第 1 メモリは、上記少なくとも 2 つの連続するサブブロック乗算計算サイクル内で同じサブブロック $A_{s,r}$ を再使用し、上記サブブロック乗算計算サイクルは、1 つのサブブロック $A_{s,r}$ および対応するサブブロック $B_{r,t}$ 上での行列乗算演算を完了させるように、 X 行 \times Y 列の上記演算ユニットによって使用される時間である、

項目 1 から 3 のいずれか一項に記載の行列乗算器。

(項目 5)

上記行列乗算器はさらに、上記演算回路に接続された第 3 メモリを含み、

上記コントローラは、上記ベクトル乗算回路および上記加算回路の演算結果を上記第 3 メモリに格納するように、 X 行 \times Y 列の上記演算ユニットを制御するように構成される、

項目 1 から 4 のいずれか一項に記載の行列乗算器。

(項目 6)

上記行列乗算器はさらに、上記第 1 メモリおよび上記第 2 メモリに接続される第 4 メモリと、上記第 3 メモリに接続される第 5 メモリとを含み、

上記コントローラはさらに、上記第 1 行列および上記第 2 行列の乗算演算を実行する前に、

上記第 4 メモリから、上記第 1 行列および上記第 2 行列のデータソースを、それぞれ上記第 1 メモリおよび上記第 2 メモリに移動させ、上記第 3 メモリから、上記計算結果を上記第 5 メモリに移動させるように、制御するように構成される、

項目 5 に記載の行列乗算器。

(項目 7)

上記ベクトル乗算回路は L 個の乗算器を含み、上記加算回路は入力数が $L + 1$ である加

算木を含む、

項目 1 から 6 のいずれか一項に記載の行列乗算器。

(項目 8)

上記第 1 メモリ、上記第 2 メモリ、上記演算回路、および上記コントローラはバスインタフェースユニットを使用して接続される、

項目 1 から 7 のいずれか一項に記載の行列乗算器。

(項目 9)

S =

(数 2 2)

$$\begin{cases} M/X, M\%X = 0 \\ \left[\frac{M}{X} \right] + 1, M\%X \neq 0 \end{cases}$$

および

R =

(数 2 3)

$$\begin{cases} K/L, K\%L = 0 \\ \left[\frac{K}{L} \right] + 1, K\%L \neq 0 \end{cases}$$

であり、

M % X = 0 のとき、計算は上記第 1 行列の第 (M + 1) 行から第 (S * X - M) 行まで実行されず、結果の値には 0 が割り当てられ、K % Y = 0 のとき、計算は上記第 1 行列の第 (K + 1) 行から第 (R * Y - K) 行まで実行されず、結果の値には 0 が割り当てられる、

項目 1 から 8 のいずれか一項に記載の行列乗算器。

(項目 10)

R =

(数 2 4)

$$\begin{cases} K/L, K\%L = 0 \\ \left[\frac{K}{L} \right] + 1, K\%L \neq 0 \end{cases}$$

および

T =

(数 2 5)

$$\begin{cases} N/Y, N\%Y = 0 \\ \left[\frac{N}{Y} \right] + 1, N\%Y \neq 0 \end{cases}$$

であり、

K % Y = 0 のとき、計算は上記第 1 行列の第 (K + 1) 列から第 (R * Y - K) 列まで実行されず、結果の値には 0 が割り当てられ、N % X = 0 のとき、計算は上記第 1 行列の第 (N + 1) 行から第 (T * X - N) 行まで実行されず、結果の値には 0 が割り当てられる、

項目 1 から 8 のいずれか一項に記載の行列乗算器。

【 国际调查报告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/CN2018/111077
A. CLASSIFICATION OF SUBJECT MATTER G06F 17/16(2006.01)i According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) G06F17; G06F7 Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) WPI, EPODOC, CNPAT, CNKI: 矩阵乘法器, 存储器, 计算, MATRIX MULTIPLICATION DEVICE, MEMORY, CALCULATING		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	CN 101086699 A (ZHEJIANG UNIVERSITY) 12 December 2007 (2007-12-12) entire document	1-10
A	CN 107315574 A (BEIJING CAMBRICON TECHNOLOGY CO., LTD.) 03 November 2017 (2017-11-03) entire document	1-10
A	CN 103902509 A (RONGCHENG DINGTONG ELECTRONIC INFORMATION TECHNOLOGY CO., LTD.) 02 July 2014 (2014-07-02) entire document	1-10
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 22 November 2018		Date of mailing of the international search report 06 December 2018
Name and mailing address of the ISA/CN State Intellectual Property Office of the P. R. China (ISA/CN) No. 6, Xitucheng Road, Jimenqiao Haidian District, Beijing 100088 China Facsimile No. (86-10)62019451		Authorized officer Telephone No.

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.

PCT/CN2018/111077

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)			Publication date (day/month/year)
CN	101086699	A	12 December 2007	CN	100465876	C	04 March 2009
CN	107315574	A	03 November 2017	WO	2017185389	A1	02 November 2017
CN	103902509	A	02 July 2014	None			

国际检索报告		国际申请号 PCT/CN2018/111077
A. 主题的分类 G06F 17/16(2006.01)i 按照国际专利分类(IPC)或者同时按照国家分类和IPC两种分类		
B. 检索领域 检索的最低限度文献(标明分类系统和分类号) G06F17; G06F7 包含在检索领域中的除最低限度文献以外的检索文献 在国际检索时查阅的电子数据库(数据库的名称, 和使用的检索词(如使用)) WPI, EPODOC, CNPAT, CNKI: 矩阵乘法器, 存储器, 计算, MATRIX MULTIPLICATION DEVICE, MEMORY, CALCULATING		
C. 相关文件		
类型*	引用文件, 必要时, 指明相关段落	相关的权利要求
A	CN 101086699 A (浙江大学) 2007年 12月 12日 (2007 - 12 - 12) 全文	1-10
A	CN 107315674 A (北京中科寒武纪科技有限公司) 2017年 11月 3日 (2017 - 11 - 03) 全文	1-10
A	CN 103902509 A (荣成市鼎通电子科技有限公司) 2014年 7月 2日 (2014 - 07 - 02) 全文	1-10
<input type="checkbox"/> 其余文件在C栏的续页中列出。		
<input checked="" type="checkbox"/> 见同族专利附件。		
* 引用文件的具体类型: “A” 认为不特别相关的表示了现有技术一般状态的文件 “E” 在国际申请日的当天或之后公布的在先申请或专利 “L” 可能对优先权要求构成怀疑的文件, 或为确定另一篇引用文件的公布日而引用的或者因其他特殊理由而引用的文件(如具体说明的) “O” 涉及口头公开、使用、展览或其他方式公开的文件 “P” 公布日先于国际申请日但迟于所要求的优先权日的文件 “T” 在申请日或优先权日之后公布, 与申请不相抵触, 但为了理解发明之理论或原理的在后文件 “X” 特别相关的文件, 单独考虑该文件, 认定要求保护的发明不是新颖的或不具有创造性 “Y” 特别相关的文件, 当该文件与另一篇或者多篇该类文件结合并且这种结合对于本领域技术人员为显而易见时, 要求保护的发明不具有创造性 “&” 同族专利的文件		
国际检索实际完成的日期 2018年 11月 22日		国际检索报告邮寄日期 2018年 12月 6日
ISA/CN的名称和邮寄地址 中华人民共和国国家知识产权局(ISA/CN) 中国北京市海淀区蓟门桥西土城路6号 100088 传真号 (86-10)62019451		受权官员 林亮亮 电话号码 86-(010)-62411890

表 PCT/ISA/210 (第2页) (2015年1月)

国际检索报告
关于同族专利的信息

国际申请号
PCT/CN2018/111077

检索报告引用的专利文件			公布日 (年/月/日)	同族专利			公布日 (年/月/日)
CN	101086699	A	2007年 12月 12日	CN	100465876	C	2009年 3月 4日
CN	107315574	A	2017年 11月 3日	WO	2017185389	A1	2017年 11月 2日
CN	103902509	A	2014年 7月 2日	无			

表 PCT/ISA/210 (同族专利附件) (2016年1月)

フロントページの続き

(81) 指定国・地域 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JO, JP, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT

(72) 発明者 リウ、フ

中華人民共和国・518129・グアンドン・シェンツェン・ロンガン・ディストリクト・バンティアン・(番地なし)・ホアウェイ・アドミニストレーション・ビルディング ホアウェイ・テクノロジー・カンパニー・リミテッド内

(72) 発明者 リアオ、ヘン

中華人民共和国・518129・グアンドン・シェンツェン・ロンガン・ディストリクト・バンティアン・(番地なし)・ホアウェイ・アドミニストレーション・ビルディング ホアウェイ・テクノロジー・カンパニー・リミテッド内

(72) 発明者 トウ、ジアジン

中華人民共和国・518129・グアンドン・シェンツェン・ロンガン・ディストリクト・バンティアン・(番地なし)・ホアウェイ・アドミニストレーション・ビルディング ホアウェイ・テクノロジー・カンパニー・リミテッド内

(72) 発明者 ユアン、ホンファイ

中華人民共和国・518129・グアンドン・シェンツェン・ロンガン・ディストリクト・バンティアン・(番地なし)・ホアウェイ・アドミニストレーション・ビルディング ホアウェイ・テクノロジー・カンパニー・リミテッド内

(72) 発明者 ラム、ホウファン

中華人民共和国・518129・グアンドン・シェンツェン・ロンガン・ディストリクト・バンティアン・(番地なし)・ホアウェイ・アドミニストレーション・ビルディング ホアウェイ・テクノロジー・カンパニー・リミテッド内

(72) 発明者 チュー、ファン

中華人民共和国・518129・グアンドン・シェンツェン・ロンガン・ディストリクト・バンティアン・(番地なし)・ホアウェイ・アドミニストレーション・ビルディング ホアウェイ・テクノロジー・カンパニー・リミテッド内

Fターム(参考) 5B056 AA05 BB71