

【特許請求の範囲】**【請求項 1】**

演算を実行する演算部と、
前記演算の演算結果の統計情報を基に特定小数点位置を決定する小数点位置決定部と、
前記演算部による候補小数点位置を用いた第 1 演算の結果である第 1 演算結果を取得し、
前記第 1 演算結果の統計情報を基に前記小数点位置決定部により決定された前記特定小数点位置を取得し、前記候補小数点位置及び前記特定小数点位置を基に、前記第 1 演算結果又は前記演算部による前記特定小数点位置を用いた第 2 演算の結果である第 2 演算結果を最終演算結果として取得する管理部と
を備えたことを特徴とする演算処理装置。

10

【請求項 2】

前記管理部は、前記特定小数点位置に一致する前記候補小数点位置が存在する場合、前記特定小数点位置に一致する前記候補小数点位置を用いた前記第 1 演算の結果である前記第 1 演算結果を前記最終演算結果とすることを特徴とする請求項 1 に記載の演算処理装置。

【請求項 3】

前記管理部は、前記特定小数点位置に一致する前記候補小数点位置が存在しない場合、前記特定小数点位置を用いた前記第 2 演算の結果である前記第 2 演算結果を取得し、前記第 2 演算結果を前記最終演算結果とすることを特徴とする請求項 1 又は 2 に記載の演算処理装置。

20

【請求項 4】

前記管理部は、予め決められた前記候補小数点位置を保持することを特徴とする請求項 1 ~ 3 のいずれか一つに記載の演算処理装置。

【請求項 5】

前記管理部は、それぞれに対応する所定演算が決められた連続する複数の処理層毎に、対応する前記所定演算を前記第 1 演算及び前記第 2 演算として前記演算部に実行させて前記最終演算結果の取得を繰り返し、且つ、特定処理層における前記最終演算結果の取得に用いた前記候補小数点位置又は前記特定小数点位置を基に前記特定処理層の次の前記処理層における前記候補小数点位置を生成することを特徴とする請求項 1 ~ 3 のいずれか一つに記載の演算処理装置。

30

【請求項 6】

前記管理部は、前記第 1 演算及び前記第 2 演算として前記演算部に所定積和演算を実行させる場合、前記所定積和演算を基に前記候補小数点位置を生成することを特徴とする請求項 1 ~ 5 のいずれか一つに記載の演算処理装置。

【請求項 7】

演算回路を有する演算処理装置の制御方法であって、
前記演算回路による候補小数点位置を用いた第 1 演算の結果である第 1 演算結果を取得し、
前記第 1 演算結果の統計情報を基に特定小数点位置を決定し、
前記候補小数点位置及び前記特定小数点位置を基に、前記第 1 演算結果又は前記演算回路による前記特定小数点位置を用いた第 2 演算の結果である第 2 演算結果を最終演算結果として取得することを特徴とする演算処理装置の制御方法。

40

【請求項 8】

演算回路による候補小数点位置を用いた第 1 演算の結果である第 1 演算結果を算出させ、
前記第 1 演算結果の統計情報を基に特定小数点位置を決定し、
前記候補小数点位置及び前記特定小数点位置を基に、前記第 1 演算結果又は前記演算回路による前記特定小数点位置を用いた第 2 演算の結果である第 2 演算結果を最終演算結果として取得する

50

処理をコンピュータに実行させることを特徴とする演算処理プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、演算処理装置、演算処理装置の制御方法及び演算処理プログラムに関する。

【背景技術】

【0002】

今日、深層学習（ディープラーニング）へのニーズが高まっている。深層学習においては、乗算、積和演算、ベクトル乗算を含む様々な演算が実行される。ところで、深層学習では、個々の演算精度への要求は、他のコンピュータ処理ほど厳密ではない。例えば、従来の信号処理等では、プログラマは極力桁あふれを発生させないようにコンピュータプログラムを開発する。一方、深層学習では、大きな値がある程度飽和することは許容される。これは、深層学習では、複数の入力データを畳み込み演算するときの係数（重み）の調整が主な処理となり、入力データのうち極端なデータは重視されないことが多いためである。また、大量のデータを繰り返し用いて係数を調整するため、一度飽和された値も、学習の進行に合わせて桁調整を行なうことで、飽和されずに係数調整に反映できるようになるためである。

10

【0003】

そこで、このような深層学習の特性を考慮し、深層学習用の演算処理装置のチップ面積の削減及び電力性能の向上等を図るため、浮動小数点数を用いずに固定小数点数による演算を用いることが考えられる。これは、浮動小数点数演算よりも固定小数点数を用いた演算の方が回路構成を簡素にできるためである。

20

【0004】

また、近年、ディープラーニング用の専用アクセラレータの開発が盛んになっている。そこで、専用アクセラレータにおける演算の面積効率を上げるためにも、固定小数点数による演算を用いることが好ましい。例えば、演算ビット数を例えば32ビットの浮動小数点数を8ビットの固定小数点数に減らして面積あたりの演算性能を向上させたハードウェアが開発されている。32ビットの浮動小数点数を8ビットの固定小数点数に減らすことで、単純に面積当たり4倍の性能が得られる。このように、十分な精度を有する実数を少ないビット数で表現する処理は量子化と呼ばれる。

30

【0005】

ただし、固定小数点数は、ダイナミックレンジが狭いため浮動小数点数より演算精度が劣化する場合がある。そのため、深層学習においても、小さな値を表現する精度、すなわち有効桁数について配慮が求められる。そこで、演算結果のビット位置の統計情報を用いて固定小数点の有効桁数を決定し小数点位置を最適化する従来技術が存在する。

【0006】

従来技術では、前のイテレーション（iteration）の統計情報を用いて次のイテレーションの小数点位置が決定され、決定された小数点位置を用いて次のイテレーションの演算が実行される。イテレーションは、ミニバッチとも呼ばれる。

【0007】

また、統計情報を用いて固定小数点における小数点位置を決定する技術として、最下位ビット位置から最上位ビット位置までの範囲を示す情報及び符号ビット位置から最下位ビットのビット位置まで範囲を示す情報を用いて小数点位置を決定する従来技術がある。また、固定小数点演算を行う技術として、指定された小数点位置を表すデータに基づいて出力された演算結果に対して丸め処理及び飽和処理を実行しつつ固定小数点演算を行う従来技術がある。

40

【先行技術文献】

【特許文献】

【0008】

【特許文献1】特開2018-124681号公報

50

【特許文献2】特開2019-74951号公報

【特許文献3】特開2009-271598号公報

【発明の概要】

【発明が解決しようとする課題】

【0009】

しかしながら、最近の深層学習のフレームワーク、特にpyTorchやchainerrで、Define by Runと呼ばれる処理方式が導入される機会が増加した。以下では、Define by RunをDbRと省略して表記する。DbRでは、ニューラルネットの構造となる計算グラフの構築が、深層学習の処理を実行しながら行われる。そして、DbRでは、最短で学習のイテレーション毎に計算グラフが変わる。そのため、過去に推定した小数点位置を記憶することが困難である。また、計算グラフが変わるとい

10

【0010】

DbRにより深層学習を行う場合、仮に前回のイテレーションにおける統計情報を用いるとしても、前回自体が存在しないあるいは前回の統計情報が実際にはいくつものイテレーションも前の情報となる場合がある。このため、DbRにより深層学習を行う場合、過去の統計情報を用いると学習が破たんするおそれがあり、過去の統計情報を用いて小数点位置を決定することは困難である。

20

【0011】

そこで、現在のレイヤの演算を実行し、その演算結果の統計情報から小数点位置を決定し、求めた小数点位置を用いて再度演算を実行する方法が考えられる。しかし、この方法では、確実に2度同じ演算を行うことになり学習時間が長くなってしまいう問題がある。

【0012】

また、最下位ビット位置から最上位ビット位置までの範囲を示す情報及び符号ビット位置から最下位ビットのビット位置まで範囲を示す情報を用いて小数点位置を決定する技術でも、過去の統計情報を用いることからDbRを用いた深層学習への適用は困難である。また、指定された小数点位置を表すデータに基づいて出力された演算結果に対して丸め処理及び飽和処理を行う従来技術では、小数点位置の決め方が考慮されておらず、DbRにより深層学習を行うことは困難である。

30

【0013】

開示の技術は、上記に鑑みてなされたものであって、Define by Runにより深層学習を行う際の固定小数点を用いた学習における学習精度を向上させつつ学習時間を短縮する演算処理装置、演算処理装置の制御方法及び演算処理プログラムを提供することを目的とする。

【課題を解決するための手段】

【0014】

本願の開示する演算処理装置、演算処理装置の制御方法及び演算処理プログラムの一つの態様において、演算部は、演算を実行する。小数点位置決定部は、前記演算の演算結果の統計情報を基に特定小数点位置を決定する。管理部は、前記演算部による候補小数点位置を用いた第1演算の結果である第1演算結果を取得し、前記第1演算結果の統計情報を基に前記小数点位置決定部により決定された前記特定小数点位置を取得し、前記候補小数点位置及び前記特定小数点位置を基に、前記第1演算結果又は前記演算部による前記特定小数点位置を用いた第2演算の結果である第2演算結果を最終演算結果として取得する。

40

【発明の効果】

【0015】

1つの側面では、本発明は、Define by Runにより深層学習を行う際の固

50

定小数点を用いた学習における学習精度を向上させつつ学習時間を短縮することができる。

【図面の簡単な説明】

【0016】

【図1】図1は、サーバの概略を表す構成図である。

【図2】図2は、ニューラルネットワークにおける深層学習の一例の図である。

【図3】図3は、DbRを説明するための図である。

【図4】図4は、演算用回路のブロック図である。

【図5】図5は、制御部の詳細を表すブロック図である。

【図6】図6は、実施例1に係る演算用回路による深層学習の処理のフローチャートである。

10

【図7】図7は、実施例2に係る演算用回路による深層学習の処理のフローチャートである。

【図8A】図8Aは、実施例3に係る演算用回路による深層学習の処理のフローチャートである。

【図8B】図8Bは、実施例3に係る演算用回路による深層学習の処理のフローチャートである。

【発明を実施するための形態】

【0017】

以下に、本願の開示する演算処理装置、演算処理装置の制御方法及び演算処理プログラムの実施例を図面に基づいて詳細に説明する。なお、以下の実施例により本願の開示する演算処理装置、演算処理装置の制御方法及び演算処理プログラムが限定されるものではない。

20

【実施例1】

【0018】

図1は、サーバの概略を表す構成図である。サーバ1は、深層学習を実行する。サーバ1は、CPU (Central Processing Unit) 2、メモリ3、及び演算用回路4を有する。CPU 2、メモリ3及び演算用回路4は、相互にPCIe (Peripheral Component Interconnect Express) バス5で接続される。

【0019】

30

CPU 2は、メモリ3に格納されたプログラムを実行しサーバ1としての各種機能を実現する。例えば、CPU 2は、PCIeバス5経由で制御信号を送信して、演算用回路4が有する制御コアを起動する。また、CPU 2は、演算で用いるデータ及び実行する演算命令を演算用回路4へ出力して、演算用回路4に演算を行わせる。

【0020】

演算用回路4は、深層学習におけるレイヤ毎の演算を実行する回路である。ここで、図2を参照して、ニューラルネットワークにおける深層学習の一例について説明する。図2は、ニューラルネットワークにおける深層学習の一例の図である。ニューラルネットワークは、例えば、画像を認識して識別するためのフォワード方向の処理と、フォワード方向の処理で使用するパラメータを決定するバックワード方向の処理を実行する。図2における上部の矢印の紙面に向かって右に進む方向がフォワード方向であり、紙面に向かって左に進む方向がバックワード方向である。

40

【0021】

図2のニューラルネットワークは、入力画像に対して、畳み込み層 (Convolution Layer) の処理及びプーリング層 (Pooling Layer) の処理を実行し、画像の特徴を抽出して画像を識別する。図2の紙面中央に記載した処理はフォワード方向の処理を表す。

【0022】

図2では、フォワード方向の処理における特徴抽出部において、入力画像に対して畳み込み層の処理及びプーリング層の処理が実行され特徴マップが生成される。その後、識別部において、特徴マップに対して全結合が行われ最終層から識別結果が出力される。畳み

50

込み層の処理は、畳み込み演算とも呼ばれる。また、プーリング層の処理は、プーリング演算とも呼ばれる。その後、識別結果は、正解データと比較され、比較結果である差分値が得られる。次に、バックワード方向の処理として、差分値からフォワード方向の畳み込み層及び全結合層における各層でのエラー及び各層での次の重みを計算する学習処理が行われる。

【 0 0 2 3 】

深層学習は、ミニバッチと呼ばれる処理の単位に区切られて実行される。ミニバッチとは、学習の対象となる入力データの集合を所定個の組に分割した複数個のデータの組み合わせである。図 2 では、N 個の画像で 1 つのミニバッチである。そして、ミニバッチ毎のフォワード方向の処理及びバックワード方向の処理の一連の処理をまとめた単位を 1 イテレーションという。

10

【 0 0 2 4 】

さらに、本実施例では、DbRにより深層学習を実行する。図 3 は、DbRを説明するための図である。Define and Runにより深層学習を実行する場合、計算グラフが固定されるため、レイヤ 5 1、5 3、5 4、5 5 と経由するか、レイヤ 5 1、5 2、5 5 と経由するかは演算を実行する前に決定される。これに対して、RbDにより演算を行う場合、演算を実行するまで計算グラフが決定されずに、レイヤ 5 1 からレイヤ 5 2 に進むかレイヤ 5 3 に進むかは確率的に決まり、演算時に計算グラフが動的に変化する。そのため、演算用回路 4 は、演算に用いる小数点位置を前もって決定することが困難である。そこで、以下の手順で演算用回路 4 は演算を実行する。

20

【 0 0 2 5 】

図 2 に戻って説明を続ける。演算用回路 4 は、深層学習中の所定数のミニバッチ毎に、各層の演算を行うとともに各層の各変数の統計情報を取得して蓄積し、深層学習に用いる変数の固定小数点位置を自動調整する。次に、演算用回路 4 の詳細について説明する。

【 0 0 2 6 】

図 4 は、演算用回路のブロック図である。図 4 に示すように、演算用回路 4 は、プロセッサ 4 0 及び命令 RAM (Random Access Memory) 4 1 及びデータ RAM 4 2 を有する。

【 0 0 2 7 】

プロセッサ 4 0 は、制御部 1 0、レジスタファイル 1 1、演算部 1 2、統計情報集約部 1 3、メモリインタフェース 1 4 及びメモリインタフェース 1 5 を有する。メモリインタフェース 1 4 は、プロセッサ 4 0 における命令 RAM 4 1 に接続するインタフェースである。また、メモリインタフェース 1 5 は、プロセッサ 4 0 におけるデータ RAM 4 2 に接続するインタフェースである。以下の説明では、プロセッサ 4 0 の各部が命令 RAM 4 1 又はデータ RAM 4 2 にアクセスする場合、メモリインタフェース 1 4 及び 1 5 の仲介を省略して説明する。

30

【 0 0 2 8 】

命令 RAM 4 1 は、CPU 2 から送信された命令を格納する記憶装置である。命令 RAM 4 1 に格納された命令は、制御部 1 0 によりフェッチされて実行される。データ RAM 4 2 は、命令で指定された演算を実行する際に使用するデータを格納する記憶装置である。データ RAM 4 2 に格納されたデータは、演算部 1 2 で実行される演算に使用される。

40

【 0 0 2 9 】

レジスタファイル 1 1 は、スカラーレジスタファイル 1 1 1、ベクタレジスタファイル 1 1 2、アキュムレータレジスタ 1 1 3、ベクタアキュムレータレジスタ 1 1 4、統計情報格納部 1 1 5 及び候補格納部 3 0 0 を有する。

【 0 0 3 0 】

スカラーレジスタファイル 1 1 1 及びベクタレジスタファイル 1 1 2 は、入力されたデータや学習処理の実行途中のデータなど演算に用いるデータを格納する。アキュムレータレジスタ 1 1 3 及びベクタアキュムレータレジスタ 1 1 4 は、演算部 1 2 が演算を実行する際に、累積などの演算を行う場合にデータを一時的に格納する。

50

【 0 0 3 1 】

統計情報格納部 1 1 5 は、統計情報集約部 1 3 が集約した統計情報を取得して格納する。統計情報は、演算結果の小数点位置に関する情報である。統計情報は、例えば、非符号となる最上位ビット位置の分布、非ゼロの最下位ビット位置の分布、非符号となる最上位ビット位置の最大値又は非ゼロの最下位ビット位置の最小値などを含む複数の情報のうちのいずれか又はその組み合わせである。

【 0 0 3 2 】

候補格納部 3 0 0 は、操作者が指定した N 個の候補小数点位置をそれぞれ用いて演算部 1 2 が実行した先行演算の N 個の演算結果を格納する。候補小数点位置及び先行演算については後で詳細に説明する。

10

【 0 0 3 3 】

次に、演算部 1 2 について説明する。演算部 1 2 は、スカラユニット 1 2 1 及びベクタユニット 1 2 2 を有する。

【 0 0 3 4 】

スカラユニット 1 2 1 は、制御部 1 0、レジスタファイル 1 1 及びメモリアンタフェース 1 5 に接続される。スカラユニット 1 2 1 は、演算器 2 1 1、統計情報取得部 2 1 2 及びデータ変換部 2 1 3 を有する。本実施例では、スカラユニット 1 2 1 は、N 個の候補小数点位置を用いた演算であり、且つ、統計情報を取得するための先行演算を実行する。また、統計情報から求められた小数点位置と一致する候補小数点位置が存在しなければ、スカラユニット 1 2 1 は、先行演算の統計情報から決定された小数点位置で演算を実行して演算結果を取得する本演算という 2 つの演算を実行する。

20

【 0 0 3 5 】

演算器 2 1 1 は、データ RAM 4 2、スカラレジスタファイル 1 1 1 及びアキュムレータレジスタ 1 1 3 が保持するデータのうちの 1 つもしくはいくつかをを用いて積和演算などの演算を実行する。この演算器 2 1 1 が演算に用いるデータが、「入力データ」の一例にあたる。演算器 2 1 1 は、先行演算及び本演算のいずれの演算でも同様の演算を実行する。演算器 2 1 1 は、演算結果を表すのに十分なビット幅を用いて演算を実行する。演算器 2 1 1 は、データ RAM 4 2、統計情報取得部 2 1 2 及びデータ変換部 2 1 3 に演算結果を出力する。

【 0 0 3 6 】

統計情報取得部 2 1 2 は、演算結果のデータの入力を演算器 2 1 1 から受ける。そして、統計情報取得部 2 1 2 は、演算結果のデータから統計情報を取得する。その後、統計情報取得部 2 1 2 は、取得した統計情報を統計情報集約部 1 3 へ出力する。ただし、統計情報取得部 2 1 2 は、本演算の場合統計情報の取得及び取得した統計情報の出力を行わなくてもよい。

30

【 0 0 3 7 】

データ変換部 2 1 3 は、演算器 2 1 1 による演算結果を取得する。次に、データ変換部 2 1 3 は、先行演算の場合、N 個の候補小数点位置の入力を制御部 1 0 から受ける。そして、候補少数点位置毎に、データ変換部 2 1 3 は、固定小数点数データを取得した各候補少数点位置で指定されたシフト量だけシフトさせる。また、データ変換部 2 1 3 は、シフトとともに、上位ビットの飽和処理および下位ビットの丸めを実行する。これにより、データ変換部 2 1 3 は、固定小数点数のデータの小数点位置を更新した N 個の演算結果を算出する。その後、データ変換部 2 1 3 は、候補小数点位置を用いた N 個の演算結果を候補格納部 3 0 0 に格納する。

40

【 0 0 3 8 】

また、データ変換部 2 1 3 は、本演算の場合、先行演算により取得された統計情報から決定された小数点位置の入力を制御部 1 0 から受ける。そして、データ変換部 2 1 3 は、固定小数点数データを取得した少数点位置で指定されたシフト量だけシフトさせる。また、データ変換部 2 1 3 は、シフトとともに、上位ビットの飽和処理および下位ビットの丸めを実行する。これにより、データ変換部 2 1 3 は、固定小数点数のデータの小数点位置

50

を更新する。データ変換部 2 1 3 は、小数点位置を更新した演算結果をスカラレジスタファイル 1 1 1 及びデータ RAM 4 2 に格納する。

【 0 0 3 9 】

ベクタユニット 1 2 2 は、制御部 1 0、レジスタファイル 1 1 及びメモリアンタフェース 1 5 に接続される。ベクタユニット 1 2 2 は、演算器 2 2 1、統計情報取得部 2 2 2 及びデータ変換部 2 2 3 の組を複数有する。本実施例では、ベクタユニット 1 2 2 は、N 個の候補小数点位置を用いた演算であり、且つ、統計情報を取得するための先行演算を実行する。また、統計情報から求められた小数点位置と一致する候補小数点位置が存在しなければ、ベクタユニット 1 2 2 は、先行演算の統計情報から決定された小数点位置で演算を実行して演算結果を取得する本演算という 2 つの演算を実行する。

10

【 0 0 4 0 】

演算器 2 2 1 は、データ RAM 4 2、ベクタレジスタファイル 1 1 2 又はベクタアキュムレータレジスタ 1 1 4 のうちの 1 つもしくはいくつかが保持するデータを用いて積和演算などの演算を実行する。演算器 2 2 1 は、演算結果を表すのに十分なビット幅を用いて演算を実行する。演算器 2 2 1 は、先行演算及び本演算のいずれの演算でも同様の演算を実行する。演算器 2 1 1 は、データ RAM 4 2、統計情報取得部 2 2 2 及びデータ変換部 2 2 3 に演算結果を出力する。

【 0 0 4 1 】

統計情報取得部 2 2 2 は、演算結果のデータの入力を演算器 2 2 1 から受ける。この時、統計情報取得部 2 2 2 は、精度を維持できる十分なビット幅で表された演算結果のデータを取得する。

20

【 0 0 4 2 】

そして、統計情報取得部 2 2 2 は、演算結果のデータから統計情報を取得する。例えば、非符号となる最上位ビット位置を取得する場合、統計情報取得部 2 2 2 は、非符号最上位ビット検出器を用いて、非符号となる最上位ビット位置の値を 1 とし、他のビット位置の値を 0 とする出力データを生成する。その後、統計情報取得部 2 2 2 は、取得した統計情報を統計情報集約部 1 3 へ出力する。ただし、統計情報取得部 2 2 2 は、本演算の場合統計情報の取得及び取得した統計情報の出力を行わなくてもよい。

【 0 0 4 3 】

データ変換部 2 2 3 は、演算器 2 2 1 による演算結果を取得する。次に、データ変換部 2 2 3 は、先行演算の場合、N 個の候補小数点位置の入力を制御部 1 0 から受ける。そして、候補少数点位置毎に、データ変換部 2 2 3 は、固定小数点数データを取得した各候補少数点位置で指定されたシフト量だけシフトさせる。また、データ変換部 2 2 3 は、シフトとともに、上位ビットの飽和処理および下位ビットの丸めを実行する。これにより、データ変換部 2 2 3 は、固定小数点数のデータの小数点位置を更新した N 個の演算結果を算出する。その後、データ変換部 2 2 3 は、候補小数点位置を用いた N 個の演算結果を候補格納部 3 0 0 に格納する。

30

【 0 0 4 4 】

データ変換部 2 2 3 は、演算器 2 2 1 による演算結果を取得する。次に、データ変換部 2 2 3 は、本演算の場合、先行演算により取得された統計情報から決定された小数点位置の入力を制御部 1 0 から受ける。そして、データ変換部 2 2 3 は、固定小数点数のデータを取得した少数点位置で指定されたシフト量だけシフトさせる。また、データ変換部 2 2 3 は、シフトとともに、上位ビットの飽和処理および下位ビットの丸めを実行する。これにより、データ変換部 2 2 3 は、固定小数点数のデータの小数点位置を更新する。データ変換部 2 2 3 は、小数点位置を更新した演算結果をアキュムレータ 1 0 3 に格納し、その後ベクタレジスタファイル 1 1 2 及びデータ RAM 4 2 に格納する。

40

【 0 0 4 5 】

統計情報集約部 1 3 は、演算器 2 2 1 による演算結果のデータから取得された統計情報の入力を統計情報取得部 2 1 2 から受ける。また、統計情報集約部 1 3 は、各演算器 2 2 1 による演算結果のデータから取得されたそれぞれの統計情報の入力を各統計情報取得部

50

2 2 2 から受ける。統計情報集約部 1 3 は、統計情報取得部 2 1 2 から取得した統計情報及び各統計情報取得部 2 2 2 から取得した各統計情報を集約して統計情報格納部 1 1 5 へ出力する。

【 0 0 4 6 】

次に、制御部 1 0 について説明する。図 5 は、制御部の詳細を表すブロック図である。ここでは、候補小数点位置が 4 個の場合を例に説明する。ただし、N は、2 以上の整数であればよい。また、N が多いほど後述する再演算の実行を回避することができるが、先行する演算の数が増えまた記憶領域も増大する。制御部 1 0 は、図 5 に示すように、統括管理部 1 0 0、小数点位置決定部 1 0 1 及び指数値変換制御部 1 0 2 を有する。

【 0 0 4 7 】

ここで、例えば、候補格納部 3 0 0 は、候補小数点位置の個数に対応する数の候補格納部 3 0 1 ~ 3 0 4 を有する。そして、候補格納部 3 0 1 ~ 3 0 4 は、演算部 1 2 によりそれぞれの候補小数点位置を用いて実行された先行演算の各演算結果をそれぞれが格納する。

【 0 0 4 8 】

統括管理部 1 0 0 は、演算部 1 2 による先行演算及び本演算の実行を管理する。統括管理部 1 0 0 は、深層学習において演算部 1 2 に演算を実行させるレイヤの情報を保持する。統括管理部 1 0 0 は、演算部 1 2 に演算を実行させるレイヤが次のレイヤに移ると、先行演算の実行を決定する。

【 0 0 4 9 】

次に、統括管理部 1 0 0 は、操作者が指定した 4 個の候補小数点位置を取得する。そして、統括管理部 1 0 0 は、取得した 4 個の候補小数点位置を指数値変換制御部 1 0 2 に通知し、演算部 1 2 に先行演算の実行を指示する。この演算部 1 2 による先行演算が、「第 1 演算」の一例にあたり、先行演算の演算結果が「第 1 演算結果」の一例にあたる。

【 0 0 5 0 】

その後、統括管理部 1 0 0 は、演算部 1 2 による先行演算の実行が完了すると、新たに算出された小数点位置を小数点位置決定部 1 0 1 から取得する。そして、統括管理部 1 0 0 は、小数点位置決定部 1 0 1 から取得した小数点位置に一致する候補小数点位置が存在するか否かを判定する。

【 0 0 5 1 】

小数点位置決定部 1 0 1 から取得した小数点位置に一致する候補小数点位置が存在する場合、統括管理部 1 0 0 は、その候補小数点位置を用いて小数点位置が更新された固定小数点数を演算結果として決定する。例えば、図 5 に示すように、統括管理部 1 0 0 は、候補格納部 3 0 1 ~ 3 0 4 が有する各候補小数点位置を用いて小数点位置が更新された演算結果の中からセレクト 3 1 0 に選択させる。その後、統括管理部 1 0 0 は、決定した演算結果をそのレイヤにおける最終演算結果としてデータ RAM 4 2 に格納する。

【 0 0 5 2 】

これに対して、小数点位置決定部 1 0 1 から取得した小数点位置に一致する候補小数点位置が存在しない場合、統括管理部 1 0 0 は、本演算の実行を決定する。そして、統括管理部 1 0 0 は、新たに決定された小数点位置の出力を指数値変換制御部 1 0 2 に指示し、演算部 1 2 に本演算の実行を指示する。本演算の演算結果はそのレイヤにおける最終演算結果としてアキュムレータレジスタ 1 1 3 を経由してデータ RAM 4 2 に格納される。この演算部 1 2 による本演算が、「第 2 演算」の一例にあたり、本演算の演算結果が「第 2 演算結果」の一例にあたる。統括管理部 1 0 0 は、以上のような先行演算及び本演算を演算部 1 2 に実行させる制御をレイヤ毎に繰り返す。

【 0 0 5 3 】

また、統括管理部 1 0 0 は、深層学習におけるイテレーションの管理も行う。例えば、所定回数のイテレーションの実行が指示された場合、統括管理部 1 0 0 は、実行されたイテレーションの回数をカウントし所定回数に達すると学習の終了を決定する。その後、統括管理部 1 0 0 は、例えば、CPU 2 に学習終了を通知して学習を終了する。この統括管

10

20

30

40

50

理部 100 が、「管理部」の一例にあたる。

【0054】

小数点位置決定部 101 は、演算部 12 による先行演算が終了すると、統計情報を統計情報格納部 115 から取得する。そして、小数点位置決定部 101 は、取得した統計情報を用いて最適な小数点位置を決定する。その後、小数点位置決定部 101 は、決定した小数点位置を指数値変換制御部 102 へ出力する。小数点位置決定部 101 は、先行計算後の小数点位置の決定処理をレイヤ毎に繰り返す。

【0055】

指数値変換制御部 102 は、各候補小数点位置を統括管理部 100 から取得する。さらに、指数値変換制御部 102 は、N 個の候補小数点位置の出力の指示を統括管理部 100 から受ける。そして、指数値変換制御部 102 は、N 個の候補小数点位置を演算部 12 に出力する。

10

【0056】

その後、演算部 12 による先行演算が完了すると、指数値変換制御部 102 は、先行演算の演算結果を用いて新たに決定された小数点位置の入力を小数点位置決定部 101 から受ける。そして、小数点位置決定部 101 が算出した小数点位置に一致する候補小数点位置が存在しない場合、指数値変換制御部 102 は、新たに決定された小数点位置の出力の指示の入力を統括管理部 100 から受ける。その後、指数値変換制御部 102 は、新たに決定された小数点位置の情報を演算部 12 に出力する。

【0057】

次に、図 6 を参照して、本実施例に係る演算用回路 4 による深層学習の処理の流れを説明する。図 6 は、実施例 1 に係る演算用回路による深層学習の処理のフローチャートである。

20

【0058】

制御部 10 の統括管理部 100 は、操作者が指定した N 個の候補小数点位置を取得する（ステップ S101）。そして、統括管理部 100 は、N 個の候補小数点位置を指数値変換制御部 102 へ出力する。

【0059】

制御部 10 の指数値変換制御部 102 は、N 個の候補小数点位置を演算部 12 へ出力する。そして、各演算器 211 及び 221 は、N 個の候補小数点位置のそれぞれを用いて、位置入力データを用いて先行演算における演算を実行する（ステップ S102）。各統計情報取得部 212 及び 222 は、対応する各演算器 211 及び 221 による演算結果から統計情報を求める。統計情報集約部 13 は、各統計情報取得部 212 及び 222 から統計情報を集約して統計情報格納部 115 に格納する。

30

【0060】

データ変換部 213 及び 223 は、演算器 211 及び 221 による演算結果を取得する。そして、データ変換部 213 及び 223 は、N 個の候補小数点位置のそれぞれを用いて小数点位置を更新する（ステップ S103）。データ変換部 213 及び 223 は、N 個の先行演算の演算結果を候補格納部 300 に格納する。

【0061】

制御部 10 の小数点位置決定部 101 は、統計情報格納部 115 に格納された統計情報を初期化する。そして、小数点位置決定部 101 は、先行演算における演算結果の統計情報を用いて新たな小数点位置を決定する（ステップ S104）。

40

【0062】

制御部 10 の統括管理部 100 は、N 個の候補小数点位置の中に小数点位置決定部 101 により決定された新たな小数点位置に一致する候補小数点位置が存在するか否かを判定する（ステップ S105）。

【0063】

新たな小数点位置に一致する候補小数点位置が存在しない場合（ステップ S105：否定）、統括管理部 100 は、新たな小数点位置の出力を指数値変換制御部 102 に指示す

50

るとともに、本演算の実行を演算部 1 2 に指示する。演算部 1 2 の各演算器 2 1 1 及び 2 2 1 は、入力データを用いて本演算における演算を実行する（ステップ S 1 0 6）。

【 0 0 6 4 】

演算部 1 2 のデータ変換部 2 1 3 及び 2 2 3 は、指数値変換制御部 1 0 2 から入力された小数点位置で、演算器 2 1 1 及び 2 2 1 による演算結果の小数点位置を更新する（ステップ S 1 0 7）。このように、演算部 1 2 は、本演算を実行する。演算部 1 2 は、本演算の演算結果をそのレイヤの演算結果とする。

【 0 0 6 5 】

これに対して、新たな小数点位置に一致する候補小数点位置が存在する場合（ステップ S 1 0 5：肯定）、統括管理部 1 0 0 は、新たな小数点位置に一致する候補小数点位置を用いて算出された固定小数点数を候補格納部 3 0 0 から取得する。そして、統括管理部 1 0 0 は、取得した固定小数点数をそのレイヤの演算結果とする（ステップ S 1 0 8）。

10

【 0 0 6 6 】

その後、制御部 1 0 の統括管理部 1 0 0 は、実行中の全てのレイヤが終了したか否かを判定する（ステップ S 1 0 9）。全てのレイヤが終了していない場合（ステップ S 1 0 9：否定）、統括管理部 1 0 0 は、次のレイヤの演算を開始させる（ステップ S 1 1 0）。その後、深層学習の処理はステップ S 1 0 1 へ戻る。

【 0 0 6 7 】

これに対して、全てのレイヤが終了した場合（ステップ S 1 0 9：肯定）、制御部 1 0 の統括管理部 1 0 0 は、学習が終了したか否かを判定する（ステップ S 1 1 1）。

20

【 0 0 6 8 】

学習が終了していない場合（ステップ S 1 1 1：否定）、統括管理部 1 0 0 は、次のイテレーションを開始する（ステップ S 1 1 2）。その後、深層学習の処理はステップ S 1 0 1 へ戻る。

【 0 0 6 9 】

これに対して、学習が終了した場合（ステップ S 1 1 1：肯定）、統括管理部 1 0 0 は、学習完了を CPU 2 に通知して学習を終了する。

【 0 0 7 0 】

以上に説明したように、本実施例に係る演算用回路は、予め決められた N 個の候補小数点位置のそれぞれを用いて先行演算を行い、N 個の先行演算の演算結果を取得する。また、演算用回路は、先行演算における演算の結果から得た統計情報を用いて、入力データを用いた演算に対する適切な小数点位置を新たに決定する。そして、新たな小数点位置と一致する候補小数点位置が存在する場合、演算用回路は、その候補小数点位置を用いて算出した先行演算の演算結果をそのレイヤの演算結果とする。これに対して、新たな小数点位置と一致する候補小数点位置が存在しない場合、演算用回路は、新たな小数点位置を用いて本演算を実行し、その本演算の演算結果をそのレイヤの演算結果とする。このように、演算用回路は、N 個の候補小数点位置を用いて先行演算を投機的に実行し、その投機的演算の演算結果が適切な小数点位置と一致する場合にはその投機的演算の演算結果をそのレイヤの演算結果とすることができる。

30

【 0 0 7 1 】

これにより、ニューラルネットの構造となる計算グラフの構築が、深層学習の処理を実行しながら行われる Define by Run により深層学習を行う際に、適切な小数点位置を決定できる。そして、固定小数点を用いた学習における学習精度を向上させることが可能となる。それに加えて、投機的演算が適切な小数点位置を用いた場合の演算と同じ結果となる場合、その投機的演算の演算結果を用いることができるため、本演算の実行を回避することができる。すなわち、投機的演算を行わずに、先行演算の統計情報を用いて適切な小数点位置を求めて本演算を実行することで 2 回の演算を行う場合に比べて、オーバーヘッドとなる演算時間を低減することができる。

40

【 0 0 7 2 】

ここで、投機的演算のため演算時間の低減の割合は確率的であるが、例えば、投機的演

50

算の演算結果がレイヤの演算結果として使用できる確率が20%である場合を考える。1回の演算時間を1基準時間とすると、この場合、本実施例に係る演算用回路の演算時間は、 1.2 基準時間 $\times 0.8 + 1$ 回 $\times (1.2 + 1)$ 基準時間 $\times 0.2$ と表される。これに対して、2回の演算を行う場合は、 2 回 $\times 1$ 基準時間と表される。すなわち、本実施例に係る演算用回路は、2回演算を行う場合に比べてオーバーヘッドを60%削減することができる。これにより、Define by Runにより深層学習を行う際に、固定小数点を用いた学習における学習精度を向上させつつ、演算のオーバーヘッドを軽減でき、学習時間を短縮することが可能となる。

【実施例2】

【0073】

次に、実施例2について説明する。本実施例に係る演算用回路4は、1つ前のレイヤで使用した小数点位置から候補小数点位置を生成することが実施例1と異なる。本実施例に係る演算用回路4も図4及び図5で表される。以下の説明では、実施例1と同様の各部の機能については説明を省略する。ここでは、候補小数点位置を4つ用いる場合で説明する。

【0074】

統括管理部100は、1つ前のレイヤである前レイヤで用いた小数点位置を取得する。例えば、前レイヤで用いた小数点位置はデータRAM42に格納される。そして、統括管理部100は、前レイヤで用いた小数点位置が指定するビット値に+1、+0、-1、-2をしたビット値を候補小数点位置として取得する。その後、統括管理部100は、候補小数点位置を用いた先行演算を演算部12に実行させる。この場合、投機的演算の演算結果を使用できる確率を向上させるために、前レイヤで用いた小数点位置を挟むように候補小数点位置を決定することが好ましい。

【0075】

その後、先行演算における演算の統計情報から小数点位置決定部101により決定される新しい小数点位置に一致する候補小数点位置がある場合、統括管理部100は、その候補小数点位置を用いた先行演算の演算結果をそのレイヤの演算結果とする。また、新しい小数点位置に一致する候補小数点位置がない場合、統括管理部100は、新しい小数点位置を用いた本演算を演算部12に実行させ、その演算結果をそのレイヤの演算結果とする。

【0076】

次に、図7を参照して、本実施例に係る演算用回路4による深層学習の処理の流れを説明する。図7は、実施例2に係る演算用回路による深層学習の処理のフローチャートである。

【0077】

制御部10の統括管理部100は、前レイヤの小数点位置を取得する(ステップS201)。

【0078】

制御部10の指数値変換制御部102は、前レイヤの小数点位置が表すビット位置を予め決められたビット数ずらすなどしてN個の候補小数点位置を生成する(ステップS202)。

【0079】

次に、指数値変換制御部102は、N個の候補小数点位置を演算部12へ出力する。そして、演算部12の各演算器211及び221は、N個の候補小数点位置のそれぞれを用いて、位置入力データを用いて先行演算における演算を実行する(ステップS203)。各統計情報取得部212及び222は、対応する各演算器211及び221による演算結果から統計情報を求める。統計情報集約部13は、各統計情報取得部212及び222から統計情報を集約して統計情報格納部115に格納する。

【0080】

データ変換部213及び223は、演算器211及び221による演算結果を取得する。そして、データ変換部213及び223は、N個の候補小数点位置のそれぞれを用いて

10

20

30

40

50

小数点位置を更新する（ステップS 2 0 4）。データ変換部 2 1 3 及び 2 2 3 は、N 個の先行演算の演算結果を候補格納部 3 0 0 に格納する。

【 0 0 8 1 】

制御部 1 0 の小数点位置決定部 1 0 1 は、統計情報格納部 1 1 5 に格納された統計情報を初期化する。そして、小数点位置決定部 1 0 1 は、先行演算における演算結果の統計情報を用いて新たな小数点位置を決定する（ステップS 2 0 5）。

【 0 0 8 2 】

制御部 1 0 の統括管理部 1 0 0 は、N 個の候補小数点位置の中に小数点位置決定部 1 0 1 により決定された新たな小数点位置に一致する候補小数点位置が存在するか否かを判定する（ステップS 2 0 6）。

10

【 0 0 8 3 】

新たな小数点位置に一致する候補小数点位置が存在しない場合（ステップS 2 0 6：否定）、統括管理部 1 0 0 は、新たな小数点位置の出力を指数値変換制御部 1 0 2 に指示するとともに、本演算の実行を演算部 1 2 に指示する。演算部 1 2 の各演算器 2 1 1 及び 2 2 1 は、入力データを用いて本演算における演算を実行する（ステップS 2 0 7）。

【 0 0 8 4 】

演算部 1 2 のデータ変換部 2 1 3 及び 2 2 3 は、指数値変換制御部 1 0 2 から入力された小数点位置で、演算器 2 1 1 及び 2 2 1 による演算結果の小数点位置を更新する（ステップS 2 0 8）。このように、演算部 1 2 は、本演算を実行する。統括管理部 1 0 0 は、本演算の演算結果をそのレイヤの演算結果とする。

20

【 0 0 8 5 】

これに対して、新たな小数点位置に一致する候補小数点位置が存在する場合（ステップS 2 0 6：肯定）、統括管理部 1 0 0 は、新たな小数点位置に一致する候補小数点位置を用いて算出された固定小数点数を候補格納部 3 0 0 から取得する。そして、統括管理部 1 0 0 は、取得した固定小数点数をそのレイヤの演算結果とする（ステップS 2 0 9）。

【 0 0 8 6 】

その後、制御部 1 0 の統括管理部 1 0 0 は、実行中の全てのレイヤが終了したか否かを判定する（ステップS 2 1 0）。全てのレイヤが終了していない場合（ステップS 2 1 0：否定）、統括管理部 1 0 0 は、次のレイヤの演算を開始させる（ステップS 2 1 1）。その後、深層学習の処理はステップS 2 0 1 へ戻る。

30

【 0 0 8 7 】

これに対して、全てのレイヤが終了した場合（ステップS 2 1 0：肯定）、制御部 1 0 の統括管理部 1 0 0 は、学習が終了したか否かを判定する（ステップS 2 1 2）。

【 0 0 8 8 】

学習が終了していない場合（ステップS 2 1 2：否定）、統括管理部 1 0 0 は、次のイテレーションを開始する（ステップS 2 1 3）。その後、深層学習の処理はステップS 2 0 1 へ戻る。

【 0 0 8 9 】

これに対して、学習が終了した場合（ステップS 2 1 2：肯定）、統括管理部 1 0 0 は、学習完了をCPU 2 に通知して学習を終了する。

40

【 0 0 9 0 】

以上に説明したように、本実施例に係る演算回路は、1 つ前のレイヤの小数点位置から候補小数点位置を生成する。深層学習におけるCNN (Convolutional Neural Network) 系のネットワークのように同じようなレイヤ構造が続く場合、小数点位置は安定する傾向にある。そこで、1 つ前のレイヤの小数点位置を用いることで、投機的演算の演算結果を使用できる確率を向上させることができる。したがって、オーバヘッドとなる演算時間をより低減することができる。これにより、Define by Runにより深層学習を行う際に、固定小数点を用いた学習における学習精度を向上させつつ、演算のオーバヘッドを軽減でき、学習時間を短縮することが可能となる。

【実施例 3】

50

【 0 0 9 1 】

次に、実施例 3 について説明する。本実施例に係る演算用回路 4 は、対象とするレイヤで実行される積和演算から小数点位置を限定して候補小数点位置を生成することが実施例 1 と異なる。すなわち、本実施例で実行される処理は、積和演算を実行するレイヤを対象とする。本実施例に係る演算用回路 4 も図 4 及び 5 で表される。以下の説明では、実施例 1 と同様の各部の機能については説明を省略する。ここでは、積和演算を実行するレイヤでの処理において候補小数点位置を 4 つ用いる場合で説明する。

【 0 0 9 2 】

統括管理部 1 0 0 は、対象とするレイヤで実行させる積和演算を K 分割する。例えば、統括管理部 1 0 0 は、次の数式 (1) で示されるように積和演算を 4 分割する。この、対象とするレイヤで実行させる積和演算が、「所定積和演算」の一例にあたる。

10

【 0 0 9 3 】

【 数 1 】

$$\sum_N XW = \sum_0^{N1} XW + \sum_{N1}^{N2} XW + \sum_{N2}^{N3} XW + \sum_{N3}^{N4} XW \dots(1)$$

【 0 0 9 4 】

ここで、数式 (1) の左辺は対象とするレイヤで実行される積和演算を表す。そして、数式 (1) の右辺が左辺の積和演算を 4 分割した積和演算である。以下では、分割後の各積和演算を「分割積和演算」と言う。

20

【 0 0 9 5 】

次に、統括管理部 1 0 0 は、K 分割後の分割積和演算毎の演算結果を演算部 1 2 に算出させる。次に、統括管理部 1 0 0 は、各分割積和演算の演算結果の中から最大値となる演算結果を選択する。

【 0 0 9 6 】

ここで、対象とするレイヤの積和演算の演算結果は、各分割積和演算の演算結果の最大値の K 倍以下となる。そこで、統括管理部 1 0 0 は、分割積和演算の演算結果のうちの最大値となる演算結果を K 倍した値を算出する。そして、統括管理部 1 0 0 は、分割積和演算の演算結果のうちの最大値となる演算結果を K 倍した値の小数点位置を上限小数点位置として求める。例えば、K 分割後の各積和演算の演算結果の最大値を V M A X とした場合、統括管理部 1 0 0 は、 $B = \log_2 (V M A X \times K)$ として上限小数点位置を求める。ここで、B は、上限小数点位置を表す。

30

【 0 0 9 7 】

対象とするレイヤの積和演算の演算結果の小数点位置は、上限小数点位置よりもビット位置が上になると考えられるので、統括管理部 1 0 0 は、上限小数点位置が示すビット位置から上のビット位置を候補小数点位置として取得する。例えば、上述した B が上限小数点位置の場合、統括管理部 1 0 0 は、B、B - 1、B - 2 及び B - 3 を 4 つの候補小数点位置として取得する。

【 0 0 9 8 】

次に、統括管理部 1 0 0 は、K 分割後の各積和演算の演算結果の総和を演算部 1 2 に算出させる。そして、統括管理部 1 0 0 は、算出された総和に対して候補小数点位置を用いたが小数点位置の変更を行わせることで先行演算を演算部 1 2 に実行させる。

40

【 0 0 9 9 】

その後、先行演算における演算の統計情報から小数点位置決定部 1 0 1 により決定される新しい小数点位置に一致する候補小数点位置がある場合、統括管理部 1 0 0 は、その候補小数点位置を用いた先行演算の演算結果をそのレイヤの演算結果とする。また、新しい小数点位置に一致する候補小数点位置がない場合、統括管理部 1 0 0 は、新しい小数点位置を用いた本演算を演算部 1 2 に実行させ、その演算結果をそのレイヤの演算結果とする。

【 0 1 0 0 】

50

次に、図 8 A 及び図 8 B を参照して、本実施例に係る演算用回路 4 による深層学習の処理の流れを説明する。図 8 A 及び図 8 B は、実施例 3 に係る演算用回路による深層学習の処理のフローチャートである。

【 0 1 0 1 】

制御部 1 0 の統括管理部 1 0 0 は、対象のレイヤで積和演算が実行されるか否かを判定する（ステップ S 3 0 1 ）。

【 0 1 0 2 】

対象のレイヤで積和演算が実行される場合（ステップ S 3 0 1 ：肯定）、統括管理部 1 0 0 は、積和演算を K 分割する（ステップ S 3 0 2 ）。

【 0 1 0 3 】

そして、統括管理部 1 0 0 は、各分割積和演算の実行を演算部 1 2 に指示する。演算部 1 2 は、分割積和演算毎に演算を実行する（ステップ S 3 0 3 ）。

【 0 1 0 4 】

統括管理部 1 0 0 は、各分割積和演算の演算結果の最大値を取得する（ステップ S 3 0 4 ）。

【 0 1 0 5 】

指数値変換制御部 1 0 2 は、分割積和演算の演算結果の最大値を K 倍した値の小数点位置を求めて上限小数点位置とする（ステップ S 3 0 5 ）。

【 0 1 0 6 】

指数値変換制御部 1 0 2 は、上限小数点位置が表すビット位置から上の位置の N 個の候補小数点位置を生成する（ステップ S 3 0 6 ）。

【 0 1 0 7 】

次に、指数値変換制御部 1 0 2 は、分割積和演算の総和の算出を演算部 1 2 に指示する。演算部 1 2 の各各演算器 2 1 1 及び 2 2 1 は、分割積和演算の総和を計算する（ステップ S 3 0 7 ）。各統計情報取得部 2 1 2 及び 2 2 2 は、対応する各演算器 2 1 1 及び 2 2 1 による演算結果から統計情報を求める。統計情報集約部 1 3 は、各統計情報取得部 2 1 2 及び 2 2 2 から統計情報を集約して統計情報格納部 1 1 5 に格納する。

【 0 1 0 8 】

データ変換部 2 1 3 及び 2 2 3 は、演算器 2 1 1 及び 2 2 1 による演算結果を取得する。そして、データ変換部 2 1 3 及び 2 2 3 は、N 個の候補小数点位置のそれぞれを用いて小数点位置を更新する。データ変換部 2 1 3 及び 2 2 3 は、算出した N 個の先行演算の演算結果を候補格納部 3 0 0 に格納する。

【 0 1 0 9 】

制御部 1 0 の小数点位置決定部 1 0 1 は、統計情報格納部 1 1 5 に格納された統計情報を初期化する。そして、小数点位置決定部 1 0 1 は、先行演算における演算結果の統計情報を用いて新たな小数点位置を決定する（ステップ S 3 0 8 ）。

【 0 1 1 0 】

制御部 1 0 の統括管理部 1 0 0 は、N 個の候補小数点位置の中に小数点位置決定部 1 0 1 により決定された新たな小数点位置に一致する候補小数点位置が存在するか否かを判定する（ステップ S 3 0 9 ）。

【 0 1 1 1 】

新たな小数点位置に一致する候補小数点位置が存在しない場合（ステップ S 3 0 9 ：否定）、統括管理部 1 0 0 は、新たな小数点位置の出力を指数値変換制御部 1 0 2 に指示するとともに、本演算の実行を演算部 1 2 に指示する。演算部 1 2 の各演算器 2 1 1 及び 2 2 1 は、入力データを用いて本演算における演算を実行する（ステップ S 3 1 0 ）。

【 0 1 1 2 】

演算部 1 2 のデータ変換部 2 1 3 及び 2 2 3 は、指数値変換制御部 1 0 2 から入力された小数点位置で、演算器 2 1 1 及び 2 2 1 による演算結果の小数点位置を更新する（ステップ S 3 1 1 ）。このように、演算部 1 2 は、本演算を実行する。統括管理部 1 0 0 は、本演算の演算結果をそのレイヤの演算結果とする。

10

20

30

40

50

【 0 1 1 3 】

一方、新たな小数点位置に一致する候補小数点位置が存在する場合（ステップ S 3 0 9 : 肯定）、統括管理部 1 0 0 は、新たな小数点位置に一致する候補小数点位置を用いて算出された固定小数点数をそのレイヤの演算結果として選択する（ステップ S 3 1 2）。

【 0 1 1 4 】

これに対して、対象のレイヤで積和演算が実行されない場合（ステップ S 3 0 1 : 否定）、統括管理部 1 0 0 は、他の手順でそのレイヤの演算を実行する（ステップ S 3 1 3）。ここで、他の手順として、例えば、実施例 1 又は 2 で説明した手順を用いることができる。

【 0 1 1 5 】

その後、制御部 1 0 の統括管理部 1 0 0 は、実行中の全てのレイヤが終了したか否かを判定する（ステップ S 3 1 4）。全てのレイヤが終了していない場合（ステップ S 3 1 4 : 否定）、統括管理部 1 0 0 は、次のレイヤの演算を開始させる（ステップ S 3 1 5）。その後、深層学習の処理はステップ S 3 0 1 へ戻る。

10

【 0 1 1 6 】

これに対して、全てのレイヤが終了した場合（ステップ S 3 1 4 : 肯定）、制御部 1 0 の統括管理部 1 0 0 は、学習が終了したか否かを判定する（ステップ S 3 1 6）。

【 0 1 1 7 】

学習が終了していない場合（ステップ S 3 1 6 : 否定）、統括管理部 1 0 0 は、次のイテレーションを開始する（ステップ S 3 1 7）。その後、深層学習の処理はステップ S 3 0 1 へ戻る。

20

【 0 1 1 8 】

これに対して、学習が終了した場合（ステップ S 3 1 6 : 肯定）、統括管理部 1 0 0 は、学習完了を CPU 2 に通知して学習を終了する。

【 0 1 1 9 】

以上に説明したように、本実施例に係る演算用回路は、対象のレイヤで実行される積和演算を分割し、分割積和演算の演算結果の最大値から上限小数点位置を求め、求めた上限小数点位置から候補小数点位置を生成する。これにより、制限された範囲の小数点位置の中から候補小数点位置を選定することができ、投機的演算の演算結果を使用できる確率を向上させることができる。したがって、オーバヘッドとなる演算時間をより低減することができる。これにより、Define by Runにより深層学習を行う際に、固定小数点を用いた学習における学習精度を向上させつつ、演算のオーバヘッドを軽減でき、学習時間を短縮することが可能となる。

30

【 符号の説明 】

【 0 1 2 0 】

- 1 サーバ
- 2 CPU
- 3 メモリ
- 4 演算用回路
- 5 PCI e バス
- 1 0 制御部
- 1 1 レジスタファイル
- 1 2 演算部
- 1 3 統計情報集約部
- 1 4 メモリインタフェース
- 1 5 メモリインタフェース
- 4 0 プロセッサ
- 4 1 命令 RAM
- 4 2 データ RAM
- 1 0 0 統括管理部

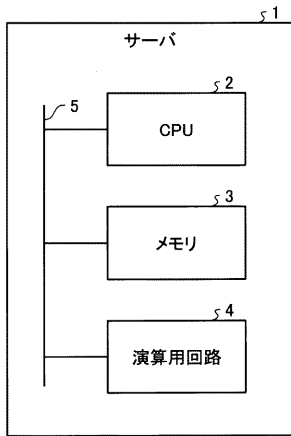
40

50

- 1 0 1 小数点位置決定部
- 1 0 2 指数値変換制御部
- 1 1 1 スカラレジスタファイル
- 1 1 2 ベクタレジスタファイル
- 1 1 3 アキュムレータレジスタ
- 1 1 4 ベクタアキュムレータレジスタ
- 1 1 5 統計情報格納部
- 1 2 1 スカラユニット
- 1 2 2 ベクタユニット
- 2 1 1 , 2 2 1 演算器
- 2 1 2 , 2 2 2 統計情報取得部
- 2 1 3 , 2 2 3 データ変換部

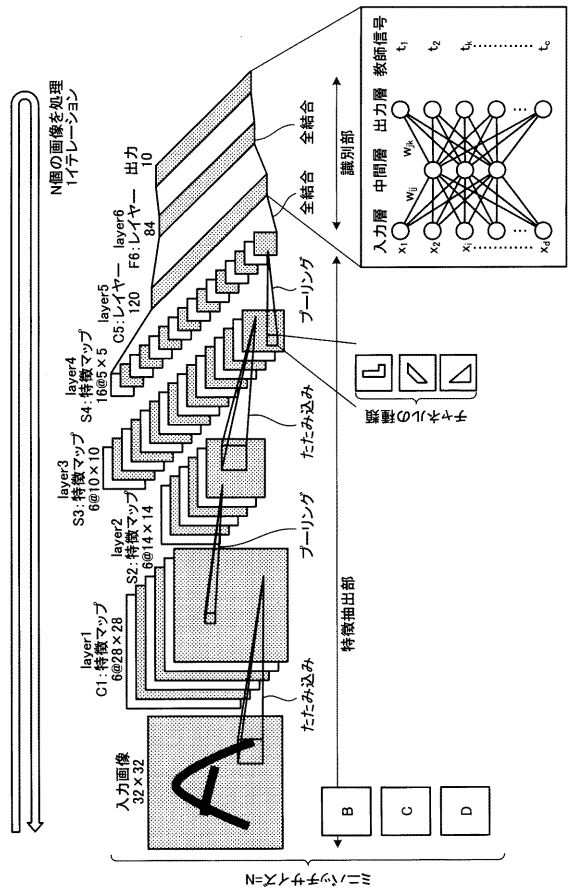
【 図 1 】

サーバの概略を表す構成図

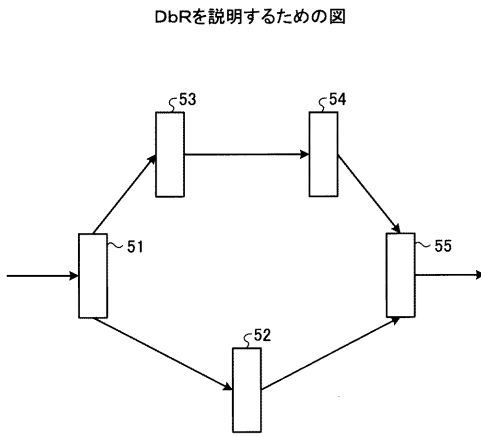


【 図 2 】

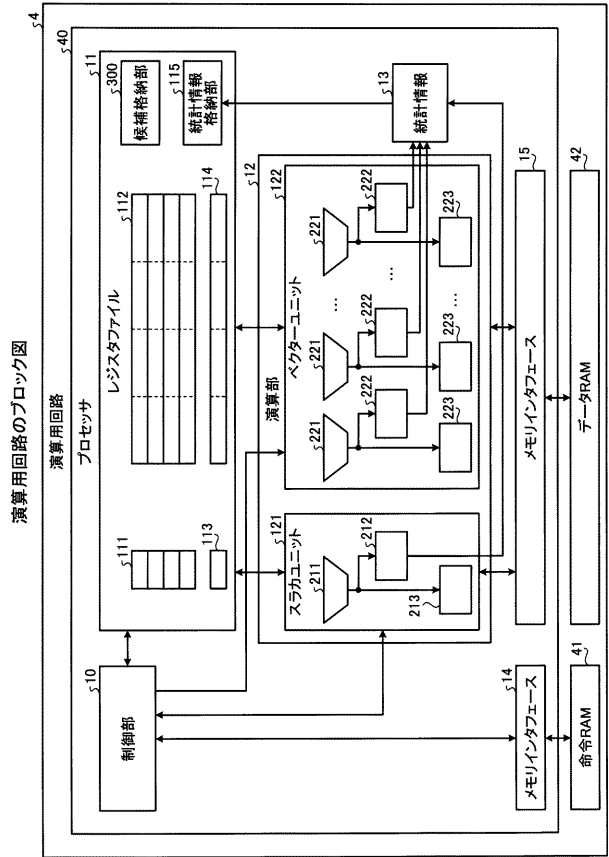
ニューラルネットワークにおける深層学習の一例の図



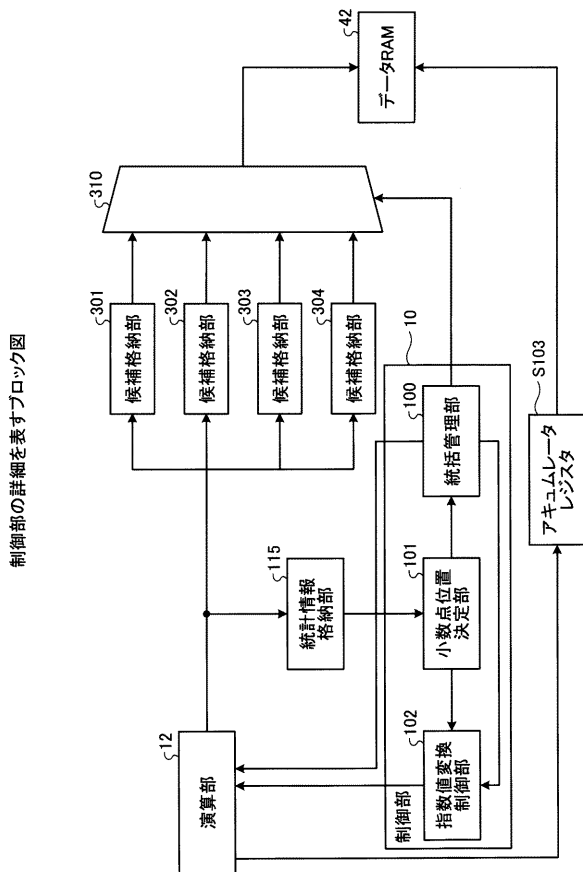
【 図 3 】



【 図 4 】

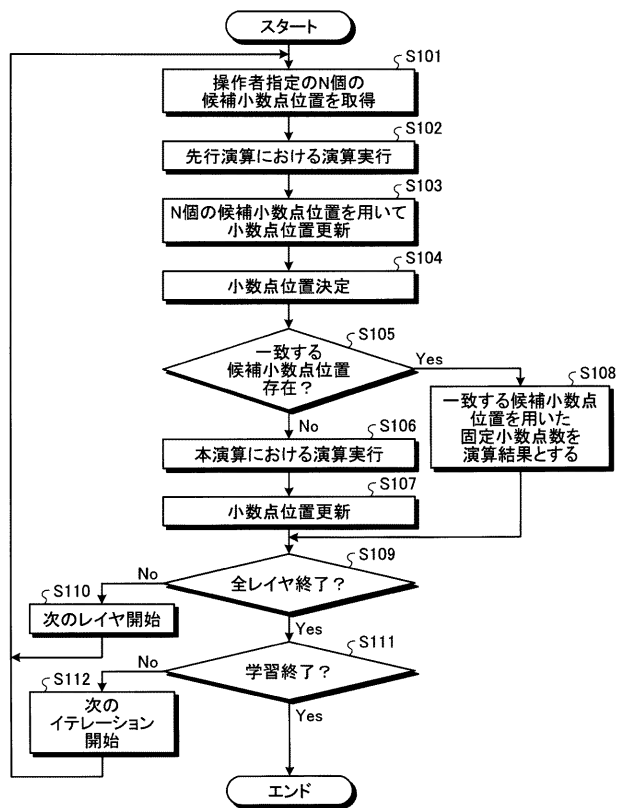


【 図 5 】



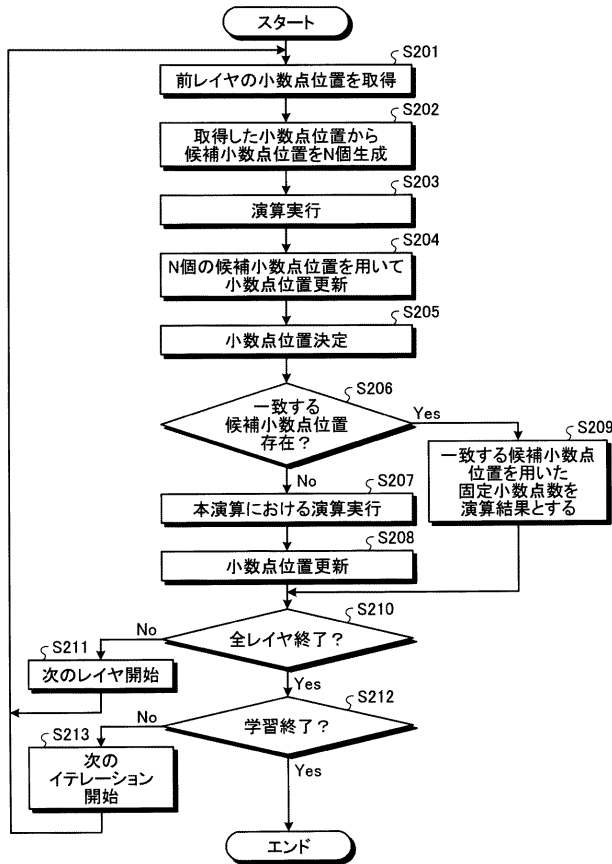
【 図 6 】

実施例1に係る演算回路による深層学習の処理のフローチャート



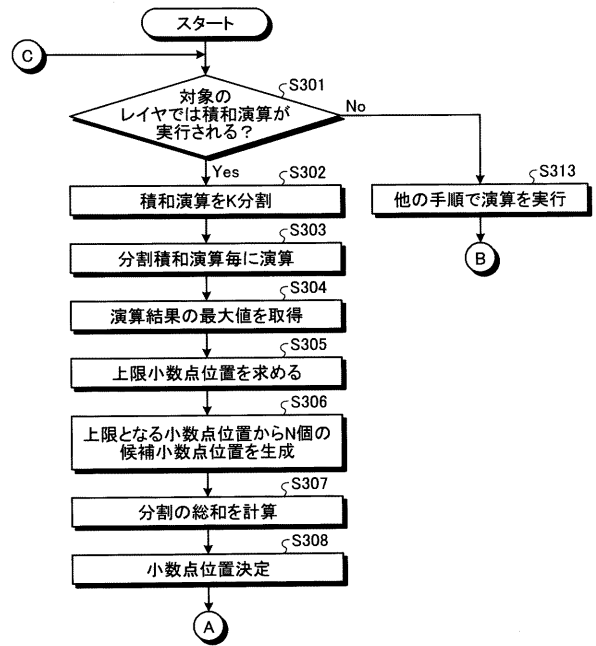
【 図 7 】

実施例2に係る演算用回路による深層学習の処理のフローチャート



【 図 8 A 】

実施例3に係る演算用回路による深層学習の処理のフローチャート



【 図 8 B 】

実施例3に係る演算用回路による深層学習の処理のフローチャート

