

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2025-505291
(P2025-505291A)

(43)公表日 令和7年2月21日(2025.2.21)

(51)国際特許分類 F I テーマコード(参考)
G 0 6 F 17/10 (2006.01) G 0 6 F 17/10 A 5 B 0 5 6

審査請求 未請求 予備審査請求 未請求 (全29頁)

(21)出願番号	特願2024-548397(P2024-548397)	(71)出願人	523062486
(86)(22)出願日	令和5年2月13日(2023.2.13)		モフェット インターナショナル カンパ ニー, リミテッド
(85)翻訳文提出日	令和6年10月15日(2024.10.15)		中華人民共和国, 香港 9 9 9 0 7 7 , クーロン, サン ポー コン, ルーク ホ ップ ストリート 2 9 , ワン フェ イ ンダストリアル ビルディング, 1 1 / フロア, ルーム 8
(86)国際出願番号	PCT/CN2023/075661	(74)代理人	100099759
(87)国際公開番号	WO2023/155748		弁理士 青木 篤
(87)国際公開日	令和5年8月24日(2023.8.24)	(74)代理人	100123582
(31)優先権主張番号	17/673,490		弁理士 三橋 真二
(32)優先日	令和4年2月16日(2022.2.16)	(74)代理人	100112357
(33)優先権主張国・地域又は機関	米国(US)		弁理士 廣瀬 繁樹
(81)指定国・地域	AP(BW,CV,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	(74)代理人	100114018
	最終頁に続く		最終頁に続く

(54)【発明の名称】 スパースニューラルネットワークのための適応テンソル畳み込みカーネル

(57)【要約】

適応テンソル演算カーネルを使用したニューラルネットワーク計算の効率化のための、コンピュータ記憶媒体にエンコードされたコンピュータプログラムを含む方法、システム、および装置が提供される。第一に、適応テンソル演算カーネルは、並列処理用の処理要素(PE)アレイに重みと入力値を分配するために、入力テンソル/重みテンソルの異なる形状に応じて形状を変形することができる。テンソル演算カーネルの形状によっては、畳み込み演算を実行するために、クラスタ間またはクラスタ内の追加の加算器が必要になる場合がある。第二に、適応テンソル演算カーネルは、全ての種類の畳み込み計算に対応するために、1×1テンソル演算モードと3×3テンソル演算モードという2つの異なるテンソル演算モードをサポートすることができる。第三に、基盤となるPEアレイは、スパースニューラルネットワークの異なる圧縮率と疎粒度に対応するために、各PE内部バッファ(レジスタファイルなど)を異なるように構成することができる。

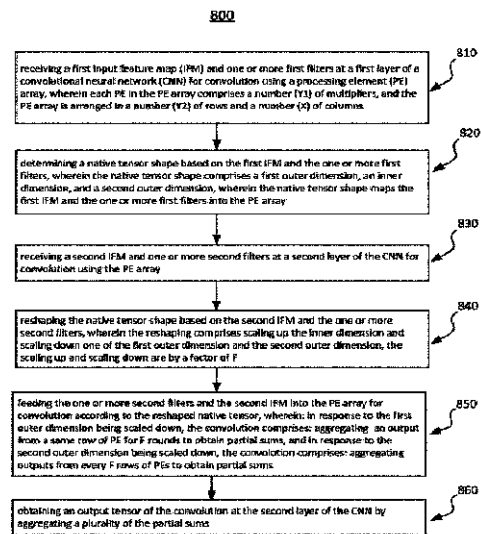


FIG. 8

【特許請求の範囲】

【請求項 1】

コンピュータで実行される方法であって、

処理要素 (PE) アレイを使用して畳み込みを行うための畳み込みニューラルネットワーク (CNN) の第 1 層において、第 1 入力特徴マップ (IFM) および 1 つまたは複数の第 1 フィルタを受信することであって、前記 PE アレイ内の各 PE は、(Y1) 個の乗算器を備え、前記 PE アレイは、(Y2) 個の行および (X) 個の列に配置される、受信することと、

前記第 1 IFM および前記 1 つまたは複数の第 1 フィルタに基づいて、ネイティブテンソル形状を決定することであって、前記ネイティブテンソル形状は、第 1 外部次元、内部次元、および第 2 外部次元を含み、前記ネイティブテンソル形状は、前記第 1 IFM および前記 1 つまたは複数の第 1 フィルタを前記 PE アレイにマッピングする、決定することと、

前記 PE アレイを使用して畳み込みを行う前記 CNN の第 2 層において、第 2 IFM および 1 つまたは複数の第 2 フィルタを受信することと、

前記第 2 IFM および前記 1 つまたは複数の第 2 フィルタに基づいて、前記ネイティブテンソル形状を変形することであって、前記変形は、前記内部次元の拡大と前記第 1 外部次元および前記第 2 外部次元のうち一方の縮小を含み、前記拡大および縮小は係数 F によって実行される、変形することと、

変形された前記ネイティブテンソルに従って、前記 1 つまたは複数の第 2 フィルタおよび前記第 2 IFM を、畳み込みのために前記 PE アレイに供給することであって、

前記第 1 外部次元が縮小された場合、PE の同一の行からの出力をフラウンドについて合計して部分和を求めることを含み、

前記第 2 外部次元が縮小された場合、前記畳み込みは、PE の全ての F 倍した行からの出力を合計して部分和を求めることを含む、供給することと、

複数の前記部分和を合計することにより、前記 CNN の前記第 2 層において、前記畳み込みの出力テンソルを取得することであって、Y1、Y2、X、および F は全て 1 より大きい整数である、取得することと、

を含む、方法。

【請求項 2】

前記 CNN の前記第 2 層は、前記 CNN の前記第 1 層の後にあり、前記第 2 IFM は、前記第 1 IFM よりも多くの入力チャンネルと、前記第 1 IFM よりも低い解像度と、を有する、請求項 1 に記載の方法。

【請求項 3】

前記 1 つまたは複数の第 2 フィルタの各々が、2 次元 (2D) カーネルの複数のチャンネルを備え、各 2D カーネルは、1 × 1 または 3 × 3 の次元を有する、請求項 1 に記載の方法。

【請求項 4】

前記変形されたネイティブテンソルに従って、前記 1 つまたは複数の第 2 フィルタを前記 PE アレイに供給することは、

前記変形されたネイティブテンソルの前記第 1 外部次元および前記内部次元に従って、前記 1 つまたは複数の第 2 フィルタを行列に変換することであって、前記 1 つまたは複数の第 2 フィルタにおける各 2D カーネルが前記 1 × 1 の次元を有する場合、前記行列の各行は、前記 1 つまたは複数の第 2 フィルタの異なる入力チャンネルからの重みを含む、変換することと、

複数の入力チャンネルが一度に同時に処理されるように、前記行列の各行の重みを PE の異なる列に分配することと、

を含む、請求項 3 に記載の方法。

【請求項 5】

10

20

30

40

50

前記変形されたネイティブテンソルに従って、前記 P E アレイに前記 1 つまたは複数の第 2 フィルタを供給することは、

前記変形されたネイティブテンソルの前記第 1 外部次元および前記内部次元に従って、前記 1 つまたは複数の第 2 フィルタを行列に変換することであって、前記 1 つまたは複数の第 2 フィルタにおける各 2 D カーネルが前記 3×3 の次元を有し、9 個の重みを含む場合、前記 9 個の重みは前記行列の同じ行に配置される、変換することと、

前記同じチャンネルからの前記重みが一度に同時に処理されるように、前記行列の前記同じ行から P E の異なる列に前記 9 個の重みを分配することと、
を含む、請求項 3 に記載の方法。

【請求項 6】

前記変形されたネイティブテンソルに従って、前記 I F M を前記 P E アレイへ供給することは、

前記変形されたネイティブテンソルの前記内部次元および前記第 2 外部次元に従って、前記 I F M を行列に変換することと、

前記行列の列に対応する前記 I F M の入力値を前記 P E の行のバッファに供給することと、

を含む、請求項 5 に記載の方法。

【請求項 7】

前記 1 つまたは複数のフィルタのチャンネルを複数のチャンネルグループに分割することであって、各チャンネルグループは 1 より大きい整数である定数のチャンネルを含む、分割することと、

前記各チャンネルグループ内の重みの固定割合が非ゼロとなるように、前記複数のチャンネルグループのそれぞれをプルーニングすることと、

をさらに含む、請求項 1 に記載の方法。

【請求項 8】

前記 P E アレイ内の各 P E に関連付けられたバッファの深さを決定することをさらに含み、

前記バッファの前記深さが前記定数より大きい場合、各 P E のプライベートメモリとして前記バッファを構成し、

前記バッファの前記深さが前記定数より小さい場合、前記 P E の前記バッファおよび隣接する P E の 1 つまたは複数のバッファを共有メモリとして結合する、

請求項 7 に記載の方法。

【請求項 9】

各 P E の前記プライベートメモリは、前記 P E 内の前記 (Y 1) 個の乗算器によって取得可能な入力値を記憶する、請求項 8 に記載の方法。

【請求項 10】

前記共有メモリが、前記 P E および前記 1 つまたは複数の近隣 P E 内の前記 (Y 1) 個の乗算器によって取得可能な入力値を記憶する、請求項 8 に記載の方法。

【請求項 11】

P E の各行は、各 P E 内の前記 (Y 1) 個の乗算器にそれぞれ対応する (Y 1) 個の加算器ツリーと結合されており、各 P E 内の各乗算器は、合計のために、対応する加算器ツリーに積算出力を送信する、請求項 1 に記載の方法。

【請求項 12】

前記 1 つまたは複数の第 2 フィルタの各々は複数の非ゼロ重みを含み、前記 1 つまたは複数の第 2 フィルタを前記 P E アレイに畳み込みのために供給することは、

各非ゼロ重みを、前記非ゼロ重みおよび対応するインデックスを含むインデックス - 値のペアとして、対応する P E の乗算器に供給することを含み、前記畳み込みは、

前記インデックスに従って、前記対応する P E のバッファから入力値を取得することと、

前記取得した値および前記非ゼロ重みを前記乗算器に送り、出力を取得することと、

10

20

30

40

50

前記対応する P E と同じ行にある他の P E の他の乗算器によって生成された出力と合計するために、対応する加算器ツリーに前記出力を送信することと、
を含む、請求項 1 に記載の方法。

【請求項 1 3】

各 P E 内の前記 (Y 1) 個の乗算器は、データを並列処理し、前記 P E アレイ内の P E は、データを並列処理する、請求項 1 に記載の方法。

【請求項 1 4】

1 つまたは複数のプロセッサと、1 つまたは複数のプロセッサに結合され、かつ、前記 1 つまたは複数のプロセッサによって実行可能な命令を備え、システムに演算を実行させるように構成された 1 つまたは複数の非一時的コンピュータ読み取り可能メモリと、を備えるシステムであって、

処理要素 (P E) アレイを使用して畳み込みを行うための畳み込みニューラルネットワーク (C N N) の第 1 層において、第 1 入力特徴マップ (I F M) および 1 つまたは複数の第 1 フィルタを受信することであって、前記 P E アレイ内の各 P E は、(Y 1) 個の乗算器を備え、前記 P E アレイは (Y 2) 個の行および (X) 個の列に配置される、受信することと、

前記第 1 I F M および前記 1 つまたは複数の第 1 フィルタに基づいて、ネイティブテンソル形状を決定することであって、前記ネイティブテンソル形状は、第 1 外部次元、内部次元、および第 2 外部次元を含み、前記ネイティブテンソル形状は、前記第 1 I F M および前記 1 つまたは複数の第 1 フィルタを前記 P E アレイにマッピングする、決定することと、

前記 P E アレイを使用して畳み込みを行う前記 C N N の第 2 層において、第 2 I F M および 1 つまたは複数の第 2 フィルタを受信することと、

前記第 2 I F M および前記 1 つまたは複数の第 2 フィルタに基づいて、前記ネイティブテンソル形状を変形することであって、前記変形は、前記内部次元の拡大と前記第 1 外部次元および第 2 外部次元のうち一方の縮小を含み、前記拡大および縮小は係数 F によって実行される、変形することと、

変形された前記ネイティブテンソルに従って、前記 1 つまたは複数の第 2 フィルタおよび前記第 2 I F M を、畳み込みのために前記 P E アレイに供給することであって、

前記第 1 外部次元が縮小された場合、前記畳み込みは、P E の同一の行からの出力を F ラウンドについて合計して部分和を求めることを含み、

前記第 2 外部次元が縮小された場合、前記畳み込みは、P E の全ての F 倍した行からの出力を合計して部分和を求めることを含む、

供給することと、

複数の前記部分和を合計することにより、前記 C N N の前記第 2 層において、前記畳み込みの出力テンソルを取得することであって、Y 1、Y 2、X、および F は全て 1 より大きい整数である、取得することと、

を含む、システム。

【請求項 1 5】

前記 C N N の前記第 2 層は、前記 C N N の前記第 1 層の後にあり、前記第 2 I F M は、前記第 1 I F M よりも多くの入力チャンネルと、前記第 1 I F M よりも低い解像度と、を有する、請求項 1 4 に記載のシステム。

【請求項 1 6】

前記演算は、さらに、

前記 1 つまたは複数のフィルタのチャンネルを複数のチャンネルグループに分割することであって、各チャンネルグループは 1 より大きい整数である定数のチャンネルを含む、分割することと、

前記複数のチャンネルグループのそれぞれにおいて、1 つのチャンネルのみが非ゼロの入力値を含み、前記各チャンネルグループの他のチャンネルは全てゼロを含むように、前記 1 つまたは複数のフィルタの各チャンネルをプルーニングすることと、

10

20

30

40

50

を含む、請求項 14 に記載のシステム。

【請求項 17】

前記演算は、

前記 P E アレイ内の各 P E に関連付けられたバッファの深さを決定することさらに含み、

前記バッファの前記深さが前記定数より大きい場合、各 P E のプライベートメモリとして前記バッファを構成し、

前記バッファの前記深さが前記定数より小さい場合、前記 P E の前記バッファおよび隣接する P E の 1 つまたは複数のバッファを共有メモリとして結合する、

請求項 16 に記載のシステム。

10

【請求項 18】

前記 1 つまたは複数の第 2 フィルタの各々は、2 次元 (2 D) カーネルの複数のチャンネルを含み、各 2 D カーネルは、 1×1 または 3×3 の次元を有する、請求項 14 に記載のシステム。

【請求項 19】

前記変形されたネイティブテンソルに従って、前記 1 つまたは複数の第 2 フィルタを前記 P E アレイに供給することは、

前記変形されたネイティブテンソルの前記第 1 外部次元および前記内部次元に従って、前記 1 つまたは複数の第 2 フィルタを行列に変換することであって、前記 1 つまたは複数の第 2 フィルタにおける各 2 D カーネルが前記 1×1 の次元を有する場合、前記行列の各行は、前記 1 つまたは複数の第 2 フィルタの異なる入力チャンネルからの重みを含む、変換することと、

20

複数の入力チャンネルが一度に同時に処理されるように、前記第 1 行列の各行の重みを P E の異なる列に分配することと、

を含む、請求項 18 に記載のシステム。

【請求項 20】

1 つまたは複数のプロセッサによって実行可能な命令を有するように構成された非一時的コンピュータ読み取り可能記憶媒体であって、前記 1 つまたは複数のプロセッサに、

処理素子 (P E) アレイを使用した畳み込みのための畳み込みニューラルネットワーク (CNN) の第 1 層において、第 1 入力特徴マップ (IFM) および 1 つまたは複数の第 1 フィルタを受信することであって、前記 P E アレイ内の各 P E は、(Y1) 個の乗算器を備え、前記 P E アレイは (Y2) 個の行および (X) 個の列に配置される、受信することと、

30

前記第 1 IFM および前記 1 つまたは複数の第 1 フィルタに基づいて、ネイティブテンソル形状を決定することであって、前記ネイティブテンソル形状は、第 1 外部次元、内部次元、および第 2 外部次元を含み、前記ネイティブテンソル形状は、前記第 1 IFM および前記 1 つまたは複数の第 1 フィルタを前記 P E アレイにマッピングする、決定することと、

前記 P E アレイを使用して畳み込みを行う前記 CNN の第 2 層において、第 2 IFM および 1 つまたは複数の第 2 フィルタを受信することと、

40

前記第 2 IFM および前記 1 つまたは複数の第 2 フィルタに基づいて、前記ネイティブテンソル形状を変形することであって、前記変形は前記内部次元の拡大と前記第 1 外部次元および第 2 外部次元のうち一方の縮小を含み、前記拡大および縮小は係数 F によって実行される、変形することと、

変形された前記ネイティブテンソルに従って、畳み込みを行うために、前記 1 つまたは複数の第 2 フィルタおよび前記第 2 IFM を前記 P E アレイに供給することであって、

前記第 1 外部次元が縮小された場合、前記畳み込みは、P E の同一の行からの出力を F ラウンドについて合計して部分和を求めることを含み、

前記第 2 外部次元が縮小された場合、前記畳み込みは、P E の全ての F 倍した行からの出力を合計して部分和を求めることを含む、

50

供給することと、

複数の前記部分和を合計することにより、前記CNNの前記第2層において前記畳み込みの出力テンソルを取得することであって、Y1、Y2、X、およびFは全て1より大きい整数である、取得することと、

を含む演算を実行させる、非一時的コンピュータ読み取り可能記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、概してニューラルネットワークの計算効率の向上に関し、特に、スパースニューラルネットワークの様々な入力テンソルの形状や演算に適応するために、ネイティブテンソル(native tensor)の次元および演算モードを動的に調整することに関する。

10

【背景技術】

【0002】

深層学習のPE(処理要素(processing element))アレイは、ほとんど全てがネイティブテンソルの次元および演算モードで定められており、通常はコンパイラに依存して、様々なネストループマッピングアプローチを使用して、様々なテンソル形状(例えば、入力特徴マップやフィルタ)や演算に対応する。PEアレイのネイティブテンソルの次元や演算モードと互換性のないテンソル形状や演算に対してPEアレイを使用して計算を行うのは明らかに非効率的である。柔軟性に欠けるPEアレイのネイティブテンソルの形状や演算モードでは、多数のゼロを持つテンソルを効率的に表現したり処理したりできないため、スパース特性を持つスパースニューラルネットワークではこの非互換性がさらに悪化する。

20

【発明の概要】

【0003】

本明細書の種々の実施形態には、ニューラルネットワークの計算において適応テンソル計算カーネルを使用するためのシステム、方法、および非一時的なコンピュータ読み取り可能媒体が含まれ得る。

【0004】

一態様によれば、ニューラルネットワークの計算における適応テンソル計算カーネルを使用するための方法は、処理要素(PE)アレイを使用した畳み込みのために、畳み込みニューラルネットワーク(CNN)の第1層において、第1入力特徴マップ(IFM)および1つまたは複数の第1フィルタを受信することであって、PEアレイ内の各PEは、(Y1)個の乗算器を備え、PEアレイは、(Y2)個の行および(X)個の列に配置される、受信することと、第1IFMおよび1つまたは複数の第1フィルタに基づいてネイティブテンソル形状を決定することであって、ネイティブテンソル形状は、第1外部次元(first outer dimension)、内部次元(inner dimension)、および第2外部次元(second outer dimension)を含み、ネイティブテンソル形状は、第1IFMおよび1つまたは複数の第1フィルタをPEアレイにマッピングする、決定することと、PEアレイを使用して畳み込みを行うCNNの第2層において、第2IFMおよび1つまたは複数の第2フィルタを受信することと、第2IFMおよび1つまたは複数の第2フィルタに基づいてネイティブテンソル形状を変形すること(reshaping)であって、その変形は、内部次元の拡大と、第1外部次元および第2外部次元のうちの一方の縮小とを含み、その拡大および縮小は係数Fによって実行される、変形することと、変形されたネイティブテンソルに従って畳み込みを行うために、1つまたは複数の第2フィルタおよび第2IFMをPEアレイに供給することであって、第1外部次元が縮小された場合、畳み込みは、PEの同一の行からの出力をFラウンドについて合計して部分和を求め、かつ、第2外部次元が縮小された場合、畳み込みは、PEの全てのF倍した行からの出力を合計して部分和を求める、供給することと、複数の部分和を合計することにより、CNNの2層目における畳み込みの出力テンソル

30

40

50

を取得することであって、 Y_1 、 Y_2 、 X 、および F は全て1より大きい整数である、取得することと、を含んでもよい。

【0005】

いくつかの実施形態において、CNNの第2層はCNNの第1層の後にあり、第2IFMは、第1IFMよりも多くの入力チャンネルと、第1IFMよりも低い解像度と、を有する。

【0006】

いくつかの実施態様において、1つまたは複数の第2フィルタの各々は、2次元(2D)カーネルの複数のチャンネルを備え、各2Dカーネルは、 1×1 (one by one) または 3×3 (three by three) の次元を有する。

10

【0007】

いくつかの実施態様において、変形されたネイティブテンソルに従ったPEアレイへの1つまたは複数の第2フィルタの供給は、変形されたネイティブテンソルの第1外部次元および内部次元に従って、1つまたは複数の第2フィルタを行列に変換することであって、1つまたは複数の第2フィルタの各2Dカーネルが 1×1 (one by one) の次元を有する場合、行列の各行は、1つまたは複数の第2フィルタの異なる入力チャンネルからの重みを含む、変換することと、複数の入力チャンネルが一度に同時に処理されるように、行列の各行の重みをPEの異なる列に分配することと、を含む。

【0008】

いくつかの実施形態において、変形されたネイティブテンソルに従ってPEアレイに1つまたは複数の第2フィルタを供給することは、変形されたネイティブテンソルの第1外部次元および内部次元に従って、1つまたは複数の第2フィルタを行列に変換することであって、1つまたは複数の第2フィルタの各2Dカーネルが 3×3 の次元を有し、9個の重みを含む場合、9個の重みが行列の同じ行に配置される、変換することと、同じチャンネルからの重みが一度に同時に処理されるように、行列の同じ行からPEの異なる列に9個の重みを分配することと、を含む。

20

【0009】

いくつかの実施例において、変形されたネイティブテンソルに従ったPEアレイへのIFMの供給は、変形されたネイティブテンソルの内部次元および第2外部次元に従って、IFMを行列に変換することと、行列の列に対応するIFMの入力値をPEの行のバッファに供給することと、を含む。

30

【0010】

いくつかの実施例において、本方法はさらに、1つまたは複数のフィルタのチャンネルを複数のチャンネルグループに分割することであって、各チャンネルグループは1より大きい整数である定数のチャンネルを含む、分割することと、複数のチャンネルグループの各々におけるいくつかのチャンネルのみが非ゼロの入力値を含み、各チャンネルグループにおける他のチャンネルは全てゼロを含むように、1つまたは複数のフィルタの各々をプルーニングすることと、を含む。プルーニング後、複数のチャンネルグループのそれぞれは、非ゼロの重みと同一の割合の重みを含む。

【0011】

いくつかの実施例において、本方法は、PEアレイ内の各PEに関連付けられたバッファの深さを決定することをさらに含んでもよく、バッファの深さが定数よりも大きい場合、各PEのプライベートメモリとしてバッファを構成し、バッファの深さが定数よりも小さい場合、PEのバッファと隣接するPEの1つまたは複数のバッファを共有メモリとして結合する。

40

【0012】

いくつかの実施例において、各PEのプライベートメモリは、PE内の(Y_1)個の乗算器によって取得可能な入力値を記憶する。

【0013】

いくつかの実施例において、共有メモリは、PEおよび1つまたは複数の隣接PE内の

50

(Y 1) 個の乗算器によって取得可能な入力値を記憶する。

【 0 0 1 4 】

いくつかの実施例において、 P E の各行は、各 P E 内の (Y 1) 個の乗算器にそれぞれ対応する (Y 1) 個の加算器ツリーと結合されており、各 P E 内の各乗算器は、合計するために、対応する加算器ツリーに積算出力を送信する。

【 0 0 1 5 】

いくつかの実施例において、 1 つまたは複数の第 2 フィルタの各々は複数の非ゼロ重みを含み、畳み込みのために 1 つまたは複数の第 2 フィルタを P E アレイに供給することは、各非ゼロ重みを、非ゼロ重みおよび対応するインデックスを含むインデックス - 値のペアとして、対応する P E の乗算器に供給することを含み、畳み込みは、インデックスに従って、対応する P E のバッファから入力値を取得することと、取得した値と非ゼロ重みを乗算器に送信して出力を取得することと、対応する P E と同じ行にある他の P E の他の乗算器によって生成された出力と合計するために、対応する加算器ツリーに出力を送信することと、を含む。

10

【 0 0 1 6 】

いくつかの実施例において、各 P E 内の (Y 1) 個の乗算器は、データを並列処理し、 P E アレイ内の P E はデータを並列処理する。

【 0 0 1 7 】

さらに他の一態様によれば、システムは、 1 つまたは複数のプロセッサと、 1 つまたは複数のプロセッサに結合され、 1 つまたは複数のプロセッサによって実行可能な命令で構成され、システムに本明細書に記載のいずれかの方法を実行させるように構成された、 1 つまたは複数の非一時的コンピュータ可読メモリと、を備えてよい。

20

【 0 0 1 8 】

さらに他の一態様によれば、非一時的なコンピュータ読み取り可能な記憶媒体は、 1 つまたは複数のプロセッサが本明細書に記載された方法のいずれかを実行するように、 1 つまたは複数のプロセッサによって実行可能な命令で構成されてよい。

【 0 0 1 9 】

本明細書に開示されているシステム、方法、および非一時的コンピュータ読み取り可能な媒体のこれらの機能およびその他の機能、ならびに構造の関連要素の動作方法および機能、ならびに一部の組み合わせおよび製造者の経済性は、以下の説明および添付の図面を参照した添付の特許請求の範囲を考慮することでより明らかにされ、これらは全て本明細書の一部を構成し、同様の参照番号は、様々な図における対応する部分を指す。ただし、図面は説明および例示のみを目的としており、本発明を限定する目的で作成されたものではないことを十分に理解されたい。

30

【 図面の簡単な説明 】

【 0 0 2 0 】

【 図 1 】 図 1 は、様々な実施形態による P E アレイにおけるニューラルネットワーク計算処理のための例示的なシステム図を示す。

【 図 2 】 図 2 は、様々な実施形態による P E アレイの例示的なアーキテクチャ図を示す。

【 図 3 】 図 3 は、様々な実施形態によるネイティブテンソル形状を使用する P E アレイにおける例示的なニューラルネットワーク計算を示す。

40

【 図 4 A 】 図 4 A は、様々な実施形態に従った適応テンソル形状を使用する P E アレイにおける例示的なニューラルネットワーク計算を示す。

【 図 4 B 】 図 4 B は、様々な実施形態に従った適応テンソル形状を使用するニューラルネットワーク計算のためのクラスタ間加算器を備える例示的な P E アレイを示す。

【 図 5 A 】 図 5 A は、様々な実施形態に従った適応テンソル形状を使用する P E アレイにおける他の例示的なニューラルネットワーク計算を示す。

【 図 5 B 】 図 5 B は、様々な実施形態に従った適応テンソル形状を使用するニューラルネットワーク計算のためのクラスタ内加算器を備えた他の例示的な P E アレイを示す。

【 図 6 A 】 図 6 A は、様々な実施形態に従った 1 x 1 テンソル演算モードを備えた例示的

50

なニューラルネットワーク計算を示す。

【図 6 B】図 6 B は、様々な実施形態に従った 3×3 テンソル演算モードによる例示的なニューラルネットワーク計算を示す。

【図 7】図 7 は、様々な実施形態に従った適応テンソル形状および PE アレイにおける 3×3 テンソル演算モードを使用するニューラルネットワーク計算の例示的な方法を説明する図である。

【図 8】図 8 は、様々な実施形態に従った適応テンソル形状を使用するニューラルネットワーク計算の例示的な方法を説明する図である。

【図 9】図 9 は、本明細書で説明される実施形態のいずれかが実装され得る例示的なコンピュータシステムを示す図である。

【発明を実施するための形態】

【0021】

本明細書で説明する実施形態は、適応テンソルの形状および演算モードを使用する PE アレイにおけるニューラルネットワーク計算のための方法、システム、装置を提供する。以下の説明では、適応テンソル計算カーネルは、異なる形状の入力特徴マップ (IMF) および重みテンソル (例えば、フィルタ) を扱うために、複数のネイティブテンソルの次元および演算モードを持つものとして説明される。入力および出力テンソルの形状および演算モードに従って、適応テンソル計算カーネル (適応ネイティブテンソルとも称する) の次元および演算モードは、PE アレイの基盤となるハードウェア資源を並列処理に完全に利用するために、動的に調整することができる。

【0022】

これらの適応テンソル計算カーネルは、3つの技術的解決策を提供することにより、ニューラルネットワーク計算における技術的課題 (背景のセクションで説明) に対処する。まず、適応テンソル計算カーネルは、入力テンソル / 重みテンソルの様々な形状に適応して形状を調整することができる。入力テンソル / 重みテンソルの様々な形状は、様々なニューラルネットワーク間に渡って存在するだけでなく、同じニューラルネットワークパイプライン内でも存在する場合がある。例えば、ニューラルネットワークの最初の数層を処理する際には、通常、テンソルは高解像度 (高さおよび幅が大きい) ではあるが、入力および出力チャンネルは少なくなるように構成され、ニューラルネットワークの最後の数層を処理する際には、テンソルは低解像度 (高さおよび幅が小さい) ではあるが、入力および出力チャンネルは多くなるように構成される。これは、ニューラルネットワークの最初の数層では、入力特徴マップからの特徴抽出に重点が置かれるのに対し、ニューラルネットワークの最後の数層では、抽出された特徴間の根本的な相関関係の学習に重点が置かれるためであると考えられる。

【0023】

第二に、適応テンソル計算カーネルは、 1×1 テンソル演算モードと 3×3 テンソル演算モードという2つの異なるテンソル演算モードをサポートすることができる。行列の乗算が関わるこれらのニューラルネットワーク層では、 1×1 カーネルの畳み込み (例えば、重みテンソル内の各カーネルが 1×1 の形状を持つ) は 1×1 テンソル演算モードにマッピングされ、その他の畳み込み (3×3 、 5×5 、 7×7 など) は 3×3 テンソル演算モードにマッピングすることができる。これらの異なるテンソル演算モードは、実行時に重みテンソルの形状に基づいて動的に決定することができる。

【0024】

第三に、基盤となる PE アレイは、スパースニューラルネットワークの異なる圧縮率と疎粒度 (sparsity granularities) をサポートするために、各 PE 内部バッファ (例えば、レジスタファイル) を異なるように構成することができる。スパースニューラルネットワークがより細かい粒度でプルーニングされる場合 (例えば、少数の入力チャンネルから1つまたは複数の非ゼロ入力チャンネルを選択し、その少数の数が閾値よりも小さい場合)、各 PE 内のレジスタファイルはプライベートメモリとして構成することができる (例えば、対応する PE のみが使用する)。スパースニューラルネット

10

20

30

40

50

トワークが粗い粒度でブルーニングされる場合（例えば、多数の入力チャンネルから1つまたは複数の非ゼロ入力チャンネルを選択し、その多数の数が閾値よりも大きい場合）、隣接するPE内の複数のレジスタファイルを、隣接するPEで共有されるマルチポートメモリとして構成することができる。

【0025】

以下の説明において、図面を参照しながら、本発明の特定の非限定的な実施態様について説明する。本明細書で開示する任意の実施態様の特定の特征および態様は、本明細書で開示する任意の他の実施態様の特定の特征および態様と併用および/または組み合わせることができる。また、このような実施態様は例示的なものであり、あくまで、本願発明の範囲内の少数の実施態様を例示するものであることを理解されたい。当業者にとって明らか

10

【0026】

図1は、様々な実施形態に従ったPEアレイにおけるニューラルネットワーク計算処理のための例示的なシステム図を示す。図1の図は、PEアレイを用いてパイプラインで実行される典型的なニューラルネットワーク計算のワークフローを示す。本開示で説明される実施形態は、図1におけるニューラルネットワーク計算の一部として、または他の適切な環境で実装されてもよい。

【0027】

ニューラルネットワーク（例えば、CNN）内の所定の（例えば、畳み込み）層において、1つまたは複数の入力特徴マップ（IFM）120は、入力ソース（例えば、例えば入力画像など）または前の層（例えば、前の層から出力されたテンソルなど）から取得され、1つまたは複数の重みテンソル110が、様々な特徴を抽出するためにIFMを畳み込むために使用されてよい。畳み込み処理は、PEアレイ160と呼ばれる処理要素（PE）のアレイ内で並列に実行されてもよい。各PEは、処理能力および記憶容量（例えば、バッファまたはキャッシュ）を有するプロセッサを指してもよい。PEは、相互接続配線を有するPEアレイ内で特定の方法で配置されてもよく、実行時に動的に再配置されなくてもよい。PEアレイは再利用され、ニューラルネットワークの異なる層において、または異なるニューラルネットワークおよび異なるユースケースに渡って演算に参与してよい。PEアレイにおける固定された内部PE配置と、潜在的に無限の種類（IFMおよび/または重みテンソルにおける）テンソル形状との間の非互換性は、通常、非効率的な資源利用と最適ではない並列処理につながる。

20

【0028】

図1を参照すると、いくつかの実施形態において、IFM120はIFMキャッシュ140に記憶され、重みテンソル110はPEアレイ160の演算のために重みキャッシュ130に記憶されてもよい。PEアレイ160は、PEの行列（例えば、 $X \times Y$ ）を含んでもよく、各PEは、並列処理用の複数の乗算器を含んでもよい。いくつかの実施形態において、IFMキャッシュ140からの各IFMは、PEアレイ160における計算を容易にするために、行列変換層150を通過してもよい。行列変換は、im2colツールを使用して、オリジナルのHWCフォーマット（Hは高さ、Wは重さ、Cはチャンネルを表す）からRSCフォーマット（Rは行、Sは列、Cはチャンネルを表す）にIFMを変形するToepplitz行列変換を含んでもよく、ここでRSCフォーマットは重さテンソルの形状に基づいて決定される。ここで、変換とは、IFM内の入力値を複製して配置し、変換されたIFMを形成することであり、これにより、変換されたIFMと重みテンソルとの間の行列乗算が、PEアレイ160内のPE間の依存性を最小限に抑えつつ、並列でPEアレイ160内で実行される。いくつかの実施形態において、PEアレイ160における並列畳み込みの各ラウンドは、複数の部分和を生成する場合があります、その部分

30

40

50

【 0 0 2 9 】

図 2 は、様々な実施形態による P E アレイの例示的なアーキテクチャ図を示す。図 2 の P E アレイにおける P E の配列は例示的なものであり、ユースケースに応じて他の方法で実装することもできる。

【 0 0 3 0 】

図 2 の左部分に示されているように、P E アレイ 2 0 0 は、P E の行列を含んでもよい。図 2 の右部分に示されているように、各 P E 2 4 0 は、複数の乗算器 (M U L ゲート) を含んでもよい。各 P E 2 4 0 内の乗算器は並列に演算してもよく、P E アレイ 2 2 0 内の P E は並列に演算してもよい。参照を容易にするため、以下の説明では、P E アレイ 2 0 0 内の P E の列 2 2 0 の数を X 、P E アレイ 2 0 0 内の P E の行 2 1 0 の数を Y_2 、各 P E 2 4 0 内の乗算器の数を Y_1 と表記する。P E 2 1 0 の各行は P E クラスタと呼ばれ、各 P E クラスタは、P E クラスタ内の乗算器によって生成された部分和を合計するために、 Y_1 個の加算器ツリー (Y_1 加算器ツリー) 2 3 0 に結合されていてもよい。即ち、P E クラスタ内の各 P E 2 4 0 における第 1 乗算器は、合計用の第 1 加算器ツリー 2 3 0 に結合され、P E クラスタ内の各 P E 2 4 0 における第 2 乗算器は、合計用の第 2 加算器ツリー 2 3 0 に結合され、以下同様である。全ての P E クラスタにわたる加算器ツリー 2 3 0 (合計 $Y_1 \times Y_2$ 加算器ツリー) を合計した結果は、合計のために加算器 2 5 0 に供給されてもよい。加算器 2 5 0 は、ネットワークオンチップ (N O C) サブシステムの一部である数値の加算を行うデジタル回路を意味する場合がある。

10

【 0 0 3 1 】

いくつかの実施形態において、P E アレイ 2 0 0 が重みを P E にブロードキャストしてもよい。スパースニューラルネットワークの場合、重みの大部分はゼロであり、従って P E にブロードキャストされる重みは全て非ゼロの重みである。非ゼロの重みは重みテンソル内の任意の位置から生じる可能性があるため、ブロードキャストされる各重みは重み値だけでなく、重み値の位置情報を示すインデックス、即ち (インデックス、重み値) のようなインデックス - 値のペアも含む可能性がある。インデックスに基づいて、各 P E 2 4 0 は I F M から対応する入力値を取得し、重み値との乗算を行うことができる。乗算結果は対応する加算器ツリーに供給することができる。図 1 に示すように、第 1 乗算器 M U L 1 は、(インデックス 1、値 1) の形式で重みを受信し、インデックス 1 に基づいて (I F M を記憶する) I B U F 2 6 0 から入力値 I F M 1 を取得し、入力値 I F M 1 と値 1 に基づいて乗算を行い、合計のために結果を加算器ツリー 1 (例えば、P E が位置する P E クラスタの Y_1 加算器ツリー 2 3 0 の第 1 加算器ツリー) に送信する。

20

30

【 0 0 3 2 】

図 3 は、様々な実施態様に従ったネイティブテンソル形状を使用する P E アレイにおける例示的なニューラルネットワーク計算を示す。図 3 の例示的な計算では、変換された重みテンソル A (3 1 0) と変換された I F M テンソル B (3 2 0) との間の行列乗算が関与し、これにより出力特徴マップ (O F M) テンソル C (3 3 0) が生成される。行列乗算では、 X および Y の次元を有する P E アレイ 3 4 0 に対応するネイティブテンソル形状が使用される。

【 0 0 3 3 】

いくつかの実施形態において、変換された重みテンソル A (3 1 0) は、R S C フォーマット (3 次元) の全ての重みテンソルを $m' \times k'$ と表記される 2 次元行列に統合することによって得られ (例えば、異なるチャンネルからの重みが同じチャンネルに再配置される)、ここで、 m' は重みテンソルの数 (通常は K と表記) によって決定される出力チャンネルの数であり、 k' は各重みテンソルの R、S、C 次元の積である (R と S は重みテンソル内の各カーネルの次元を指し、C は入力チャンネルの数を指す)。

40

【 0 0 3 4 】

いくつかの実施形態において、変換された I F M テンソル B (3 2 0) は、R S C 形式に基づく H W C 形式 (3 次元) の全ての I F M を、 $k' \times n'$ と表記される 2 次元行列に統合することによって取得することができ、ここで、 k' は依然として各重みテンソルの

50

R、S、およびC次元の積であり、 n' はIFMのHおよびW次元の積である（Hは高さ、Wは幅を指す）。行列 $m' * k'$ （重みテンソルA（3 1 0））と行列 $k' * n'$ （IFMB（3 2 0））の行列積は、 $m' * n'$ の行列としてOFMテンソルC（3 3 0）を生成することができる。

【0035】

上記の変換により、変換された重みテンソルA（3 1 0）および変換されたIFMテンソルB（3 2 0）は、並列処理のためにPEアレイ340内のPEにマッピングされ得る。PEアレイ340が Y_2 個の行のPE、 X 個の列のPEを含み、各PEが Y_1 個の乗算器を含むと仮定すると、テンソルAおよびBはPEアレイ340に以下のようにマッピングされ得る。即ち、テンソルA（3 1 0）およびテンソルB（3 2 0）の内部次元 k' はPEアレイ340の X （行）次元にマッピングされ、即ち、 $X = k' = R * S * C$ となり、テンソルAおよびテンソルBの外部次元 $m' * n'$ の乗算はPEアレイ340の Y （列）次元にマッピングされる。PEの各列には $Y_1 * Y_2$ 乗算器が含まれるため、上記のマッピングにより、 $Y = m' * n' = K * H * W$ の乗算は $Y_1 * Y_2$ 乗算器によって並列に処理されることになる。例えば、各PE内では、1つの乗算器が同じ出力チャンネルに対応する重み（例えば、全ての重みテンソルに渡る同じ位置からの重み）を処理し、即ち、 $Y_1 = K = m'$ であり、PEの各列は $H * W$ の重みを並列に処理し、即ち、 $Y_2 = H * W = n'$ である。

【0036】

上記の説明では、ワークロード（例えば、乗算のための重みと対応する入力値のペア）をPEアレイ340内のPEにマッピングするために、ネイティブテンソル形状 $m' * k' * n'$ が確定され、ここで、 $X = k'$ 、 $Y_1 * Y_2 = m' * n'$ である。即ち、PEアレイ内のPEのレイアウトに基づいて、ネイティブテンソル形状が決定される。PEアレイのレイアウトが確定すると、ネイティブテンソルの形状も確定する。全ての入力テンソル（IFMやフィルタ/重みテンソルなど）は、確定したネイティブテンソルの形状に従って変換されなければならない。しかしながら、実際のアプリケーションにおける入力テンソルは形状が様々であり、PEアレイ340におけるPEのレイアウトではなく、入力テンソルの形状に基づいて変換を行うと、最適な並列処理が可能になる。多くの場合、PEレイアウトに基づいて決定された確定したネイティブテンソル形状を使用する変換は、ワークロードをPEにマッピングするとしても、特定のPE間でいくつかの連続した依存関係（例えば、あるPEが他のPEの出力待ちになる）を引き起こす可能性がある。以下の説明では、IFMおよびフィルタの次元に基づいて決定され、同時にワークロードをPEにマッピングして並列性を最大化する適応ネイティブテンソル形状を使用する変換について説明する。

【0037】

図4Aは、様々な実施形態に従った適応テンソル形状を用いたPEアレイにおける例示的なニューラルネットワーク計算を示す。上述の通り（図3）、確定したネイティブテンソル形状 $m' * k' * n'$ （即ち、行列AおよびBをカバーする）を使用して入力テンソル（IFM）および重みテンソルを行列A（4 1 0）および行列B（4 2 0）に変換できる場合、変換されたテンソルは対応するPEアレイに分配することができる。しかしながら、実際のアプリケーションでは、乗算用のIFMや重みテンソル（例えば、CNNの異なる層における、異なるレベルのスプース化を経たテンソル）は、PEアレイに完全にマッピングできない様々な形状を有する可能性がある。確定したネイティブテンソル形状を使用してテンソルを強制的に変換すると、いくつかのPEがアイドル状態になるか、計算中にシーケンシャルな依存関係が発生する可能性がある。例えば、同じ畳み込みニューラルネットワーク（CNN）内でも、最初の数層のCNNのテンソルは高解像度（例えば、 $H * W = 64$ ）で入力チャンネル数が少ない（ $C = 16$ ）場合があり、最後の数層のCNNのテンソルは低解像度（例えば、 $H * W = 16$ ）で入力チャンネル数が多い（ $C = 64$ ）場合がある。ここで、「少ない」と「多い」は閾値に基づいて決定される。即ち、同じCNN内の畳み込み処理でも、異なる形状のテンソルが現れる可能性がある。

【 0 0 3 8 】

いくつかの実施形態において、入力テンソルおよび重みテンソルの変化する形状に対応するために、ネイティブテンソルの形状が動的に変形される場合がある。例えば、入力テンソルが、入力チャンネル数が少ない（例えば、最初の数層のCNN層において）高解像度（画素数が多い）から、入力チャンネル数が多い（例えば、最後の数層のCNN層において）低解像度に変化する場合、ネイティブテンソル形状は、それに応じて変形することができる。いくつかの実施形態において、ネイティブテンソル形状は、第1外部次元、内部次元、および第2外部次元として示される3つの次元を有する。最初の2つの次元（第1外部次元および内部次元）は、重みテンソルを行列に変換するために使用することができる。最後の2つの次元（内部次元および第2外部次元）は、IFMを行列に変換するために使用することができる。変換された行列は、重みおよび入力値がPEアレイにどのようにマッピングされるか（例えば、最適な並列性を達成するために重みおよび入力値をどのように分配するか）についての指針を提供することができる。

10

【 0 0 3 9 】

いくつかの実施形態において、以前のテンソルがマッピングおよび変換にネイティブテンソル形状 $m' * k' * n'$ を使用し、入力テンソルが以前のテンソルと比較してより低い解像度でより多くの入力チャンネルを有する場合、ネイティブテンソル形状の3つの次元は、 $m' * (F * k') * (n' / F)$ に変形することができ、ここで、 F は1より大きい整数であり、スケーリング係数を表し、最初の2つの次元（即ち、第1外部次元 m' および内部次元 $F * k'$ ）は重みテンソル行列 $4 2 0$ を表し、次の2つの次元（即ち、内部次元 $F * k'$ および第2外部次元 n' / F ）は畳み込みのためのIFMテンソル行列 $4 2 2$ を表す。即ち、ネイティブテンソル形状は、その内部次元を係数 F で拡大し、第2外部次元（IFMテンソル行列に対応する）を係数 F で縮小することができる。この変形方法は、以下の説明では $k' \& n'$ 変形と呼ばれることがある。いくつかの実施形態では、 F は2、4、8などのいずれかである。

20

【 0 0 4 0 】

いくつかの実施形態において、変形されたテンソル形状の内部次元 $F * k'$ は、重みテンソル行列 $4 2 0$ およびIFMテンソル行列 $4 2 2$ で共有され（例えば、同じ内部次元を有する）、PEアレイにおけるPEの列数に対応し、第1外部次元（例えば、重みテンソル行列 $4 2 0$ の外部次元 m' ）は、各PE内の乗算器の数に対応し、第2外部次元（例えば、IFMテンソル行列 $4 2 2$ の外部次元 n' / F ）は、PEアレイ内のPEの行の数に対応する。ここで、「対応する」とは、変換されたテンソル行列の重みおよび入力値がPEアレイにどのように分配されるかを指示するマッピング関係を指す。例えば、重みテンソル行列 $4 2 0$ の各外部次元（例えば、各列）における重みは、並列処理のために単一PE内の乗算器に分配されてもよく、IFMテンソル行列 $4 2 2$ の各外部次元（例えば、各列）における入力値は、PEアレイ内のPEの行に渡って分配されてもよい。

30

【 0 0 4 1 】

図4Aに示されているように、このように変形されたネイティブテンソル形状により、重みテンソル行列 $4 1 0$ は、その内部次元を F 倍に拡大し、その外部次元 m' を同一に保つよう変形され、それによって新たな重みテンソル行列 $4 2 0$ が形成される。即ち、重みテンソル行列の内部次元は、 $k' = R * S * C$ （例えば、 $4 1 0$ における行列A）から $F * k' = R * S * (F * C)$ （例えば、 $4 2 0$ における行列A）に変化し、その結果、新しい行列 $4 2 0$ はより多くの入力チャンネル（ C から $F * C$ ）をサポートできる。同様に、IFM行列 $4 1 2$ は、係数 F によってその外部次元を縮小し、重みテンソル行列 $4 2 0$ の縮小された内部次元と同じ方法でその内部次元を拡大し、それによって新しいIFMテンソル行列 $4 2 2$ を形成する。即ち、IFMテンソル行列の内部次元は、 $k' = R * S * C$ （例えば、 $4 1 2$ における行列B）から $F * k' = R * S * (F * C)$ （例えば、 $4 2 2$ における行列B）に変化し、IFMテンソル行列の外部次元は n' （例えば、 $4 1 2$ における行列B）から n' / F （例えば、 $4 2 2$ における行列B）に変化し、従って、新たな行列 $4 2 2$ はより少ない画素をサポートする可能性がある。従って、新たな行列 $4 2$

40

50

0 および 4 2 2 は、より少ない入力チャンネルで低解像度を有する最後のいくつかの CNN 層からのテンソルを表すのに、より適している。いくつかの実施形態において、「最初のいくつかの CNN 層」および「最後のいくつかの CNN 層」は、それぞれ、CNN 構造の最初から数えて第 1 の数の CNN 層、および CNN 構造の最後から数えて第 2 の数の CNN 層を指す場合がある。

【0042】

例として、最初の新たな CNN 層のテンソルは、高解像度 $H * W = 64$ で入力チャンネル数が少ない $C = 16$ であってよい。ここで、「少ない数」とは閾値よりも小さい数を指し、閾値は基盤となる PE アレイに従って構成されたコンパイラによって決定することができる。最初の数層の CNN 層からのこれらのテンソルのネイティブテンソル形状は、 $m' = K = 16$ 、 $k' = 1 * 1 * 16$ 、 $n' = 64$ の形状を有してよく、畳み込みが $1 * 1$ カーネルに基づいていると仮定する。畳み込みが最後の数層の CNN 層に進むと、テンソルは低解像度の $H * W = 16$ になってよく、入力チャンネル $C = 64$ (例えば、閾値よりも大きい) がより多くなると、ネイティブテンソルの形状は $m' = K = 16$ 、 $k' = 1 * 1 * 64$ 、 $n' = 16$ というように変形されてよい。

【0043】

上記の k' & n' 変形を使用してテンソルを変換した後、変換されたテンソル行列 4 2 0 および 4 2 2 を PE アレイに分散させて並列処理を行うことができる。図 4 B は、様々な実施形態に従って k' & n' 変形に基づく適応テンソル形状を使用するニューラルネットワーク計算のためのクラスタ間加算器を備えた例示的な PE アレイを示す。図 4 B に示される PE アレイを使用する並列処理スキームは、図 4 A に記述されたネイティブテンソルの k' & n' 変形に対応してよい。一貫性を保つため、PE アレイは、 Y 2 個の行と X 個の列の PE を有し、各 PE は Y 1 個の乗算器を有すると仮定する。

【0044】

k' & n' を変形することによって、重みテンソル行列と IFM テンソル行列の内部次元は F 倍に拡大され、IFM テンソル行列の外部次元は F 分の 1 に縮小される。PE アレイへの重みと入力値の分布は、重みテンソル行列の同じ行 (即ち、拡大された内部次元 / 行に沿って) からの重みと、IFM テンソル行列の同じ列 (即ち、同じく拡大された内部次元 / 列に沿って) からの入力値が、行ごとに PE に割り当てられるようにしてもよい。これは、これらの重みと入力値のペアが、PE の F 個の行に渡って分散される可能性があることを意味する。従って、PE アレイは、クラスタ間 (即ち、PE クラスタ間または行間) の加算器 4 0 0 を有し、畳み込み処理の部分和を求めるために、PE の各行によって生成された出力を加算する。各クラスタ間加算器 4 0 0 は、PE の F 個の行の Y 1 加算器ツリーからの出力を Y 1 部分和として合計することができる。これらの部分和は、畳み込みの結果として出力テンソルを構築するために合計される。このプロセスでは、部分和の合計は Y 1 * (Y 2 / F) となり、これは、出力チャンネル数 (例えば、畳み込みプロセスの出力テンソルのチャンネル数) が Y 1 であり、出力画素数が Y 2 / $F = H * W / F$ であることを意味する。

【0045】

図 5 A は、様々な実施形態に従った適応テンソル形状を使用する PE アレイにおける他の例示的なニューラルネットワーク計算を示す。上述の k' & n' の変形と比較すると、ネイティブテンソル形状もまた、重みテンソルのスパース度に基づいて動的に変形することができる。多くの実用的な用途では、畳み込み処理における重みテンソルをプルーニングまたはスパース化して、計算効率を向上させ、ニューラルネットワークのフットプリントを削減することができる。慎重にプルーニングされた重みテンソルは、ゼロ値の重みを導入することで特徴抽出の精度を犠牲にすることなく畳み込み速度を向上させ、その結果、総計算数を削減することができる (例えば、ゼロ値の重みはスキップされる)。いくつかの実施形態では、重みテンソルのプルーニングは、重みテンソルのチャンネル (フィルタとも呼ばれる) を複数のチャンネルグループに分割し、全てのチャンネルグループが同じ数のチャンネルを持つようにすることと、各チャンネルグループの少数のチャンネルのみ

を非ゼロ入力チャンネル（例えば、非ゼロ重み）として保持し、そのチャンネルグループ内の他の全てのチャンネルをゼロ（例えば、全ての重みがゼロ）にすることと、を含んでもよい。ブルーニング処理の後に、各チャンネルグループは、同じ割合の非ゼロ重みを含む。いくつかの実施形態において、ブルーニングのためのチャンネルグループのサイズ（例えば、各チャンネルグループ内のチャンネルの数）は、重みテンソル（フィルタ）の数、即ち、出力チャンネルの数に基づいて決定されてもよい。一般的に、重みテンソルのブルーニングは、2つのレベルに分類されてもよく、即ち、出力チャンネル数（例えば、重みテンソル数）が第1閾値よりも大きく、非ゼロ入力チャンネル数が第2閾値よりも小さい、重みのスパース性が高い場合と、出力チャンネル数（例えば、重みテンソル数）が第1閾値よりも小さく、非ゼロ入力チャンネル数が第2閾値よりも大きい、重みのスパース性が低い場合である。例えば、重みのスパース性が高い場合（16倍）のネイティブテンソル形状は、 $m' = K = 16$ （例えば、16個の重みテンソルまたはフィルタ）、 $k' = 3 * 3 * 4$ （例えば、各カーネルは $3 * 3$ 、1つのフィルタ内の非ゼロチャンネル数は4）、 $n' = 64$ となり、一方、重みのスパース性が低い場合（4倍）のネイティブテンソル形状は、 $m' = K = 4$ （例えば、4つの重みテンソルまたはフィルタ）、 $k' = 3 * 3 * 16$ （例えば、各カーネルは $3 * 3$ 、1つのフィルタ内の非ゼロチャンネル数は16）、 $n' = 64$ となる。

10

【0046】

いくつかの実施形態において、重みのスパース性が高から低に変化する場合、ネイティブテンソルの形状（`first_outer_dimension * inner_dimension * second_outer_dimension`と表記される）は、内部次元（重みテンソル行列とIFMテンソル行列とで共有される）を係数Fで拡大し、第1外部次元（重みテンソル行列に対応する）を係数Fで縮小することによって変形されてよい。図5Aに示されているように、元のネイティブテンソル形状 $m' * k' * n'$ は、 $(m' / F) * (F * k') * n'$ となり、重みテンソル行列510は $m' * k'$ から、次元 $(m' / F) * (F * k')$ の変形されたテンソル行列520に変化し、IFMテンソル行列512は $k' * n'$ から $(F * k') * n'$ の次元を持つ変形されたIFM行列522に変化する。この変形方法は、以下の説明では、 $k' \& m'$ 変形と呼ばれる場合がある。係数Fで拡大された内部次元は、より多くの入力チャンネル（Cから $F * C$ ）をサポートすることを示し、重みテンソル行列の縮小された外部次元は、より少ない出力チャンネル（Kから K / F ）をサポートすることを示す。

20

30

【0047】

上記の $k' \& m'$ 変形を使用してテンソルを変換した後、変換されたテンソル行列520および522をPEアレイに分散して並列処理を行うことができる。図5Bは、様々な実施形態に従った $k' \& m'$ 変形に基づく適応テンソル形状を使用するニューラルネットワーク計算のためのクラスタ間加算器を備えた他の例示的なPEアレイを示す。図5Bに示されたPEアレイを使用する並列処理スキームは、図5Aに説明されたネイティブテンソルの $k' \& m'$ 変形に対応し得る。

【0048】

一貫性を保つため、PEアレイは、Y2個の行とX個の列のPEを有し、各PEはY1個の乗算器を有すると仮定する。さらに、重み行列520とIFM行列522は、PEアレイ内のPEの(X)個の列に対応する同一の内部次元を有し、重み行列520の外部次元はPEアレイにおける各PE内の(Y1)個の乗算器に対応し、IFM行列522の外部次元はPEアレイ内のPEの(Y2)個の行に対応する。

40

【0049】

変形されたネイティブテンソルは、（重み行列520の外部次元に対応する）第1外部次元を m' / F として有するため、重みテンソル行列の各列の重みは、各PE内の $Y1 / F$ 乗算器に供給することができる。PEアレイから部分和を求めるために、Fラウンド分のY1加算器ツリーからの出力を記憶し合計するクラスタ内加算器500が実装されてもよい。ここで、「ラウンド」とは、PE内の乗算器を使用して乗算を実行するサイクルを

50

指す。各ラウンドの間、 Y_1 / F 乗算器の出力は、一時的に1つのクラスタ内加算器 500 に記憶されてもよい。Fラウンド後、クラスタ内加算器 500 には、 Y_1 加算器ツリーから収集された $F * Y_1 / F = Y_1$ 個の部分和が格納されていてもよい。これらの部分和は、畳み込みの結果として出力テンソルを構築するために合計されてもよい。このプロセスにおける部分和の合計は、 $(* Y_1 / F) * Y_2$ となり、これは出力チャンネル数（例えば、畳み込みプロセスの出力テンソルのチャンネル数）が Y_1 / F であり、出力画素数が $Y_2 = H * W$ であることを意味する。

【0050】

畳み込みニューラルネットワークの分野では、重みテンソルは、複数の2Dカーネルを含む3Dフィルタと呼ばれる場合がある。各3Dフィルタ内の2Dカーネルの数は、フィルタ内のチャンネル数と呼ばれる場合があり、各2Dカーネルは、 1×1 または 3×3 の行列であってよい。図6Aは、様々な実施形態に従った 1×1 (one by one) テンソル演算モード（即ち、 1×1 (one by one) カーネルを使用）による例示的なニューラルネットワーク計算を示し、図6Bは、様々な実施形態に従った 3×3 テンソル演算モード（即ち、 3×3 カーネルを使用）による例示的なニューラルネットワーク計算を示す。

10

【0051】

いくつかの実施形態において、一般的な行列乗算 (GEMM) および 1×1 の畳み込み演算は、 1×1 の演算モードにマッピングされてよい（例えば、 1×1 カーネルを使用する）。図6Aに示されているように、異なる入力チャンネル（または、スパース化された入力テンソル用のチャンネルグループ）からの2Dカーネル（即ち、重み）は、複数の入力チャンネルが一度に同時に処理されるように、PEの異なる列に配置されてもよく、同じ入力チャンネルからの2Dカーネルは、乗算器が一度に同時に複数の出力チャンネルを処理できるように、1つのPE内の複数の乗算器に分配されてもよい。例えば、チャンネル1 ($C = 1$) からの重み数 Y_1 およびフィルタからの $1 \sim Y_1$ カーネル（即ち、複数のフィルタの同じ入力チャンネルからの重み）を第1PEに供給し、チャンネル2 ($C = 2$) からの重み数 Y_1 およびフィルタからの $1 \sim Y_1$ カーネルを第2PEに供給することができる。このようにして、異なる入力チャンネルからの2DカーネルがPEの列に分散される。

20

【0052】

スパース化された入力テンソルを含むいくつかの実施形態において、各重みはインデックス - 値のペアとして表すことができる。インデックス - 値のペアの値は、非ゼロの重みの値であり、インデックス - 値のペアのインデックスは、非ゼロの重みのインデックスであり、これは、1つの乗算器で乗算を行うために、対応する入力値を識別するために使用することができる。いくつかの実施形態において、チャンネル数が各PEクラスタ（各行）内のPEの数よりも少ない場合、残りのPEを他のベクトル演算に使用することができる。

30

【0053】

いくつかの実施形態において、上述の 1×1 畳み込み演算以外の畳み込み演算は、1つまたは複数の 3×3 畳み込みに分解され、（例えば、 3×3 カーネルを使用して） 3×3 ネイティブ演算モードにマッピングされてもよい。図6Bに示されているように、各2Dの 3×3 カーネルは、同じ入力チャンネルからの9個の重み（ $(0, 0)$ 、 $(0, 1)$ 、 $(0, 2)$ 、 $(1, 0)$ 、 $(1, 1)$ 、 $(1, 2)$ 、 $(2, 0)$ 、 $(2, 1)$ 、 $(2, 2)$ ）を有し、それらは同時に処理するためにPEの同じ行（異なる列）に分配され得る。異なる入力チャンネルからの9個の重みは、PEの異なる行に分散させることができる。

40

【0054】

図7は、様々な実施形態によるPEのアレイにおける内部バッファの例示的なアーキテクチャ図を示す。いくつかの実施形態において、PEアレイ内の各PEは、入力値を記憶するための入力バッファ (IBUF) 722 と結合される。これらの入力値は、対応する入力値を見つけるために所定の重みインデックスに基づいてPEによって取得することが

50

できる。取得された入力値は、PEにおける乗算器内で重み値と乗算することができる。実際の実装では、IBUFの深さは通常、制限されており、1つのIBUF722は定数の入力チャンネルからの入力値のみを記憶できる。非ゼロの入力チャンネルの数も限られているため、この設計はスパース化された入力テンソルに対しては、有効に機能する。しかしながら、非ゼロの入力チャンネルの数がIBUF722の深さを超えることもよく起こる。このような場合、IBUF722はキャッシュ置換を実行して外部メモリから必要な入力値を読み込む必要があるが、これはコストが高く非効率的である。

【0055】

いくつかの実施形態において、重みテンソル（フィルタ）のスパース化の度合いに応じて、各PEのIBUFをプライベートメモリまたは共有メモリとして構成することができる。例えば、1つまたは複数の重みテンソルをスパース化することは、1つまたは複数の重みテンソルの入力チャンネルを複数のチャンネルグループに分割することであって、各チャンネルグループは、1より大きい整数である定数のチャンネルを含む、分割することと、1つまたは複数の重みテンソルをそれぞれプルーニングして、複数のチャンネルグループの各々において、複数のチャンネルグループの各々における少数のチャンネルのみが非ゼロの入力値を含み、各チャンネルグループにおける他のチャンネルは全てゼロを含むようにすることと、を含み得る。即ち、プルーニングプロセス後、各チャンネルグループは、同じ割合の非ゼロ重みを含む。スパース化の粒度は、細粒度710と粗粒度750に分類することができる。細粒度710のスパース性は、非ゼロ入力チャンネルが、定数よりも小さい数のチャンネルから選択される場合に発生し、粗粒度750のスパース性は、非ゼロ入力チャンネルが、定数よりも大きい数のチャンネルから選択される場合に発生する。例えば、重みのスパース性が15/16（16チャンネルのうちの1つが非ゼロ入力チャンネル）である場合、16個の入力チャンネルごとに1つの非ゼロ入力チャンネルを選択すること（例えば、1つのチャンネルグループには16個の入力チャンネルが含まれる）は、細粒度710のスパース性として決定されてよく、一方、64個の入力チャンネルごとに4つの非ゼロ入力チャンネルを選択すること（例えば、1つのチャンネルグループには64個の入力チャンネルが含まれる）ことは、粗粒度750のスパース性として決定されてよい。

【0056】

いくつかの実施形態において、IBUF722は、細粒度を有するスパース重みテンソル用のプライベートメモリとして、または、粗粒度を有するスパーステンソル用の共有メモリとして構成することができる。即ち、IBUF722の深さは、細粒度および粗粒度のスパース性を分類するために使用される定数と比較することができる。IBUF722の深さが定数よりも大きい場合、IBUF722は必要な入力値を記憶するのに十分であることを意味する。このように、専用プライベートメモリにより、データ取得性能が最適化される。IBUF722の深さが定数よりも小さい場合、複数の隣接するPEがそれらのIBUFを共有し、共有IBUFと表記され、これにより、それらによって取得可能な入力値を記憶する。このように、重複する入力値が削減され、全体的な記憶効率が向上する。

【0057】

図8は、様々な実施形態に従った適応テンソル形状を使用するニューラルネットワーク計算のための例示的な方法800を示す。方法800は、図1~7に記載されたデバイス、装置、またはシステムによって実行されてもよい。以下に提示される方法800の動作は、例示的なものであることを意図している。実行内容によっては、方法800は、様々な順序で、または並行して実行される追加の、より少ない、または代替のステップを含んでもよい。

【0058】

ブロック810は、処理要素（PE）アレイを使用した畳み込みのために、畳み込みニューラルネットワーク（CNN）の第1層において、第1入力特徴マップ（IFM）および1つまたは複数の第1フィルタを受信することを含み、PEアレイ内の各PEは、（Y

10

20

30

40

50

1) 個の乗算器を備え、PEアレイは、(Y2)個の行および(X)個の列に配置される。いくつかの実施形態では、PEの各行は、各PE内の(Y1)個の乗算器にそれぞれ対応する(Y1)個の加算器ツリーと結合されており、各PE内の各乗算器は、合計のために、対応する加算器ツリーに積算出力を送信する。各PE内の(Y1)個の乗算器は並列に処理データを処理し、PEアレイ内のPEは並列に処理データを処理する。

【0059】

ブロック820は、第1IFMおよび1つまたは複数の第1フィルタに基づいてネイティブテンソル形状を決定することを含み、ネイティブテンソルの形状は、第1外部次元、内部次元、および第2外部次元を含み、ネイティブテンソル形状は、第1IFMおよび1つまたは複数の第1フィルタをPEアレイにマッピングする。

10

【0060】

ブロック830は、PEアレイを使用して畳み込みを行うCNNの第2層において、第2IMFおよび1つまたは複数の第2フィルタを受信することを含む。いくつかの実施例において、CNNの第2層はCNNの第1層の後にあり、第2IMFは第1IFMよりも多くの入力チャンネルと、第1IFMよりも低い解像度と、を有する。いくつかの実施例において、1つまたは複数の第2フィルタの各々は、複数の2次元(2D)カーネルのチャンネルを備え、各2Dカーネルは1×1(one by one)または3×3(three by three)の次元を有する。

【0061】

ブロック840は、第2IMFおよび1つまたは複数の第2フィルタに基づいて、ネイティブテンソルの形状を変形することを含み、その変形は、内部次元を拡大し、第1外部次元および第2外部次元のうち的一方を縮小することを含み、拡大および縮小は、係数Fによって実行される。

20

【0062】

ブロック850は、変形されたネイティブテンソルに従って畳み込みを行うために、1つまたは複数の第2フィルタおよび第2IMFをPEアレイに供給することを含み、第1外部次元が縮小された場合、畳み込みはPEの同一の行からの出力をフラウンドについて合計して部分和を求めることを含み、第2外部次元が縮小された場合、畳み込みはPEの全てのF倍した行からの出力を合計して部分和を求めることを含む。いくつかの実施形態において、変形されたネイティブテンソルに従ってPEアレイに1つまたは複数の第2フィルタを供給することは、変形されたネイティブテンソルの第1外部次元および内部次元に従って1つまたは複数の第2フィルタを行列に変換することであって、1つまたは複数の第2フィルタの各2Dカーネルが1×1(one by one)の次元をそれぞれ有していることに応じて、行列の各行は、1つまたは複数の第2フィルタの異なる入力チャンネルからの重みを含む、変換することと、複数の入力チャンネルが一度に同時に処理されるように、行列の各行の重みをPEの異なる列に分配することと、を含む。いくつかの実施形態において、変形されたネイティブテンソルに従ってPEアレイに1つまたは複数の第2フィルタを供給することは、変形されたネイティブテンソルの第1外部次元および内部次元に従って、1つまたは複数の第2フィルタを行列に変換することであって、1つまたは複数の第2フィルタの各2Dカーネルが3×3の次元を有し、9個の重みを含む場合、9個の重みは行列の同じ行に配置される、変換することと、同じチャンネルからの重みが一度に同時に処理されるように、行列の同じ行からPEの異なる列に9個の重みを分配することと、を含む。いくつかの実施形態において、変形されたネイティブテンソルに従ったPEアレイへのIFMの供給は、変形されたネイティブテンソルの内部次元および第2外部次元に従ってIFMをマトリクスに変換することと、行列の列に対応するIFMの入力値をPEの行のバッファに供給することと、を含む。

30

40

【0063】

ブロック860は、複数の部分和を合計することにより、CNNの第2層における畳み込みの出力テンソルを取得することを含む。

【0064】

50

上記の説明において、 Y_1 、 Y_2 、 X 、および F は、全て1より大きい整数である。

【0065】

いくつかの実施形態において、方法800は、さらに、1つまたは複数のフィルタのチャンネルを複数のチャンネルグループに分割することであって、各チャンネルグループは1より大きい整数である定数のチャンネルを含む、分割することと、複数のチャンネルグループの各グループにおける少数のチャンネルのうちの一つのみが非ゼロの入力値を含み、各チャンネルグループの他のチャンネルは全てゼロを含むように、1つまたは複数のフィルタのそれぞれをブルーニングすることと、を含む。いくつかの実施形態において、方法800は、さらに、PEアレイ内の各PEに関連するバッファの深さを決定することと、バッファの深さが定数より大きい場合、各PEのプライベートメモリとしてバッファを構成することと、バッファの深さが定数より小さい場合、PEのバッファと隣接するPEの一つまたは複数のバッファを共有メモリとして結合することと、を含んでもよい。いくつかの実施形態において、各PEのプライベートメモリは、PE内の(Y_1)個の乗算器によって取得可能な入力値を記憶し、共有メモリは、PE内の(Y_1)個の乗算器および1つまたは複数の隣接するPEによって取得可能な入力値を記憶する。

10

【0066】

いくつかの実施形態において、1つまたは複数の第2フィルタの各々は複数の非ゼロ重みを含み、1つまたは複数の第2フィルタを畳み込みのためにPEアレイに供給することは、各非ゼロ重みを、対応するPEの乗算器に、非ゼロ重みと対応するインデックスを含むインデックス-値のペアとして供給することを含み、畳み込みは、インデックスに従って、対応するPEのバッファから入力値を取得することと、取得した値および非ゼロ重みを乗算器に送信して出力を取得することと、対応するPEと同じ行にある他のPEの他の乗算器によって生成された出力と合計するために、対応する加算器ツリーにその出力を送信することと、を含む。

20

【0067】

図9は、本明細書で説明されているいずれかの実施形態を実行することができる例示的な演算装置を示す。この演算装置は、図1から図8に示されているシステムおよび方法の一つまたは複数のコンポーネントを実装するために使用することができる。演算装置900は、情報を伝達するためのバス902または他の通信メカニズムと、情報を処理するためのバス902に結合された1つまたは複数のハードウェアプロセッサ704と、を備えてよい。ハードウェアプロセッサ(複数可)704は、例えば、1つまたは複数の汎用マイクロプロセッサとすることができる。

30

【0068】

演算装置900は、情報およびプロセッサ(複数可)904によって実行される命令を記憶するためにバス902に結合された、ランダムアクセスメモリ(RAM)などのメインメモリー907、キャッシュおよび/またはその他の記憶装置も備えてよい。メインメモリー907は、プロセッサ(複数可)904によって実行される命令の実行中に、一時変数またはその他の中間情報を記憶するためにも使用することができる。このような命令は、プロセッサ(複数可)904がアクセス可能な記憶媒体に記憶されると、演算装置900を、命令で規定された演算を実行するようにカスタマイズされた専用機にすることができる。メインメモリー907は、不揮発性媒体および/または揮発性媒体を含んでもよい。不揮発性媒体は、例えば、光ディスクまたは磁気ディスクを含んでもよい。揮発性媒体は、ダイナミックメモリを含んでもよい。一般的な媒体の形態には、例えば、フロッピーディスク、フレキシブルディスク、ハードディスク、半導体ドライブ、磁気テープ、その他の磁気データ記憶媒体、CD-ROM、その他の光データ記憶媒体、穿孔パターンを設けた物理的媒体、RAM、DRAM、PROM、EPROM、フラッシュ(登録商標)EPROM、NVRAM、その他のメモリチップまたはカートリッジ、またはそれらのネットワーク化されたバージョンが含まれる。

40

【0069】

演算装置900は、カスタマイズされた配線ロジック、1つまたは複数のASICまた

50

は F P G A、ファームウェアおよび / または演算装置と組み合わせることで演算装置 9 0 0 を特殊用途の機械にすることができるプログラムロジックを、本明細書で説明されている技術を実装するために使用してもよい。一実施形態によると、本明細書に記載の技術は、メインメモリー 9 0 7 に含まれる 1 つまたは複数の命令の 1 つまたは複数のシーケンスを実行するプロセッサ 9 0 4 に応答して、演算装置 9 0 0 によって実行される。このような命令は、記憶装置 9 0 9 などの他の記憶媒体からメインメモリー 9 0 7 に読み込むことができる。メインメモリー 9 0 7 に含まれる一連の命令の実行により、プロセッサ 9 0 4 が本明細書に記載された処理ステップを実行することができる。例えば、本明細書に開示されたプロセス / 方法は、メインメモリー 9 0 7 に記憶されたコンピュータプログラム命令により実行することができる。これらの命令がプロセッサ (複数可) 9 0 4 によって実行される場合、対応する図に示され、上述されたステップを実行することができる。代替の実施形態では、配線回路がソフトウェア命令の代わりに、またはソフトウェア命令と組み合わせ使用されてもよい。

10

【 0 0 7 0 】

演算装置 9 0 0 は、バス 9 0 2 に結合された通信インターフェース 9 1 0 も備える。通信インターフェース 9 1 0 は、1 つまたは複数のネットワークに接続された 1 つまたは複数のネットワークリンクに結合された双方向データ通信を提供することができる。他の例として、通信インターフェース 9 1 0 は、互換性のある L A N (または、W A N と通信する W A N コンポーネント) にデータ通信接続を提供するローカルエリアネットワーク (L A N) カードであってもよい。無線リンクを実装することもできる。

20

【 0 0 7 1 】

特定の動作の実行は、単一のマシン内だけでなく、多数のマシンに渡って展開されるプロセッサ間で分散されてもよい。いくつかの例示的な実施形態において、プロセッサまたはプロセッサ実装エンジンは単一の地理的位置 (例えば、家庭環境、オフィス環境、またはサーバーム内) に配置されてもよい。他の例示的な実施形態において、プロセッサまたはプロセッサ実装エンジンは複数の地理的位置に分散配置されてもよい。

【 0 0 7 2 】

前述の各セクションで説明されているプロセス、方法、アルゴリズムは、1 つまたは複数のコンピュータシステム、またはコンピュータハードウェアを備えるコンピュータプロセッサによって実行されるコードモジュールに実装され、完全にまたは部分的に自動化されてよい。プロセスおよびアルゴリズムは、部分的にまたは完全に、アプリケーション固有の回路に実装されてよい。

30

【 0 0 7 3 】

本明細書で開示された機能がソフトウェア機能単位で実装され、独立した製品として販売または使用される場合、それらはプロセッサで実行可能な不揮発性のコンピュータ読み取り可能な記憶媒体に記憶することができる。本明細書で開示された特定の技術的解決策 (全体または一部) または現在の技術に貢献する一態様は、ソフトウェア製品の形で具現化することができる。ソフトウェア製品は、記憶媒体に記憶され、演算装置 (パーソナルコンピュータ、サーバー、ネットワーク装置など) に本発明の実施形態の方法の全てまたはいくつかのステップを実行させる複数の命令を含んでよい。記憶媒体は、フラッシュ (登録商標) ドライブ、ポータブルハードドライブ、R O M、R A M、磁気ディスク、光ディスク、プログラムコードを記憶するために動作可能な他の媒体、またはそれらの組み合わせを含んでよい。

40

【 0 0 7 4 】

特定の実施形態はさらに、プロセッサと、プロセッサが実行可能な命令を記憶する非一時的なコンピュータ可読記憶媒体と、を備え、システムが、上述の実施形態の任意の方法におけるステップに対応する動作を実行するように構成される。特定の実施形態は、さらに、1 つまたは複数のプロセッサが実行可能な命令を含む非一時的なコンピュータ可読記憶媒体を備え、1 つまたは複数のプロセッサが、上述の実施形態の任意の方法におけるステップに対応する動作を実行するように構成される。

50

【 0 0 7 5 】

本明細書で開示されている実施形態は、クライアントとやりとりするクラウドプラットフォーム、サーバー、またはサーバーグループ（以下、総称して「サービスシステム」という）を通じて実装することができる。クライアントは、端末機器であってもよいし、プラットフォームでユーザが登録したクライアントであってもよく、端末機器は、モバイル端末、パーソナルコンピュータ（PC）、およびプラットフォームアプリケーションプログラムがインストール可能な如何なる機器であってもよい。

【 0 0 7 6 】

上記で説明した様々な機能およびプロセスは、互いに独立して使用される場合もあれば、様々な方法で組み合わせられる場合もある。全ての可能な組み合わせおよび部分的な組み合わせは、本開示の範囲に含まれることが意図されている。さらに、いくつかの実装では、特定の手法またはプロセスブロックが省略される場合がある。また、本明細書で説明した手法およびプロセスは、特定の順序に限定されるものではなく、それらに関連するブロックまたは状態は、適切な他の順序で実行することができる。例えば、説明されたブロックまたは状態は、具体的に開示された順序以外の順序で実行されてもよく、または、複数のブロックまたは状態が単一のブロックまたは状態に結合されてもよい。例示的なブロックまたは状態は、直列、並列、またはいくつかの他の方法で実行されてもよい。ブロックまたは状態は、開示された例示的な実施形態に追加されるか、またはそこから削除されてもよい。本明細書で説明される例示的なシステムおよびコンポーネントは、説明されたものとは異なる方法で構成されてもよい。例えば、開示された例示的な実施形態と比較して、要素が追加、削除、または再配置される場合がある。

【 0 0 7 7 】

本明細書で説明した例示的な方法の様々な動作は、少なくとも部分的に、アルゴリズムによって実行することができる。アルゴリズムは、メモリ（例えば、上述の非一時的なコンピュータ読み取り可能な記憶媒体）に記憶されたプログラムコードまたは命令を含んでよい。このようなアルゴリズムは、機械学習アルゴリズムを含んでよい。いくつかの実施形態において、機械学習アルゴリズムは、機能を実行するようにコンピュータに明示的にプログラムするのではなく、機能を実行する予測モデルを作成するためにトレーニングサンプルから学習することができる。

【 0 0 7 8 】

本明細書で説明されている例示的な方法の様々な動作は、関連する動作を実行するように一時的に（例えば、ソフトウェアによって）または恒久的に構成された1つまたは複数のプロセッサによって、少なくとも部分的に実行することができる。一時的にまたは恒久的に構成されたか否かに関わらず、そのようなプロセッサは、本明細書で説明されている1つまたは複数の動作または機能を実行するように動作可能に構成されたプロセッサ実装エンジンを構成することができる。

【 0 0 7 9 】

同様に、本明細書に記載される方法は、少なくとも一部がプロセッサによって実行されてもよく、特定のプロセッサまたはプロセッサ群はハードウェアの一例である。例えば、方法の少なくともいくつかの動作は、1つまたは複数のプロセッサまたはプロセッサ実装エンジンによって実行されてもよい。さらに、1つまたは複数のプロセッサは、「クラウドコンピューティング」環境内または「サービスとしてのソフトウェア」（SaaS）として、関連動作の実行をサポートするように動作してもよい。例えば、本実施例におけるいくつかの動作は、プロセッサを含む機械の一例であるコンピュータのグループによって実行され、これらの動作は、ネットワーク（例えば、インターネット）および1つまたは複数の適切なインターフェース（例えば、アプリケーションプログラムインターフェース（API））を介してアクセス可能である。

【 0 0 8 0 】

特定の動作の実行は、単一のマシン内だけでなく、複数のマシンに渡って展開されるプロセッサ間で分散されてもよい。いくつかの例示的な実施形態において、プロセッサまた

10

20

30

40

50

はプロセッサ実装エンジンは単一の地理的位置（例えば、家庭環境、オフィス環境、またはサーバーム内）に配置されてもよい。他の例示的な実施形態では、プロセッサまたはプロセッサ実装エンジンは複数の地理的位置に分散配置されてもよい。

【0081】

本明細書全体を通して、複数の例が、単一の例として説明されているコンポーネント、操作、または構造を実装する場合がある。1つまたは複数の方法の個々の動作は、個別の動作として図示および説明されているが、1つまたは複数の個々の動作は同時に実行されてもよく図示されている順序で動作を実行する必要はない。構成例における個別のコンポーネントとして示されている構造および機能は、統合された構造またはコンポーネントとして実装されてもよい。同様に、単一のコンポーネントとして提示された構造および機能は、個別のコンポーネントとして実装されてもよい。これらおよびその他のバリエーション、修正、追加、改善は、本願発明の範囲に含まれる。

10

【0082】

本書で使用される場合、「または」は、他に明確な指示がない限り、または文脈によって他に明確な指示がない限り、包括的であり排他的ではない。従って、本書では、「A、B、またはC」は、「A、B、AおよびB、AおよびC、BおよびC、または、A、B、およびC」を意味するが、他に明確な指示がない限り、または文脈によって他に明確な指示がない限り、この限りではない。さらに、「および」は、他に明確な指示がない限り、または文脈によって他に明確な指示がない限り、包括的であり個別的でもある。従って、本明細書では、「AおよびB」は「AおよびBの両方またはどちらか一方」を意味し、別段の明示的な指示または文脈による別段の指示がない限り、この限りではない。さらに、本明細書で単一の例として説明されている資源、動作、または構造については、複数の例が提供される場合がある。さらに、各種の資源、動作、エンジン、およびデータ記憶装置の境界は、ある程度任意であり、特定の動作は、特定の例示的な構成の文脈で説明されている。その他の機能の割り当ても想定されており、本開示の様々な実施形態の範囲に含めることができる。一般的に、例示的な構成における個別の資源として提示された構造および機能は、統合された構造または資源として実装することができる。同様に、単一の資源として提示された構造および機能は、個別の資源として実装することができる。これらのバリエーション、修正、追加、および改善は、添付の特許請求の範囲に示される本開示の実施形態の範囲に含まれる。従って、本明細書および図面は、制限的な意味ではなく、むしろ、例示的な意味で解釈されるべきである。

20

30

【0083】

「含む」または「備える」という用語は、後に述べられる特徴の存在を示すために使用されるが、他の特徴の追加を排除するものではない。「～できる」、「～できた」、「～かもしれない」、「～してもよい」などの条件文は、特に明記されていない限り、または文脈上、そう解釈されない限り、一般的に、特定の実施形態には特定の機能、要素および/またはステップが含まれるが、他の実施形態には含まれないことを意味する。従って、このような条件付きの表現は、機能、要素および/またはステップが1つまたは複数の実施形態に何らかの形で必要であることを意味するものではなく、また、1つまたは複数の実施形態が、これらの機能、要素および/またはステップが特定の実施形態に含まれるか、または実行されるべきであるか否かを、ユーザ入力またはプロンプトの有無にかかわらず決定するためのロジックを必ずしも含むことを意味するものではない。

40

【0084】

特定の例示的な実施形態を参照しながら発明の概要を説明したが、本開示の実施形態のより広い範囲から逸脱することなく、これらの実施形態に対して様々な修正や変更を加えることができる。本発明のこのような実施態様は、便宜上、単に「発明」という用語によって、個別にまたは集合的に本明細書で言及される場合があるが、実際には複数の開示または概念が開示されている場合であっても、本出願の範囲を任意の単一の開示または概念に自主的に限定することを意図するものではない。

【0085】

50

本明細書で説明されている実施形態は、当業者が開示された発明を実施できるように、十分詳細に説明されている。他の実施形態も利用でき、そこから派生させることも可能であり、構造的および論理的な置換や変更を、本開示の範囲から逸脱することなく行うことができる。従って、詳細な説明は限定的な意味で解釈されるべきではなく、様々な実施形態の範囲は、添付の請求の範囲のみによって定義され、そのような請求の範囲に含まれる全ての均等物とともに定義される。

【 図面 】

【 図 1 】

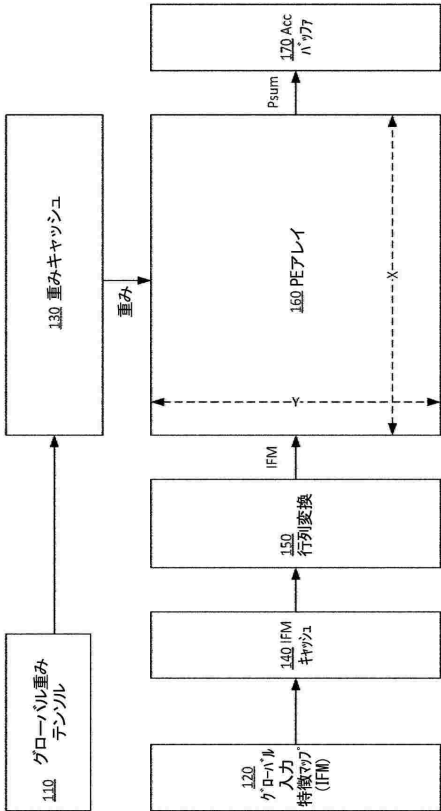


FIG. 1

【 図 2 】

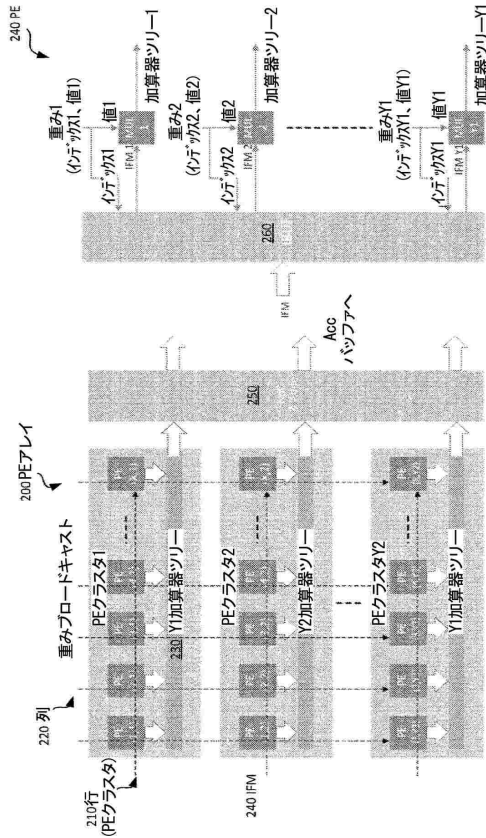


FIG. 2

10

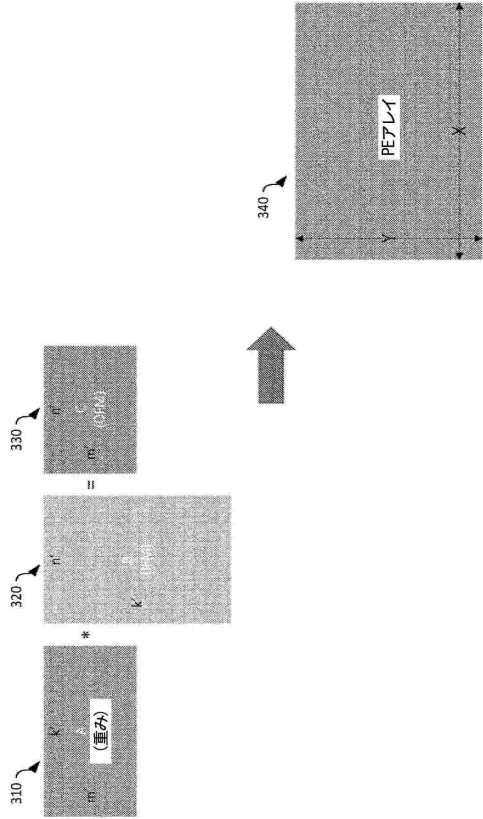
20

30

40

50

【 図 3 】



【 図 4 A 】

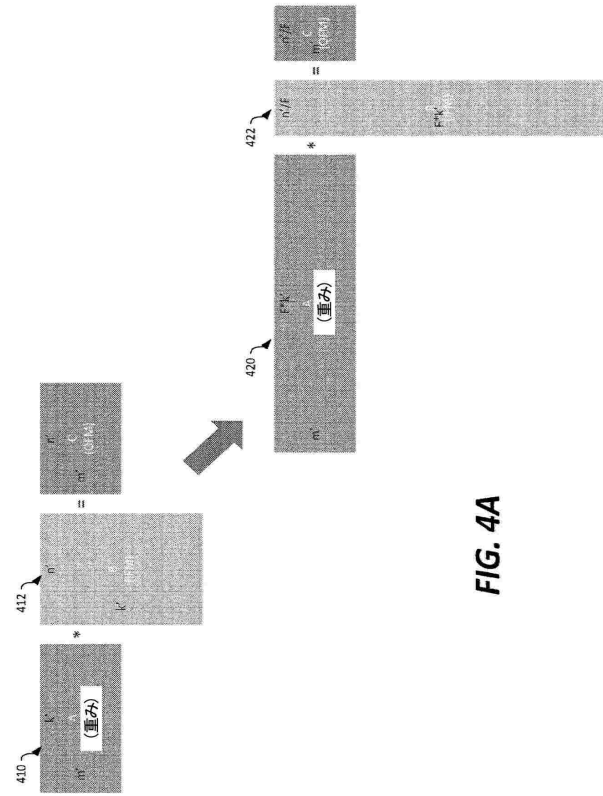


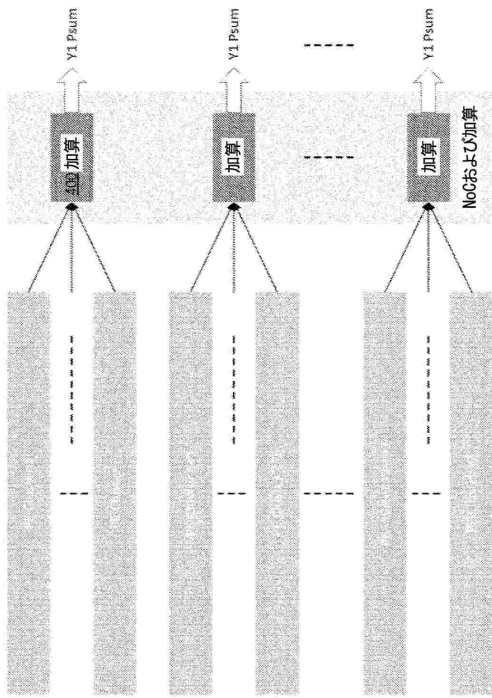
FIG. 3

FIG. 4A

10

20

【 図 4 B 】



【 図 5 A 】

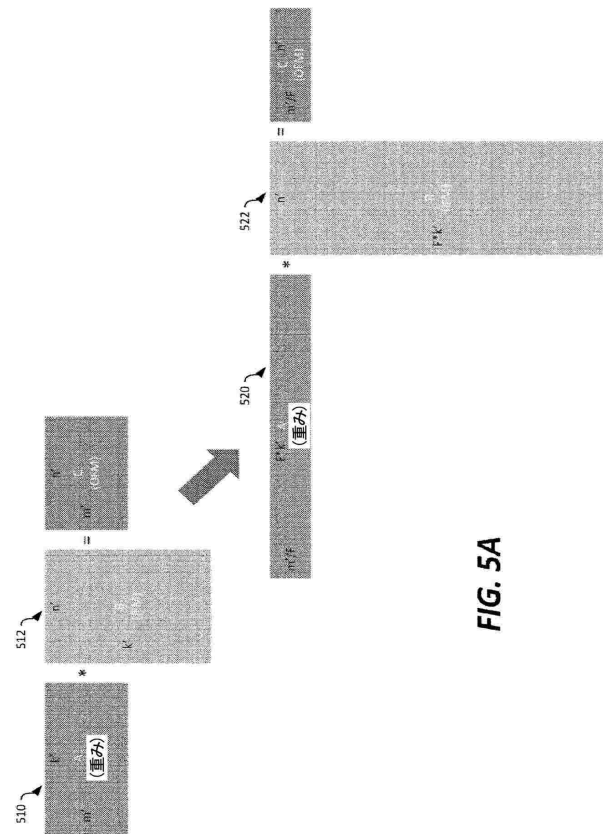


FIG. 4B

FIG. 5A

30

40

50

【 図 5 B 】

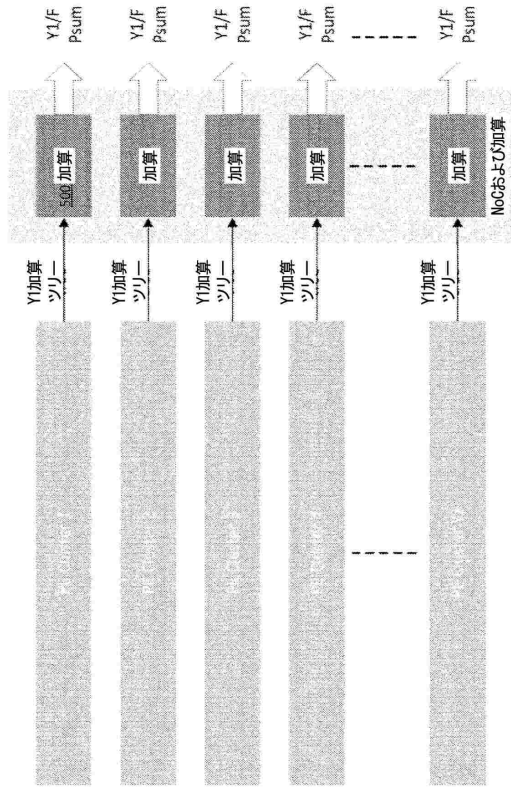


FIG. 5B

【 図 6 A 】

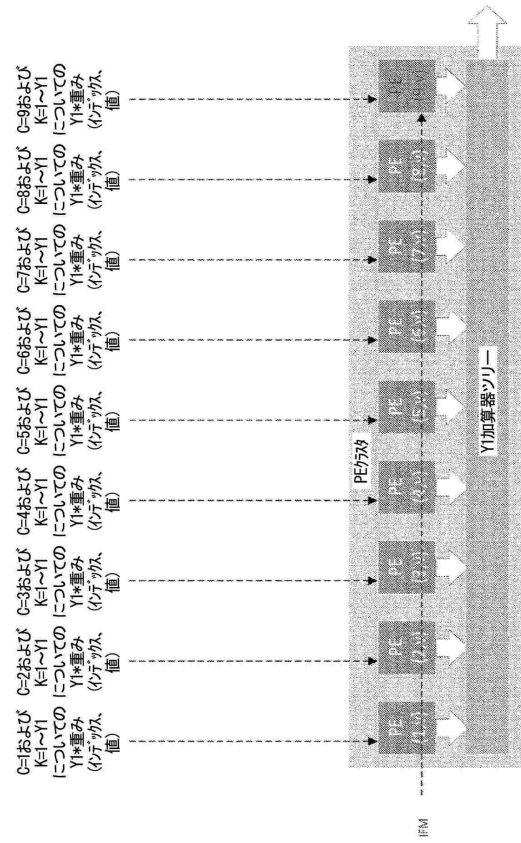


FIG. 6A

10

20

【 図 6 B 】

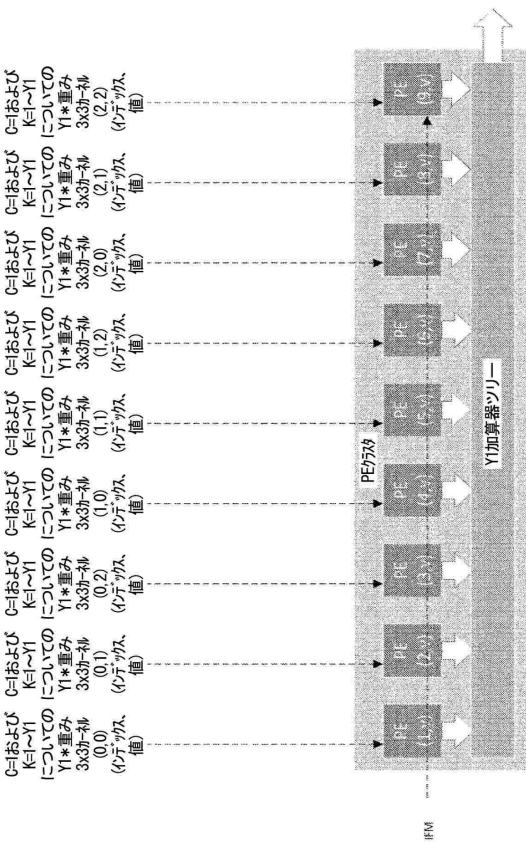


FIG. 6B

【 図 7 】

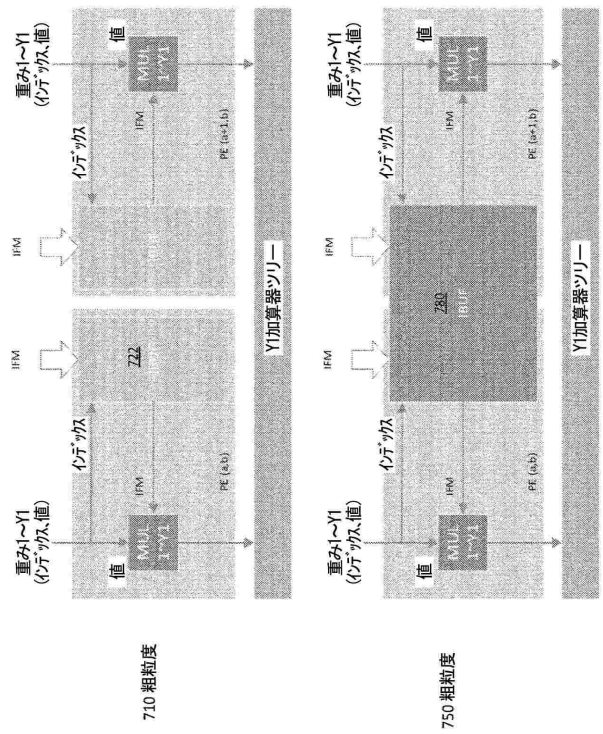


FIG. 7

30

40

50

【 図 8 】

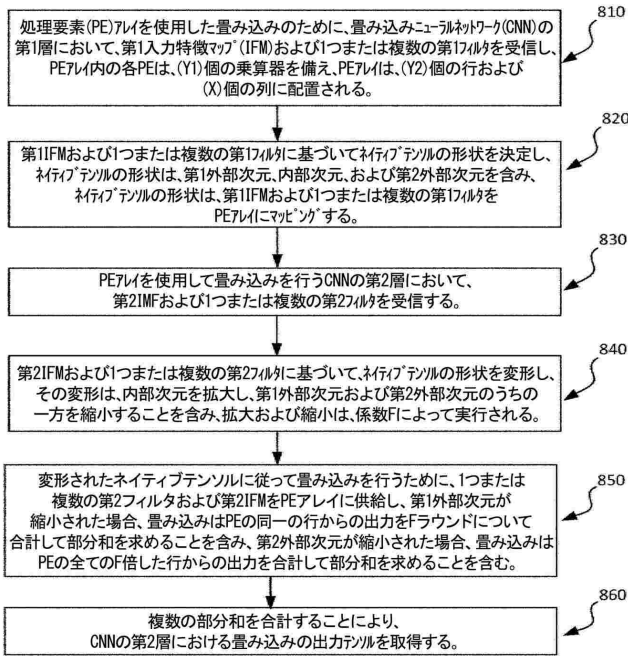


FIG. 8

【 図 9 】

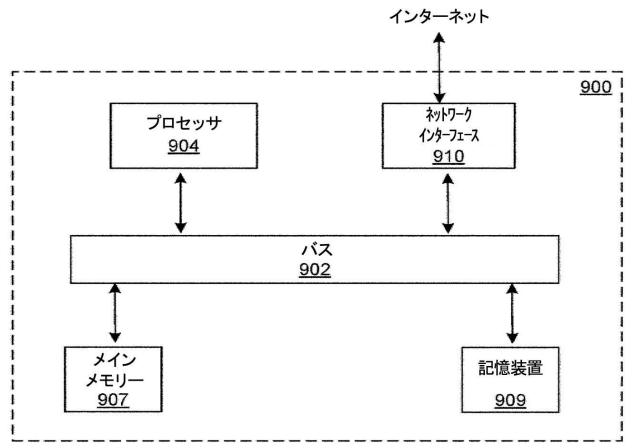


FIG. 9

10

20

30

40

50

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/CN2023/075661
A. CLASSIFICATION OF SUBJECT MATTER G06N3/04(2023.01)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) IPC: G06N,G06F 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) VEN, CNTXT, WOTXT, EPTXT, USTXT, CNKI, IEEE: processor element array, convolution, tensor, filter, neural network, kernel, shape, reshape, dimension, adaptive		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	US 2020134417 A1 (INTEL CORPORATION) 30 April 2020 (2020-04-30) description, paragraphs [0021]-[0071]	1-20
A	CN 108875958 A (GUANGZHOU YIGOU INTELLIGENT TECHNOLOGY) 23 November 2018 (2018-11-23) the whole document	1-20
A	CN 113449857 A (HUAWEI TECHNOLOGIES CO., LTD.) 28 September 2021 (2021-09-28) the whole document	1-20
A	US 2020293858 A1 (SAMSUNG ELECTRONICS CO., LTD.) 17 September 2020 (2020-09-17) the whole document	1-20
<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 "D" document cited by the applicant in the international application "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 27 April 2023		Date of mailing of the international search report 06 May 2023
Name and mailing address of the ISA/CN CHINA NATIONAL INTELLECTUAL PROPERTY ADMINISTRATION 6, Xitucheng Rd., Jimen Bridge, Haidian District, Beijing 100088, China		Authorized officer JIE,Xin Telephone No. (+86) 010-53961366

Form PCT/ISA/210 (second sheet) (July 2022)

10

20

30

40

50

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.
PCT/CN2023/075661

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)			Publication date (day/month/year)
US	2020134417	A1	30 April 2020	DE	102020131050	A1	24 June 2021
				CN	113033765	A	25 June 2021
				IN	202044041613	A	25 June 2021
				KR	20210082058	A	02 July 2021
				TW	202127324	A	16 July 2021
CN	108875958	A	23 November 2018	US	10223334	B1	05 March 2019
				US	10216704	B1	26 February 2019
				US	10169298	B1	01 January 2019
				US	10073816	B1	11 September 2018
CN	113449857	A	28 September 2021	WO	2021190127	A1	30 September 2021
				CN	115456159	A	09 December 2021
				CN	115456160	A	09 December 2021
				CN	115456161	A	09 December 2021
				EP	4123515	A1	25 January 2023
				US	2023023101	A1	26 January 2023
US	2020293858	A1	17 September 2020	IN	201941009806	A	18 September 2020
				KR	20200110165	A	23 September 2020

10

20

30

40

50

フロントページの続き

,MC,ME,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,CV,CZ,DE,DJ,DK,DM,DO,DZ,EC,EE,EG,ES,FI,GB,GD,GE,GH,GM,GT,HN,HR,HU,ID,IL,IN,IQ,IR,IS,IT,JM,JO,JP,KE,KG,KH,KN,KP,KR,KW,KZ,LA,LC,LK,LR,LS,LU,LY,MA,MD,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,ST,SV,SY,TH,TJ,TM,TN,TR,TT,TZ,UA,UG,US,UZ,VC,VN,WS,ZA,ZM,ZW

弁理士 南山 知広

(74)代理人 100153729

弁理士 森本 有一

(74)代理人 100151459

弁理士 中村 健一

(72)発明者 シアオチエン チャン

アメリカ合衆国, カリフォルニア 94022, ロス アルトス, シェアウッド アベニュー 949, スイート 200

(72)発明者 エンシュイ イェン

アメリカ合衆国, カリフォルニア 94022, ロス アルトス, シェアウッド アベニュー 949, スイート 200

(72)発明者 チーピン シアオ

アメリカ合衆国, カリフォルニア 94022, ロス アルトス, シェアウッド アベニュー 949, スイート 200

F ターム (参考) 5B056 BB26 BB38 FF01 FF02 FF05 FF07