

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7561925号
(P7561925)

(45)発行日 令和6年10月4日(2024.10.4)

(24)登録日 令和6年9月26日(2024.9.26)

(51)国際特許分類		F I	
G 0 6 F	9/38 (2018.01)	G 0 6 F	9/38 3 7 0 Z
G 0 6 F	9/30 (2018.01)	G 0 6 F	9/30 3 5 0 F
G 0 6 F	15/173 (2006.01)	G 0 6 F	15/173 6 8 3 B
G 0 6 F	15/80 (2006.01)	G 0 6 F	15/80
G 0 6 F	17/16 (2006.01)	G 0 6 F	17/16 S
請求項の数 20 外国語出願 (全18頁) 最終頁に続く			
(21)出願番号	特願2023-114361(P2023-114361)	(73)特許権者	502208397
(22)出願日	令和5年7月12日(2023.7.12)		グーグル エルエルシー
(62)分割の表示	特願2021-142529(P2021-142529))の分割		G o o g l e L L C
原出願日	平成30年5月17日(2018.5.17)		アメリカ合衆国 カリフォルニア州 9 4
(65)公開番号	特開2023-145517(P2023-145517 A)		0 4 3 マウンテン ビュー アンフィシ
(43)公開日	令和5年10月11日(2023.10.11)		アター パークウェイ 1 6 0 0
審査請求日	令和5年8月7日(2023.8.7)		1 6 0 0 Amphitheatre P
(31)優先権主張番号	62/507,771	(74)代理人	arkway 9 4 0 4 3 Mounta
(32)優先日	平成29年5月17日(2017.5.17)		in View, CA U.S.A.
(33)優先権主張国・地域又は機関	米国(US)	(72)発明者	110001195
			弁理士法人深見特許事務所
			ノリー, トーマス
			アメリカ合衆国、9 4 0 4 3 カリフォ
			ルニア州、マウンテン・ビュー、アンフ
			ィシアター・パークウェイ、1 6 0 0
			最終頁に続く

(54)【発明の名称】 専用ニューラルネットワークトレーニングチップ

(57)【特許請求の範囲】

【請求項 1】

専用ハードウェアチップを使用してニューラルネットワークをトレーニングする方法であって、

前記専用ハードウェアチップのベクトルプロセッサが、活性化入力の2次元行列を表すデータを受信することを含み、前記活性化入力の2次元行列は、ニューラルネットワークのそれぞれのネットワーク層に対する入力活性化行列の活性化入力の一部を含み、前記ベクトルプロセッサは、

2次元構造に配置された複数のベクトル処理ユニットと、

前記活性化入力の2次元行列を格納するよう構成されるベクトルレジスタとを含み、前記方法はさらに、

行列乗算ユニットが、前記2次元行列について乗算結果を生成することを含み、前記生成することは、

複数のクロックサイクルの各クロックサイクルにおいて、前記ベクトルプロセッサの前記ベクトルレジスタから前記専用ハードウェアチップの前記行列乗算ユニットに、前記活性化入力の2次元行列における活性化入力のそれぞれの行の少なくとも一部をロードすることと、

前記複数のクロックサイクルの終了前に、前記行列乗算ユニットに重み値をロードすることとを含み、前記重み値は、前記活性化入力の2次元行列における前記活性化入力に対応する形状に配置され、前記生成することはさらに、

前記行列乗算ユニットが、前記活性化入力の 2 次元行列および重み値を乗算して、前記乗算結果を得ることを含み、前記方法はさらに、

前記乗算結果に基づいて、逆伝播を通じて前記ニューラルネットワークの重み値を更新することを含む、方法。

【請求項 2】

前記行列乗算ユニットに、前記活性化入力の 2 次元行列における前記活性化入力のそれぞれの行の少なくとも前記一部をロードすることは、さらに、前記活性化入力のそれぞれの行の少なくとも前記一部を、第 1 の浮動小数点フォーマットから、前記第 1 の浮動小数点フォーマットよりも小さいビットサイズを有する第 2 の浮動小数点フォーマットに変換することを含む、請求項 1 に記載の方法。

10

【請求項 3】

前記乗算結果を前記ベクトルレジスタに格納することをさらに含む、請求項 1 または 2 に記載の方法。

【請求項 4】

前記乗算結果は 2 次元行列であり、前記乗算結果を格納することは、別のセットのクロックサイクルのクロックサイクルごとに、前記乗算結果の値のそれぞれの行を、前記行列乗算ユニットにおける先入れ先出しキューに一時的に格納することを含む、請求項 1 ~ 3 のいずれかに記載の方法。

【請求項 5】

前記ベクトルプロセッサが、前記専用ハードウェアチップのベクトルメモリから前記活性化入力の 2 次元行列を表す前記データを受信することをさらに含み、前記ベクトルメモリは、高速のプライベートメモリを前記ベクトルプロセッサに提供するように構成される、請求項 1 ~ 4 のいずれかに記載の方法。

20

【請求項 6】

前記専用ハードウェアチップの転置ユニットが、前記活性化入力の 2 次元行列の転置演算を実行することをさらに含む、請求項 1 ~ 5 のいずれかに記載の方法。

【請求項 7】

前記専用ハードウェアチップの削減ユニットが、前記乗算結果の数値の削減を実行することと、

前記専用ハードウェアチップの置換ユニットが、前記複数のベクトル処理ユニットに格納された数値を置換することとをさらに含む、請求項 1 ~ 6 のいずれかに記載の方法。

30

【請求項 8】

前記乗算結果を前記専用ハードウェアチップの高帯域幅メモリに格納することをさらに含む、請求項 1 ~ 7 のいずれかに記載の方法。

【請求項 9】

前記専用ハードウェアチップの疎計算コアによって、予め構築されたルックアップテーブルを使用して、疎な高次元データを密な低次元データにマッピングすることをさらに含む、請求項 1 ~ 8 のいずれかに記載の方法。

【請求項 10】

前記専用ハードウェアチップのチップ間相互接続を使用して、前記専用ハードウェアチップのインターフェイスまたはリソースを他の専用ハードウェアチップまたはリソースに接続することをさらに含む、請求項 1 ~ 9 のいずれかに記載の方法。

40

【請求項 11】

前記チップ間相互接続は、前記インターフェイスおよび前記専用ハードウェアチップの高帯域幅メモリを他の専用ハードウェアチップに接続する、請求項 10 に記載の方法。

【請求項 12】

前記インターフェイスは、ホストコンピュータへのホストインターフェイス、またはホストコンピュータのネットワークへの標準ネットワークインターフェイスである、請求項 10 に記載の方法。

【請求項 13】

50

前記ベクトルレジスタは 3 2 個のベクトルレジスタを含む、請求項 1 ~ 1 2 のいずれかに記載の方法。

【請求項 1 4】

前記複数のベクトル処理ユニットの各ベクトル処理ユニットは、各クロックサイクルにおいて、2つのそれぞれの算術論理ユニット（ALU）命令、それぞれのロード命令、およびそれぞれのストア命令を実行するよう構成される、請求項 1 ~ 1 3 のいずれかに記載の方法。

【請求項 1 5】

前記複数のベクトル処理ユニットにおける各ベクトル処理ユニットは、各クロックサイクルにおいてそれぞれのロード命令およびストア命令を実行するためにそれぞれのオフセットメモリアドレスを演算するよう構成される、請求項 1 4 に記載の方法。

10

【請求項 1 6】

前記複数のベクトル処理ユニットの前記 2 次元構造は、複数のレーンと、前記複数のレーンの各々のための複数のサブレーンとを備え、それぞれのベクトル処理ユニットは、前記複数のサブレーンの各々に位置し、同じレーンに位置する前記複数のベクトル処理ユニットのうちのベクトル処理ユニットは、それぞれのロード命令およびストア命令を介して互いと通信するよう構成される、請求項 1 ~ 1 5 のいずれかに記載の方法。

【請求項 1 7】

1 つ以上のコンピュータと 1 つ以上のストレージデバイスとを備えるシステムであって、前記 1 つ以上のストレージデバイスには、前記 1 つ以上のコンピュータによって実行されると前記 1 つ以上のコンピュータに専用ハードウェアチップを使用してニューラルネットワークをトレーニングするための動作を実行させるよう動作可能である命令が格納され、前記動作は、

20

前記専用ハードウェアチップのベクトルプロセッサが、活性化入力の 2 次元行列を表すデータを受信することを含み、前記活性化入力の 2 次元行列は、ニューラルネットワークのそれぞれのネットワーク層に対する入力活性化行列の活性化入力の一部を含み、前記ベクトルプロセッサは、

2 次元構造に配置された複数のベクトル処理ユニットと、

前記活性化入力の 2 次元行列を格納するよう構成されるベクトルレジスタとを含み、前記動作はさらに、

30

行列乗算ユニットが、前記 2 次元行列について乗算結果を生成することを含み、前記生成することは、

複数のクロックサイクルの各クロックサイクルにおいて、前記ベクトルプロセッサの前記ベクトルレジスタから前記専用ハードウェアチップの前記行列乗算ユニットに、前記活性化入力の 2 次元行列における活性化入力のそれぞれの行の少なくとも一部をロードすることと、

前記複数のクロックサイクルの終了前に、前記行列乗算ユニットに重み値をロードすることを含み、前記重み値は、前記活性化入力の 2 次元行列における前記活性化入力に対応する形状に配置され、前記生成することはさらに、

前記行列乗算ユニットが、前記活性化入力の 2 次元行列および重み値を乗算して、前記乗算結果を得ることを含み、前記動作はさらに、

40

前記乗算結果に基づいて、逆伝播を通じて前記ニューラルネットワークの重み値を更新することを含む、システム。

【請求項 1 8】

前記行列乗算ユニットに、前記活性化入力の 2 次元行列における前記活性化入力のそれぞれの行の少なくとも前記一部をロードすることは、さらに、前記活性化入力のそれぞれの行の少なくとも前記一部を、第 1 の浮動小数点フォーマットから、第 1 の浮動小数点フォーマットよりも小さいビットサイズを有する第 2 の浮動小数点フォーマットに変換することを含む、請求項 1 7 に記載のシステム。

【請求項 1 9】

50

1つ以上のコンピュータによって実行されると、前記1つ以上のコンピュータに、専用ハードウェアチップを使用してニューラルネットワークをトレーニングするための動作を実行させる命令を含むプログラムであって、前記動作は、

前記専用ハードウェアチップのベクトルプロセッサが、活性化入力の2次元行列を表すデータを受信することを含み、前記活性化入力の2次元行列は、ニューラルネットワークのそれぞれのネットワーク層に対する入力活性化行列の活性化入力の一部を含み、前記ベクトルプロセッサは、

2次元構造に配置された複数のベクトル処理ユニットと、

前記活性化入力の2次元行列を格納するよう構成されるベクトルレジスタとを含み、前記動作はさらに、

行列乗算ユニットが、前記2次元行列について乗算結果を生成することを含み、前記生成することは、

複数のクロックサイクルの各クロックサイクルにおいて、前記ベクトルプロセッサの前記ベクトルレジスタから前記専用ハードウェアチップの前記行列乗算ユニットに、前記活性化入力の2次元行列における活性化入力のそれぞれの行の少なくとも一部をロードすることと、

前記複数のクロックサイクルの終了前に、前記行列乗算ユニットに重み値をロードすることを含み、前記重み値は、前記活性化入力の2次元行列における前記活性化入力に対応する形状に配置され、前記生成することはさらに、

前記行列乗算ユニットが、前記活性化入力の2次元行列および重み値を乗算して、前記乗算結果を得ることを含み、前記動作はさらに、

前記乗算結果に基づいて、逆伝播を通じて前記ニューラルネットワークの重み値を更新することを含む、プログラム。

【請求項20】

前記行列乗算ユニットに、前記活性化入力の2次元行列における前記活性化入力のそれぞれの行の少なくとも前記一部をロードすることは、さらに、前記活性化入力のそれぞれの行の少なくとも前記一部を、第1の浮動小数点フォーマットから、前記第1の浮動小数点フォーマットよりも小さいビットサイズを有する第2の浮動小数点フォーマットに変換することを含む、請求項19に記載のプログラム。

【発明の詳細な説明】

【背景技術】

【0001】

背景

この明細書は、ハードウェアでのニューラルネットワーク計算の実行に関する。ニューラルネットワークは機械学習モデルであり、各々が、モデルの1つ以上の層を用いて、受信した入力に対する出力、たとえば分類などを生成する。一部のニューラルネットワークには、出力層に加えて1つ以上の隠れ層が含まれる。各隠れ層の出力は、ネットワーク内の次の層、つまりネットワークの次の隠れ層または出力層への入力として用いられる。ネットワークの各層は、それぞれのパラメータのセットの現在の値に従って、受信した入力から出力を生成する。

【発明の概要】

【0002】

概要

本明細書では、機械学習ワークロード、特にトレーニング段階に最適化されたプログラム可能な線形代数アクセラレータである専用ハードウェアチップに関する技術について説明する。

【0003】

一般に、本明細書で説明される主題の革新的な一態様は、特別なハードウェアチップで具現化することができる。

【0004】

10

20

30

40

50

この態様の他の実施形態は、各々、方法のアクションを実行するように構成された、対応するコンピュータシステム、装置、および1つ以上のコンピュータ記憶装置に記録されたコンピュータプログラムを含む。1つ以上のコンピュータのシステムが特定の操作またはアクションを実行するように構成されているとは、動作中にそのシステムにそのような操作またはアクションを実行させるソフトウェア、ファームウェア、ハードウェア、またはそれらの組み合わせがそのシステムにインストールされていることを意味する。1つ以上のコンピュータプログラムが特定の操作またはアクションを実行するように構成されるとは、1つ以上のプログラムが、データ処理装置によって実行されると、データ処理装置にそのような操作またはアクションを実行させる命令を含むことを意味する。

【0005】

10

前述の実施形態および他の実施形態は各々、以下の特徴のうちの1つ以上を単独でまたは組み合わせて任意選択で含むことができる。特に、一実施形態は、以下の特徴をすべて組み合わせて含む。

【0006】

ニューラルネットワークをトレーニングするための専用ハードウェアチップは、専用ハードウェアチップの計算動作を制御するように構成されたスカラープロセッサと、ベクトル処理ユニットの2次元配列を有するように構成されたベクトルプロセッサとを備えてもよく、ベクトル処理ユニットは、すべて、同じ命令を単一命令複数データ方式で実行し、ベクトルプロセッサのロードおよびストア命令を通して互いに通信し、専用ハードウェアチップはさらに、ベクトルプロセッサに結合され、乗算結果を得るために、少なくとも1つの2次元行列を第2の1次元ベクトルまたは2次元行列と乗算するように構成された行列乗算ユニットを備えてもよい。

20

【0007】

ベクトルメモリが、ベクトルプロセッサに高速のプライベートメモリを提供するように構成されてもよい。スカラーメモリが、スカラープロセッサに高速のプライベートメモリを提供するように構成されてもよい。転置ユニットが、行列の転置演算を実行するように構成されてもよい。削減および置換ユニットが、ベクトルアレイの異なるレーン間において、数値上で削減を実行し、数値を置換するように構成されてもよい。高帯域幅メモリが、専用ハードウェアチップのデータを記憶するように構成されてもよい。専用ハードウェアチップは、疎計算コアを備えてもよい。

30

【0008】

専用ハードウェアチップは、インターフェイスと、専用ハードウェアチップ上のインターフェイスまたはリソースを他の専用ハードウェアチップまたはリソースに接続するチップ間相互接続とを備えてもよい。

【0009】

専用ハードウェアチップは、高帯域幅メモリを備えてもよい。チップ間相互接続は、インターフェイスおよび高帯域幅メモリを他の専用ハードウェアチップに接続してもよい。インターフェイスは、ホストコンピュータへのホストインターフェイスであってもよい。インターフェイスは、ホストコンピュータのネットワークへの標準ネットワークインターフェイスであってもよい。

40

【0010】

この明細書において記載される主題は、以下の利点の1つ以上を実現するように特定の実施の形態において実現することができる。専用ハードウェアチップは、機械学習用の32ビット以下の精度計算用に最適化されながらも、従来の0次元および1次元のテンソル計算に加えて、より高次元のテンソル（つまり、2次元以上）をネイティブにサポートするプロセッサを含む。

【0011】

この明細書の主題の1つ以上の実施形態の詳細は、添付の図面および以下の詳細な説明において述べられる。主題の他の特徴、局面および利点は、詳細な説明、図面および特許請求の範囲から明らかになる。

50

【図面の簡単な説明】**【 0 0 1 2 】**

【図 1】ボード上において円形トポロジで接続された専用ハードウェアチップの例示的集合体を接続する高速接続の例示的トポロジを示す。

【図 2】ニューラルネットワークをトレーニングするための例示的な専用ハードウェアチップのハイレベル図を示す。

【図 3】コンピュータコアのハイレベルの例を示す。

【図 4】ニューラルネットワークのためにトレーニングを実行するチップのより詳細な図を示す。

【発明を実施するための形態】**【 0 0 1 3 】**

さまざまな図面の同様の参照番号および名称は、同様の要素を示す。

詳細な説明

複数の層を有するニューラルネットワークをトレーニングし、それを推論の計算に用いることができる。一般に、ニューラルネットワークの一部またはすべての層は、ニューラルネットワークのトレーニング中に調整されるパラメータを有する。たとえば、一部またはすべての層は、その層に対する、重みとも称されるパラメータの行列に、層出力の生成の一部として、その層への入力を乗算できる。行列内のパラメータの値は、ニューラルネットワークのトレーニング中に調整される。

【 0 0 1 4 】

特に、トレーニング中、トレーニングシステムは、ニューラルネットワークのトレーニング手順を実行して、ニューラルネットワークのパラメータの値を調整して、たとえば、パラメータの初期値からパラメータのトレーニングを経た値を判断する。トレーニングシステムは、逆伝播として知られる誤差の逆伝播を、最適化方法と組み合わせて用いて、ニューラルネットワークの各パラメータに関して目的関数の勾配を計算し、その勾配を用いてパラメータの値を調整する。

【 0 0 1 5 】

トレーニングされたニューラルネットワークは、順方向伝播を用いて推論を計算でき、つまり、ニューラルネットワークの層を介して入力を処理し、その入力に対するニューラルネットワーク出力を生成できる。

【 0 0 1 6 】

たとえば、入力が与えられると、ニューラルネットワークはその入力に対する推論を計算することができる。ニューラルネットワークは、ニューラルネットワークの各層を通して入力を処理することによって、この推論を計算する。いくつかの実装形態では、ニューラルネットワークの層はシーケンスで配置される。

【 0 0 1 7 】

したがって、受け取った入力から推論を計算するために、ニューラルネットワークはその入力を受け取り、それを各ニューラルネットワーク層を通してシーケンスで処理して推論を生成し、1つのニューラルネットワーク層からの出力が次のニューラルネットワーク層への入力として与えられる。ニューラルネットワーク層へのデータ入力、たとえば、ニューラルネットワークへの入力、またはシーケンス内におけるその層の下層の、あるニューラルネットワーク層への出力は、その層への活性化入力と呼ぶことができる。

【 0 0 1 8 】

いくつかの実装形態では、ニューラルネットワークの層は有向グラフで配置される。つまり、任意の特定の層が複数の入力、複数の出力、またはそれらの両方を受け取ることができる。ニューラルネットワークの層は、ある層の出力を前の層への入力として送り返すことができるように構成することもできる。

【 0 0 1 9 】

ある例示的システムは、行列乗算やその他の多次元配列の計算に最適化された高性能マルチチップテンソル計算システムである。これらの操作は、ニューラルネットワークのト

10

20

30

40

50

レーニング、およびオプションで、ニューラルネットワークを用いて推論を計算するために重要である。

【 0 0 2 0 】

ある例示的システムでは、システムがトレーニングおよび推論計算を効率的に実行するために、複数の専用チップが操作を分散するように配置される。1つの実装形態では、ボード上に4つのチップがあり、より大きなシステムでは、多くのボードがラック内で隣り合っているか、またはそうでなければ相互にデータ通信を行なう。

【 0 0 2 1 】

図1は、ボード上において円形トポロジで接続された専用ハードウェアチップ101a ~ 101dの例示的集合体を接続する高速接続の例示的トポロジを示す。各チップには2つのプロセッサ(102a ~ 102h)が含まれる。このトポロジは、1次元(1D)トラスであり、1Dトラスでは、各チップは2つの隣接チップに直接接続される。示されているように、一部の实装形態では、チップには、操作を実行するようソフトウェア命令またはファームウェア命令でプログラムされたマイクロプロセッサコアが含まれる。図1では、すべてのチップが単一のモジュール100上にある。図に示されているプロセッサ間の線は、高速データ通信リンクを表す。プロセッサは1つの集積回路基板上に有利に製造されるが、複数の基板上に製造することもできる。チップ境界を越えて、リンクは、チップ間ネットワークリンクであり、同じチップ上のプロセッサは、チップ内インターフェイスリンクを介して通信する。リンクは、一度に1つのプロセッサのみがデータを送信できる半二重リンクであってもよいし、データを両方向に同時に送信できる全二重リンクであってもよい。この例示的トポロジを用いる並列処理およびその他については、2017年2月21日に提出され、ここに引用により援用される、「PARALLEL PROCESSING OF REDUCTION AND BROADCAST OPERATIONS ON LARGE DATASETS OF NON-SCALAR DATA (非スカラーデータの大規模データセットの削減およびブロードキャスト操作の並列処理)」と題される米国特許出願第62/461,758号において詳細に説明される。

【 0 0 2 2 】

図2は、ニューラルネットワークをトレーニングするための例示的な専用ハードウェアチップのハイレベル図を示す。図示されているように、単一の専用ハードウェアチップには2つの独立したプロセッサ(202a、202b)が含まれる。各プロセッサ(202a、202b)には、2つの異なるコア:(1)コンピュートコア、たとえば非常に長い命令語(VLIW)マシン(203a、203b)、および(2)疎計算コア、つまり埋め込み層アクセラレータ(205a、205b)が含まれる。

【 0 0 2 3 】

各コア(203a、203b)は、密線形代数問題に対して最適化される。単一の非常に長い命令語が、いくつかのコンピュートコアを並列で制御する。コンピュートコアは、図3および図4を参照してより詳細に説明される。

【 0 0 2 4 】

例示的な疎計算コア(205a、205b)は、非常に疎である高次元データを密な低次元データにマッピングし、残りの層が密に充填された入力データを処理するようにする。たとえば、疎計算コアは、トレーニング中のニューラルネットワークの埋め込み層の計算を実行できる。

【 0 0 2 5 】

この疎から密へのマッピングを実行するために、疎計算コアは、埋め込みテーブルである、予め構築されたルックアップテーブルを用いる。たとえば、ユーザ入力として一連のクエリワードがある場合、各クエリワードはハッシュ識別子またはワンホットエンコードされたベクトルに変換される。識別子をテーブルインデックスとして用いて、埋め込みテーブルは対応する密なベクトルを返し、それは、次の層への入力活性化ベクトルになることができる。疎計算コアは、検索クエリワードにわたって削減操作を実行して、1つの密な活性化ベクトルを作成することもできる。疎計算コアは、効率的な疎の、分散されたル

10

20

30

40

50

ックアップを実行し、なぜならば、埋め込みテーブルが巨大であり得、専用ハードウェアチップの1つの限られた容量の高帯域幅メモリに収まらないためである。疎計算コア機能に関する詳細は、2016年2月5日に提出された「MATRIX PROCESSING APPARATUS（行

列処理装置）」と題される米国特許出願第15/016,486号に記載され、それをここに引用により援用する。

【0026】

図3は、コンピュータコア(300)のハイレベルの例を示す。コンピュータコアは、複数のコンピュータユニットを並列で制御するマシン、つまりVLIWマシンであることができる。各コンピュータコア(300)には、スカラーメモリ(304)、ベクトルメモリ(308)、スカラープロセッサ(303)、ベクトルプロセッサ(306)、および拡張ベクトルユニット(つまり、行列乗算ユニット(MXU)(313)と、転置ユニット(XU)(314)と、削減および置換ユニット(RPU)(316))が含まれる。

【0027】

例示的なスカラープロセッサは、VLIW命令のフェッチ/実行ループを実行し、コンピュータコアを制御する。命令バンドルをフェッチしてデコードした後、スカラープロセッサ自体は、スカラープロセッサ(303)およびスカラーメモリ(304)の複数のマルチビットレジスタ、つまり32の32ビットレジスタを用いて、命令バンドルのスカラーロットにおいて見いだされる命令を実行するのみである。スカラー命令セットには、アドレス計算、ロード/ストア命令、分岐命令などで使用される通常の算術演算が含まれる。残りの命令ロットは、ベクトルプロセッサ(306)または他の拡張ベクトルユニット(313、314、316)の命令をエンコードする。デコードされたベクトル命令は、ベクトルプロセッサ(306)に転送される。

【0028】

ベクトル命令とともに、スカラープロセッサ(303)は、最大3つのスカラーレジスタの値を他のプロセッサおよびユニットに転送して、演算を実行できる。スカラープロセッサは、ベクトルプロセッサから計算結果を直接取得することもできる。ただし、一部の実装形態では、例示的なチップは、ベクトルプロセッサからスカラープロセッサへの低帯域幅通信パスを有する。

【0029】

ベクトル命令ディスパッチャは、スカラープロセッサとベクトルプロセッサとの間にある。このディスパッチャは、非スカラーVLIWスロットからデコードされた命令を受け取り、それらの命令をベクトルプロセッサ(306)にブロードキャストする。ベクトルプロセッサ(306)は、同じ命令を単一命令複数データ(SIMD)方式で実行するベクトル処理ユニットの2次元配列、すなわち128×8の配列からなる。ベクトル処理ユニットは、図4を参照して詳細に説明される。

【0030】

例示的なスカラープロセッサ(303)は、小さい、高速な、プライベートスカラーメモリ(304)にアクセスする。このメモリは、はるかに大きく、低速な高帯域幅メモリ(HBM)(310)によってバックアップされる。同様に、例示的なベクトルプロセッサ(306)は、小さい、高速な、プライベートベクトルメモリ(306)にアクセスする。このメモリも、HBM(310)によってバックアップされる。ワード粒度アクセスは、スカラープロセッサ(303)とスカラーメモリ(304)との間、またはベクトルプロセッサ(306)とベクトルメモリ(308)との間で起こる。ベクトルプロセッサとベクトルメモリとの間のロードおよびストアの粒度は、128個の32ビットワードのベクトルである。ダイレクトメモリアccessは、スカラーメモリ(304)とHBM(310)との間、およびベクトルメモリ(306)とHBM(310)との間で起こる。いくつかの実装形態では、HBM(310)からプロセッサ(303、306)へのメモリ転送は、スカラーメモリまたはベクトルメモリを介してのみ実行できる。さらに、スカラーメモリとベクトルメモリとの間でダイレクトメモリ転送が行われない場合がある。

【 0 0 3 1 】

命令は、拡張ベクトルユニット演算を指定する場合がある。実行された各ベクトルユニット命令に加えて、各々が1つのレジスタ値を拡張ベクトルユニットに入力オペランドとして送ることができる2次元、つまり 128×8 のベクトルユニットがある。各拡張ベクトルユニットは、入力オペランドを受け取り、対応する演算を実行し、結果をベクトルプロセッサ(306)に返す。拡張ベクトルユニットは、図4を参照して以下に説明される。

【 0 0 3 2 】

図4は、ニューラルネットワークのためにトレーニングを実行するチップのより詳細な図を示す。上で図示し説明したように、チップには2つのコンピュータコア(480a、480b)と2つの疎計算コア(452a、452b)とが含まれる。

10

【 0 0 3 3 】

チップには、ホストコンピュータ(450)または複数のホストコンピュータへのインターフェイスを含む共有領域がある。このインターフェイスは、ホストコンピュータへのホストインターフェイス、またはホストコンピュータのネットワークへの標準ネットワークインターフェイスとすることができる。共有領域は、下部に沿って高帯域幅メモリのスタック(456a~456d)、およびインターフェイスとメモリとを接続するチップ間相互接続(448)、ならびに他のチップからのデータも有し得る。相互接続は、インターフェイスをハードウェアチップ上の計算リソースに接続することもできる。高帯域幅メモリの複数のスタック、つまり2つのスタック(456a~456b、456c~456d)が各コンピュータコア(480a、480b)に関連付けられる。

20

【 0 0 3 4 】

チップは、データを高帯域幅メモリ(456c~456d)に保存し、そのデータをベクトルメモリにおいて読込および読出し(446)、そのデータを処理する。コンピュータコア(480b)自体は、2次元に分割されたオンチップS-RAMであるベクトルメモリ(446)を含む。ベクトルメモリには、アドレスが浮動小数点数、つまり各々32ビットである128個の数値を保持するアドレス空間がある。コンピュータコア(480b)は、値を計算する計算ユニット、および計算ユニットを制御するスカラーユニットも含む。計算ユニットはベクトルプロセッサを含んでもよく、スカラーユニットはスカラープロセッサを含んでもよい。専用チップの一部を形成し得るコンピュータコアは、行列乗算ユニット、または行列、つまり 128×128 の行列の転置演算を実行する転置ユニット(422)などの別の拡張演算ユニット、ならびに削減および置換ユニットをさらに含むことができる。

30

【 0 0 3 5 】

ベクトルプロセッサ(306)は、ベクトル処理ユニットの2次元配列、すなわち 128×8 から成り、これらはすべて、同じ命令を単一命令複数データ(SIMD)方式で実行する。ベクトルプロセッサには、レーンとサブレーン、つまり128本のレーンおよび8本のサブレーンがある。レーン内において、ベクトルユニットはロード命令およびストア命令を介して互いに通信する。各ベクトルユニットは、一度に1つの4バイト値にアクセスできる。同じレーンに属さないベクトルユニットは直接通信できない。これらのベクトルユニットは、以下で説明する削減/置換ユニットを用いる必要がある。

40

【 0 0 3 6 】

計算ユニットは、ベクトル処理ユニットにおいて、浮動小数点演算および整数演算の両方に用いることができるベクトルレジスタ(440)、つまり32個のレジスタを含む。計算ユニットは、計算を実行するための2つの算術論理ユニット(ALU)(406c~406d)を含む。一方のALU(406c)は浮動小数点加算を実行し、他方のALU(406d)は浮動小数点乗算を実行する。両方のALU(406c~406d)は、シフト、マスク、比較などの他のさまざまな演算を実行できる。たとえば、コンピュータコア(480b)は、ベクトルレジスタV1と第2のベクトルレジスタV2とを加算し、結果を第3のベクトルレジスタV3に入れたい場合がある。この加算を計算するために、コンピュータコア(480b)は複数の演算を1クロックサイクルで実行する。これらのレ

50

ジスタをオペランドとして用いて、各ベクトルユニットは、クロックサイクルごとに2つのALU命令と1つのロードおよび1つのストア命令とを同時に実行できる。ロードまたはストア命令のベースアドレスは、スカラープロセッサで計算でき、ベクトルプロセッサに転送できる。各サブプレーンにおけるベクトルユニットの各々は、ストライドや特別なインデックス付きアドレスレジスタなどのさまざまな方法を用いて、それ自体のオフセットアドレスを計算できる。

【0037】

計算ユニットは、平方根や逆数などの演算を実行する拡張単項パイプライン(EUP)(416)も含む。コンピュータコア(480b)は、これらの演算を実行するのに3クロックサイクルかかり、なぜならば、それらは計算が複雑であるからである。EUP処理には1クロックサイクル以上かかるため、結果を保存する先入れ先出しのデータストレージがある。演算が終了すると、結果はFIFOに保存される。コンピュータコアは、後で別の命令を用いて、FIFOからデータを引き出し、それをベクトルレジスタに格納できる。乱数生成器(420)により、コンピュータコア(480b)はサイクルごとに複数の乱数、つまりサイクルごとに128の乱数を生成できる。

10

【0038】

上記のように、専用のハードウェアチップの一部として実装できる各プロセッサは、3つの拡張演算ユニット、つまり、行列乗算演算を実行する行列乗算ユニット(448)と、行列、すなわち 128×128 の行列の転置演算を実行する転置ユニット(422)と、削減および置換ユニット(図4において別個のユニット424、426として示される)とを有する。

20

【0039】

行列乗算ユニットは、2つの行列間で行列乗算を実行する。コンピュータコアは、乗算される行列である一連の数値を読み込む必要があるため、行列乗算ユニット(438)はデータを取り込む。図示されているように、データはベクトルレジスタ(440)から来る。各ベクトルレジスタには、 128×8 の数値、つまり32ビットの数値が含まれる。しかしながら、データを行列乗算ユニット(448)に送って、数値をより小さなビットサイズ、つまり32ビットから16ビットに変更すると、浮動小数点変換が発生する場合がある。並直列変換器(440)は、ベクトルレジスタから数値が読み取られるときに、2次元配列つまり 128×8 の行列が128個の数値のセットとして読み取られ、次の8クロックサイクルの各々ごとに行列乗算ユニット(448)に送信されることを保証する。行列乗算がその計算を完了した後、結果は非直列化され(442a、442b)、これは、結果行列が、ある数のクロックサイクルの間保持されることを意味する。たとえば、 128×8 の配列の場合、128個の数値が8クロックサイクルの各々ごとに保持され、次いでFIFOにプッシュされ、 128×8 の数値の2次元配列を1クロックサイクルで取得してベクトルレジスタ(440)に格納できる。

30

【0040】

複数すなわち128のサイクルからなる期間にわたって、重みが、行列を乗算する数値として行列乗算ユニット(448)にシフトされる。行列および重みがロードされると、コンピュータコア(480)は、数値のセット、つまり 128×8 を行列乗算ユニット(448)に送ることができる。セットの各ラインに行列を乗算して、クロックサイクルごとにある数の結果、つまり128を生成できる。コンピュータコアが行列乗算を実行している間、コンピュータコアは、前の行列の計算プロセスが完了したときに、コンピュータコアが乗算する次の行列が利用できるように、バックグラウンドにおいて次の行列になる新たな数値のセットのシフトも行う。行列乗算ユニット(448)は、「LOW MATRIX MULTIPLY UNIT COMPOSED OF MULTI-BIT CELLS(マルチビットセルで構成された低行列乗算ユニット)」と題される16113-8251001、および「MATRIX MULTIPLY UNIT WITH NUMERICS OPTIMIZED FOR NEURAL NETWORK APPLICATIONS(数値がニューラルネットワークアプリケーション向けに最適化された行列乗算ユニット)」と題される16113-8252001に、より詳細に説明され、それらの両方

40

50

をここに引用により援用する。

【 0 0 4 1 】

転置ユニットは、行列を転置する。転置ユニット (4 2 2) は、数値を取り込み、それらを転置して、レーンを横切る数が他の次元の数値と転置されるようにする。一部の実装形態では、ベクトルプロセッサは 128×8 のベクトルユニットを含む。したがって、 128×128 の行列を転置するには、完全な行列転置のために 16 個の個別の転置命令が必要である。転置が終了すると、転置された行列が利用可能になる。ただし、転置された行列をベクトルレジスタファイルに移動するには、明示的な命令が必要である。

【 0 0 4 2 】

削減 / 置換ユニット (またはユニット 4 2 4 、 4 2 6) は、置換、レーン回転、回転置換、レーン削減、置換されたレーン削減、およびセグメント化された置換されたレーン削減などのさまざまな操作をサポートすることで、クロスレーン通信の問題に対処する。図示されているように、これらの計算は別々であるが、コンピュータコアは一方もしくは他方または一方に連鎖された他方を用いることができる。削減ユニット (4 2 4) は、数値からなる各ラインにおけるすべての数値を加算し、それらの数値を置換ユニット (4 2 6) に供給する。置換ユニットは、異なるレーン間でデータを移動する。転置ユニット、削減ユニット、置換ユニット、および行列乗算ユニットは、各々、完了までに 1 クロックサイクル以上かかる。したがって、各ユニットには F I F O が関連付けられ、計算結果を F I F O にプッシュし、後で別の命令を実行して、データを F I F O からベクトルレジスタにプルできる。F I F O を用いることにより、コンピュータコアは、長々とした演算の間、複数のベクトルレジスタを予約する必要がない。図示されているように、各ユニットはベクトルレジスタ (4 4 0) からデータを取得する。

【 0 0 4 3 】

コンピュータコアは、スカラーユニットを用いて計算ユニットを制御する。スカラーユニットには 2 つの主要な機能があり、それは、(1) ループカウントおよびアドレス指定を実行すること、ならびに (2) D M A コントローラがバックグラウンドにおいて高帯域幅メモリ (4 5 6 c ~ 4 5 6 d) とベクトルメモリ (4 4 6) との間で、および次いで例示のシステムにおける他のチップへのチップ間接続 (4 4 8) にデータを移動するよう、ダイレクトメモリアドレス (D M A) 要求を生成することである。スカラーユニットは、命令メモリ (4 0 4) と、命令デコードおよび発行 (4 0 2) と、スカラーレジスタすなわち 32 ビットを含むスカラー処理ユニット (4 0 8) と、スカラーメモリ (4 1 0) と、クロックサイクルごとに 2 つの演算を実行する 2 つの A L U (4 0 6 a 、 4 0 6 b) とを含む。スカラーユニットは、オペランドと即値とをベクトル演算に渡すことができる。各命令は、ベクトルレジスタ (4 4 0) で実行される命令を含む命令バンドルとして、命令デコードおよび発行 (4 0 2) から送ることができる。各命令バンドルは非常に長い命令語 (V L I W) であり、各命令はある数のビット幅であり、ある数の命令フィールドに分割される。

【 0 0 4 4 】

チップ 4 0 0 を用いて、ニューラルネットワークのトレーニングの少なくとも一部を実行することができる。特に、ニューラルネットワークをトレーニングする場合、システムはホストインターフェイス (4 5 0) を用いてホストコンピュータからラベル付きトレーニングデータを受信する。ホストインターフェイスは、ニューラルネットワーク計算のためのパラメータを含む命令を受信することもできる。パラメータは、処理すべき層の数、各層についての対応する重み入力のセット、活性化入力の初期セット、つまり推論の計算またはトレーニングの対象となるニューラルネットワークへの入力であるトレーニングデータ、各層の対応する入力および出力サイズ、ニューラルネットワーク計算のストライド値、ならびに処理対象の層のタイプ、たとえば畳み込み層または全結合層、のうちの少なくとも 1 つ以上を含むことができる。

【 0 0 4 5 】

重み入力のセットおよび活性化入力のセットは、コンピュータコアの行列乗算ユニット

10

20

30

40

50

に送ることができる。重み入力および活性化入力を行列乗算ユニットに送る前に、システム内の他のコンポーネントが入力に対して他の計算を実行してもよい。一部の実装形態では、疎計算コアからコンピュータコアに活性化を送る方法が2つある。まず、疎計算コアは、高帯域幅メモリを介して通信を送信することができる。大量のデータの場合、疎計算コアは、ダイレクトメモリアドレス(DMA)命令を用いて活性化を高帯域幅メモリに格納でき、これにより、コンピュータコアにおいてターゲット同期フラグが更新される。コンピュータコアは、同期命令を用いてこの同期フラグを待つことができる。同期フラグがセットされると、計算コアはDMA命令を用いて、活性化を高帯域幅メモリから対応するベクトルメモリにコピーする。

【0046】

次に、疎計算コアは、通信をコンピュータコアベクトルメモリに直接送信できる。データ量が大きくない場合(つまり、コンピュータコアベクトルメモリに収まる場合)、疎計算コアは、コンピュータコアに同期フラグで通知しながら、DMA命令を用いてコンピュータコアのベクトルメモリに活性化を直接格納できる。コンピュータコアは、この同期フラグを待ったのち、活性化に依存する計算を実行することができる。

【0047】

行列乗算ユニットは、重み入力および活性化入力を処理し、出力のベクトルまたは行列をベクトル処理ユニットに与えることができる。ベクトル処理ユニットは、処理された出力のベクトルまたは行列を格納できる。たとえば、ベクトル処理ユニットは、非線形関数を行列乗算ユニットの出力に適用して、活性化された値を生成できる。いくつかの実装形態では、ベクトル処理ユニットは、正規化された値、プールされた値、またはその両方を生成する。処理された出力のベクトルは、たとえばニューラルネットワーク内の後続の層で用いるために、行列乗算ユニットへの活性化入力として用いることができる。

【0048】

トレーニングデータのバッチについての処理済み出力のベクトルが計算されると、それらの出力をラベル付きトレーニングデータの期待される出力と比較して、誤差を判断できる。その後、システムは、ネットワークをトレーニングするために、逆伝播を実行して、ニューラルネットワークを介して誤差を伝播できる。損失関数の勾配は、オンチップでベクトル処理ユニットの算術論理ユニットを用いて計算される。

【0049】

ある例示的システムでは、ニューラルネットワークを介した逆伝播を実行するために、活性化勾配が必要である。活性化勾配をコンピュータコアから疎計算コアに送るために、例示的システムでは、コンピュータコアDMA命令を用いて、ターゲット疎計算コアに同期フラグで通知しながら、活性化勾配を高帯域幅メモリに保存できる。疎計算コアは、この同期フラグを待ったのち、活性化勾配に依存する計算を実行することができる。

【0050】

行列乗算ユニットは、逆伝播のために2つの行列乗算演算を実行する。一方の行列乗算は、逆伝播誤差をネットワーク内の前の層からネットワークを通る逆方向パスに沿って重みに適用して、重みを調整してニューラルネットワークのための新たな重みを決定する。第2の行列乗算は、ニューラルネットワーク内の前の層へのフィードバックとして、元の活性化に誤差を適用する。元の活性化は、順方向パス中に生成され、逆方向パス中に用いるために保存されてもよい。計算には、浮動小数点加算、減算、および乗算を含む、ベクトル処理ユニットにおける汎用命令を用いることができる。汎用命令には、比較、シフト、マスク、および論理演算も含まれ得る。行列の乗算は非常に加速され得るが、ベクトル処理ユニットの算術論理ユニットは、サイクルあたり、コアあたり $128 \times 8 \times 2$ の演算の速度で一般的な計算を実行する。

【0051】

本明細書において記載される主題および機能的動作の実施形態は、本明細書に開示される構造およびそれらの構造的等価物を含む、デジタル電子回路系において、有形で実施されるコンピュータソフトウェアもしくはファームウェアにおいて、コンピュータハードウ

10

20

30

40

50

エアにおいて、またはそれらの1つ以上の組合せにおいて実現され得る。本明細書に記載される主題の実施形態は、1つ以上のコンピュータプログラムとして、すなわち、データ処理装置による実行のために、または、データ処理装置の動作を制御するために有形の非一時的な記憶媒体上でエンコードされたコンピュータプログラム命令の1つ以上のモジュールとして実現され得る。コンピュータ記憶媒体は、機械可読記憶装置、機械可読記憶基板、ランダムもしくはシリアルアクセスメモリデバイス、または、それらの1つ以上の組合せであり得る。代替的に、または加えて、プログラム命令は、データ処理装置による実行に対して好適な受信側装置への送信のために情報をエンコードするように生成される、たとえばマシンにより生成された電気信号、光信号、または電磁気信号などの、人為的に生成された伝搬される信号上でエンコードすることができる。

10

【0052】

「データ処理装置」という用語は、データ処理ハードウェアを指し、例としてプログラマブルプロセッサ、コンピュータ、または複数のプロセッサもしくはコンピュータを含む、データを処理するためのすべての種類の装置、デバイスおよびマシンを包含する。当該装置は、たとえばFPGA（フィールドプログラマブルゲートアレイ）またはASIC（特定用途向け集積回路）といった特定目的論理回路でもあることができるかまたはそれをさらに含み得る。当該装置は、ハードウェアに加えて、たとえばプロセッサファームウェア、プロトコルスタック、データベース管理システム、オペレーティングシステム、または、それらの1つ以上の組合せを構成するコードといった、コンピュータプログラムについて実行環境を作成するコードをオプションとして含み得る。

20

【0053】

プログラム、ソフトウェア、ソフトウェアアプリケーション、アプリケーション、モジュール、ソフトウェアモジュール、スクリプトまたはコードとも称され、または記載され得るコンピュータプログラムは、コンパイル型もしくはインタープリタ型言語、または宣言型もしくは手続き型言語を含む任意の形態のプログラミング言語で記述され得、スタンドアロンプログラムとして、または、モジュール、コンポーネント、サブルーチン、もしくは、コンピューティング環境で使用するのに好適な他のユニットとして任意の形態で展開され得る。プログラムは、ファイルシステムにおけるファイルに対応し得るが、対応する必要があるわけではない。プログラムは、当該プログラムに専用である単一のファイルにおいて、または、複数の連携ファイル（*coordinated files*）（たとえばコードの1つ以上のモジュール、サブプログラムまたは部分を格納するファイル）において、他のプログラムまたはデータ（たとえばマークアップ言語ドキュメントに格納される1つ以上のスクリプト）を保持するファイルの一部に格納され得る。コンピュータプログラムは、1つの場所に位置するかもしれないが複数の場所にわたって分散されデータ通信ネットワークによって相互接続される1つのコンピュータまたは複数のコンピュータ上で実行されるように展開され得る。

30

【0054】

本明細書に記載されるプロセスおよび論理フローは、入力データ上で動作し出力を生成することにより機能を実行するよう1つ以上のプログラマブルコンピュータが1つ以上のコンピュータプログラムを実行することによって実行され得る。本プロセスおよび論理フローは、たとえばFPGAもしくはASICといった特殊目的論理回路系によっても、または特殊目的論理回路計と1つ以上のプログラムされたコンピュータとの組み合わせによっても実行され得る。

40

【0055】

コンピュータプログラムの実行に好適であるコンピュータは、汎用マイクロプロセッサもしくは特殊目的マイクロプロセッサもしくはその両方または任意の種類の中央処理ユニットに基づき得る。一般に、中央処理ユニットは、リードオンリメモリもしくはランダムアクセスメモリまたはその両方から命令およびデータを受け取る。コンピュータの必須の要素は、命令を実行するための中央処理ユニットと、命令およびデータを格納するための1つ以上のメモリデバイスとである。中央処理ユニットおよびメモリは、特殊目的論理回

50

路系によって補足され得るか、または特殊目的論理回路に組み込まれ得る。一般に、コンピュータはさらに、たとえば磁気ディスク、光磁気ディスクまたは光ディスクといった、データを格納するための1つ以上の大容量記憶装置を含むか、当該1つ以上の大容量記憶装置からデータを受け取るかもしくは当該1つ以上の大容量記憶装置にデータを転送するよう動作可能に結合されるか、またはその両方を行う。しかしながら、コンピュータはそのような装置を有する必要はない。さらに、コンピュータはたとえば、携帯電話、携帯情報端末（PDA）、モバイルオーディオまたはビデオプレーヤ、ゲームコンソール、全地球測位システム（GPS）受信機、またはポータブル記憶装置（たとえばユニバーサルシリアルバス（USB）フラッシュドライブ）といった別のデバイスに埋め込まれ得る。

【0056】

コンピュータプログラム命令およびデータを格納するのに好適であるコンピュータ可読媒体は、例として、たとえばEPROM、EEPROMおよびフラッシュメモリデバイスといった半導体メモリデバイスを含むすべての形態の不揮発性メモリ、媒体およびメモリデバイス；たとえば内部ハードディスクまたはリムーバブルディスクといった磁気ディスク；光磁気ディスク；ならびにCD-ROMおよびDVD-ROMディスクを含む。

【0057】

ユーザとの対話を与えるために、本明細書に記載される主題の実施形態は、たとえばCRT（陰極線管）またはLCD（液晶ディスプレイ）モニタといったユーザに対して情報を表示するための表示デバイスと、たとえばマウス、トラックボールといったユーザがコンピュータに入力を与えることができるキーボードおよびポインティングデバイスとを有するコンピュータ上で実現され得る。他の種類のデバイスが、同様に、ユーザとの対話を与えるために用いられ得；たとえば、ユーザに提供されるフィードバックは、たとえば視覚フィードバック、聴覚フィードバックまたは触覚フィードバックといった任意の形態の感覚フィードバックであり得；ユーザからの入力、音響入力、音声入力、または触覚入力を含む任意の形態で受け取られ得る。加えて、コンピュータは、ユーザが使用するデバイスにドキュメントを送信しユーザが使用するデバイスからドキュメントを受信することによって、たとえば、ウェブブラウザから受信された要求に回答してユーザのデバイス上のウェブブラウザにウェブページを送信することによって、ユーザと対話し得る。また、コンピュータは、テキストメッセージまたは他の形式のメッセージを個人用デバイス、たとえばスマートフォンなどに送信し、メッセージングアプリケーションを実行し、ユーザから応答メッセージを受信することにより、ユーザと対話できる。

【0058】

本明細書に記載される主題の実施形態は、たとえばデータサーバとしてバックエンドコンポーネントを含む計算システムにおいて実現され得るか、たとえばアプリケーションサーバといったミドルウェアコンポーネントを含む計算システムにおいて実現され得るか、たとえば本明細書に記載される主題の実現例とユーザが対話することが可能であるグラフィカルユーザインターフェイス、ウェブブラウザもしくはアプリを有するクライアントコンピュータといったフロントエンドコンポーネントを含む計算システムにおいて実現され得るか、または1つ以上のそのようなバックエンドコンポーネント、ミドルウェアコンポーネントもしくはフロントエンドコンポーネントの任意の組合せの計算システムにおいて実現され得る。システムのコンポーネントは、たとえば通信ネットワークといったデジタルデータ通信の任意の形態または媒体によって相互接続され得る。通信ネットワークの例は、ローカルエリアネットワーク（LAN）およびワイドエリアネットワーク（WAN）、たとえばインターネットを含む。

【0059】

計算システムはクライアントおよびサーバを含むことができる。クライアントとサーバとは一般に互いから遠隔にあり、典型的には通信ネットワークを通じて対話する。クライアントとサーバとの関係は、それぞれのコンピュータ上で実行されるとともに互いに対してクライアント-サーバ関係を有するコンピュータプログラムによって発生する。いくつかの実施形態では、サーバは、例えば、クライアントとして振る舞うユーザデバイスと対

10

20

30

40

50

話するユーザにデータを表示し、およびそのユーザからユーザ入力を受信する目的で、データ、例えば、HTML ページをユーザデバイスに送信する。ユーザデバイスで生成されたデータ、例えば、ユーザ対話の結果は、ユーザデバイスからサーバで受信することができる。

【0060】

実施形態1は、ニューラルネットワークをトレーニングするための専用ハードウェアチップであって、専用ハードウェアチップの計算動作を制御するように構成されたスカラープロセッサと、ベクトル処理ユニットの2次元配列を有するように構成されたベクトルプロセッサとを備え、ベクトル処理ユニットは、すべて、同じ命令を単一命令複数データ方式で実行し、ベクトルプロセッサのロードおよびストア命令を通して互いに通信し、専用ハードウェアチップはさらに、ベクトルプロセッサに結合され、乗算結果を得るために、少なくとも1つの2次元行列を第2の1次元ベクトルまたは2次元行列と乗算するように構成された行列乗算ユニットを備える。

10

【0061】

実施形態2は、ベクトルプロセッサに高速のプライベートメモリを提供するように構成されたベクトルメモリをさらに備える、実施形態1の専用ハードウェアチップである。

【0062】

実施形態3は、スカラープロセッサに高速のプライベートメモリを提供するように構成されたスカラーメモリをさらに備える、実施形態1または2の専用ハードウェアチップである。

20

【0063】

実施形態4は、行列の転置演算を実行するように構成された転置ユニットをさらに備える、実施形態1～3のいずれか1つの専用ハードウェアチップである。

【0064】

実施形態5は、ベクトルアレイの異なるレーン間において、数値上で削減を実行し、数値を置換するように構成された、削減および置換ユニットをさらに備える、実施形態1～4のいずれか1つの専用ハードウェアチップである。

【0065】

実施形態6は、専用ハードウェアチップのデータを記憶するように構成された高帯域幅メモリをさらに備える、実施形態1～5のいずれか1つの専用ハードウェアチップである。

30

【0066】

実施形態7は、疎計算コアをさらに含む、実施形態1～6のいずれか1つの専用ハードウェアチップである。

【0067】

実施形態8は、インターフェイスと、専用ハードウェアチップ上のインターフェイスまたはリソースを他の専用ハードウェアチップまたはリソースに接続するチップ間相互接続とをさらに備える、実施形態1～7のいずれか1つの専用ハードウェアチップである。

【0068】

実施形態9は、複数の高帯域幅メモリをさらに備え、チップ間相互接続は、インターフェイスおよび高帯域幅メモリを他の専用ハードウェアチップに接続する、実施形態1～8のいずれか1つの専用ハードウェアチップである。

40

【0069】

実施形態10は、インターフェイスは、ホストコンピュータへのホストインターフェイスである、実施形態1～9のいずれか1つの専用ハードウェアチップである。

【0070】

実施形態11は、インターフェイスは、ホストコンピュータのネットワークへの標準ネットワークインターフェイスである、実施形態1～10のいずれか1つの専用ハードウェアチップである。

【0071】

本明細書は多くの特定の実現例の詳細を含んでいるが、これらは如何なる発明の範囲ま

50

たは請求され得るものの範囲に対する限定としても解釈されるべきではなく、特定の発明の特定の実施形態に特有の特徴であり得る記載として解釈されるべきである。別個の実施形態の文脈で本明細書において記載されるある特徴は、単一の実施形態において組合せでも実現され得る。反対に、単一の実施形態の文脈において記載されるさまざまな特徴は、複数の実施形態において別々に、または任意の好適な部分的組合せでも実現され得る。さらに、特徴は、ある組合せにおいて作用すると上で記載され、最初はそのように請求されていさえする場合もあるが、請求される組合せからの1つ以上の特徴はいくつかの場合には当該組合せから削除され得、請求される組合せは、部分的組合せまたは部分的組合せの変形例に向けられ得る。

【0072】

10

同様に、動作が図においては特定の順に示されているが、そのような動作は、望ましい結果を達成するために、示された当該特定の順もしくは連続した順で実行される必要があると理解されるべきではなく、または、すべての示された動作が実行される必要があると理解されるべきではない。ある状況においては、マルチタスキングおよび並列処理が有利であり得る。さらに、上述の実施形態におけるさまざまなシステムモジュールおよびコンポーネントの分離は、すべての実施形態においてそのような分離を必要とすると理解されるべきではなく、記載されるプログラムコンポーネントおよびシステムは一般に単一のソフトウェア製品に統合され得るかまたは複数のソフトウェア製品にパッケージ化され得ることが理解されるべきである。

【0073】

20

主題の特定の実施形態が記載された。他の実施形態は以下の請求の範囲内にある。たとえば、請求項において記載されるアクションは、異なる順で実行され得、それでも望ましい結果を達成し得る。一例として、添付の図において示されるプロセスは、望ましい結果を達成するために、示された特定の順または連続する順であることを必ずしも必要としない。ある場合においては、マルチタスキングおよび並列処理が有利であり得る。

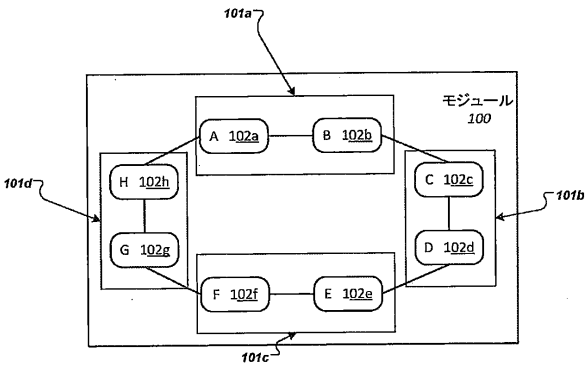
30

40

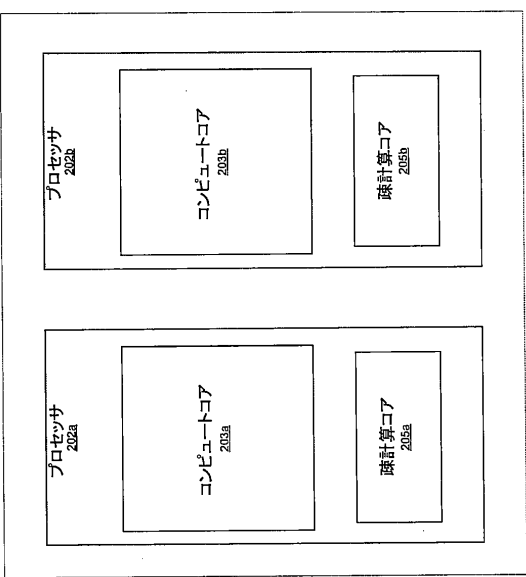
50

【図面】

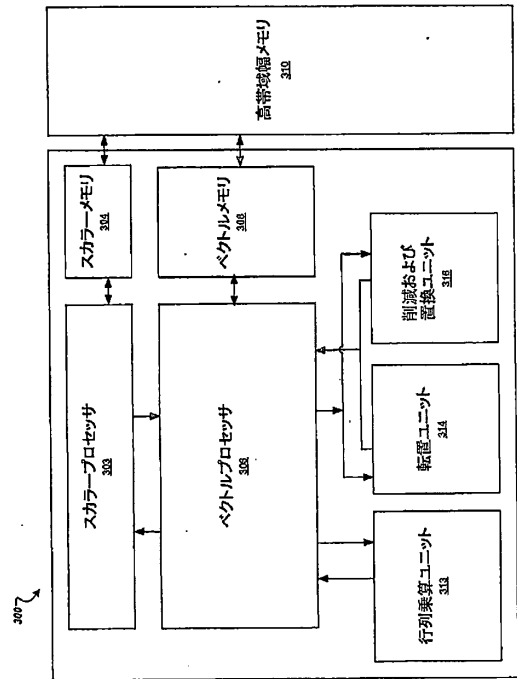
【図 1】



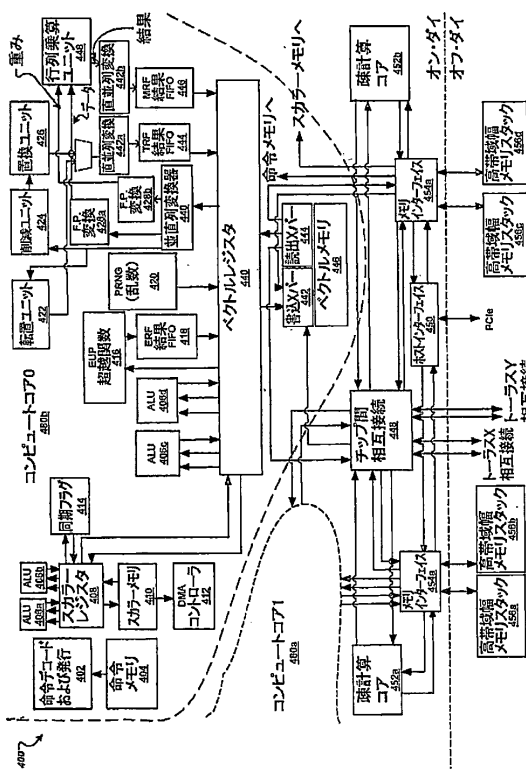
【図 2】



【図 3】



【図 4】



10

20

30

40

50

フロントページの続き

(51)国際特許分類

G 0 6 F 17/10 (2006.01)
G 0 6 N 3/063(2023.01)

F I

G 0 6 F 17/16 P
G 0 6 F 17/10 S
G 0 6 N 3/063

(72)発明者 テマム，オリビエ

アメリカ合衆国、9 4 0 4 3 カリフォルニア州、マウンテン・ビュー、アンフィシアター・パークウェイ、1 6 0 0

(72)発明者 フェルプス，アンドリュー・エバレット

アメリカ合衆国、9 4 0 4 3 カリフォルニア州、マウンテン・ビュー、アンフィシアター・パークウェイ、1 6 0 0

(72)発明者 ジョビー，ノーマン・ポール

アメリカ合衆国、9 4 0 4 3 カリフォルニア州、マウンテン・ビュー、アンフィシアター・パークウェイ、1 6 0 0

審査官 坂東 博司

(56)参考文献

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

米国特許出願公開第 2 0 1 9 / 0 3 5 4 8 6 2 (U S , A 1)

特開平 0 4 - 2 9 0 1 5 5 (J P , A)

特開 2 0 1 7 - 1 3 8 9 6 6 (J P , A)

下川 勝千，デジタルニューロコンピュータMULTINEURO，東芝レビュー 第 4 6 巻 第 1 2 号，第46巻 第12号，日本，株式会社東芝，1991年12月，pp.931~934，【ISSN】0372-0462

斎藤 康毅，ゼロから作るDeep Learning 初版，株式会社オライリー・ジャパン，第1版，日本，株式会社オライリー・ジャパン，2016年09月，pp.123-165，第1版

ISBN: 978-4-87311-758-4

(58)調査した分野 (Int.Cl.，D B 名)

G 0 6 F 9 / 3 8
G 0 6 F 9 / 3 0
G 0 6 F 1 5 / 1 7 3
G 0 6 F 1 5 / 8 0
G 0 6 F 1 7 / 1 6
G 0 6 F 1 7 / 1 0
G 0 6 N 3 / 0 6 3