

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7179853号
(P7179853)

(45)発行日 令和4年11月29日(2022.11.29)

(24)登録日 令和4年11月18日(2022.11.18)

(51)国際特許分類 F I
G 0 6 N 3/063(2006.01) G 0 6 N 3/063
G 0 6 F 15/80 (2006.01) G 0 6 F 15/80

請求項の数 15 (全51頁)

(21)出願番号	特願2020-531932(P2020-531932)	(73)特許権者	506329306 アマゾン テクノロジーズ インコーポレイテッド アメリカ合衆国 9 8 1 0 8 - 1 2 2 6 ワシントン州 シアトル ピーオー ボックス 8 1 2 2 6
(86)(22)出願日	平成30年12月10日(2018.12.10)	(74)代理人	100137969 弁理士 岡部 憲昭
(65)公表番号	特表2021-506032(P2021-506032 A)	(74)代理人	100104824 弁理士 穠場 仁
(43)公表日	令和3年2月18日(2021.2.18)	(74)代理人	100121463 弁理士 矢口 哲也
(86)国際出願番号	PCT/US2018/064777	(72)発明者	ファン, ランディ アメリカ合衆国, ワシントン州 9 8 1 0 8 - 1 2 2 6, シアトル, ピー. オー
(87)国際公開番号	WO2019/118363		
(87)国際公開日	令和1年6月20日(2019.6.20)		
審査請求日	令和2年8月1日(2020.8.1)		
(31)優先権主張番号	15/839,017		
(32)優先日	平成29年12月12日(2017.12.12)		
(33)優先権主張国・地域又は機関	米国(US)		
(31)優先権主張番号	15/839,301		
(32)優先日	平成29年12月12日(2017.12.12)		
	最終頁に続く		最終頁に続く

(54)【発明の名称】 オンチップの計算ネットワーク

(57)【特許請求の範囲】

【請求項1】

集積回路であって、

第1の処理エンジンアレイと、

ニューラルネットワークの重み値セットを記憶する複数のメモリバンクを含むオンチップメモリと、を備え、

前記複数のメモリバンクのうちの各バンクは、独立してアクセス可能であり、かつ前記複数のメモリバンクを含むオンチップメモリおよび前記第1の処理エンジンアレイは、同じダイ上にあり、

入力データを受信した時点で、前記集積回路は、前記重み値セットを用いて、前記ニューラルネットワーク用に規定されたタスクを実行するように構成され、前記タスクを実行することは、

前記オンチップメモリから前記第1の処理エンジンアレイに第1の重み値を読み取ることと、

前記入力データを前記第1の処理エンジンアレイに入力することと、

前記入力データおよび前記第1の重み値に基づいて、前記第1の処理エンジンアレイを使用して第1の中間結果を計算することと、

前記第1の中間結果を前記オンチップメモリに書き込むと同時に、前記オンチップメモリから前記第1の処理エンジンアレイに第2の重み値を読み取ることと、

前記第1の中間結果を前記第1の処理エンジンアレイに入力することと、

10

20

前記第1の中間結果および前記第2の重み値に基づいて、前記第1の処理エンジンアレイを用いて結果を計算すること、を含み、前記結果は、前記タスク実行の結果に対応する、集積回路。

【請求項2】

前記タスクを実行することは、

前記複数のメモリバンクのうちの異なるメモリバンクから2つ以上の値を同時に読み取ること、をさらに含む、請求項1に記載の集積回路。

【請求項3】

前記2つ以上の値は、重み値、入力値または中間結果を含む、請求項2に記載の集積回路。

【請求項4】

前記タスクを実行することは、

前記複数のメモリバンクのうちの第1のメモリバンクへ第1の値を書き込むことと、前記複数のメモリバンクのうちの第2のメモリバンクから第2の値を読み取ること、をさらに含む、前記第1の値は、前記第2の値が読み取られると同時に書き込まれる、請求項1～3のいずれか一項に記載の集積回路。

【請求項5】

前記第1の値および前記第2の値は、重み値、および入力値、または中間結果を含む、請求項4に記載の集積回路。

【請求項6】

前記重み値セットは、前記ニューラルネットワークの全ての重み値を含む、請求項1に記載の集積回路。

【請求項7】

第2の処理エンジンアレイをさらに備え、前記複数のメモリバンクのうちの第1のメモリバンクセットは、前記第1の処理エンジンアレイによって使用されるように構成され、前記複数のメモリバンクのうちの第2のメモリバンクセットは、前記第2の処理エンジンアレイによって使用されるように構成され、前記第1のメモリバンクセットおよび前記第2のメモリバンクセットは各々、前記ニューラルネットワークの前記重み値の一部分を含み、前記タスクを実行することは、

前記第1の処理エンジンアレイによって第2の中間結果を計算することであって、前記第1の処理エンジンアレイは、前記第1のメモリバンクセットからの重み値を用いて前記第2の中間結果を計算することと、

前記第1の処理エンジンアレイによって、前記第2のメモリバンクセットから追加の重み値を読み取ること、をさらに含む、前記第1の処理エンジンアレイは、前記第2の中間結果および前記追加の重み値を用いて前記結果を計算する、請求項1に記載の集積回路。

【請求項8】

第2のメモリバンクセットに記憶される前記ニューラルネットワークの前記重み値の前記一部分による占有は、前記第2のメモリバンクセットの全てより少なく、前記第2の処理エンジンアレイは、前記重み値セットによって占有されていない前記第2のメモリバンクセットの一部を用いて計算を実行する、請求項7に記載の集積回路。

【請求項9】

前記第1のメモリバンクセットおよび前記第2のメモリバンクセットに記憶される前記重み値の前記一部分は、前記ニューラルネットワークの全ての重み値を含む、請求項7に記載の集積回路。

【請求項10】

前記複数のメモリバンクの第1の部分は、前記結果を計算するための前記第1の中間結果を記憶するために初めに予約され、前記重み値セットは、前記ニューラルネットワークの全ての重み値より少ない重み値を含む、請求項1に記載の集積回路。

【請求項11】

前記タスクを実行することは、

10

20

30

40

50

前記中間結果を記憶するために必要なメモリの量が減少していると決定することと、別のメモリから追加の重み値セットを読み取ることと、

前記追加の重み値セットを前記複数のメモリバンクの前記第 1 の部分に記憶することであって、前記追加の重み値セットは、必要とされる前に記憶される、記憶すること、をさらに含む、請求項 1 0 に記載の集積回路。

【請求項 1 2】

前記第 1 の処理エンジンアレイは、処理エンジンセットを含み、前記処理エンジンセットのうちの各処理エンジンは、結果を、前記処理エンジンセットのうちの別の処理エンジンへ直に出力する、請求項 1 ~ 1 1 のいずれか一項に記載の集積回路。

【請求項 1 3】

前記第 1 の処理エンジンアレイのうちの各処理エンジンは、乗数累算器回路を含む、請求項 1 ~ 1 2 のいずれか一項に記載の集積回路。

【請求項 1 4】

前記ニューラルネットワークは、重み付き有向グラフから導出される複数の重み値と、前記重み付き有向グラフにおける各ノードについて計算を実行するための命令セットとを含み、前記複数の重み値は、既知の入力データを用いて前記タスクを実行することにより、先に決定されている、請求項 1 ~ 1 3 のいずれか一項に記載の集積回路。

【請求項 1 5】

方法であって、

ニューラルネットワーク処理回路のオンチップメモリであって、複数のメモリバンクを有する前記オンチップメモリに、重み値セットを記憶することであって、前記複数のメモリバンクのうちの各バンクは、独立してアクセス可能であり、前記ニューラルネットワーク処理回路は、前記オンチップメモリと同じダイ上に処理エンジンアレイを含み、前記重み値セットは、入力データの受信に先行して記憶される、記憶することと、

入力データを受信することと、

前記重み値セットを用いて、ニューラルネットワーク用に規定されるタスクを実行すること、を含み、前記タスクを実行することは、

前記オンチップメモリから前記処理エンジンアレイに第 1 の重み値を読み取ることと、

前記入力データを前記処理エンジンアレイに入力することと、

前記入力データおよび前記第 1 の重み値に基づいて、前記処理エンジンアレイを使用して中間結果を計算することと、

前記中間結果を前記オンチップメモリに書き込むと同時に、前記オンチップメモリから前記処理エンジンアレイに第 2 の重み値を読み取ることと、

前記中間結果を前記処理エンジンアレイに入力することと、前記中間結果および前記第 2 の重み値に基づいて、前記処理エンジンアレイを使用して結果を計算すること、を含み、前記結果は、前記タスク実行の結果に対応する、方法。

【発明の詳細な説明】

【背景技術】

【0001】

[0001]ニューラルネットワークは、コンピュータ技術を用いて、動物の脳を構成する生体神経回路網により実行される論理的推論を複製しようとするものである。ニューラルネットワークは、ヒトの脳の動作の力学からインスピレーションを得ている。ニューラルネットワークでは、神経細胞がノードで表され、シナプスがノード間の重み付けされた結合によって表される。重みは、入力に対する様々な応答を反映することができる。ニューラルネットワークは、層状に配置されることが可能であって、分析されるべき入力データは、ある入力層へ提供され、かつ各層の出力は、入力を次の層へ提供する。最後の層は、結果を出力することができる。重み値は、訓練を通じて、結果が既知である入力データがニューラルネットワークへ提供される間に決定されることが可能である。

【0002】

[0002]ニューラルネットワークは、計算を実行するために、中央処理装置 (CPU) を

10

20

30

40

50

用いて実装されることが可能である。しかしながら、CPUは、並列計算ではなく逐次計算用に最適化される傾向があり、よって、応答時間が長くなる可能性がある。グラフィック処理装置（GPU）は、並列計算用に最適化されているが、必ずしも、1つの計算ユニットからの結果を別の計算ユニットへ直に提供するように最適化されたものではない。多くの場合、結果は、まずメモリに書き込まなければならない。GPUは、CPUより優れた応答時間を有するものの、応答時間に遅れが出ることがある。

【0003】

【0003】専用ニューラルネットワークプロセッサは、並列の連鎖計算用に最適化された計算アレイを含む。あるニューラルネットワークプロセッサでは、計算ユニットは、結果を別の計算ユニットへ直に出力することができ、結果をメモリに書き込む必要がない。結果を、たとえばそのアレイを介して新たな計算サイクルを開始するために、メモリへ書き込む必要がある場合、結果は、計算アレイに対してローカルなメモリに記憶されることが可能である。したがって、ニューラルネットワークプロセッサは、同じ入力データに対してCPUおよびGPUのいずれよりも優れたパフォーマンスを発揮することができる。

10

【0004】

【0004】様々な例について、図面を参照して説明する。

【図面の簡単な説明】

【0005】

【図1】ニューラルネットワークの視覚モデルの一例を示す。

【図2A】順伝播重みおよび回帰重みを含むニューラルネットワークのモデルの一例を示す。

20

【図2B】異なる結合型を含むニューラルネットワークのモデルの一例を示す。

【図3A】2次元畳み込みのモデルの一例を示す。

【図3B】畳み込みニューラルネットワークのモデルの一例を示す。

【図4】ニューラルネットワークの重み値をオフチップメモリではなくオンチップで記憶する効果の一例を示す。

【図5】ニューラルネットワークを実行するために使用可能なニューラルネットワーク処理エンジンの一例を示す。

【図6A】ニューラルネットワークを実行するために使用可能なニューラル処理エンジンの一例を示す。

30

【図6B】ニューラルネットワークを実行するために使用可能なニューラル処理エンジンの一例を示す。

【図7A】ニューラルネットワークを実行するために使用可能なニューラルネットワーク処理エンジンの一例を示す。

【図7B】ニューラルネットワークを実行するために使用可能なニューラルネットワーク処理エンジンの一例を示す。

【図7C】ニューラルネットワークを実行するために使用可能なニューラルネットワーク処理エンジンの一例を示す。

【図8】複数のニューラルネットワーク処理エンジンを有するニューラルネットワークプロセッサの一例を示す。

40

【図9】ニューラルネットワークの重みが2つのニューラルネットワーク処理エンジンのメモリサブシステムに記憶されるように構成される、ニューラルネットワークプロセッサの一例を示す。

【図10】ニューラルネットワークの重みが2つのニューラルネットワーク処理エンジンのメモリサブシステムに記憶されるように構成される、ニューラルネットワークプロセッサの一例を示す。

【図11】ニューラルネットワークの全ての重みをオンチップで記憶するための別のソリューションを示す。

【図12】ニューラルネットワークプロセッサがニューラルネットワークを実行している場合にニューラルネットワークプロセッサにより実行されることが可能なプロセスの一例

50

を示す。

【図 1 3】ニューラルネットワークプロセッサがニューラルネットワークを実行している場合にニューラルネットワークプロセッサにより実行されることが可能なプロセスの一例を示す。

【図 1 4】ニューラルネットワークプロセッサを含むコンピューティングシステムにより実行されることが可能なプロセスの一例を示す。

【図 1 5】これまでに論じた様々な例を使用できるネットワークを示す。

【発明を実施するための形態】

【0006】

[0022]以下の記述では、様々な実装例について説明する。説明を目的として、実施例を完全に理解するために、特定の構成および詳細について述べる。しかしながら、当業者には、実施例が、これらの特定の詳細なしに実践され得ることも明らかであろう。さらに、実施例の説明を不明瞭にしないように、周知の特徴は、省略または簡略化されることがある。

10

【0007】

[0023]人工ニューラルネットワークは、コンピュータ技術を用いて、動物の脳を構成する生体神経回路網により実行される論理的推論を複製しようとするものである。人工ニューラルネットワーク（本明細書では、ニューラルネットワークと称する）は、人が意思決定する方法と同じ方法で意思決定することができるインテリジェントマシンの構築を目指す研究およびエンジニアリングの一分野である、人工知能（AI）分野の一部である。ニューラルネットワークは、機械学習と呼ばれる人工知能の下位分野に属する。機械学習は、コンピュータに明示的なプログラムなしに学習する能力を与えることを探る研究分野である。機械学習アルゴリズムを実装するプログラムは、プログラムがあらゆる可能性を説明するコード、および考えられる全ての挙動を記述するコードを含有する必要なしに、タスクの実行を学習することができる。

20

【0008】

[0024]ニューラルネットワークは、ヒトの脳の動作の力学から、これらの動作が理解される限りにおいて、インスピレーションを得ている。脳の様々なモデルによれば、脳の主要な計算要素は、ニューロンである。ニューロンは、幾つかの要素で互いに結合され、ニューロンに入る要素を樹状突起と称し、ニューロンを出る要素を軸索と称する。ニューロンは、樹状突起を介して信号を受け取り、信号を基に計算を実行し、かつ軸索上へ信号を出力する。入力信号および出力信号を、活性化と称する。あるニューロンの軸索は、枝を出して、複数のニューロンの樹状突起へ結合されることが可能である。軸索の枝と樹状突起との結合部を、シナプスと呼ぶ。

30

【0009】

[0025]シナプスは、シナプスを通る信号をスケールリングすることができる。スケールリング係数は、重みと称され、脳による学習を可能にする方法と考えられていて、異なる重みは、入力に対する異なる応答から生じる。学習は、重みを変更することができるが、学習を達成するためにニューロンおよびシナプスの編成が変わる必要はない。したがって、脳の静的構造は、プログラムのモデルとして使用されることが可能であり、重みは、プログラムが既の実行することを学習しているタスクを反映することができる。

40

【0010】

[0026]ニューラルネットワークは、ニューロンの計算が入力値の加重和を包含するという概念に基づいて動作する。これらの加重和は、シナプスにより実行される値のスケールリング、およびニューロンにおけるこれらの値の結合に対応する。ニューロンでは、結合された入力に対して関数演算が実行される。脳モデルにおいて、演算は、入力があるしきい値を越えるときにのみニューロンに出力を生成させる非線形関数であるように見える。したがって、類推により、ニューラルネットワークのノードは、ノードに入力される値の加重和へ非線形関数を適用することができる。

【0011】

50

[0027]図 1 は、ニューラルネットワークの視覚モデルの一例 100 を示す。本例において、モデル 100 は、入力層 104 と、隠れ層と称する場合も多い中間層 106 と、出力層 108 とを含む。各層は、幾つかのノード 102 を含む。本例において、入力層 104 のノード 102 は、隠れ層 106 の各ノード 102 へ連結されている。脳モデルにおいてシナプスと称される連結部は、重み 110 と称される。また、本例では、隠れ層 106 の各ノード 102 も、出力層の各ノード 102 との間に連結部または重み 110 を有する。入力層 104 は、入力を受信することができ、かつこれらの入力を隠れ層 106 へ伝搬することができる。ニューラルネットワークの実装は、複数の隠れ層を包含することができる。隠れ層 106 (または、複数の隠れ層) により計算される加重和は、出力層 108 へ伝搬され、出力層 108 は、最終的な出力をユーザへ提示することができる。ノード 102 の出力は、脳モデルに合わせて、活性化と呼ぶことができる。

10

【0012】

[0028]例示的なモデル 100 の各層において起こり得る計算の一例は、次の通りである。

$$y_j = f\left(\sum_{i=1}^3 W_{ij} \times x_i + b\right)$$

【0013】

[0029]上述の式において、 W_{ij} は、重みであり、 x_i は、入力活性化であり、 y_j は、出力活性化であり、 $f()$ は、非線形関数であり、 b は、バイアス項である。様々な非線形関数を用いれば、様々な目的を達成することができる。

20

【0014】

[0030]モデル 100 は、有向重み付きグラフと呼ぶことができる。有向グラフにおいて、ノードへの、またはノードからの各連結部は、方向(たとえば、ノードへ向かう、またはノードから出る方向)を示す。重み付きグラフにおいて、各連結部は、重みを有することができる。ニューラルネットワークの開発ツールは、理解およびデバッグを容易にするために、ニューラルネットワークを有向重み付きグラフとして視覚化することができる。場合によっては、これらのツールは、ニューラルネットワークを訓練しかつ訓練された重み値を出力するためにも使用可能である。ニューラルネットワークの実行は、次に、入力データに対する計算の実行に重みを用いる問題となる。

30

【0015】

[0031]4 つ以上の層(たとえば、複数の隠れ層)を有するニューラルネットワークは、深層ニューラルネットワークと呼ばれることがある。深層ニューラルネットワークは、たとえば 5 ~ 1000 超の層を有し得る。

【0016】

[0032]多層ニューラルネットワークは、より浅いネットワークを上回る複雑さおよび抽象化によって、高レベル特徴の学習が可能である。一例として、ニューラルネットワークは、画像を認識するようにと教わることが可能である。本例では、画像のピクセルをニューラルネットワークの入力層へ供給することができ、第 1 の層の出力は、線およびエッジなどの、画像における低レベル特徴の存在を示すことができる。後続層では、これらの特徴を組み合わせて、より高いレベルの特徴の存在尤度を測定することができ、つまりは、線を組み合わせて形にすることができ、これらをさらに組み合わせて形の集合にすることができる。この全ての情報を所与として、ニューラルネットワークは、高レベル特徴が特定のオブジェクトまたはシーンを表す蓋然性を出力することができる。たとえば、ニューラルネットワークは、画像に猫が含まれているか、含まれていないか、を出力することができる。

40

【0017】

[0033]ニューラルネットワークの学習フェーズは、ニューラルネットワークの訓練と称

50

される。訓練の間に、ニューラルネットワークは、タスクを実行することを教示される。タスクの学習において、重み（および場合によってはバイアスも）の値が決定される。ニューラルネットワークの基盤となるプログラム（たとえば、ノードの層編成、各層のノード間の連結、および各ノードが実行する計算）を、訓練中に変える必要はない。ニューラルネットワークは、訓練されると、訓練の間に決定された重み値を用いて結果を計算することにより、タスクを実行することができる。たとえば、ニューラルネットワークは、画像が特定のオブジェクトを包含する蓋然性、オーディオシーケンスが特定のワードを包含する蓋然性、画像におけるオブジェクトを囲むバウンディングボックス、または講じるべき対策案を出力することができる。ニューラルネットワークのプログラムの実行は、推論と称される。

10

【0018】

[0034]重みの訓練方法は、幾つか存在する。ある方法は、教師あり学習と称される。教師あり学習では、全ての訓練サンプルにラベルが付けられ、よって、各訓練サンプルをニューラルネットワークに入力すると、既知の結果が生成される。別の方法は教師なし学習と称され、訓練サンプルにラベルは付されず、訓練は、データ内の構造またはデータ内のクラスタを見つけることを目的とする。半教師あり学習は、教師あり学習と教師なし学習との中間に属する。半教師あり学習では、訓練データの部分集合にラベルが付される。ラベルなしデータを用いてクラスタ境界を画定し、かつラベル付きデータを用いてクラスタにラベルを付すことができる。

【0019】

20

[0035]ニューラルネットワークは、たとえば画像とビデオ、音声と言語、医学、ゲームプレイおよびロボット工学を含む様々なアプリケーションに使用されてきた。画像とビデオでは、画像の分類、オブジェクトの定位および検出、画像のセグメンテーションおよびアクションの認識にニューラルネットワークが使用されてきた。音声と言語では、音声認識、機械翻訳、自然言語処理およびオーディオ生成にニューラルネットワークが使用されてきた。医療分野では、ゲノミクスおよび医用画像化にニューラルネットワークが使用されてきた。ゲームプレイでは、Goなどの膨大な数の動作可能性を有するゲームを含むビデオゲームおよびボードゲームのプレイにニューラルネットワークが使用されてきた。ロボット工学では、ロボットのモーションプランニング、ビジュアルナビゲーション、制御安定化および自律車両の運転戦略にニューラルネットワークが使用されてきた。

30

【0020】

[0036]ニューラルネットワークは、多様なものが開発されてきた。ニューラルネットワークの様々な例は、順伝播型および回帰型という2つの形態に分けることができる。図2Aは、入力層204と隠れ層206との間の順伝播重み212と、出力層208における回帰重み214とを含むニューラルネットワークのモデルの一例210を示す。順伝播型ニューラルネットワークにおいて、計算は、先行層の出力に対する一連の演算であって、最終層がニューラルネットワークの出力を生成する。図2Aに示す例では、順伝播が隠れ層206によって示され、そのノード202は、入力層204のノード202の出力のみを操作する。順伝播型ニューラルネットワークにはメモリがなく、所与の入力に対する出力は、ニューラルネットワークに与えられた先行入力に関わりなく、常に同じであり得る。多層パーセプトロン(MLP)は、順伝播重みのみを有するニューラルネットワークの一種である。

40

【0021】

[0037]これに対して、回帰型ニューラルネットワークは、出力に影響する依存関係を許容し得る内部メモリを有する。回帰型ニューラルネットワークでは、中間演算の中に、内部に保存されかつ後の入力の処理に関連して他の演算への入力として使用されることが可能な値を生成できるものがある。図2Aの例では、回帰が出力層208によって示されていて、出力層208のノード202の出力は、出力層208のノード202の入力へ連結し戻される。これらのループバック連結を、回帰重み214と称し得る。長・短期記憶(LSTM)は、頻繁に使用される回帰型ニューラルネットワークの一変形例である。

50

【 0 0 2 2 】

[0038]図 2 B は、異なる結合型を含むニューラルネットワークのモデルの一例 2 2 0 を示す。このモデル例 2 2 0 において、入力層 2 0 4 および隠れ層 2 0 6 は、2 2 2 で完全に結合された層である。完全に結合された層では、全ての出力活性化が、重み付けされた入力活性化で構成される（たとえば、入力層 2 0 4 における全てのノード 2 0 2 の出力は、隠れ層 2 0 6 の入力の全てに結合される）。完全に結合された層は、大量のストレージおよび計算を要する可能性がある。多層パーセプトロンニューラルネットワークは、完全に結合されているニューラルネットワークの一種である。

【 0 0 2 3 】

[0039]アプリケーションによっては、活性化間の幾つかの結合を、たとえば、これらの結合の重みをゼロに設定することにより、出力の精度に影響を与えることなく削除することができる。結果は、図 2 B の 2 2 4 で示す、隠れ層 2 0 6 と出力層 2 0 8 との間の重みによる、疎に結合された層である。2 2 4 で示す疎に結合された層を実現できる方法の別の例に、プーリングがある。プーリングでは、ノードのクラスタの出力を、たとえば、最大値、最小値、平均値または中央値を見出すことによって組み合わせることができる。

10

【 0 0 2 4 】

[0040]ニューラルネットワークの動作効率は、幾つかの異なる方法でさらに向上され得る。たとえば、出力に寄与する重みの数は、出力を固定サイズの入力窓のみの関数とすることによって制限されることが可能である。全ての出力の計算に同じ重みセットを用いれば、さらなる効率の獲得が可能である。同じ重み値の反復使用は、重み共有と称され、重みのストレージ要件を大幅に減らすことができる。

20

【 0 0 2 5 】

[0041]ニューラルネットワーク層におけるウィンドウ処理および重み共有は、各ノードで実行される計算を畳み込みとして構造化することによって達成され得る。図 3 A は、画像処理に適用された 2 次元畳み込みのモデルの一例 3 1 0 を示す。このモデル例において、フィルタ平面 3 0 4 は、高さ R および幅 S を有する行列に配置された重みの集合である。フィルタ平面 3 0 4 は、たとえば、要素同士の乗算 3 1 2 を用いて、入力画像に適用可能であり、入力画像のデータは、入力特徴マップ 3 0 6 と称され得る。フィルタ平面 3 0 4 の高さ R および幅 S は、共に、入力特徴マップ 3 0 6 の高さ H および幅 W より小さく、よって、フィルタ平面 3 0 4 を入力特徴マップ 3 0 6 へ適用する結果、入力活性化 3 1 6 の小さい近傍が計算されることになる（たとえば、近傍を超える重みは、ゼロに設定可能である）。入力活性化 3 1 6 は、たとえば部分和の累積 3 1 4 を用いて組み合わせられ、出力特徴マップ 3 0 8 における出力活性化 3 1 8 を生成することができる。出力特徴マップ 3 0 8 は、入力特徴マップ 3 0 6 のより高レベルの抽象化を表し、高さ E および幅 F を有する。このモデル 3 1 0 では、全ての出力に関して同じ重み集合を共有することができる（たとえば、フィルタ空間は、不変である）。

30

【 0 0 2 6 】

[0042]図 3 B は、画像処理に適用された、畳み込みニューラルネットワークのモデルの一例 3 2 0 を示す。畳み込みニューラルネットワークは、複数の畳み込み層を包含することができる。畳み込みニューラルネットワークでは、各層は、入力データの（すなわち、入力特徴マップ 3 2 6 の）連続的高レベルの抽象化を生成することができる。畳み込みニューラルネットワークは、深い階層の層を採用することによって、極めて高いパフォーマンスを実現することができる。

40

【 0 0 2 7 】

[0043]図 3 B の例で示すように、畳み込みニューラルネットワークの各畳み込み層は、高次元の畳み込みから構成される。このモデル 3 2 0 において、一層の入力活性化 3 3 6 は、各々がチャンネル C と称される 2 次元入力特徴マップ 3 2 6 の集合として構造化されている。各チャンネルは、各チャンネルに 1 つのフィルタを有するフィルタスタック 3 2 4 からの特定の 2 次元フィルタで畳み込まれる。フィルタスタック 3 2 4 は、単一の 3 次元フィルタと呼ぶことができる。各点の畳み込みの結果は、全てのチャンネルに渡って合計され、

50

合わせて出力特徴マップ 3 2 8 のうちの 1 つのチャンネル M を形成する出力活性化 3 3 8 を生成する。この同じ入力上では、出力チャンネルの数に対応する追加的な 3 次元フィルタ M を用いて、追加の出力チャンネルを生成することができる。フィルタ重みの再使用を増進するために、図示の例において 1 ~ N でラベル付けされている複数の入力特徴マップ 3 2 6 をバッチ処理することができる。

【 0 0 2 8 】

[0044] 畳み込みニューラルネットワークは、5 ~ 1 0 0 0 超の範囲の層を包含することができる。例によっては、分類を目的として、畳み込み層の後へ 1 ~ 3 層などの少数の完全に結合された層を付すことができる。完全に結合された層は、入力特徴マップにフィルタを付すこともできるが、これらのフィルタのサイズは、入力特徴マップと同じである。

10

【 0 0 2 9 】

[0045] ニューラルネットワークの訓練は、オンラインで、すなわち、ニューラルネットワークが動作していてユーザが利用できる状態において生じ得る。しかしながら、より多くの場合、訓練は、オフラインで、かつニューラルネットワークが稼働される前に生じる。訓練サンプルセットは、膨大であり得、よって、訓練には数時間または数日を要することがある。オフライン訓練は、潜在的に、より正確な結果も生み出し得る。

【 0 0 3 0 】

[0046] ニューラルネットワークは、訓練されると、訓練中に決定された重みと、ネットワークの各層またはノードで実行されるべき計算を記述する命令セットとを包含する。例によっては、重みの数は、約 5 百万 ~ 1 億になることがある。例によっては、重み値を 3 2 ビットの数値を用いて表すことができるが、この場合、5 百万 ~ 1 億個の重みは、約 2 0 メガバイト (M B) ~ 4 0 0 M B の記憶容量を必要とし得る。例によっては、重みの数が僅か 1 5 0 万であることもある。

20

【 0 0 3 1 】

[0047] ニューラルネットワークの動作 (たとえば、推論の実行) は、上述のモデルが示すように、入力データまたは入力活性化をフェッチすることと、層内の各ノードについて積和演算を並列して実行することと、出力活性化を提供すること、を含む。応答時間で測定されるニューラルネットワークの最適パフォーマンスは、ハードウェアアーキテクチャが高度に並列化された計算を実行できる場合に達成可能である。汎用処理装置と呼ばれることもある中央処理装置 (C P U) は、複数のコア (たとえば、2 ~ 6 4 個またはこれを超えるコア) を有し得、かつ複数の実行スレッドの使用によって並列処理を高めることができる。しかしながら、C P U コアは、逐次処理用に最適化される傾向がある。たとえば、コアの計算エンジン (たとえば、算術論理演算装置 (A L U)) は、メモリからオペランドを取得して結果をメモリに書き込み、よって、逐次計算にはメモリ操作が必要である。この例において、各メモリ操作は、C P U の制御論理による管理を必要とすることがある。したがって、このために、かつ他の理由により、C P U は、ニューラルネットワークに関して推論を実行する場合に応答時間が遅くなる傾向がある。

30

【 0 0 3 2 】

[0048] C P U とは対照的に、グラフィック処理装置 (G P U) は、具体的には並列計算を実行するように構成される何千もの小さく効率的なコアを備えることによって並列処理を実現する。したがって、G P U は、ニューラルネットワークの実行に際して、C P U より遙かに優れたパフォーマンスを達成することができる。しかしながら、個々の G P U 計算エンジンは、本質的にやはり逐次的である可能性があり、よって、ある計算エンジンの出力を別の計算エンジンの入力へ提供するには、メモリ操作が必要である。

40

【 0 0 3 3 】

[0049] ニューラルネットワークを実行する場合、C P U および G P U の双方が遭遇し得るパフォーマンスのボトルネックは、メモリへのアクセスに存する。積和演算は、各々が重み値、入力特徴マップの活性化および部分和をフェッチする 3 つのメモリ読取りと、更新された部分和を記憶するための 1 つのメモリ書込みとを必要とし得る。最悪の事例では

50

、全てのメモリトランザクションがオフチップメモリへ、すなわちプロセッサとは別のダイ上かつ別のパッケージ内に位置決めされるメモリへ送られる。このメモリは、プロセッサメモリまたはメインメモリと称されることもあり、プロセッサが専ら、プロセッサによってアクティブに作用されているデータを一時的に記憶するためのものであり得る。プロセッサメモリには、大容量かつ低コストであるという理由で、ダイナミック・ランダム・アクセス・メモリ(DRAM)またはDRAMの変形がよく使用される。しかしながら、プロセッサメモリの読取りおよびこれへの書込みは、計算エンジンの動作より何桁分も低速である。したがって、ニューラルネットワークの速度は、オフチップメモリの待ち時間によって制限される可能性がある。

【0034】

[0050]専用ニューラルネットワークプロセッサは、ニューラルネットワークの実行に際して、CPUおよびGPUの双方より優れたパフォーマンスを達成することができる。ニューラルネットワークプロセッサは、計算エンジンが処理チェーンを形成しかつある計算エンジンから別の計算エンジンへ直にデータを渡すことができる空間アーキテクチャを採用することができる。これにより、メモリトランザクションの数を大幅に減らすことができる。例によっては、ニューラルネットワークプロセッサは、プロセッサメモリから読み取られる値を記憶し得る、かつ値をプロセッサ内の複数の計算エンジンへ配信できるオンチップバッファも含むことができる。計算エンジンは、さらに、中間結果を記憶するための小型ローカル・レジスタ・ファイル(たとえば、小型メモリ)を含むことができる。オンチップメモリ階層の保有は、メモリ待ち時間の短縮により、ニューラルネットワークの動作効率を高めることができる。

【0035】

[0051]それでもやはり、ニューラルネットワークプロセッサは、ニューラルネットワークの重み値がオフチップで記憶される場合に、メモリ帯域幅が制限されることになり得る。ニューラルネットワークプロセッサの計算行列が計算を実行できる速度は、メモリから重み値および活性化を読み取ることができる速度をすぐに超える可能性がある。たとえば、計算行列は、クロックサイクルごとに10,000回の積和演算を実行することができ、よって、サイクル当たり30,000個の入力値を必要とする。プロセッサメモリバスのクロック速度は、たとえば数千メガヘルツ(MHz)の範囲内であり得るのに対して、プロセッサのクロック速度は、ギガヘルツ(GHz)の倍数であり得る。したがって、ニューラルネットワークプロセッサの計算速度は、プロセッサメモリのデータ供給能力をすぐに追い越してしまう可能性がある。

【0036】

[0052]重み値の再使用は、メモリ帯域幅の制限を回避し得る1つの方法である。再使用は、畳み込みニューラルネットワークでは一般的であって、重み値を、たとえば平均1300回再使用することができる。後に詳述するように、重み値を頻繁に再使用するニューラルネットワークは、潜在的にメモリ帯域幅の制限を回避することができるが、代わりにプロセッサの計算速度が制限される可能性がある。

【0037】

[0053]長・短期記憶ニューラルネットワークおよび多層パーセプトロンニューラルネットワークでは、重み値の再使用率は、たとえば平均2回などと遙かに低い。

【0038】

[0054]重み値の再使用を増やすために使用される1つのソリューションが、バッチ処理である。バッチ処理は、一度に2つ以上の入力データセットをニューラルネットワークに入力することを含む。入力データセットが関連づけられる必要はない。バッチ処理では、ニューラルネットワークにたとえば10セットの入力データが提供される場合、各重みは、メモリから一旦読み取られると、その後20回(たとえば、入力データセットにつき2回)再使用が可能である。

【0039】

[0055]しかしながら、数学モデルは、ニューラルネットワークプロセッサが可能な最大

10

20

30

40

50

限のパフォーマンスを実現するには、高い再使用率が必要であることを示唆している。たとえば、一部の例は、約1000の再使用率が必要であることを示唆している。バッチ処理に際しては、たとえば一度に50～60セットの入力データを収集することが可能である場合があるが、500セットの入力データを収集すると、他の問題が生じることがある。たとえば、ニューラルネットワークのユーザは、たとえば機械翻訳または画像識別を要求する場合には迅速な応答を期待する。ニューラルネットワーク処理システムが、結果の計算を開始するまでに500もの要求の保有を待機する場合、応答時間に悪影響が及ぶ可能性がある。

【0040】

[0056]様々な実装において、ニューラルネットワーク処理システムは、ニューラルネットワークの重みをオンチップメモリに記憶することによって、メモリ帯域幅の制限を減らすことができ、よって最適効率に近づくことができる。オンチップとは、メモリが、計算行列と同じダイ上および/または同じパッケージ(たとえば、ダイの物理エンクロージャ)内にあることを意味する。ニューラルネットワークプロセッサは、中間結果を記憶するためのオンチップメモリを有し得る。様々な実装において、プロセッサのメモリサブシステムは、オンチップメモリが中間結果および重み値の双方を記憶できるように設計されることが可能である。ニューラルネットワークプロセッサは、それでもメモリバウンドである場合があるが、オンチップメモリを、たとえばオフチップメモリより10～50倍は速く読み取ることが可能であり得る。メモリの遅延をこれだけ減らせば、ニューラルネットワークの動作をプロセッサの計算速度限界まで近づけることが可能であり得る。

【0041】

[0057]事例によっては、特に小型ニューラルネットワークの場合、ニューラルネットワークの重み値を全てオンチップメモリに記憶することが可能であり得る。しかしながら、単一のモノリシックメモリを用いる場合は、単一のメモリが保有する読み取りおよび書き込みチャネルセットがたとえば1つまたは2つでしかないことがあり、よって一度に1つまたは2つの値しか読み取れないという理由で、やはりメモリ遅延が生じる場合がある。様々な実装では、1つの大型メモリの代わりに、ニューラルネットワークプロセッサに複数のメモリバンクを装備することができ、これらは個々にアクセス可能であり得る。独立してアクセス可能であることにより、同時に2つ以上のメモリバンクを読み取ることが可能であり得る。

【0042】

[0058]ニューラルネットワークの処理エンジンでは、計算行列を処理エンジンアレイとして実装することができる。ニューラルネットワークの処理エンジンは、さらに、処理エンジンアレイに対してローカルなメモリバンクセットを含むことが可能であり、ここで、ローカルとは、処理エンジンアレイに物理的に近い、かつ/または処理エンジンアレイによって直にアクセス可能であることを意味し得る。先に述べたように、ローカル・メモリ・バンクは、ニューラルネットワークの処理エンジンにより、中間結果を記憶するために使用可能である。事例によっては、特にニューラルネットワークが小型である場合、ニューラルネットワークの重み値を全てニューラルネットワーク処理エンジンのメモリバンクに記憶することもできる。これらの事例では、処理エンジンアレイは、全てのクロックサイクルにおける完全な利用を持続することが可能であり得る。

【0043】

[0059]例によっては、ニューラルネットワークの重み値が全てニューラルネットワーク処理エンジンのメモリバンクに収まるわけではない。たとえば、メモリバンクには、重み値の半分に足るスペースを有し得、残りのスペースは、入力データセットの処理中に計算される中間結果を記憶するために必要とされる。しかしながら、中間結果のサイズは、結果を計算する過程で減少する可能性がある。さらに、重み値の中には、一度使用されると不要になり得るものがある。したがって、実装によっては、計算が進行してメモリ空間が利用可能になるにつれて、ニューラルネットワークの処理エンジンは、追加の重みを利用可能なスペースにロードすることができる。事例によっては、重みは、オフチップメモリ

10

20

30

40

50

から取得されることが可能である。事例によっては、重みは、オンチップメモリから、たとえば、別のニューラルネットワーク処理エンジンのメモリバンクから、取得されることが可能である。

【0044】

[0060]実装によっては、ニューラルネットワークプロセッサは、各々が独立した処理エンジンアレイとローカル・メモリ・バンクとを有する複数のニューラルネットワーク処理エンジンで構築されることが可能である。これらの実装では、各ニューラルネットワーク処理エンジンが1つのニューラルネットワークを実行でき、よって、複数のニューラルネットワークが同時に実行されることが可能である。実装によっては、あるニューラルネットワークに関する重み値は、そのうちの1つがこのニューラルネットワークを処理するためのエンジンとして指定されている2つ以上のニューラルネットワーク処理エンジンのメモリバンクに記憶されることが可能である。指定されたニューラルネットワーク処理エンジンが、別のニューラルネットワーク処理エンジンによって記憶されている重みを必要とする場合、これらの重みは、他のニューラルネットワーク処理のメモリバンクから読み取られ、かつ指定されたニューラルネットワーク処理エンジンのメモリバンクへロードされることが可能である。他のニューラルネットワーク処理エンジンは、その固有のメモリバンクにおける残りのあらゆる利用可能なスペースを他の動作に用いることができる。

10

【0045】

[0061]実装によっては、あるニューラルネットワークプロセッサから別のニューラルネットワークプロセッサへ重みを移動させる代わりに、計算を移動させることができる。たとえば、中間結果（たとえば、ある層からの出力活性化）および状態（たとえば、計算された最後の層）は、1つのニューラルネットワーク処理エンジンから第2のニューラルネットワーク処理エンジンへコピーされることが可能であり、この場合、第2のニューラルネットワーク処理エンジンは、そのメモリバンクに、計算の継続に必要な次の重み値セットを有する。第2のニューラルネットワーク処理エンジンは、計算を再開し、かつおそらくは、計算をさらに別のニューラルネットワーク処理エンジンへ引き渡すことができる。

20

【0046】

[0062]進行中の計算の、あるニューラルネットワーク処理エンジンから別のニューラルネットワーク処理エンジンへの転送は、実装によっては、個々のニューラルネットワークプロセッサ間での転送を含む。これらの実装では、個々のニューラルネットワークプロセッサは、異なるダイ上および/または異なるパッケージ内に存在し得る。同じくこの実施例において、ニューラルネットワークプロセッサは、ホストバスまたはプロセッサバスを用いて通信することができる。ニューラルネットワーク処理エンジンが同じダイ上に存在する場合のように、中間結果および状態をコピーすると、あるニューラルネットワークプロセッサから別のニューラルネットワークプロセッサへ計算を移動させることができる。

30

【0047】

[0063]様々な実装において、あるニューラルネットワーク処理エンジンから別のニューラルネットワーク処理エンジンへの重みのコピー、ニューラルネットワーク処理エンジン間および/または物理ニューラルネットワークプロセッサチップ間での進行中の計算の移動は、ニューラルネットワークの重み値をオンチップで可能な限り多く記憶することを目的として、様々な組合せで使用されることが可能である。重み値をオンチップで有することにより、計算は、オフチップメモリの比較的長い待ち時間によって制限される代わりに、単にオンチップメモリの比較的短い待ち時間によって制限され得る。その結果、ニューラルネットワークの動作を遙かに効率的なものにすることができる。

40

【0048】

[0064]図4は、ニューラルネットワークの重み値をオフチップメモリではなくオンチップで記憶する効果の一例を示す。図4に示すグラフ400は、ルーフラインモデルと呼ばれるもののアプリケーションを示している。ルーフラインモデルは、コンピューティングシステムのパフォーマンス推定を提供するために使用可能な性能モデルである。ルーフラインモデルは、固有のハードウェア限度および最適化の潜在的利点を捉えることができる

50

。図4の例において、ルーフラインモデルは、ニューラルネットワークプロセッサの、メモリから読み取られる重み当たりの演算回数に関する性能を示すために使用されている。縦軸は、1秒間に実行できる兆単位の演算数 (t e r a o p s) を示す。横軸は、重み値当たりで実行される演算または計算の回数を示す。重み値当たりで実行される演算の回数は、重み固有の再使用 (たとえば、ニューラルネットワークの構造は、重みの再使用に繋がる)、またはバッチ処理、すなわち複数のデータセットをニューラルネットワークへ同時に、またはパイプライン式に入力すること、のいずれかを介して増加し得る。

【0049】

[0065]図4の例において、グラフ400に描かれている実線410は、重み値をオフチップメモリに記憶するニューラルネットワーク処理システムの性能の一例を示している。このようなシステムでは、重み値がプロセッサメモリに記憶されていて、ニューラルネットワークプロセッサがホストバスまたはプロセッサバス上で重み値を読み取る。重み値を別個のメモリに記憶することにより、ニューラルネットワーク処理システムは、メモリから重み値が読み取られるたびに必ず遅延を被ることになる。

10

【0050】

[0066]実線410の急勾配部402では、重み値当たりで実行される演算回数に対して、1秒間に実行できるテラ演算数がほぼ線形的に増加する。実線410の急勾配部402では、1秒当たりのテラ演算数を増やすために、所与の重みの再使用が増加されなければならない。逆に言えば、実線410の急勾配部402では、所与の再使用値において、1秒当たりのテラ演算数は、オフチップメモリから重み値を読み取れる速度によって制約される。したがって、ニューラルネットワーク処理システムは、実線410の急勾配部402においてメモリバウンドであるとされる。

20

【0051】

[0067]実線410の平坦部404では、1秒当たりのテラ演算数が最大に達している (図示の例では、1秒当たり約124兆回の演算)。実線410の平坦部404では、システムがハードウェアの最大計算速度に達していることから、重み値の再使用が増加しても、1秒当たりのテラ演算数は増加しない。実線410の平坦部404において、システムは、コンピュータバウンドであるとされる。

【0052】

[0068]したがって、このルーフラインモデルは、ニューラルネットワーク処理システムの場合、ハードウェアで可能な最大パフォーマンスを実現できることを示している。しかし、グラフ400に示すように、重みがプロセッサとは別のチップに記憶されているシステムの場合、各重み値を何度も再使用すること (図示の例では、約1600回)が必要である。先に述べたように、畳み込みニューラルネットワークは、高い重み再使用率を有し得るが、多層パーセプトロンおよび回帰型ニューラルネットワークなどの他の多くのニューラルネットワークの重み再使用率は、極めて低い。これらのタイプのニューラルネットワークの入力データセットのバッチ処理は、再使用を増やし得るが、1100個の入力データセット (各重み値に固有の再使用を2回と想定)をバッチ処理すると、入力データを提供するユーザに対する応答時間の遅れなどの他の有害な影響が生じる可能性がある。

30

【0053】

[0069]図4の例における破線420は、重み値がオンチップで記憶されるニューラルネットワーク処理システムに適用されたルーフラインモデルを示す。この事例において、重み値は、システムが入力データを受信する前、およびおそらくはシステム動作中のあらゆる時点で、ニューラルネットワークプロセッサのメモリに記憶される。重み値をオンチップで記憶することは、メモリからの重み値の読取りに対する遅延が大幅に短縮されることを意味する。破線420の急勾配部412が示すように、(本例が提供する)僅か約100回の再使用でメモリの待ち時間を克服することができる。したがって、ニューラルネットワーク処理システムは、ほとんどの時間および全ての型式のニューラルネットワークタイプについて、破線420の平坦部414で動作することができる。

40

【0054】

50

[0070]様々な実装において、ニューラルネットワークプロセッサは、重み値をプロセッサ上へ記憶することをサポートするように構築されることが可能である。図5は、ニューラルネットワークを実行するために使用可能なニューラルネットワーク処理エンジン502を示す。様々な実装において、この例示的なニューラルネットワーク処理エンジン502は、ニューラルネットワークプロセッサに含まれ得る集積回路である。ニューラルネットワークプロセッサは、追加的なニューラルネットワーク処理エンジン502を含む他のコンポーネントを有し得る集積回路である。様々な実装において、ニューラルネットワーク処理エンジン502は、メモリサブシステム504と、処理エンジンアレイ510とを含み得る。動作中(たとえば、入力データ550のセットの結果を計算しているとき)、処理エンジンアレイ510は、メモリサブシステム504から重み506および状態508の値を読み取ることができる。処理エンジンアレイ510は、計算結果を結果バッファ512へ出力することができる。事例によっては、この例示的なニューラルネットワーク処理エンジン502は、結果がメモリサブシステム504へ書き込まれる前に、計算結果に対して、活性化ブロック516を用いる活性化関数、および/またはプーリングブロック518を用いるプーリングを実行することができる。

10

【0055】

[0071]重み506は、本例では、ニューラルネットワークの重み値である。様々な実装において、重み506は、重み506の値が先に決定されていることを意味する訓練後の重みである。状態508は、本例では、計算開始時の入力データ550だけでなく、進行中の計算を反映する値も包含することができる。状態508は、たとえば、動作中のニューラルネットワークのカレント層である処理エンジンアレイ510によって決定される部分重み、および/または処理エンジンアレイ510に対する命令を包含し得、ここで、該命令は、特定の一層に関連づけられてもよい。重み506および状態508は、処理エンジンアレイ510によって作用されるべくメモリサブシステム504から読み取られ得る。

20

【0056】

[0072]様々な実装において、メモリサブシステム504は、複数のメモリバンク514を含み得る。これらの実装において、各メモリバンク514は、独立してアクセス可能であるが、これは、1つのメモリバンクの読取りが別のメモリバンクの読取りに依存しないことを意味する。同様に、あるメモリバンクへの書込みが、異なるメモリバンクへの書込みに影響を与えることも、これを制限することもない。事例によっては、各メモリバンクの読取りおよび書込みを同時に行うことができる。独立してアクセス可能なメモリバンク514を有するためには、様々な技術を用いることが可能である。たとえば、各メモリバンクは、少なくとも1つの読取りチャネルを有することができ、かつ少なくとも1つの別個の書込みチャネルを有してもよい(事例によっては、メモリバンクは、読取りチャネルと書込みチャネルとの組合せを有することがある)。これらの例において、メモリサブシステム504は、複数のメモリバンクの読取りチャネルまたは書込みチャネルへの同時アクセスを許容することができる。別の例として、メモリサブシステム504は、たとえば複数のメモリバンク514の出力間のアービトレーションによって2つ以上のメモリバンクの出力が認可されることになるように、アービトレーション論理を包含することができる。これらの例および他の例において、各メモリバンクは、全体としてはメモリサブシステム504によって管理されるものの、他のメモリバンクとは独立して動作されることが可能である。

30

40

【0057】

[0073]メモリバンク514を独立してアクセス可能とすることにより、ニューラル処理エンジン502の効率は、高まり得る。たとえば、重み506および状態508は、同時に読み取られて処理エンジンアレイ510の各口へ提供されることが可能であり、よって、処理エンジンアレイ510全体を1サイクルで起動することができる。別の例として、重み506および状態508は、中間結果がメモリサブシステム504へ書き込まれると同時に読み取られることが可能である。これに対して、単一のメモリは、重み506および状態508を処理エンジンアレイ510へなおもオフチップメモリより高速で提供で

50

きるものの、一度に提供できるものは1つの読取りまたは書込みであることがある。単一のメモリの場合、たとえば、処理エンジンアレイ510を開始できるまでに処理エンジンアレイ510の各ロウの重みを読み取るべく、複数のクロックサイクルが必要とされる可能性がある。

【0058】

[0074]様々な実装において、メモリサブシステム504は、処理エンジンアレイ510、活性化ブロック516、プーリングブロック518、およびチップ相互接続520上でメモリサブシステム504へアクセスするあらゆる外部クライアントを含む複数のクライアントへ同時にサービスを提供するように構成されることが可能である。実装によっては、複数のクライアントへサービスを提供できることは、メモリサブシステム504が少なくとも、存在するクライアントの数と同数のメモリバンクを有することを意味し得る。事例によっては、処理エンジンアレイ510の各ロウは、別々の読取りクライアントと見なされ得る。これらの事例において、重み506および状態508は、別々に記憶され、よって2つの読取りを要する場合もあれば、連結されて共に記憶され、よって1つの読取りを要する場合もある。事例によっては、処理エンジンアレイ510の各カラムは、中間値を出力することができ、よって、各カラムが別々の書込みクライアントと見なされ得る。事例によっては、処理エンジンアレイ510からの出力は、メモリバンク514に書き込まれることが可能であり、メモリバンク514は、続いて、処理エンジンアレイ510に入力データを提供することができる。メモリバンク514は、たとえば、スタティック・ランダム・アクセス・メモリ(SRAM)を用いて実装されることが可能である。

10

20

【0059】

[0075]様々な実装において、メモリサブシステム504は、制御論理を含み得る。制御論理は、たとえば、メモリバンク514の各々のアドレス空間を把握し、必要に応じて、読み取る、または書き込むべきメモリバンク514を識別し、かつ/または必要であればメモリバンク514間でデータを移動することができる。実装によっては、メモリサブシステム504は、特定のクライアントへのメモリバンクを出力するか、かつ/または特定のクライアントからどのメモリバンクの入力を受け入れるかを選択するためのマルチプレクサを包含することができる。これらの実装において、制御論理は、メモリバンク514の一部または全てが各クライアントにサービスを提供することを可能にし得る、マルチプレクサのための選択信号を生成することができる。実装によっては、メモリバンク514は、特定のクライアントへ結線で接続されることが可能である。たとえば、メモリバンク514のセットは、重み506および状態508を提供すべく処理エンジンアレイ510のロウへ結線で接続されることが可能である。これらの例において、制御論理は、たとえば、中間結果が書き込まれるメモリバンク514から、次の計算ラウンドのために中間結果が読み取られるメモリバンク514へ中間結果を移動するために、データをメモリバンク514間で移動させることができる。

30

【0060】

[0076]処理エンジンアレイ510は、ニューラル処理エンジン502の計算行列である。処理エンジンアレイ510は、たとえば、とりわけ並列積分、畳み込み、相関および/または行列乗算を実行することができる。処理エンジンアレイ510は、ロウおよびカラムに配列される複数の処理エンジン511を含み、よって、1つの処理エンジン511により出力される結果は、別の処理エンジン511へ直に入力されることが可能である。したがって、処理エンジンアレイ510の外縁上にはない処理エンジン511は、メモリサブシステム504ではなく、他の処理エンジン511から作用すべきデータを受け入れることができる。

40

【0061】

[0077]様々な例において、処理エンジンアレイ510は、シストリック実行を用い、この場合、データは、各処理エンジン511に異なる方向から定間隔で到着する。例によっては、入力データは、処理エンジンアレイ510に左から流れ込むことができ、重み値は、上部にロードされることが可能である。例によっては、重みおよび入力データは、左か

50

ら流れることができ、かつ部分和は、上から下へ流れることができる。これらの例および他の例において、積和演算は、処理エンジンアレイ 5 1 0 を介して対角波面として移動し、データは、アレイを右へ、かつ下へ移動する。制御信号は、重み 5 0 6 と同時に左側で入力されることが可能であって、計算と共に横へかつ下へ流れることができる。

【0062】

[0078]様々な実装において、処理エンジンアレイ 5 1 0 におけるカラムの数は、処理エンジンアレイ 5 1 0 の計算能力を決定し、かつロウの数は、処理エンジンアレイ 5 1 0 の最大利用を達成するために必要なメモリ帯域幅を決定する。処理エンジンアレイ 5 1 0 は、たとえば、64 のカラムおよび 256 のロウ、または他の何らかの数のカラムおよびロウを有し得る。

10

【0063】

[0079]図 5 には、処理エンジン 5 1 1 の一例が示されている。この例が示すように、処理エンジン 5 1 1 は、乗数累算器回路を包含することができる。左側からの入力には、たとえば、入力データ i および重み値 w が含まれ得、ここで、入力データは、入力データのセットから取得される値であり、重み値は、ニューラルネットワークの一層を次の層へ連結する重み値セットからの値である。入力データのセットは、たとえば、とりわけ識別またはオブジェクト認識のために提出される画像、音声認識のために提供されるオーディオクリップ、自然言語処理または機械翻訳のためのテキスト文字列、または次の移動を決定するために分析を要求するゲームのカレント状態であってもよい。例によっては、入力データおよび重み値は、次の処理エンジン 5 1 1 への入力として右側へ出力される。

20

【0064】

[0080]図示の例において、上記からの入力は、別の処理エンジン 5 1 1 から、または処理エンジンアレイ 5 1 0 による先の計算ラウンドから、のいずれかで提供される部分 p_in を包含し得る。新しい入力データセットの計算を開始する場合、処理エンジンアレイ 5 1 0 の最上ロウは、ゼロなどの p_in の固定値を受け取ることができる。本例が示すように、 i および w は互いに乗算されて結果が p_in と合計され、別の処理エンジン 5 1 1 へ入力されることが可能な新しい部分 p_out が生成される。処理エンジン 5 1 1 の実装は、他にも様々なものが可能である。

【0065】

[0081]処理エンジンアレイ 5 1 0 の最終ロウからの出力は、一時的に結果バッファ 5 1 2 に記憶されることが可能である。結果は、中間結果であり得、これは、追加計算のために処理エンジンアレイ 5 1 0 へ提供されるべくメモリバンク 5 1 4 へ書き込まれることが可能である。あるいは、結果は、最終結果であり得、これは、メモリバンク 5 1 4 へ書き込まれると、メモリサブシステム 5 0 4 からチップ相互接続 5 2 0 上で、該システムにより出力されるべく読み取られることが可能である。

30

【0066】

[0082]実装によっては、ニューラル処理エンジン 5 0 2 は、活性化ブロック 5 1 6 を含む。これらの実装において、活性化ブロック 5 1 6 は、処理エンジンアレイ 5 1 0 からの結果を組み合わせて 1 つまたは複数の出力活性化にすることができる。たとえば、畳み込みニューラルネットワークの場合、複数のチャネルからの畳み込みを合計して、単一チャネルのための出力活性化を生成することができる。他の例では、ニューラルネットワーク内の単一ノード用の出力活性化を生成するために、処理エンジンアレイ 5 1 0 内の 1 つまたは複数のカラムからの結果を累積することが必要とされる場合がある。実装によっては、活性化ブロック 5 1 6 がバイパスされてもよい。

40

【0067】

[0083]実装によっては、ニューラル処理エンジン 5 0 2 は、プーリングブロック 5 1 8 を含む。プーリングは、ニューラルネットワークの一層からのノードクラスタの出力を組み合わせることである。組み合わせられた出力は、次の層へ提供されることが可能である。組合せは、たとえば、最大値、最小値、平均値、中央値、またはノードクラスタの出力から決定される他の何らかの値を計算することを含み得る。様々な実装において、プーリン

50

グ 5 1 8 は、特定のニューラルネットワークにとっての必要性に応じて、選択的に活性化されることが可能である。

【 0 0 6 8 】

[0084]入力データ 5 5 0 は、チップ相互接続 5 2 0 を介して到着することができる。チップ相互接続 5 2 0 は、ニューラル処理エンジン 5 0 2 を、入力/出力 (I / O) デバイス、ストレージドライブまたはネットワークインタフェースから入力データ 5 5 0 を入手できるダイレクト・メモリ・アクセス (DMA) エンジンなどのニューラルネットワークプロセッサの他のコンポーネントへ接続することができる。入力データ 5 5 0 は、たとえば、文字列または数列などの 1 次元データ、または画像のピクセル値、またはオーディオ信号の経時的な周波数および振幅値のアレイなどの 2 次元データであり得る。例によつては、入力データ 5 5 0 は、たとえば、自動運転車によって使用される状況情報の場合がそうであり得るように、3 次元的であり得る。実装によっては、メモリサブシステム 5 0 4 は、入力データ 5 5 0 用に別個のバッファを包含することができる。実装によっては、入力データ 5 5 0 は、重み 5 0 6 と共にメモリバンク 5 1 4 に記憶されることが可能である。

10

【 0 0 6 9 】

[0085]様々な実装において、メモリサブシステム 5 0 4 に記憶される重み 5 0 6 は、1 つまたは複数のタスクを実行するようにニューラルネットワークを訓練することによって決定されている可能性がある。入力データ 5 5 0 は、実行すべきタスク (たとえば、画像処理、音声認識、機械翻訳、他) を示す命令を含み得る。様々な実装において、ニューラル処理エンジン 5 0 2 は、ニューラルネットワークを訓練するためではなく、推論を実行する (たとえば、タスクを実行する) ように構成される。実装によっては、ニューラル処理エンジン 5 0 2 が訓練に使用されることがあるが、おそらくは、記憶された重み 5 0 6 を更新するにはソフトウェアから支援される。

20

【 0 0 7 0 】

[0086]様々な実装において、メモリサブシステム 5 0 4 は、ニューラルネットワークの中間結果および全ての重み値の双方を記憶するに足るメモリを包含し得る。メモリサブシステム 5 0 4 は、最小でも、メモリサブシステム 5 0 4 内に中間結果を記憶するに足るメモリを有するべきであるが、多くの場合、メモリサブシステム 5 0 4 は、単に中間結果を記憶するために必要とされるものより遙かに多いメモリバンク 5 1 4 を含み得る。この追加スペースは、ニューラルネットワークの全ての重み値を記憶するために使用されることが可能である。たとえば、ニューラルネットワークは、1 5 0 万個の重みを有することがあるが、その各々が 3 2 ビットで表される場合、これにより、約 6 M B のメモリが必要となり得る。中間結果は、たとえば、最大 1 0 M B の記憶スペースを要し得る。2 0 M B のオンチップメモリは、適度なサイズであり、よつて、先の例では、重み値、中間結果、およびニューラル処理エンジン 5 0 2 が動作中に必要とし得る他のあらゆるデータに容易に対応することができる。

30

【 0 0 7 1 】

[0087]事例によっては、ニューラルネットワークの重みが全てニューラル処理エンジンのメモリに収まるわけではない。これらの事例では、重みの一部がニューラル処理エンジンのメモリに記憶され、かつ残りの重みは、計算の過程でロードされることが可能である。図 6 A ~ 6 B は、ニューラルネットワークを実行するために使用可能なニューラル処理エンジンの一例 6 0 2 を示す。様々な実装において、ニューラルネットワーク処理エンジン 6 0 2 は、メモリサブシステム 6 0 4 と、処理エンジンアレイ 6 1 0 とを含み得る。動作中 (たとえば、入力データ 6 5 0 のセットの結果を計算しているとき)、処理エンジンアレイ 6 1 0 は、メモリサブシステム 6 0 4 から重み 6 0 6 および状態 6 0 8 の値を読み取ることができ、ここで、状態 6 0 8 の値は、当初、入力データ 6 5 0 を含み得る。図示されていないが、処理エンジンアレイ 6 1 0 は、処理エンジンアレイを包含し得、これらの処理エンジンの各々が、たとえば積和演算を実行することができる。処理エンジンアレイ 6 1 0 は、計算結果を結果バッファ 6 1 2 へ出力することができる。事例によっては、この例示的なニューラルネットワーク処理エンジン 6 0 2 は、結果がメモリサブシステム

40

50

604へ書き込まれる前に、計算結果に対して、活性化ブロック616を用いる活性化関数、および/またはプーリングブロック618を用いるプーリングを実行することができる。

【0072】

[0088]動作において、ニューラルネットワーク処理エンジン602は、メモリサブシステム604のメモリバンク614から重み606および状態608を読み取り、重み606および状態608を処理エンジンアレイ610へ入力する。処理エンジンアレイ610は、中間結果を結果バッファ612へ出力することができ、中間結果は、結果バッファ612からメモリサブシステム604へ書き込まれることが可能である。複数の中間結果を組み合わせる1つの出力活性化にする場合は、活性化ブロック616が選択的に有効化され得る。また、プーリングブロック618も、ノードクラスタの出力を組み合わせる1つの出力にする場合に、選択的に起動され得る。重み606および状態608を読み取って結果を計算するサイクルは、ニューラルネットワークの最終層に到達して処理エンジンアレイ610が最終結果を出力するまで、追加的に生じ得る。

10

【0073】

[0089]様々な実装において、メモリバンク614のセットは、当初、中間結果（たとえば、状態608の値）を記憶するために予約され得る。たとえば、図6Aに示すように、メモリバンク614の半分は、即時的な結果を記憶するために予約され得る。この例において、重み606は、残りのメモリバンク614に記憶され得る。重み606は、ニューラルネットワーク処理エンジン602が入力データ650を受信するより前に記憶されることが可能であり、よって、入力データ650を受信されるとすぐに重み606を使用可能である。他の例では、メモリバンク614を、状態608の記憶と重み606の記憶との間で比率を変えて分割することができる。実装によっては、個々のメモリバンク614を状態608の記憶と重み606の記憶との間で分割することができる。

20

【0074】

[0090]図6Aの例において、メモリバンク614に記憶される重み606は、ニューラルネットワークの全ての重みより少ない。残りの重みは、ニューラルネットワーク処理エンジン602の外部に位置決めされる追加のメモリに記憶される。追加のメモリ630は、たとえば、オンチップキャッシュ内などのオンチップであり得る。あるいは、または追加的に、追加のメモリ630は、たとえば、プロセッサメモリ内などのオフチップであり得る。先に論じたように、重み606は、ニューラルネットワークの層に対応する。様々な実装において、重み606は、層に従って分割されることが可能であり、よって、最初のn個の層の重み606は、入力データ650を受信される前にメモリサブシステム604に記憶され、残りのm個の層の重み606は、追加のメモリ630に記憶される。

30

【0075】

[0091]ニューラルネットワーク処理エンジン602が入力データ650の処理を開始すると、状態608の値を記憶するために必要なメモリの量が減る可能性がある。さらに、一部のニューラルネットワークでは、重み606の中には、一度使用されるとその後は不要になり得るものがある。たとえば、先に論じたように、フィードフォワード重みは、ある層の入力値の計算に一度しか使用されないことがある。重みの中には、一度使用されると二度と必要とされないものがあることから、入力データ650のセットを処理する過程で、これらの重みを記憶するために使用されるメモリは、別の目的で使用されることが可能である。

40

【0076】

[0092]二度と必要とされない状態608の値および重み606を記憶するために先に使用されたメモリバンク614内のスペースは、追加の重み606の値を記憶するように別の目的で使われることが可能である。図6Bに示すように、計算が進むにつれて、より多くのメモリバンク614が重み606を記憶するように指定され得る。

【0077】

[0093]たとえば、あらゆるクロックサイクル上で、重み606のセットは、様々なメモ

50

リバンク 6 1 4 から読み取られることが可能であって、中間結果は、状態 6 0 8 の値として記憶され得る。メモリサブシステム 6 0 4 は、制御論理を用いて、まだ使用されたことのない重み 6 0 6 により使用されるメモリと、状態 6 0 8 の値によってその時点で占有されているスペースの量とを把握することができる。あるいは、または追加的に、ソフトウェアは、メモリサブシステム 6 0 4 におけるスペースの割当てを管理することができる。状態 6 0 8 を記憶するために必要なメモリ空間の方が少ないことから、メモリサブシステム 6 0 4 (またはソフトウェア) は、重み 6 0 6 を記憶するためにより多くのメモリ空間を割り当てることを決定し得る。さらに、様々な時点で、メモリサブシステム 6 0 4 はチップ相互接続 6 2 0 を介して、追加のメモリ 6 3 0 からより多くの重みが読み取られるように要求を送信することができる。1つの入力データセットを処理する過程で、ニューラルネットワーク処理エンジン 6 0 2 は、ニューラルネットワークの重み 6 0 6 の全てをメモリバンク 6 1 4 へ徐々に移動することができる。

10

【 0 0 7 8 】

[0094]実装によっては、1つの入力データセットの処理が完了に近づくにつれて、いまだ必要とされる重み 6 0 6 の数は、さらに少なくなっていく可能性がある。残りの重み 6 0 6 の数が減少するにつれて、重み 6 0 6 または状態 6 0 8 へのメモリの割当ては、逆転される可能性があり、よって、徐々に、状態 6 0 8 を記憶するためにより多くのメモリが蓄えられていく。計算が終わるまでには、メモリを、図 6 A に示すように再び分割することができる。

【 0 0 7 9 】

20

[0095]先に述べたように、実装によっては、重み 6 0 6 は、一度使用されると、メモリサブシステム 6 0 4 から削除され得、この場合、これらの重み 6 0 6 によって占有されていたメモリは、依然として必要とされる重み 6 0 6 を記憶することに引き渡され得る。これらの実装では、入力データセットの処理が終わりに近づくにつれて、初期の層からの重み 6 0 6 は、メモリサブシステム 6 0 4 へリロードされることが可能である。したがって、入力データの処理が完了すると、ニューラルネットワーク処理エンジン 6 0 2 は、次の入力データを処理できる態勢となり得る。

【 0 0 8 0 】

[0096]メモリサブシステム 6 0 4 に現行で存在しない重み 6 0 6 が、プロセッサメモリにある場合、プロセッサメモリから重み 6 0 6 を読み取ると、処理エンジンアレイ 6 1 0 が、重み 6 0 6 がロードされるのを待つ間にストールされる可能性がある。様々な実装において、プロセッサメモリからの読取りに由来するストールを回避するための1つの技法は、プロセッサメモリに対する読取りトランザクションを、重み 6 0 6 が必要とされる遙か前にキューイングするというものである。別の技法は、重み 6 0 6 をオンチップキャッシュにプリフェッチする。別の技法は、別のニューラルネットワーク処理エンジンのメモリを用いることであり、これについては、後に詳述する。様々な実装では、所与のシステムの要件に合わせて、これらの、および他の技法を組み合わせることができる。

30

【 0 0 8 1 】

[0097]様々な実装では、ニューラルネットワーク処理エンジンのメモリサブシステムにおいて、計算の過程で利用可能となるスペースは、別のニューラルネットワークの重みを記憶するために使用されることが可能である。この第2のニューラルネットワークは、ニューラルネットワーク処理エンジンにもともとロードされているニューラルネットワークとは異なるタスクを実行するように訓練されている場合もある。ニューラルネットワーク処理エンジンに記憶される重みを、第1のニューラルネットワークの重みから第2のニューラルネットワークの重みへ切り替えることにより、ニューラルネットワーク処理エンジンは、異なるタスクを順次実行することができる。

40

【 0 0 8 2 】

[0098]ニューラルネットワーク処理エンジンを、異なるニューラルネットワーク間で切り替えて異なるタスクを順次実行するように構成することは、所定の状況において有用であり得る。たとえば、単一のユーザからの入力、たとえば、自然言語理解、これに続く

50

機械翻訳など、複数の動作を必要とする可能性がある。ニューラルネットワークは、自然言語理解などの複雑なタスクを実行する能力を有し得るものの、1つのニューラルネットワークが、自然言語理解と共に機械翻訳をも実行することはできない場合がある。したがって、多くの場合、1つのニューラルネットワークを用いて自然言語理解が実行され、かつ別のニューラルネットワークを用いて機械翻訳が実行されることが可能である。

【0083】

[0099]ニューラルネットワークプロセッサが様々な(おそらくは無関係の)タスクに対する結果を生成する速度を最適化するために、ニューラルネットワークプロセッサは、高速コンテキストスイッチを実行するように構成されることが可能である。ニューラルネットワーク処理システムに適用されるコンテキストは、ニューラルネットワークの重みと、ニューラルネットワークが訓練されたタスクを実行するために行われる計算とを含み、1つのニューラルネットワークについて記述する。コンテキストスイッチは、ニューラルネットワーク処理エンジンを、第1のニューラルネットワークを実行するように構成された後に第2のニューラルネットワークを実行するようにプログラムし直すことを包含し得る。ニューラルネットワーク処理エンジンをプログラムし直すことは、第2のニューラルネットワークの重みをニューラルネットワーク処理エンジンのメモリにロードすることを包含し得る。

【0084】

[0100]高速コンテキストスイッチは、ニューラルネットワーク処理エンジンが、第1のニューラルネットワークを用いる計算を完了すると即時、第2のニューラルネットワークを実行できる態勢になる場合に発生し得る。即時とは、たとえば1クロックサイクルまたは数クロックサイクル内を意味し得る。様々な実装では、第1のニューラルネットワークの計算が進行し、ニューラルネットワークプロセッサのメモリにおいてメモリ空間が利用可能になるにつれて、第2のニューラルネットワークの重みをこの利用可能空間に記憶することができる。最初の計算が完了すると、第2のニューラルネットワークの重みの一部または全てがメモリに存在し、これらの重みを用いる計算を開始することができる。

【0085】

[0101]図7A~図7Cは、ニューラルネットワークを実行するために使用可能なニューラルネットワーク処理エンジンの一例702を示す。様々な実装において、ニューラルネットワーク処理エンジン702は、メモリサブシステム704と、処理エンジンアレイ710とを包含し得る。動作中(たとえば、入力データセットの結果を計算しているとき)、処理エンジンアレイ710は、メモリサブシステム704から重み706および状態708の値を読み取ることができ、ここで、状態708の値は、当初、入力データを含み得る。図示されていないが、処理エンジンアレイ710は、処理エンジンアレイを包含し得、これらの処理エンジンの各々が、たとえば積和演算を実行することができる。処理エンジンアレイ710は、計算結果を結果バッファ712へ出力することができる。事例によっては、この例示的なニューラルネットワーク処理エンジン702は、結果がメモリサブシステム704へ書き込まれる前に、計算結果に対して、活性化ブロック716を用いる活性化関数、および/またはプーリングブロック718を用いるプーリングを実行することができる。

【0086】

[0102]動作において、ニューラルネットワーク処理エンジン702は、メモリサブシステム704のメモリバンク714から重み706および状態708を読み取り、重み706および状態708を処理エンジンアレイ710へ入力する。処理エンジンアレイ710は、中間結果を結果バッファ712へ出力することができ、中間結果は、結果バッファ712からメモリサブシステム704へ書き込まれることが可能である。複数の中間結果を組み合わせて1つの出力活性化にする場合は、活性化ブロック716が選択的に有効化され得る。また、プーリングブロック718も、ノードクラスタの出力を組み合わせて1つの出力にする場合に、選択的に起動され得る。活性化716もプーリング718も有効化されない場合、処理エンジンアレイ710からの結果は、結果バッファからメモリサブシ

10

20

30

40

50

ステム 704 へ直に移動されることが可能である。重み 706 および状態 708 を読み取って結果を計算するサイクルは、ニューラルネットワークの最終層に到達して処理エンジンアレイ 710 が最終結果を出力するまで、追加的に生じ得る。事例によっては、最終結果は、メモリサブシステム 704 へ書き込まれる前に、活性化ブロック 716 および/またはプーリングブロック 718 によって作用され得る。

【0087】

[0103]図 7A に示すように、様々な実装において、メモリバンク 714 の一部は、重み 706 の値を記憶するために使用され得、異なる部分は、状態 708 の値を記憶するために使用され得る。これらの実装形態では、ニューラルネットワークの重み 706 の値は、ニューラルネットワーク処理エンジン 702 が入力データを受信する前に、重み 706 のために取っておかれたメモリバンク 714 の部分に記憶されることが可能である。事例によっては、メモリバンク 714 に記憶される重み 706 は、ニューラルネットワークの重みの全てを包含することができる。事例によっては、記憶される重み 706 は、ニューラルネットワークの重みの全てより少ない量を包含し得、残りの重みは、(先に論じたように) 計算が進行しているときに追加のメモリ 730 から読み取られる。同じくこれらの実装では、計算の過程で、中間結果(たとえば、状態 708 の値)を、メモリバンク 714 の、状態 708 の値を記憶するために予約された部分に記憶することができる。メモリバンク 714 の、重み 706 を記憶するために使用される部分と、状態 708 を記憶するために使用される部分と、任意の未使用部分とは、概して重複しない。

【0088】

[0104]図 7A の例において、第 1 の入力データ 750 が受信されるとメモリサブシステム 704 に記憶される重み 706 は、第 1 のニューラルネットワークの重み 706 である。重み 706 と、第 1 のニューラルネットワークが実行するように訓練されたタスク(以後の例では、第 1 のタスクと称する)と、おそらくは処理エンジンアレイ 710、活性化ブロック 716 および/またはプーリングブロック 718 に対する命令をも含む第 1 のニューラルネットワークは、ニューラルネットワーク処理エンジン 702 の第 1 のコンテキストを表し得る。第 1 の入力データ 750 がニューラルネットワーク処理エンジン 702 によって受信されると、第 1 の入力データ 750 に対して第 1 のタスクが実行される。

【0089】

[0105]先に論じたように、第 1 のタスクを実行するための計算が進行するにつれて、メモリサブシステム 704 内のスペースは、利用可能になることもあれば、未使用と見なされることもある。たとえば、メモリバンク 714 の、状態 708 の値を記憶するために使用される部分は、中間結果を記憶するために必要とされるメモリの量が少なくなるにつれて、小さくなる可能性がある。一部の例において、メモリサブシステム 704 における論理は、中間結果を記憶するために必要なメモリの量を決定することができる。例によっては、ソフトウェアが、中間結果を記憶するために必要なメモリの量を決定することができ、したがって、メモリバンク 714 の、状態 708 の値を記憶するために使用されている部分を減らすことができる。さらなる例として、メモリバンク 714 の、第 1 のニューラルネットワークの重み 706 を記憶するために使用される部分は、使用されたことのある重み 706 の値がもはや不要としてマーキングされることに起因して、より小さくなり得る。例によっては、メモリサブシステム 704 における論理は、使用される重みのアドレス位置を利用可能としてマーキングすることができる。例によっては、ソフトウェアは、いまだ必要とされる重み値と、使用されていて不要となった重み値とを把握することができる。

【0090】

[0106]様々な実装では、メモリサブシステム 704 におけるスペースが未使用または利用可能となるにつれて、図 7B に示すように、第 2 のニューラルネットワークの重み 707 の値を利用可能なスペースに記憶することができる。図 7B に示す例では、先に第 1 のニューラルネットワークの重み 706 を記憶するために使用された幾つかのメモリバンク 714、および先に状態 708 の値を記憶するために使用された幾つかのメモリバンク 7

14が、第2のニューラルネットワークの重み707の値を記憶するために転用されている。第2のニューラルネットワークは、第1のタスクとは異なり得る、以後第2のタスクと称するタスクを実行するように訓練されている場合がある。重み707と、第2のタスクと、おそらくは処理エンジンレイ710、活性化ブロック716および/またはプリーミングブロック718に対するコマンドとは、ニューラルネットワーク処理エンジン702の第2のコンテキストと見なされ得る。重み707は、チップ相互接続720を介して追加のメモリ730から読み取られることが可能であり、追加のメモリ730は、オンチップキャッシュ、プロセッサメモリおよび/または別のニューラルネットワーク処理エンジン702のメモリであってもよい。

【0091】

[0107]様々な実装において、第2のニューラルネットワークの重み707を記憶することは、第1の入力データ750の結果を計算する間に発生し得る。すなわち、第1の入力データ750の結果の計算が進み、メモリバンク714においてスペースが利用可能となるにつれて、ニューラルネットワーク処理エンジン702は、未使用となったスペースに第2のニューラルネットワークの重み707をロードし始めることができる。第2のニューラルネットワークの重み707の記憶は、ニューラルネットワーク処理エンジン702による第2の入力データ751の受信と同時に開始することができ、この時点で第2のタスクが実行される。事例によっては、第2の入力データ751は、ニューラルネットワーク処理エンジン702が第1の入力データ750の結果を計算する過程で受信される。事例によっては、第2のニューラルネットワークの重み707の記憶は、第2の入力データ751が受信される前または後で開始することができる。例によっては、第2の入力データ751の受信は、第2のニューラルネットワークの重み707のロードをトリガする。例によっては、第2の入力データ751の受信は、第2の入力データ751の結果の計算をトリガする。

【0092】

[0108]事例によっては、第1のニューラルネットワークの重み706を記憶するために使用されるメモリバンク714のみが、第2のニューラルネットワークの重み707を記憶するために転用される。これらの事例において、メモリバンク714の、状態708を記憶するために使用される部分は、常に、中間結果を記憶するために利用可能である。

【0093】

[0109]事例によっては、メモリバンク714を転用することは、メモリバンク714間でデータを移動することを含み得る。たとえば、メモリバンク714の、状態708を記憶するために使用される部分が小さくなるにつれて、未使用の重み706は、この部分へ移動されることが可能である。この例では、メモリバンク714の、先に第1のニューラルネットワークの重み706により占有されていた部分は、徐々に第2のニューラルネットワークの重み707によって占有されることになり得、よって、コンテキストスイッチが完了すると、メモリの割当ては、第1のニューラルネットワークの場合と同じになる。

【0094】

[0110]様々な実装において、ニューラルネットワーク処理エンジン702が第1の入力データ750の結果の計算を完了すると、ニューラルネットワーク処理エンジン702は、第2の入力データ751の計算を実行できる態勢となり得る。図7Cに示す例では、ニューラルネットワーク処理エンジン702が、第1のニューラルネットワークから第2のニューラルネットワークへのコンテキストスイッチを完了している。この例に示すように、メモリサブシステム704には第2のニューラルネットワークの重み707がロードされていて、ニューラルネットワーク処理エンジン702は、第2の入力データ751の結果の計算を開始できる態勢にある。

【0095】

[0111]様々な実装において、ニューラルネットワーク処理エンジン702は、第1の入力データ750の計算を完了する1または2サイクルのうちに、第2の入力データ751の結果の計算を開始することができる。実装によっては、ニューラルネットワーク処理工

10

20

30

40

50

エンジン 702 は、第 2 の入力データ 751 の結果の計算を、第 1 の入力データ 750 の結果が完了する前に開始することができる。たとえば、第 2 の入力データ 751 の結果を計算するための状態 708 および重み 707 が処理エンジンアレイ 710 に入力されつつある時点で、第 1 の入力データ 750 の計算の最終サイクルは、まだ処理エンジンアレイ 710 を通って進行中である場合がある。この例では、第 1 の入力データ 750 の計算および第 2 の入力データ 751 の計算が処理エンジンアレイ 710 内に同時に存在する場合がある。

【0096】

[0112] 事例によっては、いまやメモリサブシステム 704 に記憶されている重み 707 は、第 2 のニューラルネットワークの全ての重みを包含する。事例によっては、重み 707 の数は、第 2 のニューラルネットワークの重みの全てより少なく、残りの重みは、メモリバンク 714 においてスペースが利用可能となるにつれて、追加のメモリ 730 から読み取られ得る。

10

【0097】

[0113] 事例によっては、図 7C に示す例において、メモリサブシステム 704 内になおも第 1 のコンテキストに関連する幾つかのデータが存在する場合がある。たとえば、メモリバンク 714 には、第 1 のニューラルネットワークの、最後に使用されたものなどの幾つかの重み 706 の値がまだ存在している場合がある。この例では、第 2 の入力データ 751 の計算が始まると、これらの重み 706 は、上書きされることが可能である。別の例として、第 1 の入力データ 750 の計算の最終結果は、チップ相互接続 720 を介する読み出しを待機して、なおもメモリバンク 714 内に存在する場合がある。最終結果は、第 2 の入力データ 751 の計算が進行している間、メモリサブシステム 704 にとどまる場合もあれば、第 2 の入力データ 751 の計算が始まる前に読み出される場合もある。

20

【0098】

[0114] ニューラルネットワーク処理エンジン 702 が第 1 のニューラルネットワークの重み 706 値を用いて結果を計算している間に、第 2 のニューラルネットワークの重み 707 の値をロードすると、ニューラルネットワーク処理エンジン 702 が第 1 のニューラルネットワークと第 2 のニューラルネットワークとの間で高速コンテキストスイッチを実行することが有効化される。高速コンテキストスイッチは、ニューラルネットワーク処理エンジン 702 がスループットを高めかつ結果をより速く生成することを有効化し得る。

30

【0099】

[0115] ニューラルネットワークプロセッサを介してスループットを高めるための別の技法は、ニューラルネットワークプロセッサに、並行して動作できる複数のニューラルネットワーク処理エンジンを装備するというものである。ニューラルネットワークプロセッサは、たとえば、異なるニューラルネットワーク処理エンジン上で、複数の、おそらくは異なるニューラルネットワークを同時に実行することができる。別の例として、ニューラルネットワークプロセッサは、より多くの入力データセットをバッチ処理することができ、よって重み値の再使用が増加する。

【0100】

[0116] 図 8 は、複数のニューラルネットワーク処理エンジン 802a ~ 802n を有するニューラルネットワークプロセッサの一例 800 を示す。ニューラルネットワーク処理エンジン 802a ~ 802n は、各々、メモリサブシステムと、処理エンジンアレイとを含むことができ、かつニューラルネットワークが、そのためにニューラルネットワークがプログラムされたタスクを実行するために必要な計算を実行することができる。図示の例では、ニューラルネットワークプロセッサ 800 が n 個のニューラルネットワーク処理エンジン 802a ~ 802n を含む。

40

【0101】

[0117] この例示的なニューラルネットワークプロセッサ 800 は、さらに、この例では DRAM 830 を用いて実装されるプロセッサメモリと通信するための DRAM コントローラ 842a ~ 842k を含む。図示の例において、ニューラルネットワークプロセッサ

50

800は、k個のDRAMコントローラ842a~842kを含み、その各々が、独立したDRAMバンクセットと通信することができてもよい。他の例では、プロセッサメモリに他のタイプのRAM技術を使用可能である。DRAMコントローラ842a~842kは、メモリコントローラと呼ばれることもある。

【0102】

[0118]例示的なニューラルネットワークプロセッサ800は、さらに、システム内のPCIデバイス832と通信するためのペリフェラル・コンポーネント・インターコネクタ(PCI)コントローラ844a~844pを含む。本例において、PCIは、プロセッサバスの実装に使用されているバスプロトコルである。プロセッサバスは、ニューラルネットワークプロセッサ800を、たとえば、とりわけI/Oデバイス、メモリコントローラ、ストレージデバイスおよび/またはネットワーク・インタフェース・カードなどのPCIデバイス832へ接続することができる。本例において、ニューラルネットワークプロセッサ800は、p個のPCIコントローラ844a~844pを含み、その各々が、別々のルートコンプレックスを含んでもよく、かつPCIデバイス832の別々のセットと通信してもよい。他の例では、ウルトラ・バス・インターコネクタ(UPI)などの他の標準化されたバスプロトコルをホストバスに使用することができる。他の例では、独自仕様のバスプロトコルを使用することができる。

【0103】

[0119]この例示的なニューラルネットワークプロセッサ800は、さらに、ニューラルネットワーク処理エンジン802a~802nと、DRAMコントローラ842a~842kと、PCIコントローラ844a~844pとの間でデータを移動させることができるDMAエンジン846a~846dを含む。図示の例では、ニューラルネットワークプロセッサ800がd個のDMAエンジン846a~846dを含む。実装によっては、DMAエンジン846a~846dは、DRAMコントローラ842a~842kからニューラルネットワーク処理エンジン802a~802nへのデータの移動、またはPCIコントローラ844a~844pとニューラルネットワーク処理エンジン802a~802nとの間のデータの移動などの特定のタスクに割り当てられ得る。実装によっては、少なくとも1つのDMAエンジン846a~846dを、各ニューラルネットワーク処理エンジン802a~802nの専用とすることができる。実装によっては、DMAエンジン846a~846dは、1つの機能またはコンポーネントの専用とされる代わりに、プールとして扱われてもよく、よって、データを移動する必要がある度に、利用可能なDMAエンジン846a~846dが使用される。

【0104】

[0120]例示的なニューラルネットワークプロセッサ800では、様々なコンポーネントがチップ相互接続820を介して通信することができる。チップ相互接続820は、主として、ニューラルネットワークプロセッサ800のコンポーネント間でデータをルーティングするための配線を含む。事例によっては、チップ相互接続820は、データの方向を制御するためのマルチプレクサ、クロックドメイン交差を処理するためのフリップフロップ、およびタイミング論理などの、最小限の論理を包含することができる。

【0105】

[0121]図9は、ニューラルネットワークの重み906が2つのニューラルネットワーク処理エンジン902a~902bのメモリサブシステムに記憶されるように構成されている、ニューラルネットワークプロセッサの一例900を示す。ニューラルネットワークプロセッサ900は、図8に示すニューラルネットワークプロセッサ800に類似するものであり得る。図9に示すように、ニューラルネットワーク処理エンジン902a~902nは、各々、メモリサブシステム904a~904nと、処理エンジンアレイ910a~910nとを含む。様々な実装において、メモリサブシステム904a~904nは、各々、ニューラルネットワークの重みを記憶することができ、よって、ニューラルネットワーク処理エンジン902a~902nは、各々、独立してニューラルネットワークを実行することができる。

10

20

30

40

50

【 0 1 0 6 】

[0122]事例によっては、第2のニューラルネットワーク処理エンジン902bなどのニューラルネットワーク処理エンジン902a~902nのうちの1つは、ニューラルネットワーク処理エンジン902bのメモリサブシステム904bにおけるスペースの全てを必要とはしない場合がある。たとえば、ニューラルネットワーク処理エンジン902bがニューラルネットワークの重みを記憶するために必要とするスペースは、メモリサブシステム904bにおける全てのスペースに満たない場合がある。あるいは、別の例として、ニューラルネットワーク処理エンジン902bは、重み値をDRAM930のみから、または主としてDRAM930から入手するように指定されていてもよい(たとえば、ニューラルネットワーク処理エンジン902bが実行しているニューラルネットワークは、高度の重み再使用を有する)。

10

【 0 1 0 7 】

[0123]様々な実装において、第2のニューラルネットワーク処理エンジン902b内の利用可能なメモリは、第1のニューラルネットワーク処理エンジン902aにより実行されているニューラルネットワークの重み906を記憶するために使用されることが可能である。ニューラルネットワークの重み906は、たとえば、第1のニューラルネットワーク処理エンジン902aのメモリサブシステム904a内に完全に記憶するには多すぎる場合がある。先に述べたように、第1のニューラルネットワーク処理エンジン902aに収まらない重み906は、DRAM930に記憶されることが可能である。しかしながら、DRAM930から重み906を読み取ることより、1つのニューラルネットワーク処理エンジン902bから別のニューラルネットワーク処理エンジン902aへ重み906をコピーすることの方が遙かに高速で実行可能であることから、これらの重みを第2のニューラルネットワーク処理エンジン902bの利用可能なスペースに記憶することの方が遙かに効率的であり得る。

20

【 0 1 0 8 】

[0124]先に論じたように、重み906は、第1のニューラルネットワーク処理エンジン902aが入力データセットを処理するにつれて、第2のニューラルネットワーク処理エンジン902bから第1のニューラルネットワーク処理エンジン902aへコピーされることが可能である。事例によっては、使用されて不要となった重み906は、入力データの処理が完了するまで、第1のニューラルネットワーク処理エンジン902aから第2のニューラルネットワーク処理エンジン902bへ移動されてもよい。様々な実装において、DMAエンジン946a~946dは、ニューラルネットワーク処理エンジン902a~902n間でデータをコピーまたは移動させることができ、よって、ニューラルネットワーク処理エンジン902a~902nは、ニューラルネットワークの計算の実行に専念することができる。

30

【 0 1 0 9 】

[0125]第2のニューラルネットワーク処理エンジン902bのメモリサブシステム904bにおける、第1のニューラルネットワーク処理エンジン902aの重み906を記憶するために使用されないスペースはいずれも、ニューラルネットワーク処理エンジン902bがその固有の計算を実行するために使用されることが可能である。

40

【 0 1 1 0 】

[0126]実装によっては、ニューラルネットワーク処理エンジン1002a~1002n間で重みを移動させるのではなく、代わりに進行中の計算を移動させることができる。図10は、ニューラルネットワークの重み1006が2つのニューラルネットワーク処理エンジン1002a~1002bのメモリサブシステムに記憶されるように構成されている、ニューラルネットワークプロセッサの一例1000を示す。ニューラルネットワークプロセッサ1000は、図8に示すニューラルネットワークプロセッサに類似するものであり得る。図9の例の場合のように、図10の例において、ニューラルネットワークの重み1006は、おそらくは、重み1006の全てをオンチップで記憶するために必要なメモリ量が、ニューラルネットワーク処理エンジン1002a~1002nのどのメモリ量を

50

も超えるという理由で、2つのニューラルネットワーク処理エンジン1002a~1002bのメモリサブシステム1004a~1004bに記憶される。

【0111】

[0127]先に述べたように、重み1006は、ニューラルネットワークの特定の層に関連づけられ得る。したがって、第1の層セットに関する重み1006は、第1のニューラルネットワーク処理エンジン1002aに記憶されることが可能であり、残りの層に関する重み1006は、ニューラルネットワーク処理エンジン1002bに記憶されることが可能である。

【0112】

[0128]様々な実装において、重み1006に関連づけられるニューラルネットワークの入力データが受信されると、入力データは、第1のニューラルネットワーク処理エンジン1002aがニューラルネットワークの初期層の重みを有することから、第1のニューラルネットワーク処理エンジン1002aへ提供され得る。第1のニューラルネットワーク処理エンジン1002aは、入力データの処理を進めることができ、メモリサブシステム1004aから重み1006を読み取り、処理エンジンアレイ1010aで中間結果を計算しかつ中間結果をメモリサブシステム1004aに書き戻す。

10

【0113】

[0129]第1のニューラルネットワーク処理エンジン1002aが、第1のニューラルネットワーク処理エンジン1002aがその重みを有する各層についての計算を実行すると、第1のニューラルネットワーク処理エンジン1002aは、進行中の計算を第2のニューラルネットワーク処理エンジン1002bへ移動させることができる。進行中の計算を移動させることは、状態1008の値を移動させることを含み得る。状態1008は、たとえば、中間結果のセットを含み得る。例によっては、状態1008は、とりわけ、処理された最終層の識別、処理エンジンアレイ1010bの制御入力、活性化関数の状態値、および/またはプリーング関数の状態値をさらに含み得る。様々な実装において、DMAエンジン1046a~1046dは、状態1008の値を第1のニューラルネットワーク処理エンジン1002aから第2のニューラルネットワーク処理エンジン1002bへコピーすることに従事され得る。

20

【0114】

[0130]様々な実装において、状態1008が第2のニューラルネットワーク処理エンジン1002bへコピーされると、ニューラルネットワーク処理エンジン1002bは、進行中の計算を再開するための信号を受信することができる。信号は、たとえば、とりわけ、割込み、第2のニューラルネットワーク処理エンジン1002b内のレジスタへの書込み、または制御バスを介してニューラルネットワーク処理エンジン1002a~1002n間で送信される命令、の形式をとることができる。事例によっては、第2のニューラルネットワーク処理エンジン1002bは、自動的に計算を開始することができる。たとえば、中間結果は、入力データと同様にして入力されることが可能であり、よって、入力データに対して起動されるはずの第2のニューラルネットワーク処理エンジン1002b内の同じメカニズムが起動される。

30

【0115】

[0131]様々な実装において、進行中の計算を再開することは、中間層からの重み1006を、第1のニューラルネットワーク処理エンジン1002aからの中間結果と共に処理エンジンアレイ1010bへ入力することを含み得る。様々な実装において、ニューラルネットワーク処理エンジン1002bは、ニューラルネットワーク処理エンジン1002bが進行中の計算を新しい計算の開始として扱い得ること、または、ニューラルネットワーク処理エンジン1002bは、第1のニューラルネットワーク処理エンジン1002aがやめた場所を正確に捕捉できること、のいずれかを理由に、進行中の計算を再開するための特別な、または異なる挙動を必要としない。

40

【0116】

[0132]図9および10の例では、ニューラルネットワークの全ての重みが処理エンジン

50

アレイのローカルメモリにおいて保持され得なくなるのはいつであるか、に対するソリューションが提供される。ニューラルネットワーク処理エンジン間でデータを移動する際には、多少の待ち時間が発生し得るが、それでもこの待ち時間は、オフチップメモリから重みを読み取る際に発生する待ち時間より遙かに短いものである。

【0117】

[0133]図11は、重みの全てを単一の処理エンジンアレイに対してローカルに記憶し得ない場合に、ニューラルネットワークの全ての重みをオンチップで記憶するための別のソリューションを示す。図11は、2つのニューラルネットワークプロセッサ1101a~1101bを含むニューラルネットワーク処理システム1100を示す。本例において、各ニューラルネットワークプロセッサ1101a~1101bは、別々のチップであって、別々のダイを有し、かつ異なるパッケージに封入されている。ニューラルネットワークプロセッサ1101a~1101bは、ホストバスまたはプロセッサバスを用いて、互いに、かつニューラルネットワーク処理システム1100内の他のコンポーネントと通信することができる。プロセッサバスは、たとえば、PCIベースのプロトコルを用いて実装されることが可能である。

10

【0118】

[0134]様々な実装において、ニューラルネットワークプロセッサ1101a~1101bは、各々、少なくとも1つの、おそらくはそれ以上のニューラルネットワーク処理エンジン1102a~1102bを含む。先に論じたように、ニューラルネットワーク処理エンジン1102a~1102bは、各々、処理エンジンアレイ1010a~1010bと、メモリサブシステム1104a~1104bとを含むことができる。

20

【0119】

[0135]ニューラルネットワーク処理エンジン1102a~1102bと、各ニューラルネットワークプロセッサ1101a~1101b内の他のコンポーネントとの間の通信のために、ニューラルネットワークプロセッサ1101a~1101bは、チップ相互接続1020a~1020bを含むことができる。チップ相互接続1020a~1020bは、主として、ニューラルネットワークプロセッサ1101a~1101bのコンポーネント間で信号をルーティングするための配線、およびおそらくは何らかのタイミング論理も含むことができる。

【0120】

[0136]ニューラルネットワークプロセッサ1101a~1101bの他のコンポーネントは、プロセッサバスと通信するためのPCIコントローラ1144a~1144bを含むことができる。様々な実装において、ニューラルネットワークプロセッサ1101a~1101bは、DRAMおよびDMAコントローラなどの図示されていないコンポーネントを含むことができる。

30

【0121】

[0137]図11の例では、ニューラルネットワークの重み1106を2つに分割することができ、2つの部分は、同じサイズまたは異なるサイズである。ニューラルネットワークの初期層に対応する重み1106の第1のセットは、第1のニューラルネットワークプロセッサ1101aにおけるニューラルネットワーク処理エンジン1102aのメモリサブシステム1104aにロードされることが可能である。重み1106の第2のセットは、第2のニューラルネットワークプロセッサ1101bにおけるニューラルネットワーク処理エンジン1102bのメモリサブシステム1104bにロードされることが可能である。この構成では、ニューラルネットワークの入力データが受信されると、入力データは、第1のニューラルネットワークプロセッサ1101aのニューラルネットワーク処理エンジン1102aへ提供され得る。ニューラルネットワーク処理エンジン1102aは、入力データの計算を進めることができる。

40

【0122】

[0138]ニューラルネットワーク処理エンジン1102aが、ニューラルネットワーク処理エンジン1102aがその重み1106を有する層の各々についての計算を実行すると

50

、ニューラルネットワーク処理エンジン 1 1 0 2 a は、進行中の計算を第 2 のニューラルネットワークプロセッサ 1 1 0 1 b におけるニューラルネットワーク処理エンジン 1 1 0 2 b へ移動させることができる。たとえば、ニューラルネットワーク処理エンジン 1 1 0 2 a は、DMA エンジンに、計算の状態 1 1 0 8 を PCI コントローラ 1 1 4 4 a へコピーするように命令することができるが、状態 1 1 0 8 は、とりわけ、中間結果を含み得る。本例において、PCI コントローラ 1 1 4 4 a は、第 2 のニューラルネットワークプロセッサ 1 1 0 1 b への 1 つまたは複数の書込みを生成することができ、書き込みは状態 1 1 0 8 の値を含む。実装によっては、第 1 のニューラルネットワークプロセッサ 1 1 0 1 a は、ニューラルネットワーク処理エンジン 1 1 0 2 b のメモリサブシステム 1 1 0 4 b へ直に書き込むことができる。実装によっては、状態 1 1 0 8 のデータは、第 2 のニューラルネットワークプロセッサ 1 1 0 1 b 内の PCI コントローラ 1 1 4 4 b によってキャッシュされる。これらの実装において、DMA エンジンは、たとえば、状態 1 1 0 8 のデータを PCI コントローラ 1 1 4 4 b からニューラルネットワーク処理エンジン 1 1 0 2 b へ移動することができる。

10

【0123】

[0139] 状態 1 1 0 8 が第 2 のニューラルネットワークプロセッサ 1 1 0 1 b 内のニューラルネットワーク処理エンジン 1 1 0 2 b のメモリサブシステム 1 1 0 4 b へロードされると、ニューラルネットワーク処理エンジン 1 1 0 2 b は、進行中の計算を再開することができる。例によっては、第 1 のニューラルネットワークプロセッサ 1 1 0 1 a は、第 2 のニューラルネットワークプロセッサ 1 1 0 1 b への割込みを生成することにより、または第 2 のニューラルネットワークプロセッサ 1 1 0 1 b 内のレジスタへ値を書き込むことにより、第 2 のニューラルネットワーク処理エンジン 1 1 0 2 b をトリガする。ニューラルネットワーク処理エンジン 1 1 0 2 b が最終結果を計算すると、第 2 のニューラルネットワークプロセッサ 1 1 0 1 b は、結果を出力することができる。

20

【0124】

[0140] 図 1 1 の例において、ニューラルネットワークの重み 1 1 0 6 は全て、処理エンジンアレイ 1 0 1 0 a ~ 1 0 1 0 b に対してローカルに記憶され、よって、ニューラルネットワークプロセッサ 1 1 0 1 a ~ 1 1 0 1 b は、重み 1 1 0 6 を入手するためにプロセッサメモリを読み取る必要を回避することができる。状態 1 1 0 8 をニューラルネットワークプロセッサ 1 1 0 1 a ~ 1 1 0 1 b 間で転送する際には、多少の遅延が生じるものの、この遅延は、プロセッサメモリへの複数回の読取りによって引き起こされる遅延より遙かに少ないものであり得る。

30

【0125】

[0141] 様々な実装において、図 1 0 および図 1 1 に示す例は、様々な方法で組み合わせられ得る。たとえば、入力データセットの処理は、1 つのニューラルネットワーク処理エンジンにおいて開始されることが可能であって、同じニューラルネットワークプロセッサにおける第 2 のニューラルネットワーク処理エンジンへ転送されてもよい。進行中の計算は、次に、第 2 のニューラルネットワーク処理エンジンから第 3 のニューラルネットワーク処理エンジンへ転送されることが可能であり、ここで、第 3 のニューラルネットワーク処理エンジンは、異なるニューラルネットワークプロセッサ内にある。ニューラルネットワークの全ての重みをオンチップで記憶し、かつプロセッサメモリからの読取りを最小限にする、またはなくすことができるような、様々な構成が可能である。

40

【0126】

[0142] 図 1 2 は、ニューラルネットワークプロセッサがニューラルネットワークを実行している場合にニューラルネットワークプロセッサにより実行されることが可能なプロセスの一例 1 2 0 0 を示す。

【0127】

[0143] ステップ 1 2 0 2 において、プロセス 1 2 0 0 は、重み値セットをニューラルネットワーク処理回路の複数のメモリバンクに記憶することを含み、ここで、ニューラルネットワーク処理回路は、複数のメモリバンクと同じダイ上に処理エンジンアレイを含み、

50

重み値セットは、入力データの受信に先行して記憶される。様々な実装において、ニューラルネットワーク処理回路は、ニューラルネットワークプロセッサの一部である。例によっては、重み値セットは、ニューラルネットワークの全ての重み値を含む。

【0128】

[0144]ステップ1204において、プロセス1200は、入力データを受信することを含む。様々な例において、入力データは、画像、オーディオ信号、テキスト文字列、数列または他の何らかの入力であってもよい。

【0129】

[0145]ステップ1206において、プロセス1200は、重み値セットを用いて、ニューラルネットワーク用に規定されたタスクを実行することを含む。タスクは、たとえば、10 画像識別、音声認識、自然言語処理、ゲームプレイ、他を含んでもよい。重み値は、訓練を介して先に決定されている可能性があり、訓練では、たとえば、その出力結果が既知である入力データがニューラルネットワークに入力されている。

【0130】

[0146]ステップ1206は、タスクを実行するためのさらなるステップを含み得る。ステップ1208において、タスクを実行することは、複数のメモリバンクから重み値を読み取ることを含む。例によっては、重み値を読み取るとは、同時に、複数のメモリバンクのうちの第1のメモリバンクから第1の重み値を読み取り、かつ複数のメモリバンクのうちの第2のメモリバンクから第2の重み値を読み取ることを含む。

【0131】

[0147]ステップ1210において、タスクを実行することは、さらに、重み値および入力データを処理エンジンアレイへ入力することを含み得る。例によっては、処理エンジンアレイは、シストリックアレイである。例によっては、ステップ1210は、さらに、処理エンジンアレイを用いて中間結果を決定し、かつ中間結果を複数のメモリバンクにおける1つのメモリバンクに記憶することを含み得る。例によっては、中間結果は、追加の重み値が読み取られると同時に書き込まれてもよい。20

【0132】

[0148]ステップ1212において、タスクを実行することは、さらに、処理エンジンアレイを用いて結果を計算することを含み得、結果は、タスク実行の結果に対応する。

【0133】

[0149]実装によっては、ステップ1202において記憶される重み値セットは、ニューラルネットワークの全ての重み値の第1の部分を含む。これらの実装において、プロセス1200は、さらに、複数のメモリバンクが利用可能なスペースを有すると決定することを含み得る。使用可能なスペースは、より少ないスペースを要求する中間結果、および/または不要となった重み値に由来する場合がある。プロセス1200は、さらに、ニューラルネットワークの全ての重み値の第2の部分を読み取ることを含み得、第2の部分は、追加のメモリから読み取られる。例によっては、追加のメモリは、同じダイ上の第2の処理エンジンアレイに関連づけられ得る。例によっては、追加のメモリは、オフチップであってもよい。プロセス1200は、さらに、第2の部分を利用可能なスペースへ書き込むことを含み得る。重みの第2の部分は、次に、タスクの実行を継続するために使用され得る。40

【0134】

[0150]図13は、ニューラルネットワークプロセッサがニューラルネットワークを実行している場合にニューラルネットワークプロセッサにより実行されることが可能なプロセスの一例1300を示す。

【0135】

[0151]ステップ1302において、プロセス1300は、第1の重み値セットをニューラルネットワーク処理システムの第1のメモリバンクセットに記憶することを含む。

【0136】

[0152]ステップ1304において、プロセス1300は、第2の重み値セットをニュー

ラルネットワーク処理システムの第2のメモリバンクセットに記憶することを含み、第1の重み値セットおよび第2の重み値セットは、ニューラルネットワークの全ての重み値を含み、かつ第1の重み値セットおよび第2の重み値セットは、ニューラルネットワーク処理システムによる入力データの受信に先行して記憶される。

【0137】

[0153]ステップ1306において、プロセス1300は、入力データを受信することを含む。様々な例において、入力データは、画像、オーディオ信号、テキスト文字列、数列または他の何らかの入力であってもよい。

【0138】

[0154]ステップ1308において、プロセス1300は、ニューラルネットワーク用に規定されたタスクを実行することを含む。様々な例において、タスクは、入力データのタイプ（たとえば、入力画像の画像認識、オーディオ信号の音声認識、他）に相関する。

10

【0139】

[0155]ステップ1308は、タスクを実行するためのさらなるステップを含み得る。ステップ1310において、タスクを実行することは、第1の処理エンジンアレイを用いて中間結果を計算することを含み、第1の処理エンジンアレイは、第1のメモリバンクセットと同じ第1のダイ上に存在し、第1の処理エンジンアレイは、第1の重み値セットおよび入力データを用いて中間結果を計算する。中間結果を計算することは、第1のメモリバンクセットから2つ以上の重み値を同時に読み取ることを含み得る。例によっては、中間結果は、ニューラルネットワークの1つまたは複数の初期層の出力を反映する。

20

【0140】

[0156]ステップ1312において、タスクを実行することは、さらに、中間結果および第1のメモリバンクセットからの状態を第2のメモリバンクセットへコピーすることを含み、状態は、第1の処理エンジンアレイにより完了されるタスクのステージを記述する。例によっては、中間結果および状態をコピーできる態勢になると、第1の処理アレイは、入力データに対する演算を完了して、新しい入力データに対する演算を開始することができる。

【0141】

[0157]ステップ1314において、タスクを実行することは、さらに、第2の処理エンジンアレイを用いて最終結果を計算することを含み、第2の処理エンジンアレイは、第2のメモリバンクセットと同じ第2のダイ上に存在し、第2の処理エンジンアレイは、第2の重み値セット、中間結果および状態を用いて最終結果を計算し、最終結果は、タスク実行の結果に対応する。最終結果を計算することは、第2のメモリバンクセットから2つ以上の重み値を同時に読み取ることを含み得る。例によっては、最終結果は、ニューラルネットワークの1つまたは複数の最終層の出力を反映する。

30

【0142】

[0158]例によっては、第1のメモリバンクセットおよび第1の処理エンジンアレイを含む第1のダイと、第2のメモリバンクセットおよび第2の処理エンジンアレイを含む第2のダイとは、同じダイの部分である。該ダイは、ニューラルネットワークプロセッサ用であってもよく、ニューラルネットワークプロセッサは、2つ以上のニューラルネットワーク処理エンジンを含む（たとえば、ニューラルネットワーク処理エンジンは、メモリバンクセットおよび処理エンジンアレイを含む）。例によっては、ニューラルネットワークプロセッサは、内部通信ファブリックを含んでもよく、中間結果および状態は、内部通信ファブリックを介してコピーされ得る。

40

【0143】

[0159]例によっては、第1のダイおよび第2のダイは、各々、異なるニューラルネットワークプロセッサの一部であり、よって、異なるパッケージ内に存在する。これらの例において、2つのニューラルネットワークプロセッサは、ホストバスにより、中間結果および状態の一方から他方へのコピーを含む相互通信を可能にされ得る。

【0144】

50

[0160]ステップ1304において述べたように、第1の重み値セットおよび第2の重み値セットは、ニューラルネットワークの全ての重み値であり得る。これは、重み値がプロセッサメモリからは読み取られないことを意味し得る。プロセッサメモリは、第1のダイまたは第2のダイとは異なるパッケージ内に存在し、よって、読取り遅延が大きくなる場合がある。

【0145】

[0161]図14は、ニューラルネットワークプロセッサを含むコンピューティングシステムにより実行されることが可能なプロセスの一例1400を示す。

【0146】

[0162]ステップ1402において、プロセス1400は、第1のニューラルネットワークの第1の重み値セットを記憶することを含み、第1のニューラルネットワークは、第1のタスクを実行するように構成され、第1の重み値セットは、ニューラルネットワークプロセッサにおけるメモリの第1の部分に記憶される。第1の重み値セットおよび第1のタスクは、最初のコンテキストと見なされ得る。

10

【0147】

[0163]ステップ1404において、プロセス1400は、第1のタスクに関連づけられる第1の入力データを受信することを含み、第1の入力データは、第1の重み値セットの記憶後に受信される。第1のタスクは、第1の入力データを入力として用いて実行される。様々な実装において、第1の重み値セットは、入力データが受信される前にメモリに記憶される。

20

【0148】

[0164]ステップ1406において、プロセス1400は、ニューラルネットワークプロセッサを用いて第1の結果を計算することを含み、第1の結果を計算することは、第1の重み値セットおよび第1の入力データを用いることを含み、第1の結果は、第1のタスクの実行の結果に対応し、かつ第1の結果を計算する間に、中間結果がメモリの第2の部分に記憶される。メモリの第1の部分および第2の部分は、概して重ならない。

【0149】

[0165]ステップ1408において、プロセス1400は、第1の結果を計算する間に、メモリの第3の部分を識別することを含み、第3の部分は、メモリ内に未使用のスペースを含む。未使用のスペースには、以前は第1の部分（たとえば、第1の重み値セットを記憶するために使用される部分）の一部であったがもはやそうではないスペース、および/または以前は第2の部分（たとえば、中間結果を記憶するために使用される部分）の一部であったがもはやそうではないスペースが含まれ得る。第3の部分は、概して、第1の部分または第2の部分と重複しない。

30

【0150】

[0166]ステップ1410において、プロセス1400は、第2の重み値セットを第3の部分に記憶することを含み、第2の重み値セットは、第2のニューラルネットワークに対するものであり、第2のニューラルネットワークは、第2のタスクを実行するように構成される。事例によっては、第2のタスクは、第1のタスクとは異なる。第2の重み値セットは、第1の結果の計算が完了する前に記憶されることが可能である。

40

【0151】

[0167]ステップ1412において、プロセス1400は、第2のタスクに関連づけられる第2の入力データを受信することを含む。第2のタスクは、第2の入力データを入力として用いて実行される。

【0152】

[0168]ステップ1414において、プロセス1400は、ニューラルネットワークプロセッサを用いて第2の結果を計算することを含み、第2の結果を計算することは、第2の重み値セットおよび第2の入力データを用いることを含み、第2の結果の計算は、第1の結果の計算が完了した時点で発生し、かつ第2の結果は、第2のタスクの実行の結果に対応する。

50

【 0 1 5 3 】

[0169]実装によっては、プロセス1400は、さらに、メモリの第1の部分のサイズを縮小することを含み得、サイズの縮小は、第1の結果の計算の完了に必要なとされない第1の重み値セットからの重み値に対応する。必要とされない重み値には、既に使用されている重み値が含まれ得る。第1の部分のサイズを縮小することは、第1の部分の一部を未使用として指定することと、未使用の一部をもはや第1の部分の一部ではないものとして扱うことを含み得る。

【 0 1 5 4 】

[0170]実装によっては、プロセス1400は、さらに、メモリの第2の部分のサイズを縮小することを含んでもよく、該サイズは、中間結果を記憶するために必要なメモリの量に対応する。メモリの量は、第1の結果を計算する過程で減少し得る。第2の部分のサイズを縮小することは、第2の部分の一部をもはや第2の部分の一部ではないとして指定することを含んでもよい。

10

【 0 1 5 5 】

[0171]実装によっては、第1の重み値セットは、第1のニューラルネットワークの全ての重み値より少ない数の重み値を含む。これらの実装において、プロセス1400は、さらに、第1のニューラルネットワークの追加の重み値セットを、コンピューティングシステムの第2のメモリから読み取ることを含んでもよい。第2のメモリは、オフチップ・プロセッサ・メモリ、オンチップキャッシュ、および/または第2のニューラルネットワークプロセッサのメモリであり得る。プロセス1400は、さらに、追加の重み値セットをメモリの第1の部分に記憶することを含んでもよく、追加の重み値セットは、第2の重み値セットの記憶に先行して記憶される。

20

【 0 1 5 6 】

[0172]図15は、これまでに論じた様々な例を用いることができるネットワーク1500を示す。図15の例示的なネットワーク1500は、ネットワークデバイスを備えるノード、スイッチおよびルータなどの様々な異なるタイプのネットワークデバイスを含む。所定の例において、ネットワーク1500は、ポイントツーポイント・リンクを備えるスイッチドアキテクチャに基づくものであってもよい。図15に示すように、ネットワーク1500は、複数のスイッチ1104a~1104dを含み、これらは、ネットワークに配置されてもよい。事例によっては、スイッチは、Closネットワークなどの多層網に配列される。ローカル・エリア・ネットワーク(LAN)セグメント間でパケットをフィルタリングしかつ転送するネットワークデバイスを、スイッチと称する場合がある。スイッチは、概して、開放型システム間相互接続(OSI)参照モデルのデータリンク層(層2)および時としてネットワーク層(層3)において動作し、かつ幾つかのパケットプロトコルをサポートしてもよい。スイッチ1104a~1104dは、複数のノード1102a~1102hへ接続され、かつ任意の2つのノード間に複数の経路を提供してもよい。

30

【 0 1 5 7 】

[0173]ネットワーク1500は、他のサブネット、LAN、広域ネットワーク(WAN)またはインターネットなどの他のネットワーク1508との接続のための、ルータ1506と称されることもある1つまたは複数のネットワークデバイスも含んでもよい。ルータは、ヘッダおよび転送テーブルを用いてパケットを転送するための最適経路を決定し、かつインターネット制御メッセージプロトコル(ICMP)などのプロトコルを用いて相互に通信しかつ任意の2デバイス間の最良ルートを構成する。

40

【 0 1 5 8 】

[0174]例によっては、ネットワーク1500は、ケーブルネットワーク、インターネット、無線ネットワーク、セルラネットワークおよび他の専用および/または公衆ネットワークなどの多くの異なるタイプのネットワークのうちのいずれか1つまたは組合せを含んでもよい。相互接続されるスイッチ1104a~1104d、もしあればルータ1506は、スイッチファブリック、ファブリック、ネットワークファブリック、または単にネッ

50

トワークと称されることがある。コンピュータネットワークのコンテキストにおいて、本明細書における「ファブリック」および「ネットワーク」という用語は、同義で使用されてもよい。

【0159】

[0175]ノード1102a~1102hは、ユーザデバイス、サービスプロバイダ・コンピュータまたは第三者コンピュータを表す、ホストシステム、プロセッサノード、ストレージサブシステムおよびI/Oシャーシの任意の組合せであってもよい。

【0160】

[0176]ユーザデバイスは、アプリケーション1532（たとえば、ウェブブラウザまたはモバイルデバイス・アプリケーション）にアクセスするためのコンピューティングデバイスを含んでもよい。態様によっては、アプリケーション1532は、コンピューティング・リソース・サービスまたはサービスプロバイダによってホストされ、管理されかつ/または提供されてもよい。アプリケーション1532は、ユーザが、たとえばウェブコンテンツ（たとえば、ウェブページ、音楽、ビデオ、他）にアクセスするために、サービスプロバイダ・コンピュータと対話することを可能にしてもよい。ユーザデバイスは、たとえば、携帯電話、スマートフォン、携帯情報端末（PDA）、ラップトップコンピュータ、ネットブックコンピュータ、デスクトップコンピュータ、シンククライアントデバイス、タブレットコンピュータ、電子書籍（e-ブック）リーダー、ゲーム機、他などのコンピューティングデバイスであってもよい。例によっては、ユーザデバイスは、他のネットワーク1508を介してサービスプロバイダ・コンピュータと通信状態であってもよい。さらに、ユーザデバイスは、サービスプロバイダ・コンピュータによって管理され、制御される分散システムの一部であっても、他に、サービスプロバイダ・コンピュータの一部（たとえば、サービスプロバイダ・コンピュータと統合されたコンソールデバイス）であってもよい。

【0161】

[0177]図15のノード1102a~1102hは、1つまたは複数のサービスプロバイダ・コンピュータも表してもよい。1つまたは複数のサービスプロバイダ・コンピュータは、ユーザが対話し得るユーザデバイス上で実行されるように構成されるネイティブアプリケーションを提供してもよい。サービスプロバイダ・コンピュータは、例によっては、クライアントエンティティ、低遅延のデータストレージ、耐久性のあるデータストレージ、データアクセス、管理、仮想化、クラウドベースのソフトウェアソリューション、電子コンテンツのパフォーマンス管理、等々などの、但しこれらに限定されないコンピューティングリソースを提供してもよい。サービスプロバイダ・コンピュータは、また、ウェブホスティング、データベース化、コンピュータアプリケーション開発および/または実装プラットフォーム、これらの、またはこれらに類似するものの組合せをユーザへ提供するように動作可能であってもよい。実施形態によっては、サービスプロバイダ・コンピュータは、ホストされるコンピューティング環境において実装される1つまたは複数の仮想マシンとして提供されてもよい。ホストされるコンピューティング環境には、1つまたは複数の迅速にプロビジョニングおよびリリースされるコンピューティングリソースが含まれてもよい。これらのコンピューティングリソースは、コンピューティング、ネットワーキング、および/またはストレージデバイスを含んでもよい。ホストされるコンピューティング環境は、クラウドコンピューティング環境と称されることもある。サービスプロバイダ・コンピュータは、おそらくはクラスタに配置される1つまたは複数のサーバをサーバファームとして、または互いに関連づけられない個々のサーバとして含んでもよく、かつアプリケーション1532および/またはクラウドベースのソフトウェアサービスをホストしてもよい。これらのサーバは、統合された分散コンピューティング環境の一部として構成されてもよい。態様によっては、サービスプロバイダ・コンピュータは、追加的に、または代替として、たとえば、携帯電話、スマートフォン、携帯情報端末（PDA）、ラップトップコンピュータ、デスクトップコンピュータ、ネットブックコンピュータ、サーバコンピュータ、シンククライアントデバイス、タブレットコンピュータ、ゲーム機、他な

10

20

30

40

50

どのコンピューティングデバイスを含んでもよい。インスタンスによっては、サービスプロバイダ・コンピュータは、1つまたは複数の第三者コンピュータと通信してもよい。

【0162】

[0178]ある例示的な構成において、ノード1102a~1102hは、少なくとも1つのメモリ1518と、1つまたは複数の処理ユニット（または、プロセッサ1520）とを含んでもよい。プロセッサ1520は、ハードウェア、コンピュータ実行可能命令、ファームウェア、またはこれらの組み合わせで実装されてもよい。プロセッサ1520のコンピュータ実行可能命令またはファームウェア実装は、記述される様々な関数を実行するために任意の適切なプログラミング言語で書かれたコンピュータ実行可能または機械実行可能命令を含んでもよい。

10

【0163】

[0179]インスタンスによっては、ハードウェアプロセッサ1520は、シングルコアプロセッサであっても、マルチコアプロセッサであってもよい。マルチコアプロセッサは、同じプロセッサ内に複数の処理ユニットを含んでもよい。実施形態によっては、マルチコアプロセッサは、バスおよび二次または三次キャッシュなどの所定のリソースを共有してもよい。インスタンスによっては、シングルコアまたはマルチコアプロセッサにおける各コアは、複数の実行論理プロセッサ（または、実行スレッド）も含むことがある。このようなコア（たとえば、複数の論理プロセッサを備えるコア）では、実行パイプラインの幾つかのステージ、およびより低次のキャッシュをも、共有されることがある。

【0164】

20

[0180]メモリ1518は、プロセッサ1520上でロード可能かつ実行可能なプログラム命令、ならびにこれらのプログラムの実行中に発生されるデータを記憶してもよい。ノード1102a~1102hの構成およびタイプに依存して、メモリ1518は、揮発性（RAMなど）であっても、かつ/または不揮発性（ROM、フラッシュメモリ、他など）であってもよい。メモリ1518は、オペレーティングシステム1528、1つまたは複数のデータストア1530、1つまたは複数のアプリケーションプログラム1532、1つまたは複数のドライバ1534、および/または本明細書に開示している特徴を実装するためのサービスを含んでもよい。

【0165】

[0181]オペレーティングシステム1528は、タスクのスケジューリング、アプリケーションの実行、および/またはコントローラ周辺機器などのノード1102a~1102hの基本機能をサポートしてもよい。実装によっては、サービスプロバイダ・コンピュータは、1つまたは複数の仮想マシンをホストしてもよい。これらの実装において、各仮想マシンは、その独自のオペレーティングシステムを実行するように構成されてもよい。オペレーティングシステムの例としては、Unix、Linux（登録商標）、Windows、Mac OS、iOS、Androidおよびこれらに類似するものが含まれる。また、オペレーティングシステム1528は、独自仕様のオペレーティングシステムである場合もある。

30

【0166】

[0182]データストア1530は、オペレーティングシステム1528、アプリケーションプログラム1532またはドライバ1534によって使用されかつ/または作用される永続的または一時的なデータを含んでもよい。このようなデータの例としては、ウェブページ、ビデオデータ、オーディオデータ、画像、ユーザデータ、等々が含まれる。データストア1530内の情報は、実装によっては、ネットワーク1508を介してノード1102a~1102hへ提供されてもよい。事例によっては、データストア1530は、追加的に、または代替として、記憶されたアプリケーションプログラムおよび/またはドライバを含んでもよい。あるいは、または追加的に、データストア1530は、標準および/または独自のソフトウェアライブラリ、および/または標準および/または独自のアプリケーション・ユーザ・インタフェース（API）ライブラリを記憶してもよい。データストア1530に記憶される情報は、機械可読オブジェクトコード、ソースコード、イン

40

50

タープリタ型コードまたは中間コードであってもよい。

【0167】

[0183]ドライバ1534は、ノード内のコンポーネント間の通信を提供し得るプログラムを含む。たとえば、一部のドライバ1534は、オペレーティングシステム1528と、追加のストレージ1522、ネットワークデバイス1524および/またはI/Oデバイス1526との間の通信を提供してもよい。あるいは、または追加的に、一部のドライバ1534は、アプリケーションプログラム1532とオペレーティングシステム1528との間、および/またはアプリケーションプログラム1532とサービスプロバイダ・コンピュータへアクセス可能な周辺デバイスとの間の通信を提供してもよい。多くの事例において、ドライバ1534は、十分に理解されている機能を提供するドライバ(たとえば、プリンタドライバ、ディスプレイドライバ、ハードディスク・ドライバ、ソリッドステートデバイス・ドライバ)を含んでもよい。他の事例において、ドライバ1534は、独自仕様または特殊化された機能を提供してもよい。

10

【0168】

[0184]サービスプロバイダ・コンピュータまたはサーバも、追加のストレージ1522を含んでもよく、追加のストレージ1522は、取外し可能ストレージおよび/または取外し不能ストレージを含んでもよい。追加のストレージ1522としては、磁気ストレージ、光ディスク、ソリッド・ステート・ディスク、フラッシュメモリおよび/またはテープストレージが含まれてもよい。追加のストレージ1522は、ノード1102a~1102hと同じシャーシ内に収容されてもよく、または外部エンクロージャ内に存在してもよい。メモリ1518および/または追加のストレージ1522およびこれらに関連づけられるコンピュータ可読媒体は、コンピュータ可読命令、データ構造、プログラムモジュール、およびコンピューティングデバイスに関する他のデータの揮発性ストレージを提供してもよい。実装によっては、メモリ1518は、SRAM、DRAMまたはROMなどの複数の異なるタイプのメモリを含んでもよい。

20

【0169】

[0185]共に取外し可能および取外し不能であるメモリ1518および追加のストレージ1522は、コンピュータ可読記憶媒体の例である。たとえば、コンピュータ可読記憶媒体としては、情報を記憶するための方法または技術において実装される揮発性または不揮発性、取外し可能または取外し不能な媒体が含まれてもよく、情報としては、たとえば、コンピュータ可読命令、データ構造、プログラムモジュールまたは他のデータが含まれる。メモリ1518および追加のストレージ1522は、コンピュータ記憶媒体の例である。ノード1102a~1102h内に存在し得る追加のタイプのコンピュータ記憶媒体としては、PRAM、SRAM、DRAM、RAM、ROM、EEPROM、フラッシュメモリまたは他のメモリ技術、CD-ROM、DVDまたは他の光ストレージ、磁気カセット、磁気テープ、磁気ディスクストレージまたは他の磁気ストレージデバイス、ソリッド・ステート・ドライブ、または、所望される情報を記憶するために使用されることが可能でありかつノード1102a~1102hからアクセスされることが可能な他の何らかの媒体が含まれ得るが、この限りではない。コンピュータ可読媒体としては、1つの媒体タイプの複数のユニットを含む、上述の媒体タイプのいずれかの組合せも含まれる。

30

40

【0170】

[0186]あるいは、または追加的に、コンピュータ可読通信媒体は、コンピュータ可読命令、プログラムモジュール、または搬送波または他の伝送などのデータ信号内で送信される他のデータを含んでもよい。しかしながら、本明細書で使用するコンピュータ可読記憶媒体という用語は、コンピュータ可読通信媒体を包含しない。

【0171】

[0187]ノード1102a~1102hは、キーボード、マウス、ペン、音声入力デバイス、タッチ入力デバイス、ディスプレイ、スピーカ、プリンタおよびこれらに類似するものなどのI/Oデバイス1526も含んでもよい。ノード1102a~1102hは、1つまたは複数の通信チャネル1536も含んでもよい。通信チャネル1536は、ノード

50

1102a~1102hの様々なコンポーネントがそれを介して通信できる媒体を提供してもよい。1つまたは複数の通信チャネル1536は、バス、リング、スイッチングファブリックまたはネットワークの形式をとってもよい。

【0172】

[0188]ノード1102a~1102hは、ノード1102a~1102hが記憶されたデータベース、別のコンピューティングデバイスまたはサーバ、ユーザ端末、および/またはネットワーク1500上の他のデバイスと通信できるようにする、ネットワークデバイス1524も含んでもよい。

【0173】

[0189]実装によっては、ネットワークデバイス1524は、PCIベースのデバイスなどの周辺機器である。これらの実装において、ネットワークデバイス1524は、ホストデバイスと通信するためのPCIインタフェースを含む。「PCI」または「PCIベースの」という用語は、原初のPCI標準、PCI-X、アクセラレーテッド・グラフィックス・ポート(AGP)、およびPCI-Express(PCIe)を含む、バスプロトコルのPCIファミリにおける任意のプロトコル、または、本明細書において論じるPCIプロトコルを基礎とする他の任意の改良または派生プロトコルを記述するために使用される場合がある。PCIベースのプロトコルは、ローカル周辺機器などのデバイスをホストデバイスへ接続するための標準バスプロトコルである。標準バスプロトコルは、様々な製造業者がそのための仕様を規定して採用しているデータ転送プロトコルである。製造業者は、対応デバイスがバスプロトコルを実装するコンピューティングシステムと互換性があること、およびその逆を保証している。本明細書で使用するPCIベースのデバイスという用語には、不揮発性メモリエクスプレス(NVMe)を用いて通信するデバイスも含まれる。NVMeは、コンピューティングシステムへ取り付けられる不揮発性記憶媒体にPCIeを用いてアクセスするためのデバイスインタフェース仕様である。たとえば、バス・インタフェース・モジュールは、NVMeを実装してもよく、よってネットワークデバイス1524は、PCIeインタフェースを用いてコンピューティングシステムへ接続されてもよい。

【0174】

[0190]PCIベースのデバイスは、1つまたは複数の機能を含んでもよい。「機能」は、ネットワークデバイス1524により提供され得る動作を記述する。機能の例としては、とりわけ、大容量ストレージコントローラ、ネットワークコントローラ、ディスプレイコントローラ、メモリコントローラ、シリアル・バス・コントローラ、ワイヤレスコントローラ、および暗号化および復号化コントローラが含まれる。事例によっては、PCIベースのデバイスは、2つ以上の機能を含んでもよい。たとえば、PCIベースのデバイスは、大容量ストレージコントローラおよびネットワークアダプタを提供してもよい。別の例として、PCIベースのデバイスは、2つの異なるストレージリソースを制御するために2つのストレージコントローラを提供してもよい。実装によっては、PCIベースのデバイスは、最大8つの機能を有してもよい。

【0175】

[0191]実装によっては、ネットワークデバイス1524は、シングルルートI/O仮想化(SR-IOV)を含んでもよい。SR-IOVは、PCIベースのデバイスに含まれ得る拡張機能である。SR-IOVは、1つの物理リソース(たとえば、単一のネットワーク・インタフェース・コントローラ)が複数のリソース(たとえば、64個のネットワーク・インタフェース・コントローラ)として現出することを可能にする。したがって、所定の機能(たとえば、ネットワーク・インタフェース・コントローラ)を提供するPCIベースのデバイスは、PCIベースのデバイスを利用するデバイスには、同じ機能を提供する複数のデバイスであるように見えることがある。SR-IOV対応のストレージ・アダプタ・デバイスの機能は、物理機能(PF)または仮想機能(VF)として分類されてもよい。物理機能は、発見され、管理されかつ操作されることが可能な、デバイスが備える完全な機能である。物理機能は、ストレージ・アダプタ・デバイスを構成または制

10

20

30

40

50

御するために使用されることが可能な構成リソースを有する。物理機能は、非仮想化デバイスが有すると思われる同じ構成アドレス空間およびメモリアドレス空間を含む。1つの物理機能は、それに関連づけられる幾つかの仮想機能を有する場合がある。仮想機能は、物理機能に類似するものであるが、概して構成リソースが不足し得る軽量の機能であり、概して、その基礎を成す物理機能の構成によって制御される。物理機能および/または仮想機能は、各々、ホストデバイス上で実行される個々の実行スレッド（たとえば、仮想マシンなど）へ割り当てられてもよい。

【0176】

[0192]一連の機能の一部として、様々な例では、ノード1102a~1102hのうち
の1つまたはそれ以上が、これまでにその様々な例について述べたニューラルネットワー
ク処理エンジンおよび/またはニューラルネットワークプロセッサを包含することができ
る。これらの例において、こうして装備されたノード1102a~1102hは、機械学
習タスクおよびニューラルネットワーク処理タスク用に最適化されることが可能である。

10

【0177】

[0193]本明細書に記述しているモジュールは、ソフトウェアモジュール、ハードウェア
モジュール、またはこれらの適切な組合せであってもよい。モジュールがソフトウェアモ
ジュールであれば、モジュールは、非一時的コンピュータ可読媒体上に具現され、かつ本
明細書に記述しているコンピュータシステムのいずれかにおけるプロセッサによって処理
されることが可能である。記述しているプロセスおよびアーキテクチャが、リアルタイム
、またはユーザ対話に先行する非同期モードのいずれかで実行可能であることは、留意さ
れるべきである。モジュールは、図15において提案されている方法で構成されてもよく
、かつ/または本明細書に記載の機能は、別個のモジュールとして存在する1つまたは複
数のモジュールによって提供されることが可能であり、かつ/または本明細書に記載のモ
ジュール機能は、複数のモジュールに分散されることが可能である。

20

【0178】

[0194]したがって、本明細書本文および図面は、限定を意味するものではなく、例示的
なものとなされるべきである。しかしながら、これらに対して、特許請求の範囲に記載
されている本開示のより広い精神および範囲を逸脱することなく、様々な改変および変更
が行われ得ることは、明らかであろう。

【0179】

[0195]開示する技術の例示的な実施形態は、以下の条項によって記述することができる
。

30

1. ニューラルネットワーク処理のための集積回路であって、

ニューラルネットワークの重み値セットを記憶する複数のメモリバンクであって、重み
値セットは、ニューラルネットワークの全ての重み値を含み、重み値セットは、結果が既
知である入力データを用いて先に決定されていて、複数のメモリバンクからの各バンクは
、独立してアクセス可能である、複数のメモリバンクと、

第1の処理エンジンアレイであって、各処理エンジンは、乗数累算器回路を含み、第1
の処理エンジンアレイは、複数のメモリバンクと同じダイ上にある、第1の処理エンジン
アレイと、を備え、

40

入力データを受信した時点で、集積回路は、重み値セットを用いて、ニューラルネット
ワークが実行するように訓練されたタスクを実行するように構成され、タスクは、結果が
既知である入力データによって規定され、かつ、タスクを実行することは、

複数のメモリバンクから重み値を読み取ることと、

重み値および入力データを第1の処理エンジンアレイに入力することであって、第1の
処理エンジンアレイにおける各処理エンジンは、複数の重み値からの1つの重み値と、複
数の入力データからの1つの入力値とを用いて加重和を計算することと、

結果を計算すること、を含み、結果を計算することは、第1の処理エンジンアレイから
の出力を累積することを含み、結果は、タスク実行の結果に対応する、集積回路。

2. 第2の処理エンジンアレイをさらに備え、第1のメモリバンクセットは、第1の処

50

理エンジンアレイによって使用されるように構成され、第2のメモリバンクセットは、第2の処理エンジンアレイによって使用されるように構成され、かつ重み値セットは、第1のメモリバンクセットおよび第2のメモリバンクセットに記憶される、条項1に記載の集積回路。

3. オフチップメモリとの通信を可能にするメモリコントローラと、
 ホストバスとの通信を可能にするバス・インタフェース・コントローラと、
 集積回路のコンポーネント間でデータを移動するように構成される管理コントローラと、
 複数のメモリバンク、メモリコントローラ、バス・インタフェース・コントローラ、および管理コントローラ間の通信を可能にする通信ファブリックと、をさらに備える、条項1～2に記載の集積回路。

10

4. 集積回路であって、
 第1の処理エンジンアレイと、
 ニューラルネットワークの重み値セットを記憶する複数のメモリバンクと、を備え、複数のメモリバンクからの各バンクは、独立してアクセス可能であり、かつ複数のメモリバンクおよび第1の処理エンジンアレイは、同じダイ上にあり、

入力データを受信した時点で、集積回路は、重み値セットを用いて、ニューラルネットワーク用に規定されたタスクを実行するように構成され、かつ、タスクを実行することは、複数のメモリバンクから重み値を読み取ることと、

重み値および入力データを第1の処理エンジンアレイに入力することと、
 第1の処理エンジンアレイを用いて結果を計算すること、を含み、結果は、タスク実行の結果に対応する、集積回路。

20

5. タスクを実行することは、さらに、
 複数のメモリバンクのうちの異なるメモリバンクから2つ以上の値を同時に読み取ることを含む、条項4に記載の集積回路。

6. 2つ以上の値には、重み値、入力値または中間結果が含まれる、条項5に記載の集積回路。

7. タスクを実行することは、さらに、
 複数のメモリバンクのうちの第1のメモリバンクへ第1の値を書き込むことと、
 複数のメモリバンクのうちの第2のメモリバンクから第2の値を読み取ること、を含み、第1の値は、第2の値が読み取られると同時に書き込まれる、条項4～6のいずれかに記載の集積回路。

30

8. 第1の値および第2の値は、重み値、および入力値、または中間結果を含む、条項7に記載の集積回路。

9. 重み値セットは、ニューラルネットワークの全ての重み値を含む、条項4～8のいずれかに記載の集積回路。

10. 第2の処理エンジンアレイをさらに備え、複数のメモリバンクのうちの第1のメモリバンクセットは、第1の処理エンジンアレイによって使用されるように構成され、複数のメモリバンクのうちの第2のメモリバンクセットは、第2の処理エンジンアレイによって使用されるように構成され、第1のメモリバンクセットおよび第2のメモリバンクセットは各々、重み値セットの一部を含み、かつ、タスクを実行することは、さらに、

40

第1の処理エンジンアレイによって中間結果を計算することであって、第1の処理エンジンアレイは、中間結果を、第1のメモリバンクセットからの重み値を用いて計算することと、

第1の処理エンジンアレイによって、第2のメモリバンクセットから追加の重み値を読み取ることであって、第1の処理エンジンアレイは、中間結果および追加の重み値を用いて結果を計算すること、を含む、条項4～9のいずれかに記載の集積回路。

11. 重み値セットによる占有は、第2のメモリバンクセットの全てより少なく、第2の処理エンジンアレイは、重み値セットによって占有されていない第2のメモリバンクセットの一部を用いて計算を実行する、条項10に記載の集積回路。

12. 第1のメモリバンクセットおよび第2のメモリバンクセットに記憶される重み値

50

セットの部分は、ニューラルネットワークの全ての重み値を含む、条項 10 に記載の集積回路。

13. 複数のメモリバンクの第 1 の部分は、結果を計算するための中間結果を記憶するために予約され、かつ重み値セットは、ニューラルネットワークの全ての重み値より少ない重み値を含む、条項 4 ~ 12 のいずれかに記載の集積回路。

14. タスクを実行することは、さらに、
中間結果を記憶するために必要なメモリの量が減少していると決定することと、
別のメモリから追加の重み値セットを読み取ることと、
追加の重み値セットを複数のメモリバンクの第 1 の部分に記憶することを含み、追加の重み値セットは、結果の計算に必要とされる前に記憶される、条項 13 に記載の集積回路。

15. 第 1 の処理エンジンアレイは、1 つの処理エンジンセットを含み、処理エンジンセットのうちの各処理エンジンは、結果を、処理エンジンセットのうちの別の処理エンジンへ直に出力する、条項 4 ~ 14 のいずれかに記載の集積回路。

16. 第 1 の処理エンジンアレイのうちの各処理エンジンは、乗数累算器回路を含む、条項 4 ~ 15 のいずれかに記載の集積回路。

17. ニューラルネットワークは、重み付き有向グラフから導出される複数の重み値と、重み付き有向グラフにおける各ノードについて計算を実行するための命令セットとを含み、複数の重み値は、既知の入力データを用いてタスクを実行することにより、先に決定されている、条項 4 ~ 16 のいずれかに記載の集積回路。

18. 方法であって、
ニューラルネットワーク処理回路の複数のメモリバンクに、1 つの重み値セットを記憶することであって、ニューラルネットワーク処理回路は、複数のメモリバンクと同じダイ上に 1 つの処理エンジンアレイを含み、かつ重み値セットは、入力データの受信に先行して記憶されることと、

入力データを受信することと、
重み値セットを用いて、ニューラルネットワーク用に規定されるタスクを実行すること、
を含み、タスクを実行することは、

複数のメモリバンクから重み値を読み取ることと、
重み値および入力データを処理エンジンアレイに入力することと、
処理エンジンアレイを用いて結果を計算すること、を含み、結果は、タスク実行の結果に対応する、方法。

19. 重み値セットは、ニューラルネットワークの全ての重み値を含む、条項 18 に記載の方法。

20. 重み値セットは、ニューラルネットワークの全ての重み値の第 1 の部分を含み、かつ、さらに、

複数のメモリバンクが利用可能なスペースを有すると決定することと、
ニューラルネットワークの全ての重み値の第 2 の部分を読み取ることであって、第 2 の部分は、追加のメモリから読み取られることと、

第 2 の部分を使用可能なスペースへ書き込むこと、を含む、条項 18 または 19 に記載の方法。

21. 追加のメモリは、同じダイ上の第 2 の処理エンジンアレイに関連づけられる、請求項 18 ~ 20 のいずれかに記載の方法。

22. 追加のメモリは、オフチップである、請求項 18 ~ 21 のいずれかに記載の方法。

23. 重み値を読み取るとは、同時に、複数のメモリバンクのうちの第 1 のメモリバンクから第 1 の重み値を読み取り、かつ複数のメモリバンクのうちの第 2 のメモリバンクから第 2 の重み値を読み取ることを含む、条項 18 ~ 22 のいずれかに記載の方法。

24. 処理エンジンアレイを用いて中間結果を決定することと、
中間結果を、複数のメモリバンクのうちの 1 つのメモリバンクに記憶すること、をさらに含み、中間結果は、重み値が読み取られると同時に記憶される、条項 18 ~ 23 のいずれかに記載の方法。

10

20

30

40

50

25. ニューラルネットワーク処理システムであって、

第1の処理エンジンアレイであって、各処理エンジンが乗数累算器回路を含む、第1の処理エンジンアレイと、

第1の重み値セットを記憶する第1のメモリバンクセットであって、第1の重み値セットは、結果が既知である入力データを用いて先に決定されており、第1のメモリバンクセットのうちの各バンクは、独立してアクセス可能であり、かつ第1の処理エンジンアレイおよび第1のメモリバンクセットは、同じ第1のダイ上にある、第1のメモリバンクセットと、

第2の処理エンジンアレイであって、各処理エンジンが乗数累算器回路を含む、第2の処理エンジンアレイと、

第2の重み値セットを記憶する第2のメモリバンクセットであって、第2の重み値セットは、結果が既知である入力データを用いて先に決定されており、第2のメモリバンクセットのうちの各バンクは、独立してアクセス可能であり、かつ第2の処理エンジンアレイおよび第2のメモリバンクセットは、同じ第2のダイ上にある、第2のメモリバンクセットと、を備え、第1の重み値セットおよび第2の重み値セットは、ニューラルネットワークの全ての重み値を含み、かつ第1の重み値セットおよび第2の重み値セットは、ニューラルネットワーク処理システムによる入力データの受信に先行して記憶され、

入力データを受信した時点で、ニューラルネットワーク処理システムは、ニューラルネットワークが実行するように訓練されたタスクを実行するように構成され、タスクは、結果が既知である入力データによって規定され、かつ、タスクを実行することは、

第1の処理エンジンアレイを用いて中間結果を計算することであって、第1の処理エンジンアレイにおける各処理エンジンは、第1の重み値セットからの1つの重み値および入力データからの1つの入力値を用いて加重和を計算することと、

第1のメモリバンクセットからの中間結果および状態を、第2のメモリバンクセットへコピーすることであって、状態は、第1の処理エンジンアレイにより完了されるタスクのステージを記述することと、

第2の処理エンジンアレイを用いて最終結果を計算すること、を含み、第2の処理エンジンアレイは、状態を用いてタスクの状態を決定し、第2の処理エンジンアレイにおける各処理エンジンは、第2の重み値セットからの1つの重み値、および中間結果からの1つの値を用いて加重和を計算し、かつ最終結果は、タスク実行の結果に対応する、ニューラルネットワーク処理システム。

26. 第1のダイおよび第2のダイは、同じダイの部分であり、ダイは、通信ファブリックと、ダイレクト・メモリ・アクセス・コントローラとを含み、かつダイレクト・メモリ・アクセス・コントローラは、中間結果および状態を通信ファブリックを介してコピーする、条項25に記載のニューラルネットワーク処理システム。

27. 第1のダイは、第1のパッケージ内にあり、かつ第2のダイは、第2のパッケージ内にあり、かつ、

第1のダイおよび第2のダイに連結されるバスをさらに備え、中間結果および状態は、バスを介してコピーされる、条項25または26に記載のニューラルネットワーク処理システム。

28. ニューラルネットワーク処理システムであって、

第1の処理エンジンアレイと、

第1の重み値セットを記憶する第1のメモリバンクセットであって、第1の処理エンジンアレイおよび第1のメモリバンクセットは、同じ第1のダイ上にある、第1のメモリバンクセットと、

第2の処理エンジンアレイと、

第2の重み値セットを記憶する第2のメモリバンクセットであって、第2の処理エンジンアレイおよび第2のメモリバンクセットは、同じ第2のダイ上にある、第2のメモリバンクセットと、を備え、第1の重み値セットおよび第2の重み値セットは、ニューラルネットワークの全ての重み値を含み、かつ第1の重み値セットおよび第2の重み値セットは

10

20

30

40

50

、ニューラルネットワーク処理システムによる入力データの受信に先行して記憶され、
入力データを受信した時点で、ニューラルネットワーク処理システムは、ニューラルネットワーク用に規定されるタスクを実行するように構成され、かつ、タスクを実行することは、

第1の処理エンジンアレイを用いて中間結果を計算することであって、第1の処理エンジンアレイは、第1の重み値セットおよび入力データを用いて中間結果を計算することと、
第1のメモリバンクセットからの中間結果および状態を、第2のメモリバンクセットへコピーすることであって、状態は、第1の処理エンジンアレイにより完了されるタスクのステージを記述することと、

第2の処理エンジンアレイを用いて最終結果を計算すること、を含み、第2の処理エンジンアレイは、最終結果を、第2の重み値セット、中間結果および状態を用いて計算し、かつ最終結果は、タスク実行の結果に対応する、ニューラルネットワーク処理システム。

10

29. 第1のダイおよび第2のダイは、同じダイの部分である、条項28に記載のニューラルネットワーク処理システム。

30. 第1のダイおよび第2のダイは、異なるダイである、条項28または29に記載のニューラルネットワーク処理システム。

31. バスをさらに備え、第1のダイおよび第2のダイは各々、バスに連結され、かつ中間結果および状態は、バスを介してコピーされる、条項30に記載のニューラルネットワーク処理システム。

32. タスクを実行することは、さらに、

20

第1のメモリバンクセットまたは第2のメモリバンクセットのうちの異なるメモリバンクから、2つ以上の値を同時に読み取ることを含む、条項28～31のいずれかに記載のニューラルネットワーク処理システム。

33. 2つ以上の値には、重み値、入力値または中間結果が含まれる、請求項32に記載のニューラルネットワーク処理システム。

34. タスクを実行することは、さらに、

第1のメモリバンクセットまたは第2のメモリバンクセットのうちの第1のメモリバンクへ第1の値を書き込むことと、

第1のメモリバンクセットまたは第2のメモリバンクセットのうちの第2のメモリバンクから第2の値を読み取ること、を含み、第1の値は、第2の値が読み取られると同時に書き込まれる、条項28～33のいずれかに記載のニューラルネットワーク処理システム。

30

35. 第1の値または第2の値には、重み値、入力値または中間結果が含まれる、請求項34に記載のニューラルネットワーク処理システム。

36. タスクを実行する間、プロセッサメモリから重み値は読み取られず、プロセッサメモリは、第1のダイまたは第2のダイとは異なるダイ上にある、条項28～35のいずれかに記載のニューラルネットワーク処理システム。

37. 第1の処理エンジンアレイおよび第2の処理エンジンアレイのうちの各処理エンジンは、乗数累算器回路を含む、条項28～36のいずれかに記載のニューラルネットワーク処理システム。

38. ニューラルネットワークは、重み付き有向グラフから導出される複数の重み値と、重み付き有向グラフにおける各ノードについて計算を実行するための命令セットとを含み、複数の重み値は、既知の入力データを用いてタスクを実行することにより、先に決定されている、条項28～37のいずれかに記載のニューラルネットワーク処理システム。

40

39. ニューラルネットワークを処理する方法であって、

第1の重み値セットを、ニューラルネットワーク処理システムの第1のメモリバンクセットに記憶することと、

第2の重み値セットを、ニューラルネットワーク処理システムの第2のメモリバンクセットに記憶することであって、第1の重み値セットおよび第2の重み値セットは、ニューラルネットワークの全ての重み値を含み、かつ第1の重み値セットおよび第2の重み値セットは、ニューラルネットワーク処理システムによる入力データの受信に先行して記憶さ

50

れることと、

入力データを受信することと、

ニューラルネットワーク用に規定されるタスクを実行すること、を含み、タスクを実行することは、

第1の処理エンジンアレイを用いて中間結果を計算することであって、第1の処理エンジンアレイは、第1のメモリバンクセットと同じ第1のダイ上にあり、かつ第1の処理エンジンアレイは、第1の重み値セットおよび入力データを用いて中間結果を計算することと、

第1のメモリバンクセットからの中間結果および状態を、第2のメモリバンクセットへコピーすることであって、状態は、第1の処理エンジンアレイにより完了されるタスクのステージを記述することと、

10

第2の処理エンジンアレイを用いて最終結果を計算すること、を含み、第2の処理エンジンアレイは、第2のメモリバンクセットと同じ第2のダイ上にあり、第2の処理エンジンアレイは、最終結果を、第2の重み値セット、中間結果および状態を用いて計算し、かつ最終結果は、タスク実行の結果に対応する、方法。

40．第1のダイおよび第2のダイは、同じダイの部分である、条項39に記載の方法。

41．第1のダイおよび第2のダイは、異なるパッケージ内にある、請求項39に記載の方法。

42．プロセッサメモリから重み値は読み取られず、プロセッサメモリは、第1のダイまたは第2のダイのいずれかとは異なるパッケージ内にある、条項39～41のいずれかに記載の方法。

20

43．第1のメモリバンクセットから2つ以上の値を同時に読み取ることと、

第2のメモリバンクセットから2つ以上の値を同時に読み取ること、をさらに含む、条項39～42のいずれかに記載の方法。

44．第1のメモリバンクセットから第1の値を読み取ると同時に、第1のメモリバンクセットへ第2の値を書き込むことと、

第2のメモリバンクセットから第3の値を読み取ると同時に、第2のメモリバンクセットへ第4の値を書き込むこと、をさらに含む、条項39～43のいずれかに記載の方法。

45．ニューラルネットワークを処理するための集積回路であって、

処理エンジンアレイと、

30

第1のニューラルネットワークの第1の重み値セットを記憶する複数のメモリバンクと、を備え、第1のニューラルネットワークは、第1のタスクを実行するように構成され、第1の重み値セットは、結果が既知である入力データを用いて第1のタスクを実行することにより、先に決定されており、複数のメモリバンクからの各バンクは、独立してアクセス可能であり、かつ複数のメモリバンクおよび処理エンジンアレイは、同じダイ上にあり、集積回路は、

第1のタスクに関連づけられる第1の入力データを受信するように構成され、第1の入力データは、複数のメモリバンクに記憶され、かつ集積回路は、

第1の重み値セットおよび第1の入力データを用いて第1の結果を計算するように構成され、第1の結果は、処理エンジンアレイを用いて計算され、かつ第1の結果は、第1のタスク実行の結果に対応し、かつ集積回路は、

40

第1の結果を計算する間に、複数のメモリバンクが使用可能なスペースを有することを決定し、

利用可能なスペースに第2の重み値セットを記憶するように構成され、第2の重み値セットは、第2のニューラルネットワーク用であり、第2のニューラルネットワークは、第2のタスクを実行するように構成され、第2のタスクは、第1のタスクとは異なり、かつ、第2の重み値セットを記憶することは、集積回路が第1の入力データに基づく第1のタスクの実行を完了する前に、第2のタスクを実行するように集積回路を準備し、かつ集積回路は、

第2のタスクに関連づけられる第2の入力データを受信するように構成され、第2の入

50

カデータは、複数のメモリバンクに記憶され、かつ集積回路は、

第2の重みセットおよび第2の入力データを用いて第2の結果を計算するように構成され、第2の結果は、処理エンジンアレイを用いて計算され、第2の結果の計算は、第1の結果の計算が完了した時点で発生し、かつ第2の結果は、第2のタスク実行の結果に対応する、集積回路。

46．複数のメモリバンクが利用可能なスペースを有すると決定することは、

第1の結果を計算するための中間結果を記憶するために必要なスペースの量を決定することと、

スペースの量が減少したと決定すること、を含む、条項45に記載の集積回路。

47．複数のメモリバンクが利用可能なスペースを有すると決定することは、

第1の結果を計算するために既に使用されている1つまたは複数の重み値を識別することと、

1つまたは複数の重み値によって占有されるスペースを利用可能としてマーキングすること、を含む、条項45または46に記載の集積回路。

48．第1の重み値セットおよび第1のタスクは、第1のコンテキストに関連づけられる、条項45～47のいずれかに記載の集積回路。

49．第2の重み値セットおよび第2のタスクは、第2のコンテキストを構成する、条項45～48のいずれかに記載の集積回路。

50．集積回路であって、

処理エンジンアレイと、

第1のニューラルネットワークの第1の重み値セットを記憶する複数のメモリバンクと、を備え、第1のニューラルネットワークは、第1のタスクを実行するように構成され、複数のメモリバンクからの各バンクは、独立してアクセス可能であり、かつ複数のメモリバンクおよび処理エンジンアレイは、同じダイ上にあり、

集積回路は、

第1のタスクに関連づけられる第1の入力データを受信し、

第1の重み値セットおよび第1の入力データを用いて第1の結果を計算するように構成され、第1の結果は、処理エンジンアレイを用いて計算され、かつ第1の結果は、第1のタスク実行の結果に対応し、かつ集積回路は、

複数のメモリバンクに第2の重み値セットを記憶するように構成され、第2の重み値セットは、第2のニューラルネットワーク用であり、第2のニューラルネットワークは、第2のタスクを実行するように構成され、かつ第2の重み値セットを記憶することは、第1の結果を計算する間に発生し、かつ集積回路は、

第2のタスクに関連づけられる第2の入力データを受信し、かつ、

第2の重み値セットおよび第2の入力データを用いて第2の結果を計算するように構成され、第2の結果は、処理エンジンアレイを用いて計算され、第2の結果の計算は、第1の結果の計算が完了した時点で発生し、かつ第2の結果は、第2のタスク実行の結果に対応する、集積回路。

51．第1の重み値セットは、第1の入力データの受信に先行して、複数のメモリバンクに記憶される、条項50に記載の集積回路。

52．第1の重み値セットは、第1のニューラルネットワークの全ての重み値を含む、条項50または51に記載の集積回路。

53．第2の結果の計算の完了に先行して、第2のニューラルネットワークの全ての重み値が複数のメモリバンクに記憶される、条項50～52のいずれかに記載の集積回路。

54．第2のニューラルネットワークの追加の重み値セットを記憶するようにさらに構成され、追加の重み値セットは、第2の結果を計算する間に記憶される、条項50～53のいずれかに記載の集積回路。

55．複数のメモリバンクにおける、第1の結果を計算するための中間結果を記憶するために必要なスペースの量を決定し、

第1の結果を計算する間にスペースの量が減少したと決定し、かつ、

10

20

30

40

50

複数のメモリバンクにおける、もはや中間結果を記憶するために使用されないスペースを用いて、第2の重み値セットを記憶するように、さらに構成される、条項50～54のいずれかに記載の集積回路。

56. 第1の重み値セットのうちの、第1の結果の計算の完了に必要なとされない1つまたは複数の重み値を識別し、かつ、

複数のメモリバンクにおける、1つまたは複数の重み値により占有されるスペースを用いて第2の重み値セットを記憶するように、さらに構成される、条項50～55のいずれかに記載の集積回路。

57. 複数のメモリバンクのうちの異なるメモリバンクから2つ以上の値を同時に読み取るように、さらに構成される、条項50～56のいずれかに記載の集積回路。

10

58. 2つ以上の値には、重み値、入力値または中間結果が含まれる、条項50～57のいずれかに記載の集積回路。

59. 複数のメモリバンクのうちの第1のメモリバンクへ第1の値を書き込み、かつ、複数のメモリバンクのうちの第2のメモリバンクから第2の値を読み取るように、さらに構成され、第1の値は、第2の値が読み取られると同時に書き込まれる、条項50～58のいずれかに記載の集積回路。

60. 第1の値または第2の値は、重み値、入力値、または中間結果を含む、条項59に記載の集積回路。

61. 処理エンジンアレイは、1つの処理エンジンセットを含み、処理エンジンセットのうちの各処理エンジンは、結果を、処理エンジンセットのうちの別の処理エンジンへ直

20

に出力する、条項50～60のいずれかに記載の集積回路。

62. 処理エンジンアレイのうちの各処理エンジンは、乗数累算器回路を含む、条項50～61のいずれかに記載の集積回路。

63. 各ニューラルネットワークは、個々の重み付き有向グラフから導出される個々の複数の重み値と、重み付き有向グラフにおける各ノードについて計算を実行するための個々の命令セットとを含み、複数の重み値は、既知の入力データを用いて個々のタスクを実行することにより、先に決定されている、条項50～62のいずれかに記載の集積回路。

64. コンピュータ実行方法であって、

ニューラルネットワークプロセッサを含むコンピューティングシステムにより、第1のニューラルネットワークの第1の重み値セットを記憶することであって、第1のニューラルネットワークは、第1のタスクを実行するように構成され、かつ第1の重み値セットは、ニューラルネットワークプロセッサにおけるメモリの第1の部分に記憶されることと、

30

第1のタスクに関連づけられる第1の入力データを受信することであって、第1の入力データは、第1の重み値セットの記憶後に受信されることと、

ニューラルネットワークプロセッサを用いて第1の結果を計算することであって、第1の結果を計算することは、第1の重み値セットおよび第1の入力データを用いることを含み、第1の結果は、第1のタスク実行の結果に対応し、かつ第1の結果を計算する間に、メモリの第2の部分に中間結果が記憶されることと、

第1の結果を計算する間に、メモリの第3の部分に識別することであって、第3の部分は、メモリにおける未使用のスペースを含むことと、

40

第3の部分に第2の重み値セットを記憶することであって、第2の重み値セットは、第2のニューラルネットワーク用であり、第2のニューラルネットワークは、第2のタスクを実行するように構成されることと、

第2のタスクに関連づけられる第2の入力データを受信することと、

ニューラルネットワークプロセッサを用いて第2の結果を計算すること、を含み、第2の結果を計算することは、第2の重みセットおよび第2の入力データを用いることを含み、第2の結果の計算は、第1の結果の計算が完了した時点で発生し、かつ第2の結果は、第2のタスク実行の結果に対応する、コンピュータ実装方法。

65. メモリの第1の部分のサイズを縮小することをさらに含み、サイズの縮小は、第1の結果の計算の完了に必要なとされない第1の重み値セットからの重み値に対応する、条

50

項 6 4 に記載のコンピュータ実装方法。

6 6 . メモリの第 2 の部分のサイズを縮小することをさらに含み、サイズは、中間結果を記憶するために必要なメモリの量に対応する、条項 6 4 または 6 5 に記載のコンピュータ実装方法。

6 7 . 第 1 の部分、第 2 の部分および第 3 の部分は重複しない、請求項 6 4 ~ 6 6 のいずれかに記載のコンピュータ実装方法。

6 8 . 第 1 の重み値セットは、第 1 のニューラルネットワークの全ての重み値より少ない重み値を含み、かつ、

コンピューティングシステムの第 2 のメモリから、第 1 のニューラルネットワークの追加の重み値セットを読み取ることと、

追加の重み値セットをメモリの第 1 の部分に記憶すること、をさらに含み、追加の重み値セットは、第 2 の重み値セットの記憶に先行して記憶される、条項 6 4 ~ 6 7 のいずれかに記載のコンピュータ実装方法。

【 0 1 8 0 】

[0196]他の変形例も、本開示の精神の範囲内である。したがって、開示している技術に対しては、様々な改変および代替構造が可能であるが、図面に示しかつこれまでに詳述したものは、その所定の例示的な実施形態である。しかしながら、本開示を、開示した特定の 1 つまたは複数の形式に限定する意図はなく、逆に、添付の特許請求の範囲に規定されている本開示の精神および範囲に含まれる全ての改変、代替構造および同等物をカバーすることが意図される点は、理解されるべきである。

【 0 1 8 1 】

[0197]開示される実施形態を説明する文脈における（特に、添付の特許請求の範囲の文脈における）不定冠詞および定冠詞および類似の指示語の使用は、本明細書における別段の指摘または文脈による明らかな矛盾のない限り、単数および複数の双方を包含するものと解釈されるべきである。「備える」、「有する」、「含む」および「包含する」という用語は、別段の記載のない限り、制約のない（すなわち、「を含むが、これに限定されない」を意味する）用語として解釈されるべきである。「接続される」という用語は、介在する何かが存在するとしても、部分的または完全に、内部に含まれる、付着される、または互いに接合される、と解釈されるべきである。本明細書における値範囲の列挙は、本明細書において別段の指摘のない限り、単に、その範囲に含まれる別個の各値に個々に言及する簡略法として機能することが意図されるものであって、別個の値は各々、本明細書において個々に列挙されたものであるかのように明細書に組み込まれる。本明細書に記述している全ての方法は、本明細書における別段の指摘または文脈による別段の明らかな矛盾のない限り、任意の適切な順序で実行されることが可能である。本明細書において述べる全ての例または例示的な言語（たとえば、「など」）の使用は、単に、本開示の実施形態をよりよく明らかにすることを意図したものであり、別段の主張のない限り、本開示の範囲を限定するものではない。明細書における文言は、いずれも、請求の範囲に記載されていない要素が本開示の実施に不可欠であることを示すものとして解釈されるべきではない。

【 0 1 8 2 】

[0198]「X、YまたはZのうちの少なくとも1つ」という言い回しなどの選言的言語の意図は、別段で具体的な明記のない限り、一般的に使用される文脈において、あるアイテム、用語、他がX、YまたはZのいずれか、またはこれらの任意の組合せ（たとえば、X、Yおよび/またはZ）であり得ることを提示するものと理解されることにある。したがって、このような選言的言語は、概して、所定の実施形態が少なくとも1つのX、少なくとも1つのY、または少なくとも1つのZを、各々が存在するために必要とすることを含意するものではなく、また含意すべきではない。

【 0 1 8 3 】

[0199]本明細書には、発明者らが知るところの本開示を実施するための最良の形態を含む、本開示の様々な実施形態が記述されている。これまでの説明を読めば、一般的な当業者には、これらの実施形態の変形が明らかになり得る。本発明者らは、当業者がこのよう

10

20

30

40

50

な変形例を適宜使用することを予期するものであり、かつ発明者らの意図は、本開示が本明細書における具体的な記述以外の方法で実施されることにある。したがって、本開示は、適用法で許可される、添付の特許請求の範囲に列挙された主題の全ての変更および同等物を含む。さらに、これまでに述べた要素の、その全ての可能な変形における組合せは、いずれも、本明細書における別段の指摘または文脈による別段の明らかな矛盾のない限り、本開示に包含される。

【図面】

【図 1】

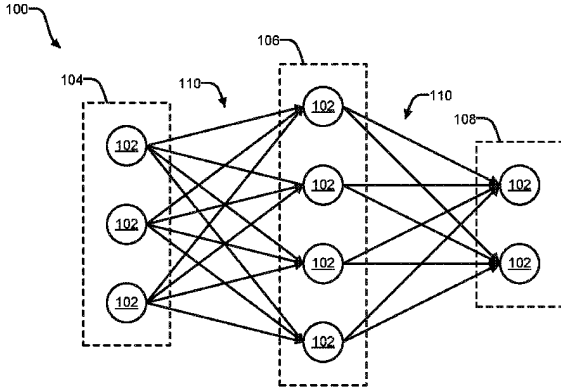


FIG. 1

【図 2 A】

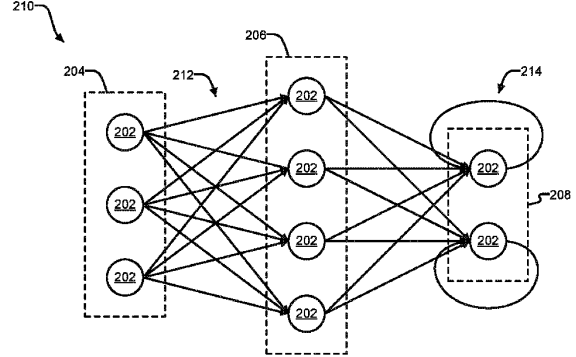


FIG. 2A

10

20

【図 2 B】

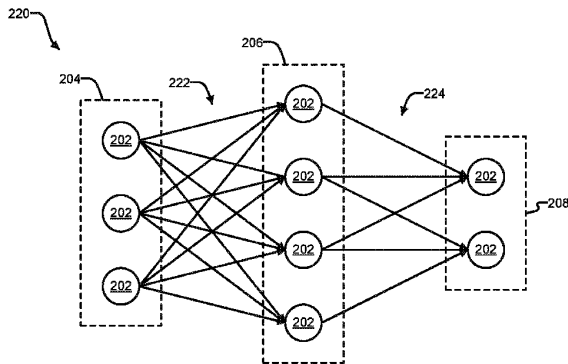


FIG. 2B

【図 3 A】

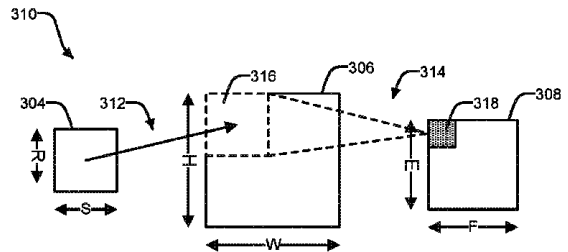


FIG. 3A

30

40

50

【図 3 B】

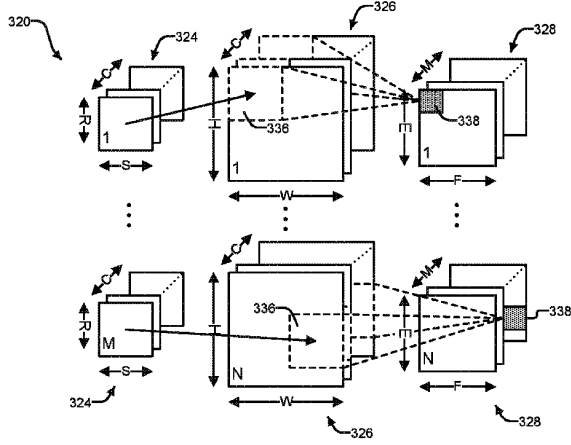


FIG. 3B

【図 4】

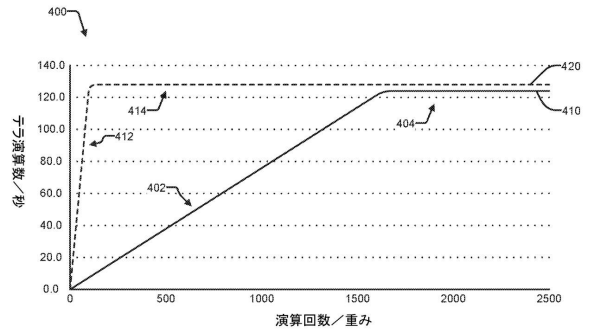


FIG. 4

10

【図 5】

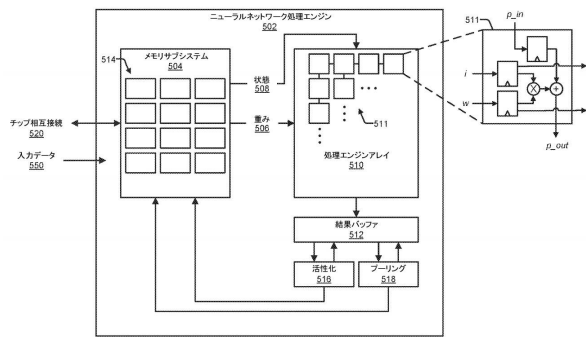


FIG. 5

【図 6 A】

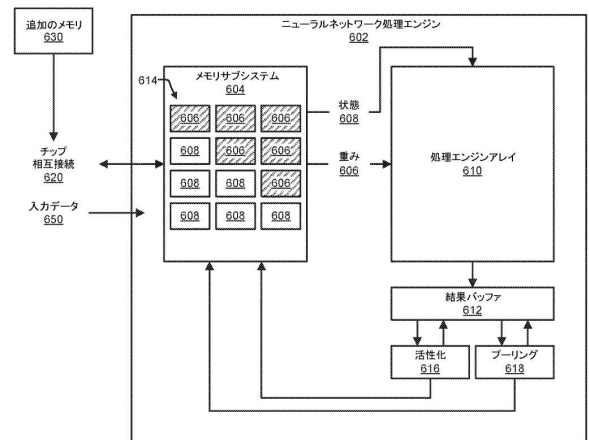


FIG. 6A

20

30

40

50

【図 6 B】

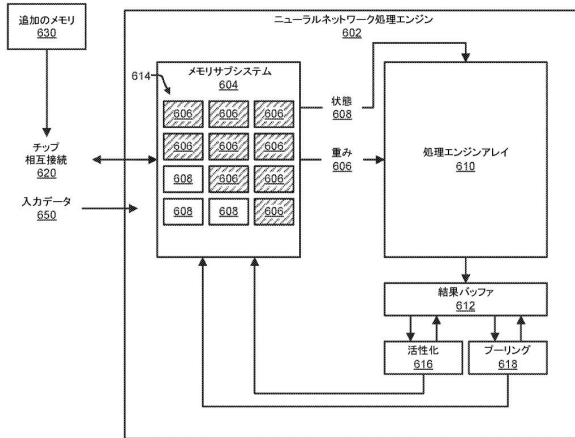


FIG. 6B

【図 7 A】

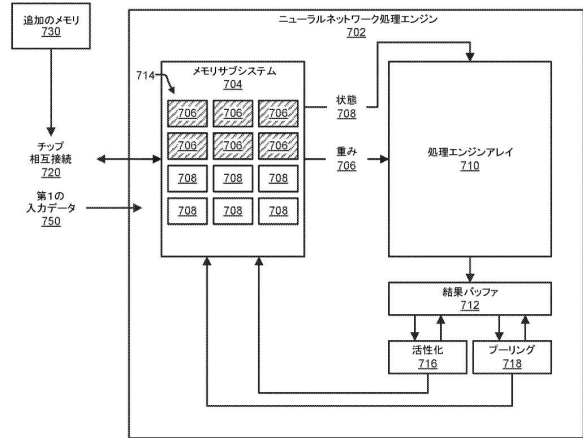


FIG. 7A

【図 7 B】

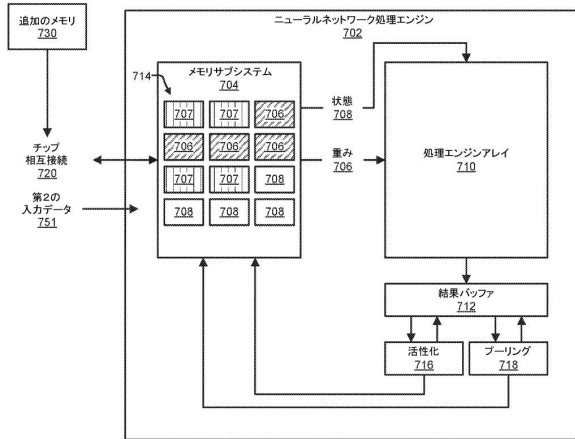


FIG. 7B

【図 7 C】

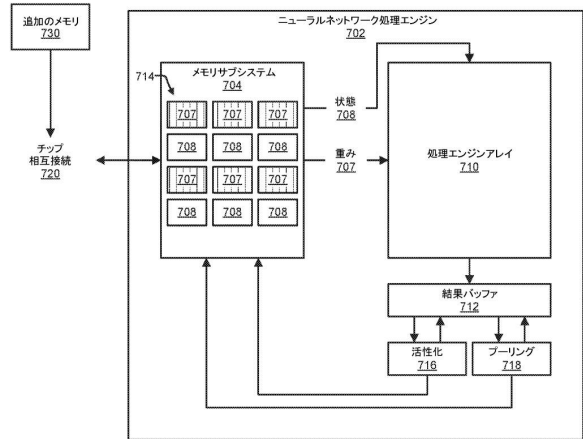


FIG. 7C

10

20

30

40

50

【図 8】

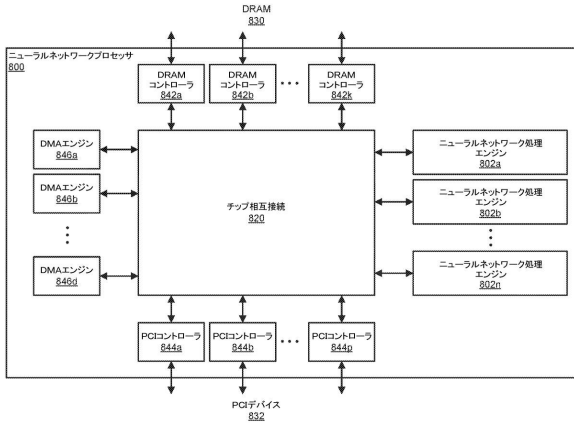


FIG. 8

【図 9】

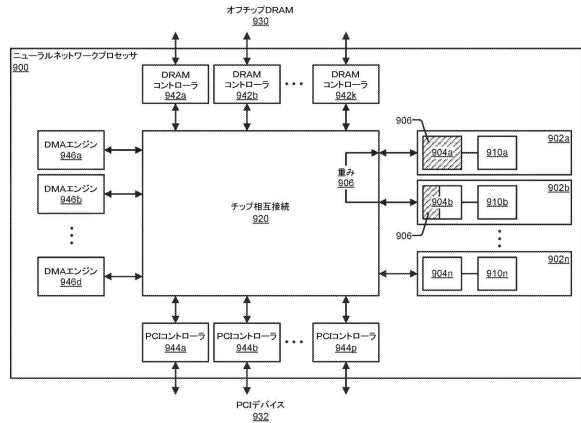


FIG. 9

【図 10】

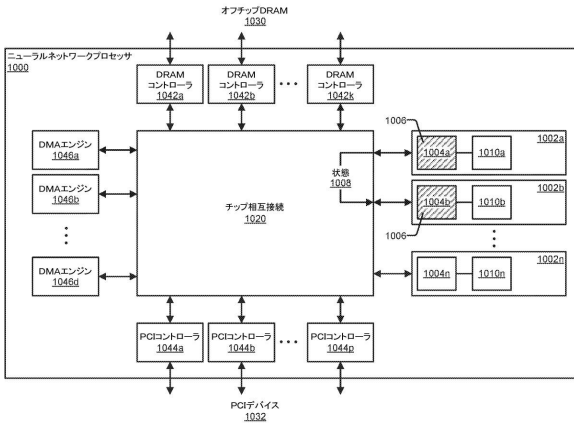


FIG. 10

【図 11】

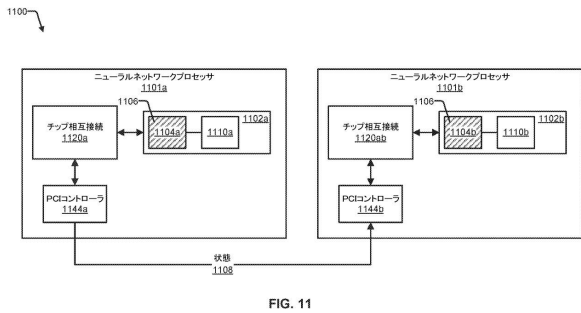


FIG. 11

10

20

30

40

50

【 図 1 2 】

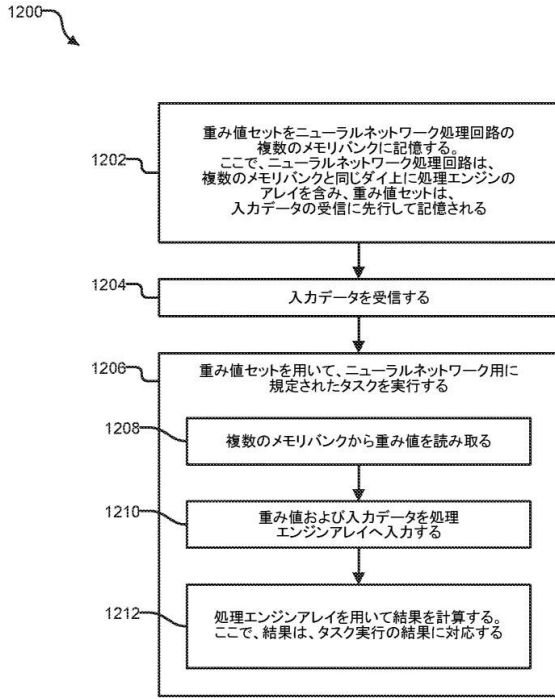


FIG. 12

【 図 1 3 】

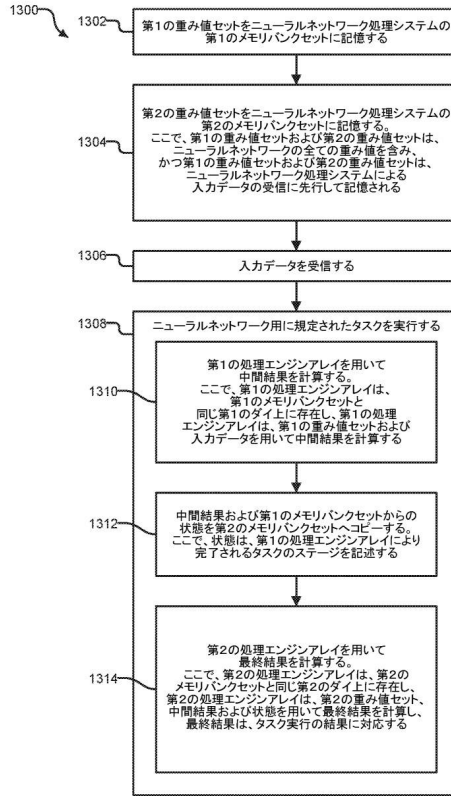


FIG. 13

【 図 1 4 】

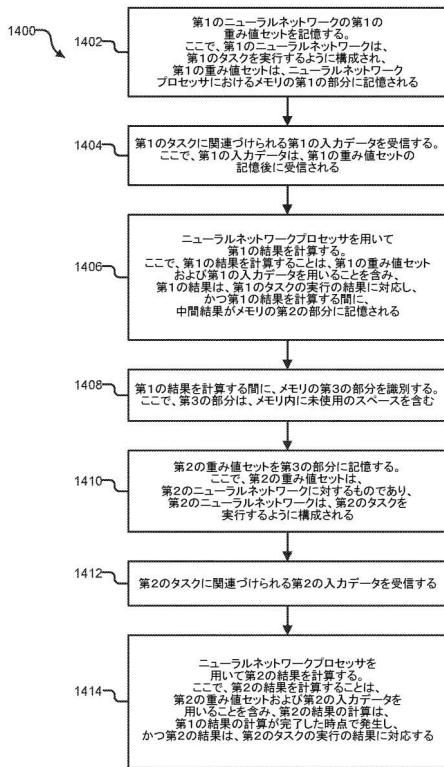


FIG. 14

【 図 1 5 】

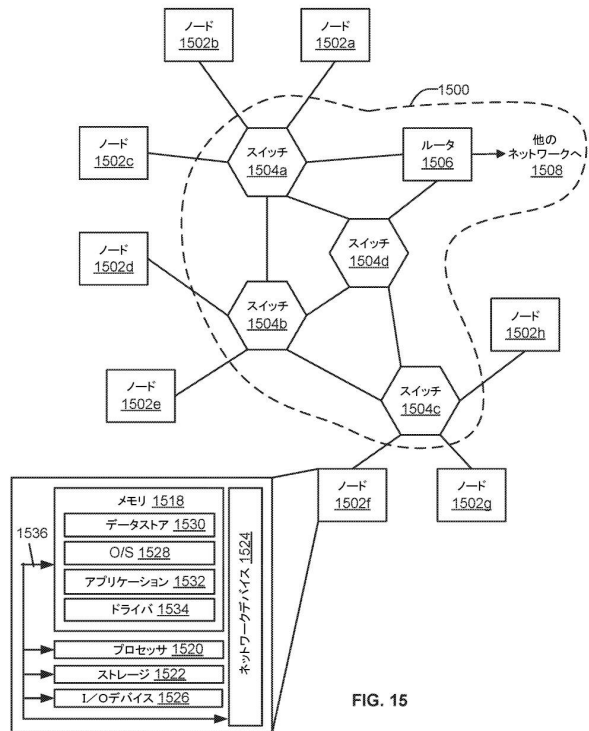


FIG. 15

10

20

30

40

50

フロントページの続き

(33)優先権主張国・地域又は機関

米国(US)

(31)優先権主張番号 15/839,157

(32)優先日 平成29年12月12日(2017.12.12)

(33)優先権主張国・地域又は機関

米国(US)

前置審査

. ボックス 8 1 2 2 6

(72)発明者 ディアマント, ロン

アメリカ合衆国, ワシントン州 9 8 1 0 8 - 1 2 2 6 , シアトル, ピー . オー . ボックス 8 1 2
2 6

(72)発明者 ゼジュダ, ジンドリック

アメリカ合衆国, ワシントン州 9 8 1 0 8 - 1 2 2 6 , シアトル, ピー . オー . ボックス 8 1 2
2 6

(72)発明者 ボルコヴィッチ, ドラジェン

アメリカ合衆国, ワシントン州 9 8 1 0 8 - 1 2 2 6 , シアトル, ピー . オー . ボックス 8 1 2
2 6

審査官 武田 広太郎

(56)参考文献 特表2013-529342(JP, A)

特開2002-117389(JP, A)

特開2013-178294(JP, A)

(58)調査した分野 (Int.Cl., DB名)

G 0 6 N 3 / 0 6 3

G 0 6 F 1 5 / 8 0