

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7232599号
(P7232599)

(45)発行日 令和5年3月3日(2023.3.3)

(24)登録日 令和5年2月22日(2023.2.22)

(51)国際特許分類 F I
G 0 6 N 3/082(2023.01) G 0 6 N 3/082

請求項の数 11 外国語出願 (全11頁)

(21)出願番号	特願2018-165782(P2018-165782)	(73)特許権者	507186492 ピバンテ コーポレーション アメリカ合衆国, カリフォルニア州 9 5 0 0 2, サン ノゼ, ゴールド ストリ ート 2 1 5 0, スイート 2 0 0
(22)出願日	平成30年9月5日(2018.9.5)	(74)代理人	100121083 弁理士 青木 宏義
(65)公開番号	特開2019-49977(P2019-49977A)	(74)代理人	100138391 弁理士 天田 昌行
(43)公開日	平成31年3月28日(2019.3.28)	(74)代理人	100074099 弁理士 大菅 義之
審査請求日	令和3年8月25日(2021.8.25)	(74)代理人	100106851 弁理士 野村 泰久
(31)優先権主張番号	15/699,438	(72)発明者	シン ウォン アメリカ合衆国, カリフォルニア州 9
(32)優先日	平成29年9月8日(2017.9.8)		最終頁に続く
(33)優先権主張国・地域又は機関	米国(US)		

(54)【発明の名称】 畳み込みニューラルネットワークのための刈り込みと再学習法

(57)【特許請求の範囲】

【請求項1】

(a) コンピューティングデバイスによって、ニューラルネットワークを学習し、
(b) 前記コンピューティングデバイスによって、前記ニューラルネットワークを刈り込み、

(c) 前記コンピューティングデバイスによって、前記(b)を実行した後に、再学習の期間中に前記(b)の刈り込みに従って制約を課すことなく前記ニューラルネットワークを前記再学習し、

(d) 前記コンピューティングデバイスによって、前記(c)の再学習が、後続繰り返しにおいて、前の繰り返しでの前記(b)の刈り込みに従う制約なしに実行されるように、前記(b)および前記(c)を1回以上繰り返す、
ことを含む方法であって、

前記(b)が、前記コンピューティングデバイスによって、残りの結合の数が、前記ニューラルネットワーク内の可能な結合の数に関する刈り込み率となるように、除去されていない前記ニューラルネットワークの前記残りの結合よりも低い重みを有する前記ニューラルネットワーク内の結合を除去することによって、前記ニューラルネットワークを刈り込むことを含む、

前記刈り込み率は所定の値であり、

前記方法が、

(e) 前記コンピューティングデバイスによって、前記(d)の各繰り返しの前に前記刈

り込み率を増加させる、
ことを更に含む、
 方法。

【請求項 2】

目標率と等しい前記刈り込み率で前記 (d) が実行されるまで、前記 (d) 及び前記 (e) を実行することを、更に含む、請求項 1 に記載の方法。

【請求項 3】

(f) 前記 (d) および前記 (e) の実行の後に、前記 (d) が前記目標率に等しい前記刈り込み率で実行されるまで、前記コンピューティングデバイスによって、前記 (b) の最後の繰り返しでゼロに重み付けするように制約されて刈り込まれた結合を用いて前記ニューラルネットワークを再学習し、

前記コンピューティングデバイスによって、前記 (a) から前記 (f) の各繰り返しで前記目標率を増加しながら、前記 (a) から前記 (f) を 1 回以上繰り返す、

ことを更に含む、請求項 2 に記載の方法。

【請求項 4】

前記 (e) の後で前記 (f) を実行する前に、前記ニューラルネットワークの精度が安定な閾値条件を達成するまで、前記刈り込み率を前記目標率に固定して前記 (d) を繰り返すことを、更に含む請求項 3 に記載の方法。

【請求項 5】

前記 (d) を実行することは、前記ニューラルネットワークの精度が前記 (a) を実行することによって達成される前記ニューラルネットワークの精度よりも小さく、閾値量の範囲内になるまで、前記 (b) および前記 (c) を実行することを含む、請求項 3 に記載の方法。

【請求項 6】

前記ニューラルネットワークは、畳み込みニューラルネットワーク (CNN) である、請求項 1 に記載の方法。

【請求項 7】

1 つ以上の処理装置と 1 つ以上のメモリデバイスとを含むシステムであって、前記 1 つ以上のメモリデバイスは、

(a) 畳み込みニューラルネットワーク (CNN) であるニューラルネットワークを学習し、

(b) 前記ニューラルネットワークを刈り込み、および、

(c) 前記 (b) を実行した後で、再学習の期間中に前記 (b) の刈り込みに従って前記ニューラルネットワークに制約を課すことなく、前記ニューラルネットワークを前記再学習し、

(d) 前記 (c) の再学習が、次の繰り返しにおいて前の繰り返しにおける前記 (b) の刈り込みに従って制約されることなしに実行されるように、前記 (b) および前記 (c) を 1 回以上繰り返す、

処理を、前記 1 つ以上の処理装置に実行させるのに有効である実行可能コードを格納し、前記実行可能コードは、

残りの結合の数が、前記ニューラルネットワーク内の可能な結合の数に関する刈り込み率となるように、除去されていない前記ニューラルネットワークの前記残りの結合よりも低い重みを有する前記ニューラルネットワーク内の結合を除去することによって、前記 (b) において刈り込みを実行する、処理を、前記 1 つ以上の処理装置に実行させるのに更に有効であり、

前記刈り込み率は所定の値であり、

前記実行可能コードは、

(e) 前記刈り込み率を前記 (d) の各繰り返しの後に増加させる、

処理を、前記 1 つ以上の処理装置に実行させるのに更に有効である、

システム。

10

20

30

40

50

【請求項 8】

前記実行可能コードは、目標率と等しい前記刈り込み率で前記 (d) が実行されるまで、前記 1 つ以上の処理装置に前記 (d) および前記 (e) を実行させるのに更に有効である、請求項 7 に記載のシステム。

【請求項 9】

前記実行可能コードは、

(f) 前記 (d) および前記 (e) の実行の後に、前記 (d) が前記目標率に等しい前記刈り込み率で実行されるまで、前記 (b) の最後の繰り返しでゼロに重み付けするように制約されて刈り込まれた結合を用いて前記ニューラルネットワークを再学習し、

前記ニューラルネットワークの精度が安定な閾値条件を満たすまで、前記 (a) から前記 (f) の各繰り返しで前記目標率を増加しながら、1 回以上前記 (a) から前記 (f) を繰り返す、

処理を、前記 1 つ以上の処理装置に実行させるのに更に有効である、請求項 8 に記載のシステム。

【請求項 10】

前記実行可能コードは、

前記 (e) の後で前記 (f) を実行する前に、前記ニューラルネットワークの精度が安定な閾値条件を達成するまで、前記刈り込み率を前記目標率に固定して前記 (d) を繰り返す、

処理を、前記 1 つ以上の処理装置に実行させるのに更に有効である、請求項 9 に記載のシステム。

【請求項 11】

前記実行可能コードは、

前記ニューラルネットワークの精度が前記 (a) を実行することによって達成される前記ニューラルネットワークの精度よりも小さく、閾値量の範囲内になるまで、前記 (b) および前記 (c) を実行することによって、前記 (d) を実行する、

処理を、前記 1 つ以上の処理装置に実行させるのに更に有効である、請求項 9 に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、畳み込みニューラルネットワークなどのニューラルネットワークの学習を実行するシステムおよび方法に関する。

【背景技術】

【0002】

畳み込みニューラルネットワーク (CNN) を使用する多くのマシンラーニングアプリケーションは、非常に高い計算およびメモリ帯域幅を必要とする。計算負荷を軽減する 1 つの方法は、係数をゼロに刈り込み、係数がゼロである場合に計算をスキップすることである。様々な既存のソフトウェアおよびハードウェアの最適化技術は、ゼロに評価された係数の利点を利用する。1 つの例は、疎畳み込みニューラルネットワーク (Baoyuan Liu、Min Wang¹、Hassan Foroosh¹、Marshall Tappen、および Marianna Pensky) およびディープニューラルネットワーク圧縮および効率的推論エンジン (Song Han CVA group、Stanford University) に記載される疎行列乗算技術であり、これらの文献の両方は、その全体が参照により本明細書に組み込まれる。

【0003】

CNN は、閾値よりも弱いいくつかの結合荷重をゼロにクランプすることで刈り込むことができる。精度は、刈り込みによって大きく影響される。刈り込まれた結合を維持しながら、刈り込まれたモデル精度を復元するための特別な再学習方法が必要である。刈り込みの従来アプローチでは、再学習の間、刈り込まれた結合荷重に無効化マスク (disabling mask) が適用される。刈り込みと再学習の複数回の繰り返しは、一般に、精度の更な

10

20

30

40

50

る向上に役立てることができる。しかし、各繰り返しに対して、再学習する前に無効化マスクが固定される。

【0004】

本明細書に開示されたシステムおよび方法は、高いレベルの精度を達成しながらもゼロ係数の数を増加させるためにCNNを刈り込みする改善されたアプローチを提供する。

【図面の簡単な説明】

【0005】

本発明の利点が容易に理解されるように、添付図面に示される特定の実施形態を参照することによって、上で簡潔に説明された本発明のより詳細な説明が提供されるだろう。これらの図面は、本発明の典型的な実施形態のみを示すものであり、したがって、その範囲を限定するものとみなされるべきではないことを理解されたく、本発明は、添付の図面の使用を通じてさらなる具体性および詳細とともに記載および説明されるだろう。

【図1】本発明の一実施形態に係るニューラルネットワークを刈り込む方法の処理フロー図である。

【図2】本発明の一実施形態に係る図2の方法を繰り返し実行するための方法の処理フロー図である。

【図3】図1の刈り込み方法を示すプロットである。

【図4】本発明の方法を実施するためのコンピューティングデバイスの概略ブロック図である。

【発明を実施するための形態】

【0006】

本明細書の図面に概して記述および図示されているような本発明の構成要素は、多種多様な異なる構成で配置および設計できることが容易に理解されるであろう。したがって、図面に表された本発明の実施形態の以下のより詳細な説明は、請求されるような本発明の範囲を限定するものではなく、単に本発明による現在考えられる実施形態の特定の例を表すに過ぎない。現在記述されている実施形態は、同様の部材が全体を通じて同様の数字で示されている図面を参照することで最もよく理解されるであろう。

【0007】

本発明に係る実施形態は、装置、方法、またはコンピュータプログラム製品として具体化され得る。したがって、本発明は、完全にハードウェアの実施形態、完全にソフトウェアの実施形態（ファームウェア、常駐ソフトウェア、マイクロコードなどを含む）、または、本明細書ではすべて「モジュール」または「システム」として参照され得るソフトウェアおよびハードウェアの態様を組み合わせた実施形態の形態をとってよい。さらに、本発明は、媒体に埋め込まれたコンピュータ使用可能プログラムコードを有する表現の任意の有形の媒体に具体化されたコンピュータプログラム製品の形態を取ってよい。

【0008】

非一時的媒体を含む、1以上のコンピュータ使用可能媒体またはコンピュータ可読媒体の任意の組み合わせが利用され得る。例えば、コンピュータ可読媒体は、ポータブルコンピュータディスク、ハードディスク、ランダムアクセスメモリ（RAM）デバイス、読み出し専用メモリ（ROM）デバイス、消去可能プログラマブル読み出し専用メモリ（EPROMまたはフラッシュメモリ）デバイス、ポータブルコンパクトディスク読み取り専用メモリ（CDROM）、光記憶装置、および磁気記憶装置の1つ以上を含む。選ばれた実施形態では、コンピュータ可読媒体は、命令実行システム、装置、またはデバイスによって使用される、又は関連するプログラムを、含む、格納する、通信する、伝播する、または転送することができる任意の非一時的媒体を含み得る。

【0009】

本発明の動作を実行するためのコンピュータプログラムコードは、Java（登録商標）、Smalltalk、またはC++などのオブジェクト指向プログラミング言語、および"C"プログラミング言語または同様のプログラミング言語などの従来の手続き型プログラミング言語を含む、1つ以上のプログラミング言語の任意の組み合わせで記載され得

10

20

30

40

50

る。プログラムコードは、全体がスタンドアローンのソフトウェアパッケージとしてコンピュータシステム上で、全体がスタンドアローンのハードウェアユニット上で、一部がコンピュータからある程度の距離離れたリモートコンピュータ上で、または全体がリモートコンピュータ或いはサーバ上で実行され得る。後者のシナリオでは、リモートコンピュータは、ローカルエリアネットワーク（LAN）またはワイドエリアネットワーク（WAN）を含む任意のタイプのネットワークを介してコンピュータに接続されてもよく、または、（例えば、インターネットサービスプロバイダを使用してインターネットを介して）外部コンピュータに接続されてもよい。

【0010】

本発明は、本発明の実施形態に係る方法、装置（システム）およびコンピュータプログラム製品のフローチャート図および/またはブロック図を参照して以下に記述される。フローチャート図および/またはブロック図の各ブロック、ならびにフローチャート図および/またはブロック図のブロックの組み合わせは、コンピュータプログラム命令またはコードによって実装できることが理解されよう。これらのコンピュータプログラム命令を、汎用コンピュータ、専用コンピュータ、または他のプログラム可能データ処理装置のプロセッサに提供して、マシンを製造することもでき、こうすることで、その命令が、コンピュータまたは他のプログラム可能データ処理装置のプロセッサを介して実行され、フローチャートおよび/またはブロック図のブロックまたは複数のブロックで指定された機能/動作を実装するための手段を作成する。

【0011】

これらのコンピュータプログラム命令はまた、コンピュータ可読媒体に格納された命令が、フローチャート並びに/或いはブロック図のブロックまたは複数のブロックで指定された機能/動作を実装する命令手段を含む製造品を製造するように、コンピュータまたは他のプログラム可能データ処理装置に特定の仕方でも機能するように指示することができる非一時的コンピュータ可読媒体に格納されてもよい。

【0012】

コンピュータプログラム命令はまた、コンピュータまたは他のプログラム可能装置上で実行される命令が、フローチャート並びに/或いはブロック図のブロックまたは複数のブロックで指定される機能/動作を実装するためのプロセスを提供するように、コンピュータまたは他のプログラム可能データ処理装置にロードされて、コンピュータまたは他のプログラム可能装置上で一連の動作ステップを実行させて、コンピュータ実装プロセスを生成してもよい。

【0013】

図1を参照すると、本明細書で開示されるシステムおよび方法は、高精度を達成しつつも、ゼロ係数の数を増加させるために畳み込みニューラルネットワーク（CNN）の学習に刈り込みを組み込む改良されたアプローチを提供する。

【0014】

CNNの応用例は、2016年8月11日に出願された「ZERO COEFFICIENT SKIPPING CONVOLUTION NEURAL NETWORK ENGINE」と題する米国特許出願第62/373,518号、2017年8月8日に提出された「ZERO COEFFICIENT SKIPPING CONVOLUTION NEURAL NETWORK ENGINE」と題する米国特許出願第15/671,829号、および、2017年8月8日に提出された「ZERO COEFFICIENT SKIPPING CONVOLUTION NEURAL NETWORK ENGINE」と題する米国特許出願第15/671,860号に記載されており、これらの全ては、その全体が参照により本明細書に組み込まれる。

【0015】

図示された方法100は、ゼロ係数の数を増やすためにCNNの結合を刈り込みもするが、一方で、方法100は、CNN、他の種類のニューラルネットワーク、または他の種類のマシンラーニングモデルの学習のために使用されてよい。以下の記載では、CNNを参照する。しかし、開示された方法は、任意の種類のニューラルネットワークでの使用に

10

20

30

40

50

適用可能であると理解されるべきである。

【0016】

方法100は、CNNモデルのフル学習(full training)を実行する(102)ことを含み得る。フル学習は、当技術分野で知られているCNNのための任意の学習方法を含み得る。当技術分野で知られているように、学習(102)は、CNNモデルの結合の重みを設定するために、学習データを入力として取り、学習アルゴリズムを適用することを含み得る。フル学習の期間中に、CNNモデルの任意の結合の重みは、学習データに基づいて所与の入力データに対して正確な出力を生成するように変更され得る。当技術分野で知られているように、学習データは、エントリを含み、各エントリは、1つ以上の入力と1つ以上の所望の出力とを含む。学習アルゴリズムは、1つ以上の入力を受信したときに、モデルが1つ以上の所望の出力を与えるように、CNNの結合の重みを選択する。ステップ102は、例えば、前の繰り返しに対して閾値の割合よりも小さく変化する(例えば、0.01~2パーセントの間の値で、例えば1パーセントが許容可能な値である)といったように、CNNモデルの精度が安定するまで学習アルゴリズムを繰り返すことによってなど、1以上のフル再学習の繰り返しを含み得る。

10

【0017】

方法100は、CNNモデルを刈り込む(104)ことを含み得る。これは、ゼロに近い結合の重みをゼロにクランプすることを含み得る。刈り込むステップ104は、所定の刈り込み率で実行され得る。刈り込み率は、ゼロに設定する結合の割合である。刈り込み率Pは、初期値(例えば50%から60%の間の値)に設定され得る。したがって、刈り込みステップ104は、ステップ102の後で最も低い重みを有するPパーセントの結合について重みをゼロにクランプすることを含み得る。言い換えると、保持される1-Pパーセントの結合は、ゼロにクランプされるPパーセントの結合のすべての刈り込み前の重みよりも、高い重みをそれぞれが有することになる。

20

【0018】

方法100は、刈り込みの後にCNNモデルのフル再学習を実行する(106)ことを含み得る。フル再学習は、刈り込みステップに従ってCNNモデルを制約しない。しかしながら、ステップ106の開始時のCNNモデルは、ステップ104由来の刈り込まれたモデルである。これにより、刈り込みされた状態にとどまる制約なしに、刈り込みされたモデルに向かってバイアスが与えられる。したがって、刈り込みされた結合の重みは、再学習(106)の期間中に増加することが許容されることになる。

30

【0019】

ステップ106は、CNNモデルの精度が復元されるべきと分かる(108)まで、すなわち、ステップ102に続くCNNモデルの精度の或る閾値パーセンテージ(1)内の精度を達成された状態になるまで、繰り返され得る。

【0020】

方法100は、刈り込み率Pが目標の刈り込み率 P_T に等しいかどうかを評価することを含み得る。目標の刈り込み率の例は、50%から90%の範囲であってよい。しかしながら、何が妥当な目標の刈り込み率であるかは、用途に大きく依存する。そうでない場合、刈り込み率Pは、 $P + \Delta P$ に増加されるか、又は、 P_T に設定され得(いずれか小さい方)、ここで ΔP はインクリメント量であるとする。インクリメント ΔP は、定数であってよく、又は、非線形関数に従ってステップ112の各繰り返しごとに変化してもよい。方法100は、その後、ステップ112で設定された値に等しいPで、CNNモデルを再び刈り込むことによって、ステップ104に進み得る。次いで、方法は、上記のステップ104から進み得る。

40

【0021】

CNNモデルが、いったん目標刈り込み率 P_T に刈り込み104されており、閾値精度内に学習106されると、方法は、ステップ114、116、および118に進み得、ここで、CNNモデルは、その精度が安定した状態にあると分かる(116)まで、1回以上、再学習(114)され、および刈り込まれる(118)。安定性は、ステップ114

50

の1回の繰り返し後の精度が、前の繰り返しから閾値差（例えば0.01%～2%の間の値、例えば1%の値が許容されることが多い）内である場合に、達成されたと判定され得る。

【0022】

方法100は、その後、CNNモデルのマスクされた再学習を実行する（120）ことを含み得る。マスクされた再学習では、非ゼロ結合のみが値を変更することを許可される。換言すれば、ステップ120では、選択された結合、すなわち、最後の刈り込み（ステップ104または118）の期間中に刈り込まれなかった結合のみが変更可能である。

【0023】

マスクされた再学習120は、当技術分野で知られている任意のマスクされたCNN再学習アプローチによる無効化マスクを用いた学習を含み得る。ステップ120は、CNNモデルの精度が安定な状態にあると分かる（122）まで、1回以上繰り返され得る。安定性は、ステップ116と同じ方法で決定され得る。

10

【0024】

その後、方法100は終了し得、それからCNNモデルは、生産データを処理するために使用され得る。

【0025】

図2を参照すると、いくつかの実施形態では、方法100（本明細書ではPFP-S：刈り込み（prune）、フル再学習（full retrain）、刈り込み（prune）、および選択された結合の再学習（selected connection retraining）として参照される）は、図示された方法200の内容の範囲内で繰り返され得る。

20

【0026】

方法200は、CNNモデルが安定化された状態にあることが分かる（204）まで（ステップ116参照）、PFP-S方法100を実行する（202）ことを含み得る。PFP-S方法100が繰り返される前に、刈り込み目標 P_T は、例えば、固定のまたは可変の減少量分など、増加され得る。刈り込み率 P の初期値は、PFP-S方法100の以前の実行由来の P_T の値に、またはそれよりも小さい値に設定され得る。

【0027】

その後、方法200の繰り返し由来のCNNモデルが使用され得る。これは、方法200の最後の繰り返しによって生成されるCNNモデルであり得る。刈り込み目標 P_T の増加が、許容できない量でCNNモデルの精度を低下させる場合、CNNモデルの最新バージョンではなく、CNNモデルの前のバージョンが、ステップ208で復元され得、生産データを処理するために使用され得る。

30

【0028】

図3は、方法100がどのように動作し得るかの例を示す。横軸は刈り込み率 P を表す。縦軸はCNNモデルの精度を表す。点300は、ピーク精度（“ピーク%”）を有する方法100の初期学習ステップ102に続くCNNモデルの状態を表す。図3に示すように、CNNモデルは刈り込みされ、その結果、刈り込み率が増加し、また、割合（Peak - 0）分精度が減少する。CNNモデルは、その後、各学習の後で精度がPeak - 1に達するように、繰り返し学習され（106）、および刈り込まれる（104）。学習ステップ106の期間中、刈り込みされた結合が非ゼロ値に変更されることが許されるので、刈り込み率は減少することに留意されたい。

40

【0029】

点302は、 P が刈り込み目標 P_T に到達した後のCNNモデルの状態を示す。その後、ステップ118の最後の繰り返しの後の点304によって示される状態にCNNモデルになる点において、CNNモデルの精度が安定すると分かる（116）まで、再学習（114）及び刈り込み（118）が実行され得る。その後、CNNモデルはマスクを用いて再学習され（120）、刈り込み目標 P_T における刈り込み率を維持しながらモデルの精度を向上させる。その後、CNNモデルは点306になる。その後、方法100は、方法200に従って、異なる刈り込み目標 P_T で繰り返されてもよく、または、CNNモデル

50

が生産データを処理するために使用されてもよい。

【0030】

図4は、例示的なコンピューティングデバイス400を示すブロック図である。コンピューティングデバイス400は、本明細書で論じられるものなどの様々な手順を実行するために使用され得る。コンピューティングデバイス400は、デスクトップコンピュータ、ノートブックコンピュータ、サーバコンピュータ、ハンドヘルドコンピュータ、タブレットコンピュータなどの多種多様なコンピューティングデバイスのいずれかとすることができる。

【0031】

コンピューティングデバイス400は、1つ以上のプロセッサ402、1つ以上のメモリデバイス404、1つ以上のインタフェース406、1つ以上の大容量記憶デバイス408、1つ以上の入力/出力(I/O)デバイス410、および表示装置430を含み、これらのすべてがバス412に接続される。1つの(または複数の)プロセッサ402は、1つ以上のプロセッサまたはコントローラを含み、これらは1つの(または複数の)メモリデバイス404および/または1つの(または複数の)大容量記憶デバイス408に記憶された命令を実行する。1つの(または複数の)プロセッサ402はまた、キャッシュメモリなどの様々な種類のコンピュータ可読媒体を含み得る。

10

【0032】

1つの(または複数の)メモリデバイス404は、揮発性メモリ(例えば、ランダムアクセスメモリ(RAM)414)および/または不揮発性メモリ(例えば、読み出し専用メモリ(ROM)416)などの様々なコンピュータ可読媒体を含む。メモリデバイス404はまた、フラッシュメモリなどの書き換え可能なROMを含み得る。

20

【0033】

1つの(または複数の)大容量記憶デバイス408は、磁気テープ、磁気ディスク、光ディスク、ソリッドステートメモリ(例えば、フラッシュメモリ)などの様々なコンピュータ可読媒体を含む。図4に示すように、具体的な大容量記憶装置はハードディスクドライブ424である。1つの(または複数の)大容量記憶デバイス408には様々なドライブが含まれ、様々なコンピュータ可読媒体からの読み取りおよび/または書き込みを可能にしてもよい。1つの(または複数の)大容量記憶デバイス408は、リムーバブルメディア426および/または非リムーバブルメディアを含む。

30

【0034】

1つの(または複数の)I/Oデバイス410は、データおよび/または他の情報がコンピューティングデバイス400に入力されるか、またはコンピューティングデバイス400から取り出されることを可能にする様々なデバイスを含む。例示的な1つの(または複数の)I/Oデバイス410は、カーソル制御デバイス、キーボード、キーパッド、マイク、モニタまたは他の表示装置、スピーカ、プリンタ、ネットワークインタフェースカード、モデム、レンズ、およびCCDまたは他の画像キャプチャ装置などを含む。

【0035】

表示装置430は、コンピューティングデバイス400の1人または複数のユーザに情報を表示することができる任意の種類デバイスを含む。表示装置430の例は、モニタ、表示端末、およびビデオ投影装置などを含む。

40

【0036】

グラフィックスプロセッシングユニット(GPU)432は、1つの(または複数の)プロセッサ402および/または表示装置430に接続されてよい。GPUは、コンピュータ生成画像をレンダリングし、他のグラフィカル処理を実行するように動作可能であってよい。GPUは、1つの(または複数の)プロセッサ402などの汎用プロセッサの機能の一部または全部を含み得る。GPUは、グラフィックス処理に特有の追加の機能を含んでもよい。GPUは、座標変換、シェーディング、テクスチャリング、ラスターライズ、およびコンピュータ生成画像のレンダリングに役立つ他の機能に関連するハードコードされたグラフィックス機能、および/またはハードワイヤードのグラフィックス機能を含み

50

得る。

【0037】

1つの（または複数の）インタフェース406は、コンピューティングデバイス400が他のシステム、デバイス、またはコンピューティング環境と情報をやり取りすることを可能にする様々なインタフェースを含む。例示的な1つの（または複数の）インタフェース406は、ローカルエリアネットワーク（LAN）、ワイドエリアネットワーク（WAN）、無線ネットワーク、およびインターネットへのインタフェースなどの任意の数の様々なネットワークインタフェース420を含む。他の1つ（または複数の）インタフェースは、ユーザインタフェース418と周辺デバイスインタフェース422とを含む。1つの（または複数の）インタフェース406はまた、1つ以上のユーザインタフェース素子418を含んでもよい。1つの（または複数の）インタフェース406はまた、例えば、プリンタ、ポインティングデバイス（マウス、トラックパッドなど）、およびキーボードなどのインタフェースのような、1つ以上の周辺インタフェースを含んでもよい。

10

【0038】

バス412は、1つの（または複数の）プロセッサ402、1つの（または複数の）メモリデバイス404、1つの（または複数の）インタフェース406、1つの（または複数の）大容量記憶デバイス408、および1つの（または複数の）I/Oデバイス410が相互に通信することを可能にし、バス412に接続された他のデバイスまたは構成要素との通信も可能にする。バス412は、システムバス、PCIバス、IEEE1394バス、およびUSBバスなどのいくつかの種類バス構造の1つ以上を表す。

20

【0039】

説明のために、プログラムおよび他の実行可能なプログラム構成要素は、本明細書では個別のブロックとして示されているが、そのようなプログラムおよび構成要素は、様々な時間で、コンピューティングデバイス400の異なる記憶構成要素内に存在し得、1つの（または複数の）プロセッサ402によって実行され得ることが理解される。代わりに、本明細書に記載のシステムおよび手順は、ハードウェアで、またはハードウェア、ソフトウェア、および/またはファームウェアの組み合わせで実装されることができる。例えば、1つ以上の特定用途向け集積回路（ASIC）は、本明細書に記載されるシステムおよび手順の1つ以上を実行するようにプログラムすることができる。

【0040】

本発明は、その趣旨または本質的な特徴から逸脱することなく、他の特定の形態で実施することができる。記載された実施形態は、すべての点において、例示的なものであり、限定的なものではないとみなされるべきである。したがって、本発明の範囲は、前述の説明によってではなく、添付の特許請求の範囲によって示される。特許請求の範囲と均等の意味および範囲内に入るすべての変更は、その範囲内に含まれるべきである。

30

【0041】

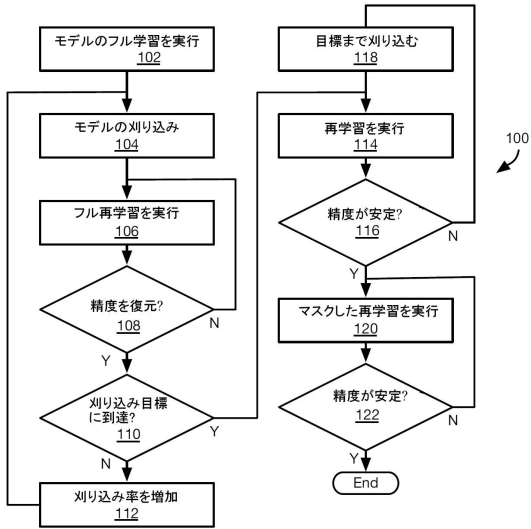
請求されているものは以下に列挙される。

40

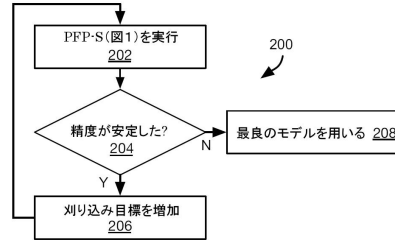
50

【図面】

【図1】



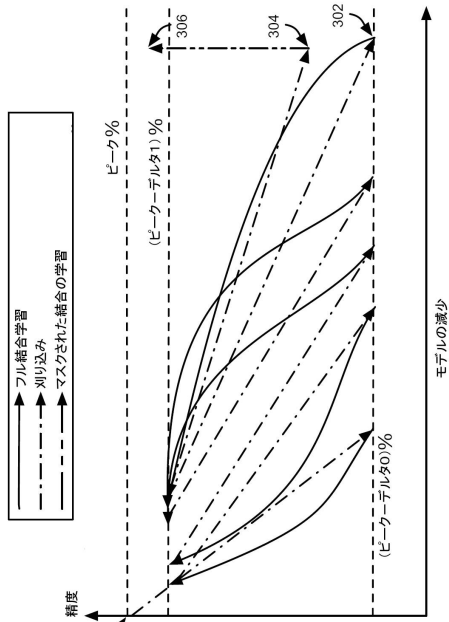
【図2】



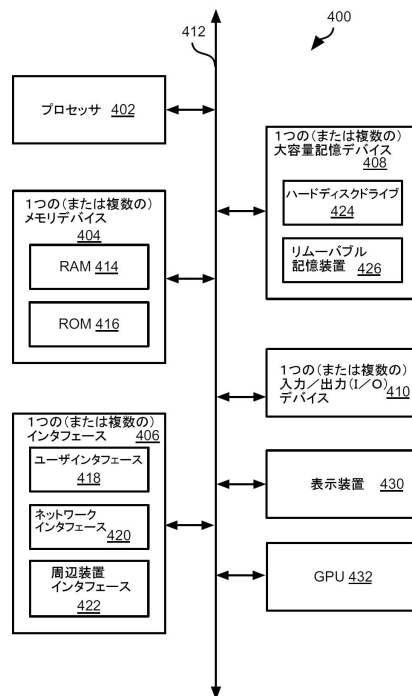
10

20

【図3】



【図4】



30

40

50

フロントページの続き

- 5 0 0 2 , サ ン ノ ゼ , ゴ ー ル ド ス ト リ ー ト 2 1 5 0 , ス イ ー ト 2 0 0
(72)発明者 シ ャ ン ハ ン リ ン
ア メ リ カ 合 衆 国 , カ リ フ ォ ル ニ ア 州 9 5 0 0 2 , サ ン ノ ゼ , ゴ ー ル ド ス ト リ ー ト 2 1 5 0 ,
ス イ ー ト 2 0 0
審 査 官 石 川 亮
(56)参考文献 特 開 2 0 1 8 - 1 5 2 0 0 0 (J P , A)
米 国 特 許 出 願 公 開 第 2 0 1 6 / 0 3 0 7 0 9 8 (U S , A 1)
特 開 2 0 1 1 - 0 5 4 2 0 0 (J P , A)
HU Hengyuan et al. , Network Trimming: A Data-Driven Neuron Pruning Approach towards
Efficient Deep Architectures , arXiv.org [online] , 2016年 , pp.1-9 , [2022年7月28日検索],
イ ン タ ー ネ ッ ト URL: <https://arxiv.org/abs/1607.03250>
(58)調査した分野 (Int.Cl. , D B 名)
G 0 6 N 3 / 0 2 - 3 / 1 0