

(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2021-5358

(P2021-5358A)

(43) 公開日 令和3年1月14日(2021.1.14)

(51) Int.Cl.			F I			テーマコード (参考)		
G06F	17/16	(2006.01)	G06F	17/16		K	5B056	
G06N	3/04	(2006.01)	G06N	3/04				
G06N	3/02	(2006.01)	G06N	3/02				
G06N	3/10	(2006.01)	G06N	3/10				
G06N	3/06	(2006.01)	G06N	3/06				

審査請求 有 請求項の数 20 O L 外国語出願 (全 23 頁)

(21) 出願番号 特願2020-2746 (P2020-2746)
 (22) 出願日 令和2年1月10日 (2020.1.10)
 (31) 優先権主張番号 201910559362.7
 (32) 優先日 令和1年6月26日 (2019.6.26)
 (33) 優先権主張国・地域又は機関 中国 (CN)

(71) 出願人 514322098
 ベイジン バイドウ ネットコム サイエ
 ンス アンド テクノロジー カンパニー
 リミテッド
 中華人民共和国 ペキン 100085,
 ハイディアン ディストリクト, シャ
 ンディ テンス ストリート, 10番,
 バイドウ キャンパス 2階
 (74) 代理人 110001508
 特許業務法人 津国

最終頁に続く

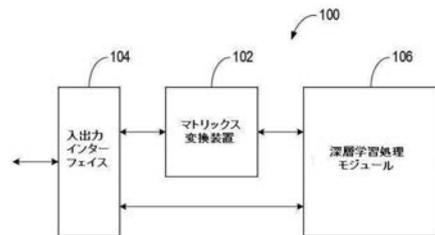
(54) 【発明の名称】 マトリックスを変換するための装置及び方法、データ処理システム

(57) 【要約】 (修正有)

【課題】 マトリックスを変換するための装置、方法及びデータ処理システムを提供する。

【解決手段】 データ処理システム 100 において、マトリックス装置は、マトリックスデータを受信し、マトリックスデータに対して第一の循環シフトを行うことによって、第一のデータを生成するように配置されている第一のシフトユニットと、第一のデータの中の各行のデータを、当該行のデータの中の各データと異なる配列の配列順で書き込むように配置されているキャッシュユニットと、キャッシュユニットから第二のデータを読み取り、第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリックスデータを生成するように配置されている第二のシフトユニットとを含む。

【選択図】 図 1



【特許請求の範囲】

【請求項 1】

マトリックスを変換するための装置であって、

マトリックスデータを受信し、前記マトリックスデータに対して第一の循環シフトを行うことによって、第一のデータを生成するように配置されている第一のシフトユニットと

、
前記第一のデータの中の各行のデータを、前記行のデータの中の各データと異なる配列の配列順で書き込むことによって、前記第一のデータを第二のデータとして記憶するように配置されているキャッシュユニットと、

前記キャッシュユニットから前記第二のデータを読み取り、前記第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリックスデータを生成するように配置されている第二のシフトユニットと、を含むことを特徴とする装置。

10

【請求項 2】

前記キャッシュユニットは、それぞれが複数の記憶アドレスを有する複数の記憶ユニットグループを含み、前記第一のデータの中の各行のデータを、異なる記憶ユニットグループの異なる記憶アドレスにそれぞれ書き込むことによって、前記第一のデータを前記第二のデータとして記憶するように更に配置されていることを特徴とする請求項 1 に記載の装置。

【請求項 3】

前記第二のシフトユニットは、

前記第二のデータの中の、異なる記憶ユニットグループの同じ記憶アドレスに記憶されているデータをそれぞれ読み取ることによって、前記第二のデータの中の対応する行のデータとし、

前記第二のデータの中の各行のデータに対して前記第二の循環シフトを行うことによって、前記変換後のマトリックスデータの中の対応する行のデータを生成するように更に配置されていることを特徴とする請求項 2 に記載の装置。

20

【請求項 4】

前記第一のシフトユニットは、前記マトリックスデータの中の第 i 行のデータを $(i - 1)$ 桁右循環シフトすることによって、前記第一のデータを生成し、 i は、自然数であるように更に配置されており、

前記第二のシフトユニットは、前記第二のデータの中の第 i 行のデータを $(i - 1)$ 桁左循環シフトすることによって、前記変換後のマトリックスデータを生成するように配置されていることを特徴とする請求項 3 に記載の装置。

30

【請求項 5】

前記マトリックスデータで表されるマトリックスは、 n 行及び m 列を含み、 n 及び m は、それぞれ自然数であり、

前記キャッシュユニットは、

前記第一のデータの中の第 1 行のデータの中の m 個の列のデータの中の第 j データを、前記複数の記憶ユニットグループの中の第 j 記憶ユニットグループの第 j 記憶アドレスにそれぞれ書き込み、 j は、1 以上且つ m 以下である自然数であり、

40

前記第一のデータの中の第 i 行のデータの中の m 個の列のデータの中の第 j データを、前記複数の記憶ユニットグループの中の第 j 記憶ユニットグループの第一の記憶アドレス及び第二の記憶アドレスにそれぞれ書き込み、 i は、2 以上且つ n 以下である自然数であり、

j が 1 以上且つ $i - 1$ 以下であるとき、前記第一の記憶アドレスは、第 $m+j-i+1$ 記憶アドレスであり、

j が i 以上且つ m 以下であるとき、前記第二の記憶アドレスは、第 $j-i+1$ 記憶アドレスであるように更に配置されていることを特徴とする請求項 4 に記載の装置。

【請求項 6】

前記第一のシフトユニットは、

前記マトリックスデータで表されるマトリックスを、 p 行を含む第一のマトリックス及

50

び q 行を含む第二のマトリックス、又は、 p 列を含む第一のマトリックス及び q 列を含む第二のマトリックスに分割し、 p 及び q は、それぞれ自然数であり、

前記第一のマトリックスのマトリックスデータに対して前記第一の循環シフトを行うことによって、前記第一のマトリックスの前記第一のデータを生成し、

前記第二のマトリックスのマトリックスデータに対して前記第一の循環シフトを行うことによって、前記第二のマトリックスの前記第一のデータを生成するように更に配置されていることを特徴とする請求項 1 に記載の装置。

【請求項 7】

前記キャッシュユニットは、

前記第一のマトリックスの前記第一のデータを、第一の記憶アドレスをスタートアドレスとして前記キャッシュユニットに書き込むことによって、前記第一のマトリックスの前記第一のデータを前記第一のマトリックスの前記第二のデータとして記憶し、

前記第二のマトリックスの前記第一のデータを、第 $k + 1$ 記憶アドレスをスタートアドレスとして前記キャッシュユニットに書き込むことによって、前記第二のマトリックスの前記第一のデータを前記第二のマトリックスの前記第二のデータとして記憶し、 k は、 p 以上の自然数であるように更に配置されていることを特徴とする請求項 6 に記載の装置。

【請求項 8】

前記第二のシフトユニットは、

前記キャッシュユニットから前記第一のマトリックスの前記第二のデータを読み取り、前記第一のマトリックスの前記第二のデータに対して前記第二の循環シフトを行うことによって、変換後の第一のマトリックスデータを生成し、

前記キャッシュユニットから前記第二のマトリックスの前記第二のデータを読み取り、前記第二のマトリックスの前記第二のデータに対して前記第二の循環シフトを行うことによって、変換後の第二のマトリックスデータを生成し、

前記変換後の第一のマトリックスデータ及び前記変換後の第二のマトリックスデータを結合することによって、前記変換後のマトリックスデータを生成するように更に配置されていることを特徴とする請求項 7 に記載の装置。

【請求項 9】

前記キャッシュユニットは、複数の記憶ユニットグループを含み、前記複数の記憶ユニットグループの中の各 s 個の記憶ユニットグループが 1 グループの記憶ユニットグループに分割され、各グループの記憶ユニットグループは、複数の記憶アドレスを含み、前記複数の記憶アドレスの中の各 t 個の記憶アドレスが 1 グループの記憶アドレスに分割され、

前記キャッシュユニットは、前記第一のデータの中の各 t 行のデータの中の複数のグループのデータを、異なるグループの記憶ユニットグループの異なるグループの記憶アドレスにそれぞれ書き込むことによって、前記第一のデータを前記第二のデータとして記憶し、前記複数のグループのデータの中の各グループのデータは、 $s \times t$ 個のデータを含み、 s 及び t は、自然数であるように配置されていることを特徴とする請求項 1 に記載の装置。

【請求項 10】

前記第二のシフトユニットは、

前記第二のデータの中の、異なるグループの記憶ユニットグループの同じグループの記憶アドレスに記憶されている各グループのデータをそれぞれ読み取り、前記第二のデータの中の相応する行のデータとし、

前記第二のデータの中の各行のデータに対して前記第二の循環シフトを行うことによって、前記変換後のマトリックスデータの中の相応する行のデータを生成するように更に配置されていることを特徴とする請求項 9 に記載の装置。

【請求項 11】

前記第一のシフトユニットは、前記マトリックスデータの中の各列のデータの中の各 s 個のデータを 1 グループのデータに分割し、前記マトリックスデータの中の第 i 行のデータの中の各グループのデータを $(i-1) \times s$ 桁右循環シフトすることによって、前記第一の

10

20

30

40

50

データを生成し、 i 及び s は、自然数であるように更に配置されており、

前記第二のシフトユニットは、前記第二のデータの中の第 i 行のデータの中の各グループのデータを $(i-1) \times s$ 桁左循環シフトするように更に配置されていることを特徴とする請求項10に記載の装置。

【請求項12】

前記変換後のマトリクスデータで表されるマトリクスは、前記マトリクスデータで表されるマトリクスの転置マトリクスであることを特徴とする請求項1～11の何れか一項に記載の装置。

【請求項13】

請求項1～12の何れか一項に記載の装置と、

前記装置に電氣的に結合され、前記マトリクスデータを前記装置に送るように配置されている入出力インターフェイスと、

前記装置に電氣的に結合され、深層学習モデルに基づいて前記変換後のマトリクスデータに対して処理を行う深層学習処理モジュールとを備えることを特徴とするデータ処理システム。

【請求項14】

前記深層学習処理モジュールは、前記処理の結果を他のマトリクスデータとして前記装置に送るように更に配置されており、

前記装置は、前記他のマトリクスデータに基づいて変換後の他のマトリクスデータを生成し、前記変換後の他のマトリクスデータを前記入出力インターフェイスに送るように更に配置されていることを特徴とする請求項13に記載のデータ処理システム。

【請求項15】

前記深層学習処理モジュールは、前記入出力インターフェイスに電氣的に結合され、前記処理の結果を前記入出力インターフェイスに送るように更に配置されていることを特徴とする請求項13に記載のデータ処理システム。

【請求項16】

マトリクスを変換するための方法であって、

マトリクスデータを受信し、前記マトリクスデータに対して第一の循環シフトを行うことによって、第一のデータを生成することと、

前記第一のデータの中の各行のデータを、当該行のデータの中の各データと異なる配列の配列順でキャッシュユニットに書き込むことによって、前記キャッシュユニットにおいて前記第一のデータを第二のデータとして記憶することと、

前記キャッシュユニットから前記第二のデータを読み取り、前記第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリクスデータを生成することと、を含むことを特徴とする方法。

【請求項17】

前記キャッシュユニットは、それぞれが複数の記憶アドレスを有する複数の記憶ユニットグループを含み、

前記第一のデータを前記第二のデータとして記憶することは、前記第一のデータの中の各行のデータを、異なる記憶ユニットグループの異なる記憶アドレスにそれぞれ書き込むことを含むことを特徴とする請求項16に記載の方法。

【請求項18】

前記変換後のマトリクスデータを生成することは、

前記第二のデータの中の、異なる記憶ユニットグループの同じ記憶アドレスに記憶されているデータをそれぞれ読み取ることによって、前記第二のデータの中の相応する行のデータとすることと、

前記第二のデータの中の各行のデータに対して前記第二の循環シフトを行うことによって、前記変換後のマトリクスデータの中の相応する行のデータを生成することを含むことを特徴とする請求項17に記載の方法。

【請求項19】

10

20

30

40

50

前記マトリックスデータの中の第*i*列のデータを(*i* - 1)桁右循環シフトすることによって、前記第一のデータを生成し、*i*は、自然数であり、

前記第二のデータの中の第*i*列のデータを(*i* - 1)桁左循環シフトすることによって、前記変換後のマトリックスデータを生成することを特徴とする請求項18に記載の方法。

【請求項20】

前記マトリックスデータで表されるマトリックスは、*n*行及び*m*列を含み、*n*及び*m*は、それぞれ自然数であり、

前記第一のデータを前記第二のデータとして記憶することは、

前記第一のデータの中の第1行のデータの中の*m*個の列のデータの中の第*j*データを、前記複数の記憶ユニットグループの中の第*j*記憶ユニットグループの第*j*記憶アドレスにそれぞれ書き込み、*j*は、1以上且つ*m*以下である自然数であることと、

10

前記第一のデータの中の第*i*行のデータの中の*m*個の列のデータの中の第*j*データを、前記複数の記憶ユニットグループの中の第*j*記憶ユニットグループの第一の記憶アドレス及び第二の記憶アドレスにそれぞれ書き込み、*i*は、2以上且つ*n*以下である自然数であり、

*j*が1以上且つ*i* - 1以下であるとき、前記第一の記憶アドレスは、第*m*+*j*-*i*+1記憶アドレスであり、*j*が*i*以上且つ*m*以下であるとき、前記第二の記憶アドレスは、第*j*-*i*+1記憶アドレスであることとを含むことを特徴とする請求項19に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

20

本願の実施形態は、主にマトリックスデータ処理の分野に関し、具体的には、マトリックスを変換するための装置、方法、及び当該装置を含むデータ処理システムに関する。

【背景技術】

【0002】

深層学習アクセラレータにおいては、入力データの最も核心的なものとしての特徴、モデルパラメータ、中間結果等のデータは、何れもマトリックス方式で表され、例えば、一次元マトリックス、二次元マトリックス及び多次元マトリックスである。中央処理装置(CPU)は、通常、クロックサイクルごとに処理する単位が1つの数であるが、深層学習アクセラレータの処理能力は、伝統的なCPUの数十倍乃至数百倍である。故に、深層学習アクセラレータは、クロックサイクルごとに処理する単位が1つの数ではなく、1つの多次元マトリックスであり、例えば、入力マトリックスの1つのマトリックスブロック又はサブマトリックスである。

30

【0003】

通常、深層学習モデルは、1つの複雑な計算プロセスであり、より良い性能を得るために、マトリックスのようなデータに対して幾つかのフォーマットの変化及び整理を行う必要がある。マトリックスの変換においては、マトリックス転置は、最も頻繁に使われている操作であり、このような操作では、マトリックス内部のデータを変えないが、マトリックスの次元及びデータの順番を調整することがある。

【0004】

深層学習アクセラレータのハードウェアの設計においては、性能が高く、コストが低いマトリックスの変換案が期待されている。

40

【発明の概要】

【0005】

本願の実施形態は、マトリックスを変換するための装置、方法、及び当該装置を含むデータ処理システムを提供する。これらの装置、方法、及びデータ処理システムは、深層学習アクセラレータにおいて、多次元のマトリックスの変換案を提供することができ、ストリーム処理、簡単な回路、低い消費電力等の特徴を有する。

【0006】

本願の第一の態様は、マトリックスを変換するための装置を提供する。当該装置は、マトリックスデータを受信し、前記マトリックスデータに対して第一の循環シフトを行うこ

50

とによって、第一のデータを生成するように配置されている第一のシフトユニットと、前記第一のデータの中の各行のデータを、前記行のデータの中の各データと異なる配列の配列順で書き込むことによって、前記第一のデータを第二のデータとして記憶するように配置されているキャッシュユニットと、前記キャッシュユニットから前記第二のデータを読み取り、前記第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリックスデータを生成するように配置されている第二のシフトユニットとを含む。

【0007】

本願の第二の態様は、データ処理システムを提供する。当該データ処理システムは、上述した装置と、前記装置に電氣的に結合され、前記マトリックスデータを前記装置に送るように配置されている入出力インターフェイスと、前記装置に電氣的に結合され、深層学習モデルに基づいて前記変換後のマトリックスデータに対して処理を行うように配置されている深層学習処理モジュールとを備える。

10

【0008】

本願の第三の態様は、マトリックスを変換するために用いられる方法を提供する。当該方法は、マトリックスデータを受信し、前記マトリックスデータに対して第一の循環シフトを行うことによって、第一のデータを生成することと、前記第一のデータの中の各行のデータを、当該行のデータの中の各データと異なる配列の配列順でキャッシュユニットに書き込むことによって、前記キャッシュユニットにおいて前記第一のデータを第二のデータとして記憶することと、前記キャッシュユニットから前記第二のデータを読み取り、前記第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリックスデータを生成することを含む。

20

【0009】

発明の内容の部分に記載の内容は、本願の実施形態の肝心又は重要な特徴を制限せず、本願の範囲も制限しない。本願の他の特徴は、以下の説明により理解しやすくなる。

【図面の簡単な説明】

【0010】

図面とともに以下の詳しい説明を参照することによって、本願の各実施形態の上述した及び他の特徴、メリットと態様は、より明らかになる。図面においては、同じ又は類似する図面符号は、同じ又は類似する部材を示す。

【図1】本願の実施形態に係るデータ処理システムを示すブロック図である。

30

【図2A】マトリックスの変換を示す模式図である。

【図2B】マトリックスの変換を示す模式図である。

【図2C】マトリックスの変換を示す模式図である。

【図3】本願の実施形態に係るマトリックス変換装置を示すブロック図である。

【図4】本願の実施形態に係るキャッシュユニットを示す模式図である。

【図5】本願の実施形態に係るマトリックスの変換のプロセスを示す模式図である。

【図6】本願の実施形態に係るマトリックスの変換のプロセスを示す模式図である。

【図7】本願の実施形態に係るマトリックスの変換のプロセスを示す模式図である。

【図8】本願の実施形態に係るマトリックスの変換のプロセスを示す模式図である。

【図9】本願の他の一実施形態に係るマトリックスの変換のプロセスを示す模式図である

40

。【図10】本願の他の一実施形態に係るマトリックスの変換のプロセスを示す模式図である。

【図11】本願の実施形態に係るマトリックスを変換するために用いられる方法を示すフローチャート図である。

【発明を実施するための形態】

【0011】

以下、図面を参照しながら、本願の実施形態をより詳しく説明する。図面は、本願の幾つかの実施形態を示しているが、本願は、様々な態様で実現することができ、ここで説明した実施形態に限定されないことが理解される。逆に、本願をより詳しく及び完全に理解

50

するために、これらの実施形態を提供する。本願の図面及び実施形態は、例示するためのものであり、本願の保護範囲を制限するためのものではない。

【0012】

本願の実施形態の記載においては、用語「含む」及びその類似する用語は、開放的な含みであり、即ち、「含むが、限定されない」と理解すべきである。用語「基づく」は、「少なくとも部分的に基づく」と理解すべきである。用語「一つの実施形態」又は「当該実施形態」は、「少なくとも一つの実施形態」と理解すべきである。用語「第一」、「第二」等は、異なる又は同じ対象を示すことができる。後述は、他の明確な及び隠れた定義を含む可能性がある。

【0013】

上述したように、深層学習アクセラレータの入力は、通常、マトリックスデータであり、深層学習モデルの計算プロセスにおいては、マトリックスデータに対して変換を行う必要がある。一つの伝統的な方案においては、汎用プロセッサ(CPU又はARM)により表示されるソフトウェア前処理又はコプロセッシングによって、マトリックスの次元拡張、転置等の操作を実現する。しかしながら、このような方案のマトリックスの変換性能が比較的悪く、ユーザーのコードが増えたので、不便である。また、もう一つの伝統的な方案においては、レジスタ行列によりマトリックスの転置を実現する。これは、一つの $N \times N$ のマトリックスを縦方向でレジスタ行列に移し入れ、そして横方向でレジスタ行列から移し出すことによって、マトリックスの転置を実現するものである。しかし、このような方案は、大量のレジスタを消費する必要があり、マトリックスの次元拡張をサポートすることができず、三次元以上のマトリックスの転置をサポートすることができず、フレキシビリティが比較的悪い。

【0014】

深層学習アクセラレータのハードウェア設計は、FPGA回路設計又はASICチップ設計等の形式を含む。性能が高く、コストが低いマトリックス変換回路を如何に実現し、二次元のマトリックスの転置、三次元のマトリックスの転置、マトリックスの次元拡張及び転置等の様々なマトリックスの変換を如何に柔軟にサポートするかは、一つの難点である。

【0015】

本願の実施形態は、マトリックスを変換するために用いられる装置を提供する。第一のシフトユニットにより、マトリックスデータを受信し、前記マトリックスデータに対して第一の循環シフトを行うことによって、第一のデータを生成する。キャッシュユニットにより、前記第一のデータの中の各行のデータを、前記行のデータの中の各データと異なる配列の配列順で前記キャッシュユニットに書き込むことによって、前記第一のデータを第二のデータとして記憶する。第二のシフトユニットにより、前記キャッシュユニットから前記第二のデータを読み取り、前記第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリックスデータを生成する。このような方式では、本願の実施形態により深層学習アクセラレータにおいて多次元のマトリックスの変換方案を実現し、二次元のマトリックスの転置、三次元のマトリックスの転置、マトリックスの次元拡張及び転置等の様々なマトリックスの変換をサポートし、ストリーム処理、簡単な回路、低い消費電力等の特徴を有する。

【0016】

以下、図面を参照しながら、様々な実施形態を用いて本願を詳しく説明する。

【0017】

図1は、本願の一つの実施形態によるデータ処理システムを示すブロック図である。図1に示すデータ処理システム100は、深層学習アクセラレータで実現することができる。図1に示すように、データ処理システム100は、マトリックス変換装置102、入出力インターフェイス104及び深層学習処理モジュール106を含んでも良い。

【0018】

入出力インターフェイス104は、データ処理システム100の外部からマトリックス

10

20

30

40

50

データを受信し、データ処理システム 100 の外部へデータ処理システム 100 により処理された結果を送信するように配置されている。幾つかの実施形態においては、入出力インターフェイス 104 は、外部のメモリからマトリクスデータを読み取り、処理結果を示すデータを外部のメモリに書き込むように配置されている。幾つかの実施形態においては、入出力インターフェイス 104 は、マトリクスブロックデータ又はサブマトリクスデータを受信するように配置されている。また、入出力インターフェイス 104 は、受信されたマトリクスデータをマトリクス変換装置 102 へ送信するように更に配置されている。

【0019】

マトリクス変換装置 102 は、入出力インターフェイス 104 に電氣的に結合される。マトリクス変換装置 102 は、マトリクスデータを受信し、マトリクスデータに対して変換を行うことによって、変換後のマトリクスデータを生成するように配置されている。幾つかの実施形態においては、マトリクス変換装置 102 は、マトリクスデータに対してフォーマットの変換を行うように配置されている。幾つかの実施形態においては、マトリクス変換装置 102 は、二次元のマトリクスの転置、三次元のマトリクスの転置、マトリクスの次元拡張及び転置等のような様々なマトリクスの変換を実行するように配置されている。幾つかの実施形態においては、マトリクス変換装置 102 は、計算処理に用いられるように、変換後のマトリクスデータを深層学習処理モジュール 106 に送信するように更に配置されている。マトリクス変換装置 102 の詳しい配置については後に説明する。

【0020】

深層学習処理モジュール 106 は、マトリクス変換装置 102 に電氣的に結合される。深層学習処理モジュール 106 は、深層学習モデルに基づいて変換後のマトリクスデータに対して処理を行うように配置されている。幾つかの実施形態においては、深層学習処理モジュール 106 は、様々なマトリクス計算、ベクトル計算、非線形計算等の処理を実行するように配置されている。幾つかの実施形態においては、深層学習処理モジュール 106 は、当該分野の既知の深層学習モデルに基づく知的財産権 (IP) の核とするように配置されている。

【0021】

幾つかの実施形態においては、深層学習処理モジュール 106 は、処理結果を他のマトリクスデータとしてマトリクス変換装置 102 に送信するように更に配置されている。幾つかの実施形態においては、マトリクス変換装置 102 は、他のマトリクスデータに基づいて変換後の他のマトリクスデータを生成し、変換後の他のマトリクスデータを入出力インターフェイス 104 に送信することによって、変換後の他のマトリクスデータを外部の装置に出力するように更に配置されている。

【0022】

幾つかの実施形態においては、深層学習処理モジュール 106 は、入出力インターフェイス 104 に直接に電氣的に結合され、処理結果を入出力インターフェイス 104 に直接に送信することによって、処理結果を外部の装置に出力するように更に配置されている。幾つかの実施形態においては、深層学習処理モジュール 106 は、入出力インターフェイス 104 からマトリクスデータを直接に受信し、マトリクスデータに対して処理を行うように更に配置されている。

【0023】

幾つかの実施形態においては、データ処理システム 100 は、流れの方式で操作を行い、マトリクス変換装置 102 により目下のマトリクスデータに対して変換を行うプロセスにおいて、入出力インターフェイス 104 が次の 1 つのマトリクスデータを受信することができ、深層学習処理モジュール 106 が前の 1 つの変換後のマトリクスデータに対して処理を行うことができる。

【0024】

図 2A ~ 図 2C は、マトリクス変換装置により実行されるマトリクス変換を示す模

10

20

30

40

50

式図である。通常、二次元のマトリックスは、 $[n,m]$ で表され、 n は、一番目の次元のサイズを表し、 m は、二番目の次元のサイズを表す。また、三次元のマトリックスは、 $[n,m,k]$ で表され、順に類推する。

【0025】

図2Aは、(a)マトリックス $[4,2]$ 及び(b)マトリックス $[2,4]$ を示しており、なお、左側のマトリックス $[4,2]$ は、右側のマトリックス $[2,4]$ に転置される。マトリックスデータの中の(0,0)データは、第一の行の第一の列のデータを表し、(2,1)データは、第三の行の第二の列のデータを表し、このように類推する。図2Aは、マトリックス $[n,m]$ の二次元のマトリックスの転置を示している。

【0026】

図2Bは、(a)マトリックス $[4,3,2]$ 及び(b)マトリックス $[3,4,2]$ を示しており、なお、左側のマトリックス $[4,3,2]$ は、右側のマトリックス $[3,4,2]$ に転置される。図2Bは、マトリックス $[n,m,k]$ の三次元のマトリックスの転置を示している。

【0027】

図2Cは、(a) $[4,4]$ マトリックス及び(b) $[2,4,2]$ マトリックスを示しており、なお、左側の $[4,4]$ マトリックスは、まず、次元拡張されて $[4,2,2]$ になり、その後、次元拡張されたマトリックスは、右側の $[2,4,2]$ マトリックスに転置される。図2Cは、二次元のマトリックスが三次元のマトリックスに拡張され、三次元のマトリックスが更に転置される変換を示している。

【0028】

図3は、本願の実施形態によるマトリックス変換装置を示すブロック図である。図3に示すように、マトリックス変換装置102は、第一のシフトユニット202、キャッシュユニット204及び第二のシフトユニット206を含む。

【0029】

第一のシフトユニット202は、マトリックスデータを受信し、マトリックスデータに対して第一の循環シフトを行うことによって、第一のデータを生成するように配置されている。幾つかの実施形態においては、第一のシフトユニット202は、マトリックスデータの中の各行のデータに対して第一の循環シフトをそれぞれ行うように配置されている。幾つかの実施形態においては、第一のシフトユニット202は、マトリックスデータの第 i 行のデータを $(i-1)$ 桁右循環シフトすることによって、第一のデータを生成するように配置されている。

【0030】

キャッシュユニット204は、第一のシフトユニット202に電氣的に結合され、第一のデータをキャッシュユニット204に書き込むことによって、第一のデータを第二のデータとして記憶するように配置されている。キャッシュユニット204は、第一のデータの中の各行のデータを、前記行のデータの中の各データと異なる配列の配列順でキャッシュユニット204に書き込む。幾つかの実施形態においては、キャッシュユニット204は、記憶コントローラ及び記憶ユニットグループを含み、記憶コントローラは、第一のデータを記憶ユニットグループに書き込むように制御する。幾つかの実施形態においては、キャッシュユニット204は、それぞれが複数の記憶アドレスを有する複数の記憶ユニットグループを含み、キャッシュユニット204は、第一のデータの中の各行のデータを、異なる記憶ユニットグループの異なる記憶アドレスにそれぞれ書き込むことによって、第一のデータを第二のデータとして記憶するように更に配置されている。第一のデータは、自身の元々の行及び列の配列順でキャッシュユニット204に書き込むのではないので、キャッシュユニット204に記憶されている第二のデータは、第一のデータの配列順又はフォーマットの変更後のデータであることを理解すべきである。しかしながら、第一のデータに比べ、第二のデータの中の各データの内容は、変化がない。後に、第一のデータをキャッシュユニットに書き込むことを詳しく説明する。

【0031】

第二のシフトユニット206は、キャッシュユニット204に電氣的に結合され、キャ

10

20

30

40

50

キャッシュユニット 204 から第二のデータを読み取り、第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリクスデータを生成するように配置されている。幾つかの実施形態においては、第二のシフトユニット 206 は、第二のデータの中の、異なる記憶ユニットグループの同じ記憶アドレスに記憶されているデータをそれぞれ読み取ることによって、第二のデータの中の相応する行のデータとするように配置されている。幾つかの実施形態においては、第二のシフトユニット 206 は、第二のデータの中の各行のデータに対して第二の循環シフトを行うことによって、変換後のマトリクスデータの中の相応する行のデータを生成するように配置されている。幾つかの実施形態においては、第二のシフトユニット 206 は、第二のデータの中の第 i 行のデータを $(i-1)$ 桁左循環シフトすることによって、変換後のマトリクスデータを生成するように配置されている。また、第二のシフトユニット 206 は、変換後のマトリクスデータを図 1 の深層学習処理モジュール 106 に送信するように更に配置されている。

10

【0032】

図 4 は、本願の実施形態によるキャッシュユニットを示す模式図である。図 4 に示すように、キャッシュユニット 204 は、複数の記憶ユニットグループ 302 を含む。

【0033】

図 4 に示すように、幾つかの実施形態においては、複数の記憶ユニットグループ 302 は、 x 軸の方向に沿って順に配列し、 x 軸は、第 i 記憶ユニットグループ 302 を表す。また、それぞれの記憶ユニットグループ 302 は、複数の記憶アドレスを含む。説明の便宜上、幾つかの実施形態においては、それぞれの記憶ユニットグループ 302 の複数の記憶アドレスは、 x 軸に直交する y 軸の方向に沿って配列するように示されており、 y 軸は、それぞれの記憶ユニットグループ 302 の第 j 記憶アドレスを示す。幾つかの実施形態においては、記憶アドレスは、深さと呼ばれても良い。幾つかの実施形態においては、記憶アドレス $[i, j]$ は、第 i 記憶ユニットグループ 302 の第 j 記憶アドレスを示す。ここでは、 i 及び j は、それぞれ自然数である。

20

【0034】

幾つかの実施形態においては、記憶ユニットグループ 302 は、スタティックランダムアクセスメモリ (SRAM) グループである。後述においては、記憶ユニットグループ 302 を SRAM グループとする例示を説明する。記憶ユニットグループ 302 は、SRAM グループと限らず、他のタイプの記憶ユニットの集合を採用することもできる。

30

【0035】

図 5 は、本願の実施形態によるマトリクス変換のプロセスを示す模式図である。図 5 に示される左側のデータは、マトリクスデータの中の第一の行のデータ $(0,0)$ 、 $(0,1)$ 、 $(0,2)$ 、 $(0,3)$ 、 $(0,4)$ 、 $(0,5)$ 、 $(0,6)$ 、 $(0,7)$ である。幾つかの実施形態においては、第一のシフトユニット 202 は、マトリクスデータの中から当該第一の行のデータを読み取り、当該第一の行のデータに対して $(1-1=0)$ 桁右循環シフトし、即ち、実質的に当該第一の行のデータに対して右循環シフトしない。キャッシュユニット 204 は、第一の対角線の記憶アドレスで当該第一の行のデータをキャッシュユニット 204 に書き込むように、当該第一の行のデータの書き込みを制御する。幾つかの実施形態においては、図 5 に示すように、キャッシュユニット 204 は、当該第一の行のデータを複数の記憶ユニットグループの記憶アドレス $[1,1]$ 、 $[2,2]$ 、 $[3,3]$ 、 $[4,4]$ 、 $[5,5]$ 、 $[6,6]$ 、 $[7,7]$ 、 $[8,8]$ にそれぞれ書き込む。このように、マトリクスデータの中の第一の行のデータは、自身の元々のフォーマット又は配列順ではなく、第一の対角線の記憶アドレスでキャッシュユニット 204 に記憶されている。

40

【0036】

図 6 は、本願の実施形態によるマトリクスの変換のプロセスを示す模式図である。図 6 に示される左側のデータは、マトリクスデータの中の第二の行のデータである。幾つかの実施形態においては、第一のシフトユニット 202 は、マトリクスデータの中から当該第二の行のデータを読み取り、当該第二の行のデータに対して $(2-1=1)$ 桁右循環シフトすることによって、1 桁右循環シフトした後の第二の行のデータ $(1,7)$ 、 $(1,0)$ 、

50

(1,1)、(1,2)、(1,3)、(1,4)、(1,5)、(1,6)を獲得する。キャッシュユニット204は、右循環シフトした後の第二の行のデータの書き込みを制御することによって、第二の対角線の記憶アドレスで、右循環シフトした後の第二の行のデータをキャッシュユニット204に書き込む。幾つかの実施形態においては、図6に示すように、キャッシュユニット204は、右循環シフトした後の第二の行のデータを、複数の記憶ユニットグループの記憶アドレス[1,8]、[2,1]、[3,2]、[4,3]、[5,4]、[6,5]、[7,6]、[8,7]にそれぞれ書き込む。このように、マトリクスデータの中の右循環シフトした後の第二の行のデータは、自身の元々のフォーマット又は配列順ではなく、第二の対角線の記憶アドレスでキャッシュユニット204に記憶されている。

【0037】

図7は、本願の実施形態によるマトリクスの変換のプロセスを示す模式図である。図7に示される左側のデータは、マトリクスデータの中の第三の行のデータである。幾つかの実施形態においては、第一のシフトユニット202は、マトリクスデータの中から当該第三の行のデータを読み取り、当該第三の行のデータに対して(3-1=2)桁右循環シフトすることによって、2桁右循環シフトした後の第三の行のデータ(2,6)、(2,7)、(2,0)、(2,1)、(2,2)、(2,3)、(2,4)、(2,5)を獲得する。キャッシュユニット204は、右循環シフトした後の第三の行のデータの書き込みを制御することによって、第三の対角線の記憶アドレスで、右循環シフトした後の第三の行のデータをキャッシュユニット204に書き込む。幾つかの実施形態においては、図7に示すように、キャッシュユニット204は、右循環シフトした後の第三の行のデータを、複数の記憶ユニットグループの記憶アドレス[1,7]、[2,8]、[3,1]、[4,2]、[5,3]、[6,4]、[7,5]、[8,6]にそれぞれ書き込む。このように、マトリクスデータの中の右循環シフトした後の第三の行のデータは、自身の元々のフォーマット又は配列順ではなく、第三の対角線の記憶アドレスでキャッシュユニット204に記憶されている。

【0038】

幾つかの実施形態においては、これによって類推し、第一のシフトユニット202は、マトリクスデータの中から第*i*行のデータを読み取り、第*i*行のデータに対して(*i*-1)桁右循環シフトする。キャッシュユニット204は、第*i*の対角線の記憶アドレスで右循環シフトした後の第*i*行のデータをキャッシュユニット204に書き込むように、右循環シフトした後の第*i*行のデータの書き込みを制御する。ここでは、仮に、マトリクスデータで表されるマトリクスは、*n*行及び*m*列を含むとする。なお、*n*及び*m*は、それぞれ自然数である。幾つかの実施形態においては、キャッシュユニット204は、右循環シフトした後の第*i*行のデータの中の*m*個の列のデータの中の第*j*データを、第*j*記憶ユニットグループの第(*m*+*j*-*i*+1)記憶アドレス(*j*は、1以上且つ*i*-1以下である)及び第(*j*-*i*+1)記憶アドレス(*j*は、*i*以上且つ*m*以下である)にそれぞれ書き込み、なお、*i*は、2以上且つ*n*以下である。また、*i*が1に等しい場合、上述したように、第一の行のデータの中の*m*個の列のデータの中の第*j*データを第*j*記憶ユニットグループの第*j*記憶アドレスにそれぞれ書き込む。このように、キャッシュユニット204は、右循環シフトした後の各行のデータの中の各データを、相応する対角線の記憶アドレスで異なる記憶ユニットグループの異なる記憶アドレスにそれぞれ書き込むように制御する。これにより、第一の循環シフトの後のマトリクスデータを、自身の元々のフォーマット又は配列順と異なるフォーマット又は配列順でキャッシュユニットに記憶するので、読み取ることで変換後のマトリクスを生成することができる。

【0039】

図8は、本願の実施形態によるマトリクスの変換のプロセスを示す模式図である。図8の左側の部分は、キャッシュユニット204に記憶されている第二のデータの一部を示している。幾つかの実施形態においては、第二のシフトユニット206は、キャッシュユニット204から異なる記憶ユニットグループの第一の記憶アドレスに記憶されているデータ(0,0)、(1,0)、(2,0)、(3,0)、(4,0)、(5,0)、(6,0)、(7,0)を順に読み取って第二のデータの中の第一の行のデータとする。そして、第二のシフトユニット

10

20

30

40

50

206は、第二のデータの中の第一の行のデータに対して(1-1=0)桁左循環シフトし、即ち、第二のデータの中の第一の行のデータに対して実質的に左循環シフトしない。第二のシフトユニット206は、第二のデータの中の第一の行のデータを変換後のマトリックスデータの中の第一の行のデータとして出力する。

【0040】

幾つかの実施形態においては、第二のシフトユニット206は、キャッシュユニット204から異なる記憶ユニットグループの第二の記憶アドレスに記憶されているデータ(7,1)、(0,1)、(1,1)、(2,1)、(3,1)、(4,1)、(5,1)、(6,1)を順に読み取り、第二のデータの中の第二の行のデータとする。第二のシフトユニット206は、第二のデータの中の第二の行のデータに対して(2-1=1)桁左循環シフトし、左循環シフトした後の第二の行のデータ(0,1)、(1,1)、(2,1)、(3,1)、(4,1)、(5,1)、(6,1)、(7,1)を生成する。第二のシフトユニット206は、第二のデータの中の、左循環シフトした後の第二の行のデータを変換後のマトリックスデータの中の第二の行のデータとして出力する。

10

【0041】

幾つかの実施形態においては、第二のシフトユニット206は、キャッシュユニット204から異なる記憶ユニットグループの第三の記憶アドレスに記憶されているデータ(6,2)、(7,2)、(0,2)、(1,2)、(2,2)、(3,2)、(4,2)、(5,2)を順に読み取り、第二のデータの中の第三の行のデータとする。第二のシフトユニット206は、第二のデータの中の第三の行のデータに対して(3-1=2)桁左循環シフトし、左循環シフトした後の第三の行のデータ(0,2)、(1,2)、(2,2)、(3,2)、(4,2)、(5,2)、(6,2)、(7,2)を生成する。第二のシフトユニット206は、第二のデータの中の、左循環シフトした後の第三の行のデータを変換後のマトリックスデータの中の第三の行のデータとして出力する。

20

【0042】

幾つかの実施形態においては、これによって類推し、第二のシフトユニット206は、キャッシュユニット204から第二のデータの中の、異なるユニットグループの同じ記憶アドレスに記憶されているデータを順に読み取り、第二のデータの中の相応する行のデータとする。そして、第二のシフトユニット206は、第二のデータの中の第*i*行のデータに対して(*i*-1)桁左循環シフトすることによって、変換後のマトリックスデータの中の相応する行のデータを生成する。このように、マトリックス変換装置102は、転置後のマトリックスデータを出力する。

30

【0043】

幾つかの実施形態においては、マトリックス[n, m]のnがm以下である場合、上述した方式でマトリックスデータに対して第一の循環シフトを行うことによって第一のデータを生成し、第一のデータが書き込まれることを制御することによって第二のデータとして記憶し、第二のデータの中の各行のデータを読み取り、各行のデータに対して第二の循環シフトを行うことによって、変換後のマトリックスを生成する。例えば、nが3に等しく且つmが8に等しい場合、マトリックスデータの中の各行のデータは、図7に示すように、右循環シフトされ、書き込まれる。幾つかの実施形態においては、第二のシフトユニット206は、異なる記憶ユニットグループの、第一の記憶アドレスに記憶されているデータ(0,0)、(1,0)、(2,0)、(dummy)、(dummy)、(dummy)、(dummy)、(dummy)、第二の記憶アドレスに記憶されているデータ(dummy)、(0,1)、(1,1)、(2,1)、(dummy)、(dummy)、(dummy)、(dummy)、・・・第八の記憶アドレスに記憶されているデータ(1,7)、(2,7)、(dummy)、(dummy)、(dummy)、(dummy)、(dummy)、(0,7)をそれぞれ読み取り、第二のデータの中の相応する行のデータとし、なお、dummyは、ダミーデータであり、マトリックス変換装置102から出力される時に省略される。第二のシフトユニット206は、第二のデータの中の第*i*行のデータをそれぞれ(*i*-1)桁左循環シフトすることによって、第一の行のデータ(0,0)、(1,0)、(2,0)、第二の行のデータ(0,1)、(1,1)、(2,1)・・・及び第八の行のデータ(0,7)、(

40

50

1,7)、(2,7)を含む変換後のマトリクスデータを出力する。このように、マトリクス変換装置102は、マトリクス[3,8]の転置マトリクス[8,3]を出力する。

【0044】

幾つかの実施形態においては、 n が m より大きい場合、第一のシフトユニット202は、マトリクス $[n,m]$ をマトリクス $[m,m]$ とマトリクス $[n-m, m]$ に分割し、上述した方式でマトリクス $[m,m]$ に対して変換を行い、変換後の第一のマトリクスを出力する。マトリクス $[n-m, m]$ に対し、第一の循環シフトの後のマトリクスデータを、複数の記憶ユニットグループの第 $(m+1)$ 記憶アドレスからキャッシュユニット204に書き込み、上述した方式で、マトリクス $[n-m, m]$ に対して変換を行い、なお、第二のシフトユニットは、第 $(m+1)$ 記憶アドレスからデータを読み取り、読み取られたデータに対して第二の循環シフトを行った後に変換後の第二のマトリクスを出力する。そして、第一のマトリクス及び第二のマトリクスを組み合わせることで転置後のマトリクスを生成する。

10

【0045】

幾つかの実施形態においては、 m がキャッシュユニット204の総データの幅 k より大きい場合、第一のシフトユニット202は、マトリクス $[n, m]$ をマトリクス $[n,k]$ 及びマトリクス $[n,m-k]$ に分割し、上述した方式で、マトリクス $[n,k]$ に対して変換を行い、変換後の第一のマトリクスを出力する。マトリクス $[n,m-k]$ に対し、第一の循環シフトの後のマトリクスデータを、記憶ユニットグループの第 $(k+1)$ 記憶アドレスからキャッシュユニット204に書き込み、上述した方式で、マトリクス $[n,m-k]$ に対して変換を行い、なお、第二のシフトユニットは、第 $(k+1)$ 記憶アドレスから読み取り、読み取られたデータに対して第二の循環シフトを行った後に変換後の第二のマトリクスを出力する。そして、第一のマトリクス及び第二のマトリクスを組み合わせることで転置後のマトリクスを生成する。

20

【0046】

幾つかの実施形態においては、上述した、対角線の方式で書き込んで、行方向でデータを読み取る方式と逆に、まず、第一のデータの中の各行のデータを異なる記憶ユニットグループの同じアドレスに記憶し、その後、対応する対角線の方式でキャッシュユニットに記憶されているデータを読み取り、変換後のマトリクスを生成しても良い。

【0047】

幾つかの実施形態においては、第一のデータの中の各行のデータの、キャッシュユニットにおける書き込みアドレスを好ましく設定し、及び/又は読み取り方式を変えることによって、次元の変換、データ位置の変換、データの入れ替え等を含む望ましいマトリクスの変換を行うことができる。

30

【0048】

図9及び図10は、本願のもう一つの実施形態によるマトリクス変換のプロセスを示す模式図である。図9及び図10は、マトリクス $[n, m]$ が $[n, p,q]$ に拡張され、その後、マトリクス $[n, p,q]$ に対して変換を行うプロセスの例示を示している。

【0049】

幾つかの実施形態においては、複数の記憶ユニットグループ302の中の各 s 個の記憶ユニットグループは、1グループの記憶ユニットグループに分けられ、各グループの記憶ユニットグループは、複数の記憶アドレスを含む。幾つかの実施形態においては、複数の記憶アドレスの中の各 t 個の記憶アドレスが1グループの記憶アドレスに分けられる。幾つかの実施形態においては、キャッシュユニット204は、第一のデータの中の各 t 行のデータの中の複数のグループのデータを異なるグループの記憶ユニットグループの異なるグループの記憶アドレスにそれぞれ書き込み、第一のデータを第二のデータとして記憶するように更に配置されており、なお、複数のグループのデータの中の各グループのデータは、 $s \times t$ 個のデータを含み、なお、 s 及び t は、自然数である。

40

【0050】

幾つかの実施形態においては、第一のシフトユニット202は、マトリクスデータの中の各行のデータの中の各 s 個のデータを1グループのデータに分け、マトリクスデー

50

タの中の第*i*行のデータの中の各グループのデータを $(i-1) \times s$ 桁右循環シフトし、第一のデータを生成するように更に配置されている。

【 0 0 5 1 】

幾つかの実施形態においては、第二のシフトユニット 2 0 6 は、第二のデータの中の、異なるグループの記憶ユニットグループの同じグループの記憶アドレスに記憶されている各グループのデータをそれぞれ順に取り、第二のデータの中の相応する行のデータとするように更に配置されている。幾つかの実施形態においては、第二のシフトユニット 2 0 6 は、第二のデータの中の各行のデータに対して第二の循環シフトを行うことによって、変換後のマトリクスデータの中の相応する行のデータを生成する。幾つかの実施形態においては、第二のシフトユニット 2 0 6 は、第二のデータの中の第*i*行のデータの中の各グループのデータを $(i-1) \times s$ 桁左循環シフトするよう、即ち、 $(i-1)$ グループであるように更に配置されている。

10

【 0 0 5 2 】

図 9 及び図 1 0 に示される例示においては、入力マトリクスは、例えば、二次元のマトリクス[4,8]であり、マトリクス変換装置 1 0 2 は、まず、二次元のマトリクス[4,8]を三次元のマトリクス[4,4,2]に拡張する。幾つかの実施形態においては、各行のデータの中のそれぞれの二つのデータが 1 グループのデータに分けられ、それぞれの二つの記憶ユニットグループが 1 グループの記憶ユニットグループに分けられ、異なるグループの記憶ユニットグループの同じ記憶アドレスが 1 グループの記憶アドレスに分けられる。

20

【 0 0 5 3 】

図 9 に示すように、第一のシフトユニット 2 0 2 は、第一の行のデータの中の 4 グループのデータを $(1-1=0) \times 2$ 桁右循環シフトし、即ち、実質的に右循環シフトを行わない。図 9 に示すように、キャッシュユニット 2 0 4 は、第一の行のデータの中の 4 グループのデータ{(0,0)、(0,1)}、{(0,2)、(0,3)}、{(0,4)、(0,5)}、{(0,6)、(0,7)}を、第一及び第二の記憶ユニットグループの第一の記憶アドレス、第三及び第四の記憶ユニットグループの第二の記憶アドレス、第五及び第六の記憶ユニットグループの第三の記憶アドレス、第七及び第八の記憶ユニットグループの第四の記憶アドレスにそれぞれ書き込む。

30

【 0 0 5 4 】

図 1 0 に示すように、第一のシフトユニット 2 0 2 は、第二の行のデータの中の 4 グループのデータ{(1,0)、(1,1)}、{(1,2)、(1,3)}、{(1,4)、(1,5)}、{(1,6)、(1,7)}を $(2-1=1) \times 2$ 桁右循環シフトし、即ち、2 桁又は 1 つのグループ右循環シフトし、右循環シフトした第二の行のデータ{(1,6)、(1,7)}、{(1,0)、(1,1)}、{(1,2)、(1,3)}、{(1,4)、(1,5)}を生成する。図 1 0 に示すように、キャッシュユニット 2 0 4 は、当該右循環シフトした第二の行のデータを第一及び第二の記憶ユニットグループの第四の記憶アドレス、第三及び第四の記憶ユニットグループの第一の記憶アドレス、第五及び第六記憶ユニットグループの第二の記憶アドレス、第七及び第八の記憶ユニットグループの第四の記憶アドレスにそれぞれ記憶する。

40

【 0 0 5 5 】

これによって類推し、図 5 ~ 図 8 を参照しながら説明したマトリクス変換方式と類似的にマトリクス[4,4,2]に対して第一の循環シフト、書き込み、読み取り及び第二の循環シフトを行うことによって、マトリクス[4,8]を拡張して変換し、変換後のマトリクス[4,4,2]を生成する。

【 0 0 5 6 】

本願の実施形態により、深層学習アクセラレータにおいては、深層学習処理モジュールに対してマトリクスを変換するために用いられる装置を実現し、当該装置は、チップの外から入力された、マトリクスブロック又はサブマトリクス等のようなマトリクスデータに対して変換を行い、変換後のマトリクスデータを深層学習処理モジュールに送信し、計算処理に用いられることができる。本願の実施形態の装置は、多次元のマトリッ

50

クスの変換方案を実現することができ、二次元のマトリックス転置、三次元のマトリックス転置、マトリックスの次元拡張及び転置等の様々なマトリックスの変換をサポートし、ストリーム処理、簡単な回路、低い消費電力等の特徴を有する。

【0057】

図11は、本願の実施形態によるマトリックスを変換するために用いられる方法を示すフローチャート図である。図11に示すように、マトリックスを変換するために用いられる方法400は、ブロック402～ブロック406を含む。

【0058】

ブロック402においては、マトリックスデータを受信し、マトリックスデータに対して第一の循環シフトを行うことによって、第一のデータを生成する。幾つかの実施形態においては、マトリックスデータの中の第*i*行のデータを(*i*-1)桁右循環シフトすることによって、第一のデータを生成し、*i*は、自然数である。

【0059】

ブロック404においては、第一のデータの中の各行のデータを、当該行のデータの中の各データの配列順と異なる配列順でキャッシュユニットに書き込むことによって、キャッシュユニットにおいて第一のデータを第二のデータとして記憶する。幾つかの実施形態においては、キャッシュユニットは、それぞれが複数の記憶アドレスを含む複数の記憶ユニットグループを含む。幾つかの実施形態においては、第一のデータの中の各行のデータを異なる記憶ユニットグループの異なる記憶アドレスにそれぞれ書き込む。

【0060】

幾つかの実施形態においては、マトリックスデータで表されるマトリックスは、*n*行及び*m*列を含み、なお、*n*及び*m*は、それぞれ自然数である。幾つかの実施形態においては、第一のデータの中の第1行のデータの中の*m*個の列のデータの中の第*j*データを複数の記憶ユニットグループの中の第*j*記憶ユニットグループの第*j*記憶アドレスにそれぞれ書き込み、なお、*j*は、1以上且つ*m*以下の自然数である。また、第一のデータの中の第*i*行のデータの中の*m*個の列のデータの中の第*j*データを、複数の記憶ユニットグループの中の第*j*記憶ユニットグループの第(*m*+*j*-*i*+1)記憶アドレス(*j*は、1以上且つ*i*-1以下である)及び第(*j*-*i*+1)記憶アドレス(*j*は、*i*以上且つ*m*以下である)にそれぞれ書き込み、なお、*i*は、2以上且つ*n*以下の自然数である。

【0061】

ブロック406においては、キャッシュユニットから第二のデータを読み取り、第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリックスデータを生成する。幾つかの実施形態においては、第二のデータの中の、異なる記憶ユニットグループの同じ記憶アドレスに記憶されているデータをそれぞれ読み取り、第二のデータの中の相応する行のデータとする。幾つかの実施形態においては、第二のデータの中の各行のデータに対して第二の循環シフトを行うことによって、変換後のマトリックスデータの中の相応する行のデータを生成する。幾つかの実施形態においては、第二のデータの中の第*i*行のデータを(*i*-1)桁左循環シフトすることによって、変換後のマトリックスデータを生成する。

【0062】

本願の実施形態を実施するための方法は、システムオンチップ(SoC)アーキテクチャに適用するために、1つ又は複数のプログラミング言語の任意の組み合わせを用いてプログラミングすることができるかと理解すべきである。また、特定の順序でそれぞれの操作を説明したが、理想な結果が得られるために、このような操作は、示された特定の順序又は順に実行することが求められ、又は、全ての図示の操作は、実行する必要があると理解すべきである。一定の環境においては、マルチタスク及びパラレル処理が有利である可能性がある。

【0063】

上述した記載は、若干の具体的な実現の細部を含むが、これらは、本願の範囲に対する制限であると解釈されるべきではない。1つの実施形態の上下の文脈に記載の幾つかの特

10

20

30

40

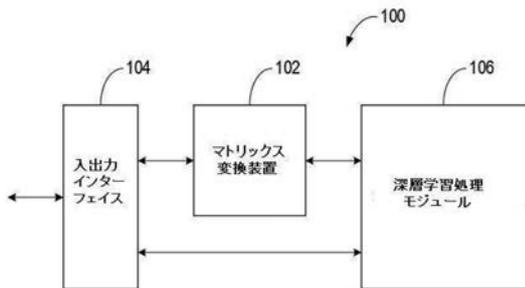
50

徴は、1つの実現において組合わせて実現されることもできる。逆に、1つの実現の上下の文脈に記載の様々な特徴は、複数の実現において単独又は任意の適切なサブ組み合わせの方式で実現されることもできる。

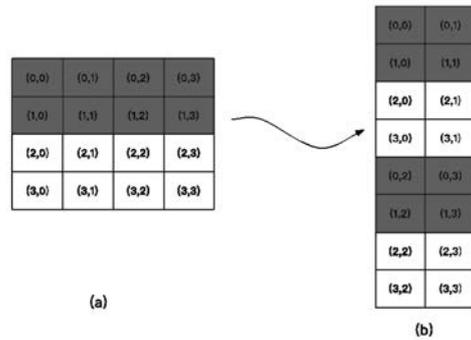
【0064】

構造の特徴及び/又は方法のロジック動作に特定した言語を用いて本願の主題を記載したが、特許請求の範囲に限定される主題は、上述した記載の特定の特徴又は動作に限らないと理解すべきである。逆に、上述した記載の特定の特徴及び動作は、特許請求の範囲を実現する例示の形式に過ぎない。

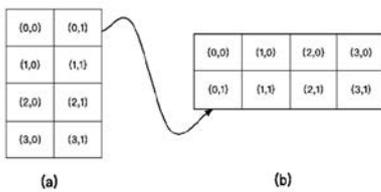
【図1】



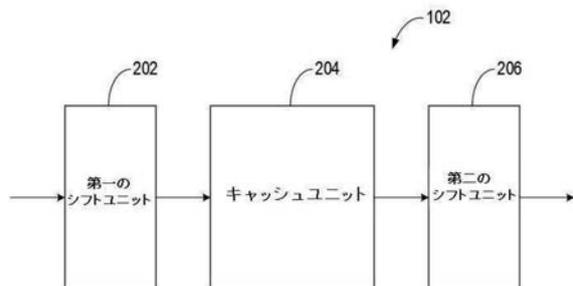
【図2C】



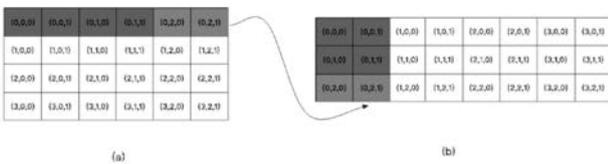
【図2A】



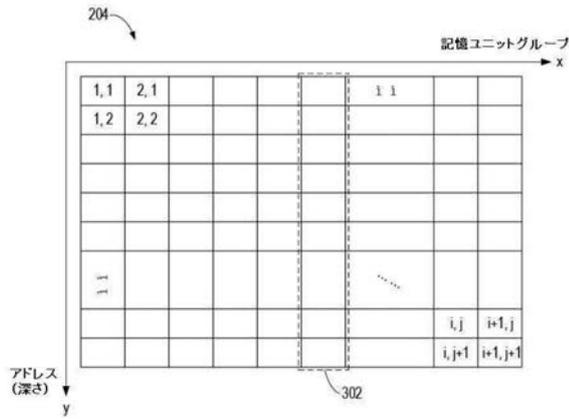
【図3】



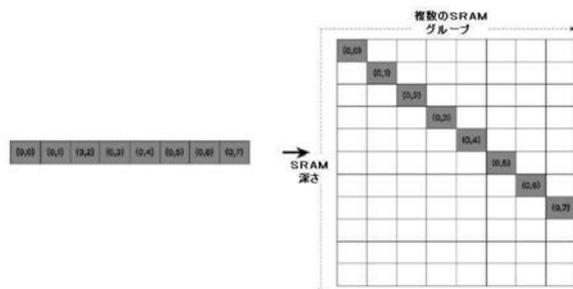
【図2B】



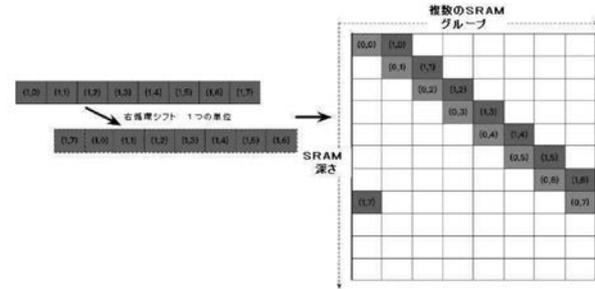
【 図 4 】



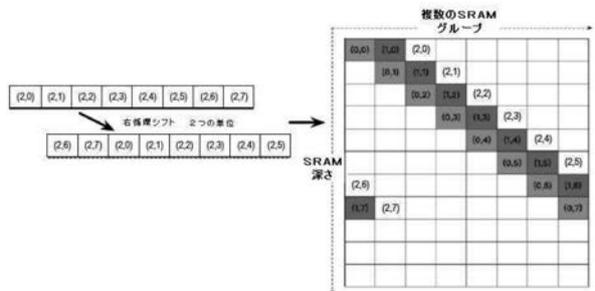
【 図 5 】



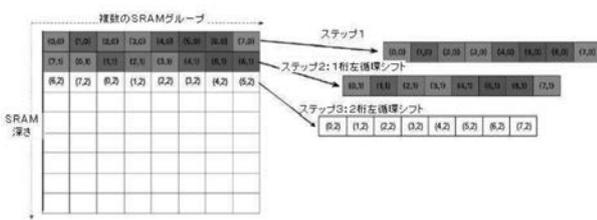
【 図 6 】



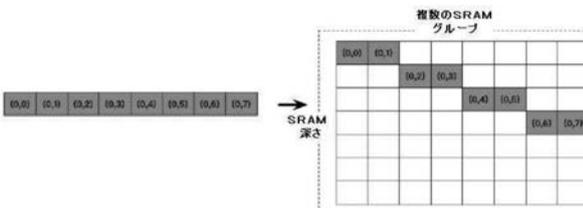
【 図 7 】



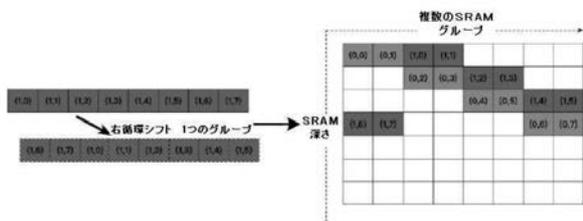
【 図 8 】



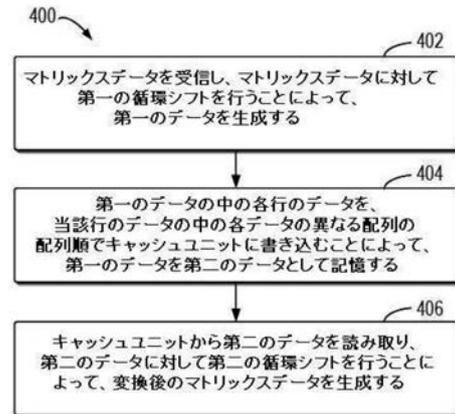
【 図 9 】



【 図 10 】



【 図 11 】



【手続補正書】

【提出日】令和2年8月27日(2020.8.27)

【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】

マトリクスを変換するための装置であって、

マトリクスデータを受信し、前記マトリクスデータに対して第一の循環シフトを行うことによって、第一のデータを生成するように配置されている第一のシフトユニットと

、
前記第一のデータの中の各行のデータを、前記行のデータの中の各データと異なる配列の配列順で書き込むことによって、前記第一のデータを第二のデータとして記憶するように配置されているキャッシュユニットと、

前記キャッシュユニットから前記第二のデータを読み取り、前記第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリクスデータを生成するように配置されている第二のシフトユニットと、を含むことを特徴とする装置。

【請求項2】

前記キャッシュユニットは、それぞれが複数の記憶アドレスを有する複数の記憶ユニットグループを含み、前記第一のデータの中の各行のデータを、異なる記憶ユニットグループの異なる記憶アドレスにそれぞれ書き込むことによって、前記第一のデータを前記第二のデータとして記憶するように更に配置されていることを特徴とする請求項1に記載の装置。

【請求項3】

前記第二のシフトユニットは、

前記第二のデータの中の、異なる記憶ユニットグループの同じ記憶アドレスに記憶されているデータをそれぞれ読み取ることによって、前記第二のデータの中の相応する行のデータとし、

前記第二のデータの中の各行のデータに対して前記第二の循環シフトを行うことによって、前記変換後のマトリクスデータの中の相応する行のデータを生成するように更に配置されていることを特徴とする請求項2に記載の装置。

【請求項4】

前記第一のシフトユニットは、前記マトリクスデータの中の第*i*行のデータを(*i* - 1)桁右循環シフトすることによって、前記第一のデータを生成し、*i*は、自然数であるように更に配置されており、

前記第二のシフトユニットは、前記第二のデータの中の第*i*行のデータを(*i* - 1)桁左循環シフトすることによって、前記変換後のマトリクスデータを生成するように配置されていることを特徴とする請求項3に記載の装置。

【請求項5】

前記マトリクスデータで表されるマトリクスは、*n*行及び*m*列を含み、*n*及び*m*は、それぞれ自然数であり、

前記キャッシュユニットは、

前記第一のデータの中の第1行のデータの中の*m*個の列のデータの中の第*j*データを、前記複数の記憶ユニットグループの中の第*j*記憶ユニットグループの第*j*記憶アドレスにそれぞれ書き込み、*j*は、1以上且つ*m*以下である自然数であり、

前記第一のデータの中の第*i*行のデータの中の*m*個の列のデータの中の第*j*データを、前記複数の記憶ユニットグループの中の第*j*記憶ユニットグループの第一の記憶アドレス及び第二の記憶アドレスにそれぞれ書き込み、*i*は、2以上且つ*n*以下である自然数であり、

j が 1 以上且つ $i - 1$ 以下であるとき、前記第一の記憶アドレスは、第 $m+j-i+1$ 記憶アドレスであり、

j が i 以上且つ m 以下であるとき、前記第二の記憶アドレスは、第 $j-i+1$ 記憶アドレスであるように更に配置されていることを特徴とする請求項 4 に記載の装置。

【請求項 6】

前記第一のシフトユニットは、

前記マトリクスデータで表されるマトリクスを、 p 行を含む第一のマトリクス及び q 行を含む第二のマトリクス、又は、 p 列を含む第一のマトリクス及び q 列を含む第二のマトリクスに分割し、 p 及び q は、それぞれ自然数であり、

前記第一のマトリクスのマトリクスデータに対して前記第一の循環シフトを行うことによって、前記第一のマトリクスの前記第一のデータを生成し、

前記第二のマトリクスのマトリクスデータに対して前記第一の循環シフトを行うことによって、前記第二のマトリクスの前記第一のデータを生成するように更に配置されていることを特徴とする請求項 1 に記載の装置。

【請求項 7】

前記キャッシュユニットは、

前記第一のマトリクスの前記第一のデータを、第一の記憶アドレスをスタートアドレスとして前記キャッシュユニットに書き込むことによって、前記第一のマトリクスの前記第一のデータを前記第一のマトリクスの前記第二のデータとして記憶し、

前記第二のマトリクスの前記第一のデータを、第 $k+1$ 記憶アドレスをスタートアドレスとして前記キャッシュユニットに書き込むことによって、前記第二のマトリクスの前記第一のデータを前記第二のマトリクスの前記第二のデータとして記憶し、 k は、 p 以上の自然数であるように更に配置されていることを特徴とする請求項 6 に記載の装置。

【請求項 8】

前記第二のシフトユニットは、

前記キャッシュユニットから前記第一のマトリクスの前記第二のデータを読み取り、前記第一のマトリクスの前記第二のデータに対して前記第二の循環シフトを行うことによって、変換後の第一のマトリクスデータを生成し、

前記キャッシュユニットから前記第二のマトリクスの前記第二のデータを読み取り、前記第二のマトリクスの前記第二のデータに対して前記第二の循環シフトを行うことによって、変換後の第二のマトリクスデータを生成し、

前記変換後の第一のマトリクスデータ及び前記変換後の第二のマトリクスデータを結合することによって、前記変換後のマトリクスデータを生成するように更に配置されていることを特徴とする請求項 7 に記載の装置。

【請求項 9】

前記キャッシュユニットは、複数の記憶ユニットグループを含み、前記複数の記憶ユニットグループの中の各 s 個の記憶ユニットグループが 1 グループの記憶ユニットグループに分割され、各グループの記憶ユニットグループは、複数の記憶アドレスを含み、前記複数の記憶アドレスの中の各 t 個の記憶アドレスが 1 グループの記憶アドレスに分割され、

前記キャッシュユニットは、前記第一のデータの中の各 t 行のデータの中の複数のグループのデータを、異なるグループの記憶ユニットグループの異なるグループの記憶アドレスにそれぞれ書き込むことによって、前記第一のデータを前記第二のデータとして記憶し、前記複数のグループのデータの中の各グループのデータは、 $s \times t$ 個のデータを含み、 s 及び t は、自然数であるように配置されていることを特徴とする請求項 1 に記載の装置。

【請求項 10】

前記第二のシフトユニットは、

前記第二のデータの中の、異なるグループの記憶ユニットグループの同じグループの記憶アドレスに記憶されている各グループのデータをそれぞれ読み取り、前記第二のデータの中の相応する行のデータとし、

前記第二のデータの中の各行のデータに対して前記第二の循環シフトを行うことによって、前記変換後のマトリクスデータの中の相応する行のデータを生成するように更に配置されていることを特徴とする請求項 9 に記載の装置。

【請求項 11】

前記第一のシフトユニットは、前記マトリクスデータの中の各列のデータの中の各 s 個のデータを 1 グループのデータに分割し、前記マトリクスデータの中の第 i 行のデータの中の各グループのデータを $(i-1) \times s$ 桁右循環シフトすることによって、前記第一のデータを生成し、 i 及び s は、自然数であるように更に配置されており、

前記第二のシフトユニットは、前記第二のデータの中の第 i 行のデータの中の各グループのデータを $(i-1) \times s$ 桁左循環シフトするように更に配置されていることを特徴とする請求項 10 に記載の装置。

【請求項 12】

前記変換後のマトリクスデータで表されるマトリクスは、前記マトリクスデータで表されるマトリクスの転置マトリクスであることを特徴とする請求項 1 ~ 11 の何れか一項に記載の装置。

【請求項 13】

請求項 1 ~ 12 の何れか一項に記載の装置と、

前記装置に電氣的に結合され、前記マトリクスデータを前記装置に送るように配置されている入出力インターフェイスと、

前記装置に電氣的に結合され、深層学習モデルに基づいて前記変換後のマトリクスデータに対して処理を行う深層学習処理モジュールとを備えることを特徴とするデータ処理システム。

【請求項 14】

前記深層学習処理モジュールは、前記処理の結果を他のマトリクスデータとして前記装置に送るように更に配置されており、

前記装置は、前記他のマトリクスデータに基づいて変換後の他のマトリクスデータを生成し、前記変換後の他のマトリクスデータを前記入出力インターフェイスに送るように更に配置されていることを特徴とする請求項 13 に記載のデータ処理システム。

【請求項 15】

前記深層学習処理モジュールは、前記入出力インターフェイスに電氣的に結合され、前記処理の結果を前記入出力インターフェイスに送るように更に配置されていることを特徴とする請求項 13 に記載のデータ処理システム。

【請求項 16】

マトリクスを変換するための方法であって、

マトリクスデータを受信し、前記マトリクスデータに対して第一の循環シフトを行うことによって、第一のデータを生成することと、

前記第一のデータの中の各行のデータを、当該行のデータの中の各データと異なる配列の配列順でキャッシュユニットに書き込むことによって、前記キャッシュユニットにおいて前記第一のデータを第二のデータとして記憶することと、

前記キャッシュユニットから前記第二のデータを読み取り、前記第二のデータに対して第二の循環シフトを行うことによって、変換後のマトリクスデータを生成することと、を含むことを特徴とする方法。

【請求項 17】

前記キャッシュユニットは、それぞれが複数の記憶アドレスを有する複数の記憶ユニットグループを含み、

前記第一のデータを前記第二のデータとして記憶することは、前記第一のデータの中の各行のデータを、異なる記憶ユニットグループの異なる記憶アドレスにそれぞれ書き込むことを含むことを特徴とする請求項 16 に記載の方法。

【請求項 18】

前記変換後のマトリクスデータを生成することは、

前記第二のデータの中の、異なる記憶ユニットグループの同じ記憶アドレスに記憶されているデータをそれぞれ読み取ることによって、前記第二のデータの中の相応する行のデータとすることと、

前記第二のデータの中の各行のデータに対して前記第二の循環シフトを行うことによって、前記変換後のマトリクスデータの中の相応する行のデータを生成することを含むことを特徴とする請求項 17 に記載の方法。

【請求項 19】

前記マトリクスデータの中の第 i 列のデータを $(i - 1)$ 桁右循環シフトすることによって、前記第一のデータを生成し、 i は、自然数であり、

前記第二のデータの中の第 i 列のデータを $(i - 1)$ 桁左循環シフトすることによって、前記変換後のマトリクスデータを生成することを特徴とする請求項 18 に記載の方法。

【請求項 20】

前記マトリクスデータで表されるマトリクスは、 n 行及び m 列を含み、 n 及び m は、それぞれ自然数であり、

前記第一のデータを前記第二のデータとして記憶することは、

前記第一のデータの中の第 1 行のデータの中の m 個の列のデータの中の第 j データを、前記複数の記憶ユニットグループの中の第 j 記憶ユニットグループの第 j 記憶アドレスにそれぞれ書き込み、 j は、1 以上且つ m 以下である自然数であることと、

前記第一のデータの中の第 i 行のデータの中の m 個の列のデータの中の第 j データを、前記複数の記憶ユニットグループの中の第 j 記憶ユニットグループの第一の記憶アドレス及び第二の記憶アドレスにそれぞれ書き込み、 i は、2 以上且つ n 以下である自然数であり、

j が 1 以上且つ $i - 1$ 以下であるとき、前記第一の記憶アドレスは、第 $m + j - i + 1$ 記憶アドレスであり、 j が i 以上且つ m 以下であるとき、前記第二の記憶アドレスは、第 $j - i + 1$ 記憶アドレスであることとを含むことを特徴とする請求項 19 に記載の方法。

【請求項 21】

コンピュータープログラムであって、

プロセッサにより実行されると、請求項 16 ~ 20 の何れか一項に記載の方法を実行させるコンピュータープログラム。

フロントページの続き

(72)発明者 シャオチャン・ゴン

中華人民共和国、 Beijing 100085、ハイディアンのディストリクト、シャンディ・10
ティーエイチ・ストリート、ナンバー10、バドウ・キャンパス 2/エフ

Fターム(参考) 5B056 BB42

【外国語明細書】

2021005358000001.pdf