

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2018-18220
(P2018-18220A)

(43) 公開日 平成30年2月1日(2018. 2. 1)

(51) Int. Cl.		F I		テーマコード (参考)
G06N	3/10	(2006.01)	G06N 3/10	
G06F	9/50	(2006.01)	G06F 9/46	4 6 5 C
G06F	9/48	(2006.01)	G06F 9/46	4 5 7

審査請求 未請求 請求項の数 10 O L (全 35 頁)

(21) 出願番号	特願2016-146731 (P2016-146731)	(71) 出願人	000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番1号
(22) 出願日	平成28年7月26日 (2016. 7. 26)	(74) 代理人	100113608 弁理士 平川 明
		(74) 代理人	100105407 弁理士 高田 大輔
		(72) 発明者	山崎 雅文 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	田原 司睦 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

最終頁に続く

(54) 【発明の名称】 並列情報処理装置、情報処理方法、およびプログラム

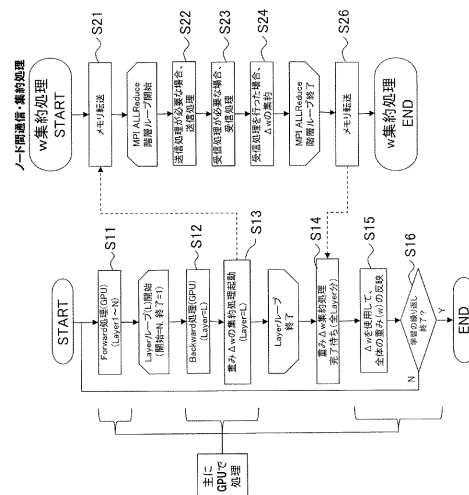
(57) 【要約】

【課題】 ノード間並列による深層学習において、係数演算に使用する係数の勾配情報を深層学習に反映する処理の時間を短縮する。

【解決手段】 各ノードの演算部は、処理対象のデータに対する係数による演算処理を実行し、演算処理の結果を基に係数の変化量を算出し、算出した係数の変化量を処理部に転送するとともに、係数の変化量を並列情報処理装置内の他のノードとの間で授受する処理の実行を処理部に要求する。各ノードの処理部は、演算部から転送された係数の変化量を並列情報処理装置の他のノードに送信するとともに他のノードで算出された係数の変化量を受信する通信処理と、自ノードの演算部から転送された係数の変化量と他のノードで算出された係数の変化量とを積算する集約処理とを実行する。

そして、演算部および処理部の少なくとも一方が積算された係数の変化量を基に次回以降の演算処理で使用される係数を更新する。

【選択図】 図7



【特許請求の範囲】**【請求項 1】**

演算部と処理部とを有するノードを複数備えた並列情報処理装置において、
それぞれのノードの演算部は、処理対象のデータに対する係数による演算処理を実行し、前記演算処理の結果を基に前記係数の変化量を算出し、算出した係数の変化量を前記処理部に転送するとともに、前記係数の変化量を前記並列情報処理装置内の他のノードとの間で授受する処理の実行を前記処理部に要求し、

前記それぞれのノードの処理部は、前記演算部から転送された係数の変化量を前記並列情報処理装置の他のノードに送信するとともに前記他のノードで算出された係数の変化量を受信する通信処理と、前記演算部から転送された係数の変化量と前記他のノードで算出された係数の変化量とを積算する集約処理とを実行し、

前記演算部および前記処理部の少なくとも一方が前記積算された係数の変化量を基に次回以降の演算処理で使用される係数を更新する並列情報処理装置。

10

【請求項 2】

前記演算処理は、所定順序で実行される複数階層の層別処理を有し、それぞれの階層の層別処理はそれぞれの階層の前の階層から入力されるデータに前記係数による演算を実行して次の階層に出力する処理であり、

前記演算部は、それぞれの階層での前記層別処理の結果を基に前記それぞれの階層での前記係数の変化量を算出し、前記算出した係数の変化量を前記処理部に転送し、

前記処理部は、前記それぞれの階層での前記係数の変化量に対する前記集約処理を 2 以上並列に実行する請求項 1 に記載の並列情報処理装置。

20

【請求項 3】

前記係数は、前記複数階層のそれぞれにおいて複数使用され、係数列を形成し、

前記演算部は、前記複数階層のそれぞれの前記係数列を複数の部分列に分割して部分列ごとに前記変化量を前記処理部に転送するとともに、前記部分列ごとに前記授受する処理の実行を前記処理部に要求する請求項 2 に記載の並列情報処理装置。

【請求項 4】

前記演算部は、前記複数階層のうち、前記演算処理の実行の順序が早い階層の係数の変化量を優先して前記処理部に転送し、前記授受する処理の実行を要求する請求項 2 または 3 に記載の並列情報処理装置。

30

【請求項 5】

前記処理部は、前記複数階層のうち前記演算処理の実行の順序が早い階層の係数を優先して前記演算部に前記次回以降の演算処理で使用される係数を更新させる請求項 2 から 4 のいずれか 1 項に記載の並列情報処理装置。

【請求項 6】

前記演算部は、前記複数階層の層別処理を前記所定順に繰り返して実行し、前記複数階層のうち、実行の順序が先の階層で使用される係数に対して前記積算された変化量を基に前記次回以降の演算処理で使用される係数が更新された場合には、実行の順序が後の階層で使用される係数に対する前記積算された変化量の反映を待たないで、次の演算処理における前記実行順が先の階層の層別処理を開始する請求項 2 から 5 のいずれか 1 項に記載の並列情報処理装置。

40

【請求項 7】

前記演算処理と前記積算された変化量を基に前記次回以降の演算処理で使用される係数を更新する処理とが複数回繰り返して実行される場合に、前記演算部は、現在の演算処理による変化量を基に前記次回以降の演算処理で使用される係数が更新される前に次の演算処理を開始し、前記現在の演算処理による変化量を基に次々回の演算処理で使用される係数が更新される請求項 2 から 6 のいずれか 1 項に記載の並列情報処理装置。

【請求項 8】

前記係数を格納するための 2 組以上の記憶部を有し、

前記演算部は第 1 の記憶部に格納した係数を用いて前記演算処理を実行し、前記演算処

50

理による変化量を基に更新した係数を第 2 の記憶部に格納する請求項 1 から 7 のいずれか 1 項に記載の並列情報処理装置。

【請求項 9】

演算部と処理部とを有するノードを複数備えた並列情報処理装置における情報処理方法であって、

それぞれのノードの演算部は、処理対象のデータに対する係数による演算処理を実行し、前記演算処理の結果を基に前記係数の変化量を算出し、算出した係数の変化量を前記処理部に転送するとともに、前記係数の変化量を前記並列情報処理装置内の他のノードとの間で授受する処理の実行を前記処理部に要求し、

前記それぞれのノードの処理部は、前記演算部から転送された係数の変化量を前記並列情報処理装置の他のノードに送信するとともに前記他のノードで算出された係数の変化量を受信する通信処理と、前記演算部から転送された係数の変化量と前記他のノードで算出された係数の変化量とを積算する集約処理とを実行し、

前記演算部および前記処理部の少なくとも一方が前記積算された係数の変化量を基に次回以降の演算処理で使用される係数を更新する情報処理方法。

【請求項 10】

演算部と処理部とを有するノードを複数備えた並列情報処理装置に実行させるためのプログラムであり、

それぞれのノードの演算部に、処理対象のデータに対する係数による演算処理を実行し、前記演算処理の結果を基に前記係数の変化量を算出し、算出した係数の変化量を前記処理部に転送するとともに、前記係数の変化量を前記並列情報処理装置内の他のノードとの間で授受する処理の実行を前記処理部に要求することを実行させるプログラムと、

前記それぞれのノードの処理部に、前記演算部から転送された係数の変化量を前記並列情報処理装置の他のノードに送信するとともに前記他のノードで算出された係数の変化量を受信する通信処理と、前記演算部から転送された係数の変化量と前記他のノードで算出された係数の変化量とを積算する集約処理とを実行させるプログラムとを含み、

前記演算部および前記処理部の少なくとも一方に前記積算された係数の変化量を基に次回以降の演算処理で使用される係数を更新させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、並列情報処理装置、情報処理方法、およびプログラムに関する。

【背景技術】

【0002】

近年、深層学習 (deep learning、DL) に関する研究が盛んである。例えば、画像や音声、文章等に対し、その内容の認識・理解といった研究領域が例示される。このような研究領域の具体的な応用 (アプリケーション) としては、携帯電話における通信時の音声認識、ネットワーク上の検索、大量のログ情報からの異常検出、さらには自動運転などが例示される。このような応用のプロジェクトは、実際に動き始めており、今後もさらに広い分野への応用が進むと考えられる。

【0003】

ところで、深層学習が導入されたシステムにおいては、学習処理は、膨大なデータを繰り返し学習させる手法が例示される。したがって、この学習処理には膨大な計算量が費やされる。例えば、画像識別等の分野では、学習用としてラベル付きの百万枚を超える静止画を、繰り返し学習する。このために、GPU のような積和演算のような学習処理で多用される演算を高速に演算可能な演算部品 (以下、演算部品) を利用したシステム、あるいは、演算部品を含むノードを複数組み合わせたクラスタ環境が利用される。すなわち、学習処理には、GPU のような演算部品の利用は有効であり、さらに複数の演算部品で処理を分散し実行させることで処理の高速化が可能である。複数の演算部品で処理を分散し実行させる方法としては、ノード内並列と、ノード間並列が考えられる。

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2010-020445号公報

【特許文献2】特開2012-022558号公報

【特許文献3】特開2005-182785号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

上述のように、深層学習についてはこれまで、ノード内に複数のGPU等の演算部品を実装し、並列に処理を行うことで、ノード内並列による高速化が行われてきた。一方、演算部品が実装されているノードを複数組み合わせたノード間並列による成果は少ない。

10

【0006】

これまでノード間並列による成果が少ない理由としては、ノード数が増加するに従いノード間をまたいだ深層学習において、係数演算に使用する係数情報のノード間での集約処理、および集約された結果を深層学習に反映する処理に時間がかかることが想定できる。つまり、ノード数の増加による演算能力の向上が実行速度の増加に十分に寄与していないことが想定できる。

【0007】

深層学習では、処理対象のデータに対する係数による演算処理と、演算処理の結果を係数に反映する処理が繰り返し実行される。そこで、1つの側面では、本実施の形態は、演算部品が実装されたノードを組み合わせることで係数演算を並列に実行するとき、係数演算に使用する係数情報のノード間での処理の時間を短縮することを目的とする。

20

【課題を解決するための手段】

【0008】

本発明の一側面は、並列情報処理装置によって例示される。すなわち、本並列情報処理装置は、演算部と処理部とを有するノードを複数備える。それぞれのノードの演算部は、処理対象のデータに対する係数による演算処理を実行し、演算処理の結果を基に係数の変化量を算出し、算出した係数の変化量を処理部に転送するとともに、係数の変化量を並列情報処理装置内の他のノードとの間で授受する処理の実行を処理部に要求する。

30

【0009】

それぞれのノードの処理部は、演算部から転送された係数の変化量を並列情報処理装置の他のノードに送信するとともに他のノードで算出された係数の変化量を受信する通信処理と、自ノードの演算部から転送された係数の変化量と他のノードで算出された係数の変化量とを積算する集約処理とを実行する。

そして、演算部および処理部の少なくとも一方が積算された係数の変化量を基に次回以降の演算処理で使用される係数を更新する。

【発明の効果】

【0010】

本並列情報処理装置によれば、演算部品が実装されたノードを組み合わせることで係数演算を並列に実行するとき、係数演算に使用する係数情報のノード間での処理の時間を短縮することができる。

40

【図面の簡単な説明】

【0011】

【図1】ニューラルネットワークの処理を例示する図である。

【図2】フォワード方向の処理とバックワード方向の処理を例示する図である。

【図3】並列情報処理装置の構成図を例示する図である。

【図4】比較例による処理を示す図である。

【図5】比較例による処理を例示するタイムチャートである。

【図6】実施形態1の処理を例示するタイムチャートである。

50

- 【図 7】実施形態 1 の計算ノードの処理を例示するフローチャートである。
- 【図 8】実施形態 1 の計算ノードにおけるデータフローを例示する図である。
- 【図 9】実施形態 2 の計算ノードの処理を例示するフローチャートである。
- 【図 10】実施形態 2 の計算ノードにおけるデータフローを例示する図である。
- 【図 11】実施形態 3 の処理を例示するタイムチャートである。
- 【図 12】実施形態 3 の計算ノードの処理を例示するフローチャートである。
- 【図 13】分割重みの反映処理を起動する処理の詳細を例示するフローチャートである。
- 【図 14】キュー情報を例示する図である。
- 【図 15】実施形態 4 の処理を例示するタイムチャートである。
- 【図 16】学習処理後のメモリ転送において、層 1、2 が層 3 よりも優先される処理例のタイムチャートである。
- 【図 17】実施形態 4 の学習処理を例示するフローチャートである。
- 【図 18】実施形態 4 の処理の起動を例示するフローチャートである。
- 【図 19】実施形態 5 の処理のタイムチャートを実施形態 4 と対比して例示する図である。
- 【図 20】実施形態 5 における学習処理結果を集約する集約処理を例示するフローチャートである。
- 【図 21】実施形態 6 のタイムチャートを実施形態 4 と対比して例示する図である。
- 【図 22】実施形態 6 における集約処理および反映処理を例示するフローチャートである。

10

20

【発明を実施するための形態】

【0012】

以下、図面を参照して一実施形態に係る並列情報処理装置について説明する。

< 深層学習の処理例 >

【0013】

図 1 に、ニューラルネットワークの処理を例示する。ニューラルネットワークは、画像を認識し、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（バックワードプロパゲーションともいう）を実行する。

【0014】

図 1 のニューラルネットワークは、入力画像に対して、畳み込み演算を実行する畳み込み層（Convolution Layer）の処理と、サブサンプリング層（sub sampling Layer）の処理とを実行し、画像の特徴を抽出し、画像を識別する。すなわち、図 1 では、フォワード方向の処理が例示されている。

30

【0015】

フォワード方向の処理は、入力画像に対して、畳み込み層の処理と、サブサンプリング層の処理を繰り返し実行する特徴抽出部の処理と、識別結果を出力する識別部の処理を含む。特徴抽出部は、入力画像に対して、畳み込み層の処理と、サブサンプリング層の処理を繰り返し実行することで、間引かれた画像を抽出する。畳み込み層の処理は、畳み込み演算ともいう。畳み込み演算は、例えば、 N 個 \times N 個の画素を有する画像の情報（第 $N - 1$ 層）に対して、例えば、 $a \times b$ 個の重み $w_{a,b}$ ($a, b = 0, \dots, m - 1$) のフィルタによる畳み込み演算を実行することで、次の層（第 N 層）の画像の情報を作る。サブサンプリング層の処理は、画像間引き処理であり、プーリング演算ともいう。

40

【0016】

畳み込み層およびサブサンプリング層での演算の入力画像および演算の出力画像はフィーチャマップとも呼ばれる。図 1 の例では、フィーチャマップは、例えば、画像のチャンネル数、あるいは、RGB 等の色に対応して 1 つのニューロン層で複数作成されている。

【0017】

図 2 に、フォワード方向の認識処理および識別処理とともに、バックワード方向の処理を例示する。本実施の形態では、フォワード方向の処理とバックワード方向の処理を

50

合わせて学習処理と呼ぶ。図2のニューラルネットワークにおいても、フォワード方向の認識処理は、入力画像に畳み込み演算を実行する畳み込み層、間引き処理を実行するサブサンプリング層（図2ではpoolingと記載）によって実行される。また、識別結果を出力する識別処理は、全結合層（図2では、Fully connectedと記載）によって実行される。フォワード方向の畳み込み層とサブサンプリング層とを1つのニューロン層という。また、フォワード方向の全結合層も1つのニューロン層ということができる。

【0018】

フォワード方向の処理の結果は、正解値と比較され、比較結果である差分値がエラーとして出力される。エラーは、バックワード方向に各ニューロン層によって処理される。バックワード方向の処理は、全結合層のエラーから、バックワード方向に順次、各ニューロン層でのエラーの評価関数（ERROR）および各ニューロン層での次の重みを計算する処理である。図2では、現在の重みとして、畳み込み層（1層）における1つの重み w_i と、全結合層（1層）における1つの重み w_j が例示されている。また、次の重みとして、畳み込み層（1層）における1つの重み w_{i+1} と、全結合層（1層）における1つの重み w_{j+1} が例示されている。

【0019】

勾配降下法によるニューラルネットワークの学習処理においては、エラーの評価関数（ERROR）の勾配と、学習係数イータの積が重み w の変化量（例えば、現在の重み w_t と次の重み w_{t+1} の差分値）となる。すなわち、深層学習においては、フォワード方向に各ニューロン層の処理が実行され、バックワード方向に、各ニューロン層でのエラーの評価関数（ERROR）が伝搬される。そして、各ニューロン層は、バックワード方向に伝搬するエラーの評価関数（ERROR）から、エラーの評価関数（ERROR）の勾配を求める。そして、各ニューロン層は、エラーの評価関数（ERROR）が小さくする方向でのエラーの評価関数（ERROR）の勾配と、学習係数イータの積から重み w_t の変化量（勾配情報ともいう）を算出し、次の重み w_{t+1} を求める。ここで、現在の重みを w_t で表し、次の演算で使用される重みを w_{t+1} で表した。また、図1で説明したように、学習処理において、重み w は1以上の成分を有する係数列（ベクトル）である。

【0020】

このようにして、バックワード方向に順次、各ニューロン層で、エラーの評価関数（ERROR）を小さくする方向に重みを変化させるための変化量が求められる。そして、バックワード方向に順次伝搬するエラーの評価関数（ERROR）と重み w の変化量が算出され、最終的に、入力層に最も近い層の重み w の変化量が算出される。重み w_t の変化量は、各層において、次の重み w_{t+1} に反映され、次の学習処理に使用される。なお、以下の説明においては、並列演算処理装置における学習処理の時間の短縮について説明するが、学習処理自体のアルゴリズムの詳細は省略する。

<構成>

【0021】

図3に、並列情報処理装置1の構成図を例示する。並列情報処理装置1は、計算ノード10-1、10-2、10-3、10-4等を有する。計算ノード10-1、10-2、10-3、10-4等は、ノード間高速ネットワーク20で接続される。以下、計算ノード10-1等を総称する場合には、単に計算ノード10という。本実施の形態において計算ノード10の数に限定がある訳ではない。並列情報処理装置1は、本実施形態の情報処理方法を実行する。

【0022】

計算ノード10は、Central Processing Unit（CPU11）とメモリ12とGraphics Processing Unit（GPU13）とメモリ14を有する。CPU11とGPU13とはバス15によって接続される。さらにバス15を介して、CPU11とGPU13とは、ノード間インターフェース（ノード間IF16）に接続される。計算ノード10はノードの一例である。

10

20

30

40

50

【0023】

CPU11は、メモリ12に実行可能に展開されたコンピュータプログラムにしたがって、計算ノード10の処理、例えば、他の計算ノード10との通信処理、あるいは、GPU13を制御し、管理する処理を実行する。CPU11は、MPU(Microprocessor)、プロセッサとも呼ばれる。CPU11は、単一のプロセッサに限定される訳ではなく、マルチプロセッサ構成であってもよい。また、単一のソケットで接続される単一のCPU11がマルチコア構成を有していても良い。上記CPU11の少なくとも一部の処理は、CPU11以外のプロセッサ、例えば、GPU13で実行されてもよい。CPU11は、処理部の一例である。メモリ12は、CPU11で実行されるコンピュータプログラム、CPU11が処理するデータを格納する。

10

【0024】

GPU13は、例えば、高速のVRAM、高速の演算器を複数搭載し、積和演算機能等を高速に実行する。GPU13は、メモリ14に実行可能に展開されたコンピュータプログラムにしたがって、計算ノード10の処理のうち、例えば、学習処理を実行する。GPU13は、演算部の一例である。メモリ14は、GPU13で実行されるコンピュータプログラム、GPU13が処理するデータを格納する。

【0025】

上記CPU11およびGPU13の少なくとも一部の処理は、例えば、Digital Signal Processor(DSP)、数値演算プロセッサ、ベクトルプロセッサ、画像処理プロセッサ等の専用プロセッサで行われても良い。また、上記各部の少なくとも一部の処理は、集積回路(IC)、その他のデジタル回路で実行されてもよい。また、上記各部の少なくとも一部にアナログ回路が含まれても良い。集積回路は、LSI, Application Specific Integrated Circuit(ASIC), プログラマブルロジックデバイス(PLD)を含む。PLDは、例えば、Field-Programmable Gate Array(FPGA)を含む。

20

【0026】

すなわち、CPU11あるいはGPU13の処理の少なくとも一部は、プロセッサと集積回路との組み合わせであっても良い。組み合わせは、例えば、マイクロコントローラ(MCU), SoC(System-on-a-chip), システムLSI, チップセットなどと呼ばれる。

30

【0027】

BUS15は、CPU11およびGPU13の例えば内部バスに接続され、CPU11およびGPU13を相互に接続する。また、BUS15は、CPU11およびGPU13をノード間IF16に接続する。BUS15は、例えば、PCI-Expressの規格に従うバスである。

【0028】

ノード間IF16は、ノード間高速ネットワーク20を介して計算ノード10同士を接続するインターフェースである。ノード間高速ネットワーク20は、例えば、クロスバー、インターコネクタ等とも呼ばれる。なお、ノード間高速ネットワーク20は、どのようなネットワーク構成であってもよい。例えば、ノード間高速ネットワーク20は、トラス構造のメッシュであってもよいし、Local Area Network(LAN)のようなバス型のネットワークであってもよい。

40

<複数ノードによる学習処理>

【0029】

学習処理では、まず、フォワード方向の処理が、各ニューロン層に対して、それぞれのニューロン層が持つ重みパラメータ(w)を用いて、バッチ単位で実行され、次に、バックワード方向の処理が各ニューロン層に対して順次実行される。ここで、バッチ単位とは、学習処理の対象をまとめた処理の単位である。例えば、ニューラルネットワークが画像の認識を行う場合に、バッチ単位として、数十枚から数千枚分の画像のデータが学習処理に用いられ、画像の認識と、正解判定が繰り返し実行される。

50

【0030】

図3に例示した複数の計算ノード10がバッチ内の画像データを分担して処理することで、学習処理が並列に実行される。一度のバッチ単位での学習処理の結果としては、重みパラメータ(w)の変化量(Δw)が算出される。図1で述べたように、重みパラメータ(w)は、1以上の成分を有するベクトルである。以下、重みパラメータ(w)は、単に重み(w)ともいう。上述のように、重み(w)の変化量(Δw)は、エラーの評価関数(ERRO)を小さくする方向に算出される。各計算ノード10は、次のバッチ処理に向けて、自身のバッチ単位での重み(w)の変化量(Δw)の計算結果と、他の計算ノード10でのバッチ単位での重み(w)の変化量(Δw)の計算結果と相互に授受し、相互の計算結果を積算する。重み(w)の変化量(Δw)の計算ノード10相互の積算処理を集約処理ともいう。そして、各計算ノード10は、相互の計算結果を集約処理した変化量(Δw)を用いて、重み(w)の更新処理を行う。各層の重み(w)を集約処理された変化量(Δw)を用いて更新することを、集約処理された変化量(Δw)を重み(w)に反映する、ともいう。

10

【0031】

3以上のノード数の計算ノード10が相互に計算結果を授受する場合、計算ノード10の1対1の通信が複数回実行される。例えば、計算ノード10-1、10-2、10-3、10-4が相互に情報をバタフライ方式(Recursive Doubling)で授受する場合、まず、1回目の授受で、計算ノード10-1と計算ノード10-2が情報を授受し、計算ノード10-3と計算ノード10-4が情報を授受する。次に、2回目の授受で、計算ノード10-1と計算ノード10-3が情報を授受し、計算ノード10-2と計算ノード10-4が情報を授受する。以上の2回の情報の授受によって、計算ノード10-1、10-2、10-3、10-4の間での情報の授受が完了する。

20

【0032】

本実施の形態で、ノード間通信アルゴリズムはRecursive Doublingに限定される訳ではない。例えば、ノード間通信アルゴリズムとして、Reduce+Broadcast(Bcast)、Reduce_scatter+Allgather等の方式を用いてもよい。このようなノード間通信処理は、MPI AllReduce処理として、コンピュータプログラムが提供されている。なお、以下の実施の形態説明では、MPI AllReduce処理が実装された計算ノード10を用いて説明するが、計算ノード10間の通信処理がMPI AllReduce処理に限定される訳ではない。また、計算ノード10間の通信処理が実行されるネットワーク構成に限定がある訳ではなく、どのようなネットワーク構成が用いられてもよい。

30

<比較例>

【0033】

比較例では、図2に例示したニューラルネットワークに含まれる各ニューロン層(例えば、ニューロン層1からN)が1つの計算ノード10内に構築される。つまり、比較例では、各ニューロン層の処理は、計算ノード10のコンピュータプログラムによって実行される。なお、以下の説明で用いる図中には、ニューロン層NをLayer Nのように記述する。

40

【0034】

図4に、比較例による処理を示す。比較例では、それぞれの計算ノード10が図2に例示したフォワード処理およびバックワード処理を実行する。また、比較例では、計算ノード10は、フォワード方向の処理をすべてのニューロン層(ニューロン層1からN)において順次実行する(S301)。次に、計算ノード10は、バックワード方向の処理をすべてのニューロン層(ニューロン層Nから1)において順次実行する(S302)。

【0035】

各計算ノード10は、各ニューロン層1~Nにおける重み(w)の変化量(Δw)を相互に転送し、相互に転送した演算結果(各ニューロン層1~Nにおける重みwの変化量 Δw)を積算する。上述のように、それぞれの計算ノード10において計算された演算結果

50

をそれぞれの計算ノード10で積算することを集約するともいう(S303)。そして、各計算ノードは、集約した各ニューロン層1~Nにおける重み(w)の変化量(Δw)を各層の重み(w)に反映する(S304)。そして、計算ノード10は、学習処理の繰り返しを終了するか否かを判定する(S305)。ここで、計算ノード10は、未学習のバッチが存在する場合には、処理をS301に戻し、次のバッチでの学習処理を実行する(S305でNO)。一方、計算ノード10は、すべてのバッチで学習した場合には、処理を終了する(S305でYES)。

【0036】

図5は、比較例による処理を例示するタイムチャートである。図5では、比較のため、単一ノードでの処理も例示されている。図5の左側に例示したように、単一ノードでの処理は、バッチ単位での学習処理、重み(w)の更新処理、バッチ単位での学習処理の繰り返しとなる。

10

【0037】

一方、図5の右側に例示したように、複数ノードでは、バッチ単位での学習処理が、計算ノード10の数だけ並列で実行可能である。しかしながら、それぞれの計算ノード10は、バッチ単位での学習処理が終了すると、重み(w)の変化量(Δw)をノード間通信で授受し、集約した後に、それぞれの計算ノード10での重み(w)を更新することになる。したがって、比較例の処理では、計算ノード10の数が増加しても、ノード間通信・集約処理、更新処理の時間が増加し、計算ノード数の増加による学習処理の時間短縮効果が十分に発揮されない結果となる。

20

<実施形態1>

【0038】

図6は、実施形態1の処理を例示するタイムチャートである。ところで、計算ノード10の構成要素のうち、GPU13は、グラフィックス処理で用いる積和演算を高速に実行する。したがって、GPU13は、学習処理で主体となる重み(w)による演算を高速に実行可能である。しかしながら、学習処理、ノード間通信・集約処理、反映処理を演算部が主体となって処理すると、処理手順としては、図4のフローチャートと同様であり、重み(w)の変化量(Δw)をノード間通信で授受し、集約処理、反映処理を実行する時間が無視できない。

【0039】

30

そこで、実施形態1の並列情報処理装置1は、演算部(GPU13)と処理部(CPU11)を備えた計算ノード10を複数備え、学習処理を演算部(GPU13)で行い、ノード間通信、集約処理、反映処理は処理部(CPU11)で行う。

(1)学習処理

【0040】

学習処理は、主にGPU13で実行される。学習処理は、ニューロン層毎にフォワード処理とバックワード処理(ニューロン層の処理の順番はフォワード処理の逆)を順に行う。複数の計算ノード10がバッチ内の画像データを分担して処理することで、並列に学習処理が実行される。図6では、ニューロン層として、ニューロン層1(LAYER1)から4(LAYER4)が例示されている。ニューロン層1から4は、複数階層の一例である。各ニューロン層1から4におけるフォワード処理およびバックワード処理は、層別処理の一例である。また、各ニューロン層1から4におけるフォワード処理およびバックワード処理は、それぞれの階層の前の階層から入力されるデータに係数による演算を実行して次の階層に出力する処理の一例である。フォワード処理がニューロン層1から4の順に実行され、バックワード処理がニューロン層4から1の順に実行されることは、所定順序の一例である。

40

(2)メモリ転送(GPU13からCPU11への転送)

【0041】

演算部(GPU13)は、学習処理の各ニューロン層で計算された重み(w)の変化量(Δw)を学習処理が終わったニューロン層ごとに順次、処理部(CPU11)へメモリ

50

転送する。これによって、演算部（GPU13）は、ニューロン層ごとに、ノード間通信・集約処理、反映処理を処理部（CPU11）に開始させる。ニューロン層ごとにノード間通信・集約処理、反映処理を開始することで、次のバッチ単位での学習処理の開始を早め、高速化が実現される。

【0042】

具体的には、各計算ノード10において各層のバックワード処理が終わる毎に、演算部（GPU13）に割り当てられた学習処理用のスレッドはメモリ転送を起動するためのキューを発行する。キューは要求と呼ぶこともできる。メモリ転送（GPU13からCPU11への転送）用処理スレッドは、キューを受けると転送対象のデータをGPU13からCPU11へ転送し、最後に集約処理のキューをCPU11に発行する。図6では、ニューロン層として、ニューロン層4（LAYER4）から層1（LAYER1）のバックワード処理で、重みの変化量としてそれぞれ、 $WL4-1$ 、 $WL3$ 、 $WL2$ 、 $WL1$ が算出されている。

10

（3）集約処理および（4）ノード間通信

【0043】

予め、指定数（1個から数十個）が用意されている集約処理用スレッドは、キューを受けると、まず、ノード間通信処理のためのキューを発行する。ノード間通信処理用スレッドは、ノード間通信処理のためのキューを受けるとノード間通信の Message Passing Interface（MPI）リクエストを、ノンブロッキング通信を指定してMPI通信プログラムに投入する。リクエストに対応する通信が完了した時点で、MPI通信プログラムから集約処理用スレッドへ通信完了が通知され、集約処理用スレッドにしたがい集約処理が実行される。集約処理には多数回の演算が実行されるため、集約処理は複数のスレッドを並列で実行することで高速化を実現する。すなわち、計算ノード10に複数のCPU11が搭載される場合には、複数のスレッドを並列で実行することで、CPU11による並列処理が実行される。また、単一のCPU11がマルチコアを有する場合も同様である。

20

【0044】

図6では、第1回目のノード間通信において、例えば、ニューロン層4（LAYER4）については、ノード間通信用スレッドは、 $WL4-1$ を他ノードに送信し、 $WL4-2$ を他ノードから受信する。そして、集約処理用のスレッド1は、 $WL4-1$ と $WL4-2$ を積算し、集約処理を実行する。集約処理によって $WL4-1 + WL4-2$ が得られる。

30

【0045】

次に、第1回目のノード間通信において、例えば、ニューロン層4（LAYER4）については、ノード間通信用スレッドは、 $WL4-1 + WL4-2$ を他ノードに送信し、 $WL4-3 + WL4-4$ を他ノードから受信する。そして、集約処理のスレッド1は、 $WL4-1 + WL4-2$ と $WL4-3 + WL4-4$ を積算し、集約処理を実行する。図6のスレッド1から3は、一例として、それぞれの階層での係数の変化量に対する集約処理を2以上並列に実行する。

（5）メモリ転送（CPU11からGPU13への転送）

40

【0046】

他の全ノードと情報を授受するための回数分のノード間通信と集約処理が完了すると、CPU11は、メモリ転送（CPU11からGPU13への転送）処理のキューを発行する。メモリ転送処理用スレッドがキューを受けて、メモリ転送（CPU11からGPU13への転送）を実行する。

（6）反映処理

【0047】

各層のメモリ転送（CPU11からGPU13への転送）が完了すると、主にGPU13側での反映処理が、メモリ転送が完了したニューロン層から順に実行される。

【0048】

50

図7は、実施形態1の計算ノード10の処理を例示するフローチャートである。図で左側のフローチャートは、主としてGPU13が実行する学習処理と反映処理を例示する。また、右側のフローチャートは、主としてCPU11が実行するノード間通信・集約処理を例示する。図7の処理では、まず、GPU11がニューロン層（例えば、ニューロン層1からN）について、フォワード処理を実行する（S11）。

【0049】

フォワード処理は、図1に例示したように、入力データと重み（ w ）とによる演算処理である。演算処理は、例えば、入力データの要素 $x(i, j)$ と $a \times b$ 個の重み $w_{a,b}$ （ $a, b = 0, \dots, m-1$ ）のフィルタによる畳み込み演算、サブサンプリング層のプーリング演算、全結合層の演算等である。S11の処理は、処理対象のデータに対する係数による演算処理の一例である。

10

【0050】

次に、GPU13は、バックワード方向にニューロン層Nから1のループ（LAYERループ（L）、開始=N、終了=1）の中で、S12、S13の処理を実行する。S12の処理では、GPU13は、バックワード方向に各ニューロン層（L）において上位の層（L+1）でのエラーの評価関数（ERROR）から当該ニューロン層（L）におけるエラーの評価関数（ERROR）を求める。そして、GPU13は、当該ニューロン層（L）のエラーの評価関数（ERROR）に基づいて当該ニューロン層（L）のエラーの評価関数（ERROR）を減少させる方向の重み（ w ）の変化量（ Δw ）を求める。S12の処理は、演算処理の結果を基に係数の変化量を算出することの一例である。S12の処理は、それぞれの階層での層別処理の結果を基にそれぞれの階層での係数の変化量を算出することの一例でもある。

20

【0051】

また、S13の処理は、CPU11に対して重みの変化量（ Δw ）の集約処理の起動を要求する処理である。S13の処理により、GPU13は、S12で求めた当該ニューロン層（L）について計算された重み（ w ）の変化量（ Δw ）をCPU11にメモリ転送するとともに、集約処理を実行するCPU11のスレッドにキューを登録する（S13）。したがって、実施形態1では、各ニューロン層（L）でバックワード処理が終了するごとに、CPU11に対して、重み（ w ）の変化量（ Δw ）の集約処理の起動が要求される。S13の処理は、算出した係数の変化量を処理部に転送するとともに、係数の変化量を並列情報処理装置内の他のノードとの間で授受する処理の実行を処理部に要求することの一例である。S13の処理は、算出した係数の変化量を処理部に転送することの一例でもある。

30

【0052】

以降、GPU13は、CPU11からの重み（ w ）の変化量（ Δw ）の集約処理の完了を全ニューロン層数待つ（S14）。そして、CPU11で集約処理された各ニューロン層（L）の重み（ w ）の変化量（ Δw ）が、CPU11からGPU13にメモリ転送される。そして、全レイヤの集約処理が完了すると、GPU13は、集約処理された変化量（ Δw ）を各層の重み（ w ）に反映する（S15）。すなわち、GPU13は、次のバッチのフォワード処理およびバックワード処理で使用される各層の重み（ w ）を更新する。S15の処理は、演算部が積算された係数の変化量を基に次回以降の演算処理で使用される係数を更新することの一例である。

40

【0053】

そして、GPU13は、学習の終わりか否かを判定する（S16）。学習の終わりとは、例えば、計算ノード10について用意されたすべてのバッチが終了する場合である。計算ノード10について用意された未学習のバッチが残っている場合には、GPU113は、処理をS11に戻し、次のバッチを実行する。

【0054】

S13の処理によって、集約処理の起動が要求されると、キューがCPU11のスレッドに登録され、キューが順次処理される。CPU11は、まず、メモリ転送を実行し、G

50

GPU 13で計算されたニューロン層Lの重み(w)の変化量(Δw)を取得する(S 2 1)。そして、ニューロン層Lの重み(w)の変化量(Δw)を他の計算ノード10との間で授受する。上述のように、本実施の形態では、ノード間のデータ交換の処理として、MPI仕様のAllReduceアルゴリズムが用いられる。ただし、本実施の形態のノード間のデータ交換の処理がAllReduceアルゴリズムに限定される訳ではない。図7において、CPU 11は、MPI AllReduceの階層ループにおいて、S 2 2からS 2 4の処理を繰り返し実行する。

【0055】

例えば、ノード数が4であって(計算ノード10-1~10-4)、Recursive Doublingの場合には、以下の処理が実行される。計算ノード10-1と10-2の組と計算ノード10-3と10-4の組のそれぞれ組でCPU 11が、S 2 2からS 2 4の処理を実行する。すなわち、自ノードで計算されたニューロン層Lの重み(w)の変化量(Δw)を相手ノードに送信する(S 2 2)。S 2 2の処理は、演算部から転送された係数の変化量を並列情報処理装置の他のノードに送信することの一例である。

10

【0056】

また、CPU 11は相手ノードで計算されたニューロン層Lの重み(w)の変化量(Δw)を受信する(S 2 3)。S 2 3の処理は、他のノードで算出された係数の変化量を受信することの一例である。したがって、S 2 2とS 2 3の処理は通信処理の一例である。

【0057】

そして、CPU 11は自ノードで計算されたニューロン層Lの重み(w)の変化量(Δw)と相手ノードで計算されたニューロン層Lの重み(w)の変化量(Δw)を積算する(S 2 4)。S 2 4の処理は、演算部から転送された係数の変化量と他のノードで算出された係数の変化量とを積算する集約処理の一例である。

20

【0058】

さらに、計算ノード10-1と10-3の組と計算ノード10-2と10-4の組のそれぞれの組でCPU 11が、S 2 2からS 2 4の処理を実行する。この処理によって、計算ノード10-1~10-4の間でニューロン層Lの重み(w)の変化量(Δw)が集約される。ニューロン層Lの重み(w)の変化量(Δw)が集約されると、CPU 11は、集約されたニューロン層Lの重み(w)の変化量(Δw)をメモリ転送し、GPU 13に戻す(S 2 6)。計算ノード10は、S 2 1からS 2 6の処理をキューの蓄積順にすべてのニューロン層Lについて繰り返し実行する。

30

【0059】

図8に、実施形態1の計算ノード10におけるデータフローを例示する。計算ノード10においては、まず、GPU 13による学習処理では、GPU 13による演算結果がGPU 13のメモリ14に格納される(矢印A 1)。上述のように演算結果は、ニューロン層Lの重み(w)の変化量(Δw)である。

【0060】

次に、ノード間通信処理が実行される。まず、GPU 13とCPU 11との間のメモリ転送が実行され、メモリ14に格納されたニューロン層Lの重み(w)の変化量(Δw)がCPU 11のメモリ12へ転送される(矢印A 2-1)。ここでは、メモリ12に格納された重み(w)の変化量を w_1 とする。そして、ノード間IFを介して、メモリ12に格納された重み(w)の変化量(w_1)が他の計算ノード10に送信される(矢印A 2-2)。一方、ノード間IFを介して、他の計算ノード10で計算されたニューロン層Lの重み(w)の変化量(w_2)が当該計算ノード10で受信される(矢印A 2-3)。

40

【0061】

さらに集約処理が実行される(矢印A 3)。集約処理では、CPU 11は、メモリ12のデータ(変化量 w_1 と w_2)を加算する。ここでは、加算結果は集約された重みの変化量として w_2 に保持されるとする。ノード数が3以上の場合には、矢印A 2-2からA 3がノード間通信のアルゴリズムで実行される回数だけ繰り返される。

50

【 0 0 6 2 】

そして、GPU 1 1 は、メモリ転送で GPU 1 3 に集約されたニューロン層 L の重み (w) の変化量 (Δw) を転送する (矢印 A 5 - 1)。転送先の GPU 1 3 は、転送された重みの変化量を変化量 (Δw) に保存する。そして、GPU 1 3 は、集約された層 L の重み (w) の変化量 (Δw) を使用して、重み (w) を更新する (A 5 - 2)。

【 0 0 6 3 】

以上述べたように、実施形態 1 の並列情報処理装置 1 は、複数の計算ノード 1 0 がそれぞれのバッチによって、入力データに対する重み (w) の演算を複数のニューロン層について実行するため、重み (w) の学習処理が並列に実行される。そして、並列に実行された学習処理によって得られた重み (w) の変化量 (Δw) を複数の計算ノード 1 0 間で集約し、各ニューロン層についてのすべての計算ノード 1 0 のバッチの結果を反映した重み (w) を各計算ノード 1 0 が取得する。

10

【 0 0 6 4 】

このような処理において、各計算ノード 1 0 は、GPU 1 3 が各ニューロン層の学習処理を順次実行する。すなわち、GPU 1 3 は、フォワード方向にニューロン層 1 からニューロン層 N に対して、重み (w) による演算を実行する。次に、GPU 1 3 は、バックワード方向にニューロン層 N からニューロン層 1 について、各ニューロン層 L の重み (w) の変化量 (Δw) を計算する処理を実行する。そして、各ニューロン層 L の重み (w) の変化量 (Δw) の計算が終了するごとに、GPU 1 3 は、計算した変化量 (Δw) を CPU 1 1 にメモリ転送するとともに、集約処理のキューを CPU 1 1 のスレッドに発行し、集約処理を依頼する。

20

【 0 0 6 5 】

以上述べたように、積和演算等の重み (w) による演算を高速に実行可能な GPU 1 3 が学習処理を複数の計算ノード 1 0 において並列に実行し、CPU 1 1 が重みの変化量 (Δw) のメモリ転送、ノード間通信、および集約処理を実行する。したがって、GPU 1 3 は、CPU 1 1 との連携により、もっぱら学習処理を実行すればよく、GPU 1 3 の演算性能が発揮されやすい。

【 0 0 6 6 】

また、CPU 1 1 は、集約処理の依頼を受けると、キューの順に、ノード間通信を実行する。例えば、CPU 1 1 は、ALLReduce アルゴリズムにより、自ノードで計算した重み (w) の変化量 (Δw) を他の計算ノード 1 0 に送信するとともに、他の計算ノード 1 0 で得られた計算結果を受信する。そして、CPU 1 1 は、ニューロン層ごとに、順次重み (w) の変化量 (Δw) を集約する。したがって、比較例で例示した図 4 のようにバックワード方向の処理がすべてのニューロン層について完了した後に重み (w) の変化量 (Δw) の集約処理を実行する場合と比較して、各層の集約処理が早期に開始される。例えば、CPU 1 1 がマルチコア構成の場合には、図 6 のように、集約処理を複数のスレッドに分けて、異なるニューロン層の集約処理を割り当てることで、複数のニューロン層の集約処理が並列して実行される。

30

【 0 0 6 7 】

また、あるニューロン層 L の集約処理を実行中に、他のニューロン層 L + 1 のノード間通信が並列して実行可能である。また、メモリ転送用のスレッドがニューロン層 L の集約処理結果を GPU 1 3 にメモリ転送中に、集約処理の複数のスレッドが複数のレイヤ L + 1、L + 2、L + 3 に対して集約処理とノード間通信処理とを並列に実行できる。図 5 に例示した比較例では、バッチ単位で全ニューロン層について学習処理を実行し、全ニューロン層について集約処理を実行し、全ニューロン層について次の学習処理を実行する。このような比較例の処理に対して、実施形態 1 の計算ノード 1 0 は、少なくとも集約処理の処理時間が短縮される。また、次のバッチにおけるフォワード方向の処理の開始を早めることができる。

40

< 実施形態 2 >

【 0 0 6 8 】

50

図9および図10により実施形態2に係る並列情報処理装置1について説明する。実施形態2の並列情報処理装置1は、図6に例示した「(6)反映処理」をニューロン層単位でCPU11が実行する。そして、CPU11は、ニューロン層単位での反映処理の後に、(5)メモリ転送(CPU11からGPU13)を実行する。実施形態2の他の構成および作用は実施形態1と同様である。そこで、実施形態2の並列情報処理装置1の構成要素のうち、実施形態1と同一の構成要素については、同一の符号を付してその説明を省略する。

【0069】

図9に、実施形態2の計算ノード10の処理を例示するフローチャートである。図9の処理は、変化量(w)を重み(w)に反映する処理がGPU13ではなく、CPU11によって実行される点で図7と相違する。例えば、図9では、ノード間通信・集約処理において、S25の処理が追加されている。

10

【0070】

まず、GPU13は、学習処理によって計算された変化量(w)を重み(w)に反映する処理を起動する(S13A)。このとき、メモリ転送処理によってGPU13からCPU11に当該ニューロン層の重み(w)の変化量(w)が送信される点は図7と同様である。すると、GPU13は、キューの優先順に変化量(w)のメモリ転送(S21)、および集約処理を実行する(S22-S24)。そして、MPI ALL Reduce階層ループが終了すると、CPU11は、集約処理されたあるニューロン層Lの重みの変化量(w)を重み(w)に反映する(S25)。S25の処理は、処理部が積算された係数の変化量を基に次回以降の演算処理で使用される係数を更新することの一例である。

20

【0071】

そして、CPU11は、変化量(w)が反映された重み(w)をメモリ転送でGPU13に送信する(S26A)。すると、GPU13は、変化量(w)が反映された重み(w)をメモリ転送で受信し、メモリ14に保存する(S14A)。そして、GPU13は、未学習のバッチが残っている場合には(S16でN)、次のバッチの学習を実行する。

【0072】

図10に、実施形態2の計算ノード10におけるデータフローを例示する。図10の処理は、学習処理(矢印A1)、ノード間通信処理(A2-2、A2-3)、集約処理(矢印A3)までは、図8と同一である。ただし、ノード間通信処理前のメモリ転送(矢印A2-1)において、CPU11は、GPU13から重みの変化量(w)とともに重み(w)を受信し、 w_1 としてメモリ12に格納する。

30

【0073】

そして、CPU11は、重みの変化量(w)の集約処理の後、集約された重みの変化量(w)を重み w に反映し、重み w_1 としてメモリ12に格納する(矢印A5-3)。そして、CPU11は、重みの変化量(w)が反映された重み(w_1)をメモリ転送でGPUに転送し、重み(w)としてメモリ14に保存する(矢印A5-4)。

【0074】

以上述べたように、実施形態2では、CPU11が変化量(w)を重み(w)に反映する処理を実行する。この構成および手順により、GPU13は重みの変化量(w)の演算により専念することが可能となる。また反映処理のスレッドは、集約処理と同様にCPU11のコア数に応じて並列処理することで、学習処理の高速処理が可能となる。

40

<実施形態3>

【0075】

図11から図13により実施形態3の並列情報処理装置1について説明する。上記実施形態1では、CPU11が学習結果のノード間通信・集約処理を実行する際に、各ニューロン層単位で処理を分割した。すなわち、CPU11は、1つのニューロン層について学習結果のノード間通信・集約処理を個別に実行し、それぞれのニューロン層の重みの変化

50

量 (w) が集約されるごとに、GPU 13 にメモリ転送した。また、実施形態 2 では、CPU 11 が重みの変化量 (w) を重み (w) に反映し、GPU 13 にメモリ転送した。しかし、実施形態 1、2 の処理でも、1 つのニューロン層が大きなパラメータ数の重みを持つ場合には転送処理に時間がかかり、マルチコアの CPU 11 が複数スレッドによって並列処理を実行する構成を有していても、並列化の効果が発揮されない場合がある。そこで、実施形態 3 では、GPU 13 および CPU 11 は、ノード間通信スレッド、複数の集約処理スレッド、および反映処理スレッドの実行単位をニューロン層単位よりも細かく分割して処理する。このような手順により、計算ノード 10 は、各処理をパイプライン化し、高速化する。

【0076】

例えば、あるニューロン層 L の重み (w) が $w = (p_1, p_2, \dots, p_X)$ のようなパラメータ列であるとする。パラメータ列は、係数列の一例である。つまり、ニューロン層 L の重み (w) は、複数使用され、係数列を形成する。そして、学習処理の結果、重みの変化量は $w = (p_1, p_2, \dots, p_X)$ のような多数のパラメータ列として計算されるとする。このような場合に、GPU 13 は w を部分列に区切り、 $w_1 = (p_1, p_2, \dots, p_{X1})$ 、 $w_2 = (p_{X1+1}, \dots, p_{X2})$ 、 $w_3 = (p_{X2+1}, \dots, p_{X3})$ 、 \dots 、 $w_x = (p_{X1}, \dots, p_X)$ のように分割する。

【0077】

図 11 は、実施形態 3 の処理を例示するタイムチャートである。なお、図 11 では、実施形態 3 の処理が適用される前のタイムチャート(「適用前」)が、実施形態 3 の処理が適用された場合のタイムチャートとともに例示されている。適用前の例(図 11 の上側)では、ニューロン層 N に対するバックワード処理の終了後、GPU 13 から CPU 11 へのメモリ転送が実行され、その後、スレッド 1 による集約処理が 2 回のノード間データ通信(例えば、ALLReduce アルゴリズム)とともに実行されている。

【0078】

一方、適用後の例(図 11 の下側)では、ニューロン層 N に対するバックワード処理の終了後、GPU 13 は、学習処理で計算した重みの変化量 (w 、パラメータ列)を w_1 、 w_2 、 w_3 、 w_4 の部分列に分割し、CPU 11 にメモリ転送する。

【0079】

CPU 11 はメモリ転送で分割された変化量 w_1 、 w_2 、 w_3 、 w_4 を取得し、集約処理用のスレッド 1 から 3 により順次集約処理を起動する。例えば、スレッド 1 が分割された変化量 (w_1) を受け取ると、まず、ノード間通信処理のスレッドを起動する。ノード間通信処理のスレッドは、分割された変化量 (w_1) を他の計算ノード 10 - 2 に送信するとともに、計算ノード 10 - 2 からニューロン層 N の分割された変化量 w_1 を受信する。今、自ノードと他ノードとで変化量 w_1 を区別するため、自ノードで計算されたものを $w_1 - 1$ とし、計算ノード 10 - 2 で計算されたものを $w_1 - 2$ とする。スレッド 1 は、自ノードで計算され、分割された変化量 ($w_1 - 1$) と、ノード間通信処理で得られた他ノードで計算された変化量 ($w_1 - 2$) とを積算し、計算ノード 10 - 2 との間で集約処理を実行する。このとき、スレッド 1 の集約処理と並行してスレッド 2 は、分割された変化量 (w_2) について、ノード間通信処理のスレッドを起動しており、スレッド 2 もスレッド 1 と同様に、ノード間通信処理と集約処理をパイプラインで実行する。スレッド 3 も、スレッド 1、2 と同様に、ノード間通信処理と集約処理をパイプラインで実行する。

【0080】

スレッド 1 は、自ノードで計算された重みの変化量 ($w_1 - 1$) と他ノード計算された重みの変化量 ($w_1 - 2$) との間の集約処理が完了すると、再びノード間通信処理のスレッドを起動し、計算ノード 10 - 3 との間で、集約処理を実行する。また、スレッド 2、3 についても、1 回目の集約処理が終了すると、スレッド 1 と同様に、再びノード間通信処理のスレッドを起動し、計算ノード 10 - 3 との間で、集約処理を実行する。

10

20

30

40

50

【0081】

そして、例えば、スレッド1が分割された変化量 (w_1) について、他のすべての計算ノード10との間で集約処理を完了すると、メモリ転送スレッドを起動する。メモリ転送スレッドにより、CPU11は、集約された変化量 (w_1) をGPU13に転送する。スレッド2、スレッド3も同様である。

【0082】

また、スレッド1は、分割された変化量 (w_1) についてメモリ転送スレッドのキューを発行すると、分割された次の変化量 (w_4) について、分割された変化量 (w_1) と同様の処理を実行する。このようにして、例えば、CPU11が複数、例えば、5つのコアを有している場合には、CPU11はスレッド1から3およびメモリ転送スレッド、およびノード間通信スレッドを並行して実行できる。したがって、例えば、ある分割された変化量 (w_k) についてのノード間通信の時間が、別の分割された変化量 (w_j) についての集約処理の時間に実行できる。また、仮に、あるニューロン層Lの重み (w_L) のパラメータ数が他の層よりも多いものであっても、GPU13およびCPU11は重み (w_L) に含まれるパラメータを複数部分に分割し、複数スレッドで並行して処理できる。

【0083】

図12は、実施形態3の計算ノード10の処理を例示するフローチャートである。図12の処理は、反映処理の起動と反映処理待ちにおいて、図9の処理と相違する。すなわち、実施形態3においては、図11で説明したように、GPU13は、ニューロン層のループにおいて(ニューロン層1からN)、各ニューロン層Lの重みの変化量 (w_L) を複数の部分に分割する (w_{Lk} 、kは分割された部分列に対応する数)。そして、GPU13は、メモリ転送を行い、各部分列ごとに集約処理、反映処理を起動する(S13B)。そして、ニューロン層のループの終了後、GPU13は、分割された重みの変化量 (w_{Lk}) の反映処理の完了待ちとなる(S14B)。そして、すべてのニューロン層のすべての分割された重みの変化量 (w_{Lk}) についての反映処理が終了すると、GPU13は、学習の繰り返しの終了か否かを判定し、未学習のバッチがある場合に、処理をS11に戻し、次のバッチの学習を実行する。

【0084】

なお、図12の処理フローは、図9を変形したもので、CPU11が重みの変化量 (w_{Lk}) を基に重み (w_{Lk}) を更新する反映処理を実行する。しかし、図7に例示したように、CPU11が重みの変化量 (w_{Lk}) をメモリ転送でGPU13に転送し、GPU13が反映処理を実行してもよい。

【0085】

図13は、実施形態3におけるGPU13によって分割重み (w_{Lk}) の反映処理を起動する処理(図12の13A)の詳細を例示するフローチャートである。この処理では、GPU13は、レイヤLの重み (w_L) のk番目の分割重みの部分列 (w_{Lk}) と重みの変化量 (w_{Lk}) のメモリ転送を起動する(S13B1)。S13B1の処理は、係数列を複数の部分列に分割して部分列ごとに変化量を処理部に転送することの一例である。

【0086】

次に、GPU13は、分割された重みの部分列 (w_{Lk}) の変化量 (w_{Lk}) の集約処理、および、重みの部分列 (w_{Lk}) への反映処理をスレッド S_n ($n = 1 \sim N$) のキューに登録する(S13B2)。S13B2の処理は、部分列ごとに授受する処理の実行を処理部に要求することの一例である。

【0087】

以上述べたように、本実施形態の並列情報処理装置1は、複数のスレッドによって、メモリ転送(GPU13からCPU11)、ノード間通信、集約および反映処理、メモリ転送(CPU11からGPU13)を実行できる。さらに、実施形態3では、GPU13はニューロン層Lの重みのパラメータ列 (w_L) を複数の部分列 (w_{Lk} 、 $k = 1, 2, 3, \dots$) に分ける。そして、GPU13はそれぞれの重みの変化量の部分列 (w_{Lk}

10

20

30

40

50

、 $k = 1, 2, 3, \dots$) ごとに、メモリ転送、集約および反映を起動する。すると、CPU 11は、重みの変化量の部分列 (w_{Lk} , $k = 1, 2, 3, \dots$) ごとに、メモリ転送 (GPU 13からCPU 11)、集約および反映、メモリ転送 (CPU 11からGPU 13) を実行する。したがって、ニューロン層の重み (w) に含まれるパラメータ数が多い場合であっても、メモリ転送、ノード間通信、集約処理のパイプラインを形成し、例えば、ノード間通信処理に要する時間 (またはその一部) を集約処理の時間で隠すことができる。なお、重みのパラメータ列 (w_L) は、係数列の一例である。

< 実施形態 4 >

【0088】

図14から図18により実施形態4を説明する。上記実施形態1から実施形態3では、例えば、学習処理の終了順にニューロン層ごとのデータがメモリ転送され、ノード間通信処理、集約処理、反映処理が実行された。実施形態4では、各スレッドは、ニューロン層のうち、最も階層が低い層、すなわち、図2の入力画像が入力される層 (例えば、ニューロン層1)) の優先順位を高くし、階層が上がるほど優先順位が低くなるようにキューの発行を制御する。このような処理によって、1つのバッチが終了前にすでに、階層が低いニューロン層の重み (w) に対して、変化量 (w) が反映されている場合には、階層が低いニューロン層における次のバッチの開始を可能とする。

10

【0089】

図14は、Reduce処理で用いられるキュー情報を例示する図である。キュー情報は、キュー情報を発行する処理 (前処理、キュー情報発行スレッドともいう) から発行され、後続処理 (キュー処理スレッドともいう) によって処理される。図14では、前処理として、処理A-1、処理A-2等が例示されている。また、後続処理として処理B-1、処理B-2が例示されている。

20

【0090】

図14の例では、前処理 (キュー発行スレッド) は、処理が終わる毎に後続処理のキューを登録する。後続処理 (キュー処理スレッド) は、処理が要求されているキューが存在しない場合は何もしない。一方、処理が要求されているキューが存在する場合、後続処理 (キュー処理スレッド) は、要求された処理を実行し、処理が終了すると処理完了フラグ情報を更新する。処理完了フラグ情報は、例えば、完了した処理数 (または未完了の処理数) のカウンタである。なお、ある前処理が、それ以前に実行される前処理 (例えば、処理A-1、処理A-2) に依存する場合には、処理を行う前に、依存する前処理の完了を確認してから処理を開始する。

30

【0091】

以上のようにして、後続処理 (キュー処理スレッド) は、登録されたキューの順に処理を実行する。以下、実施形態4では、登録されるキューの順を所定の優先順序で優先する制御、具体的には、ニューロン層のうち、階層の低いニューロン層を優先して処理を実行する制御手順を例示する。

【0092】

図15は、実施形態4の処理を例示するタイムチャートである。図15では、ニューロン層として、ニューロン層1から4が想定されている。ただし、実施形態4のニューロン層が4つのニューロン層に限定される訳ではない。バックワード方向の処理がニューロン層4から1までの順でそれぞれ終了すると、この終了順にメモリ転送処理が起動され、ノード間通信処理、集約処理が実行される。さらに、各ニューロン層の集約処理が完了後、メモリ転送 (CPU 11からGPU 13) が実行される。

40

【0093】

ところで、図15の例では、ニューロン層1の集約された重みの変化量がCPU 11からGPU 13にメモリ転送可能となったときに、まだ、ニューロン層2についても、集約された変化量のメモリ転送が起動されていない。例えば、ニューロン層2のメモリ転送処理 (CPU 11からGPU 13) は、キューが登録された状態で未実行の状態となっている。実施形態4では、このような場合にニューロン層1の集約処理が終了すると、集約処

50

理用のスレッドは、ニューロン層 2 よりもニューロン層 1 のメモリ転送を優先する。すなわち、CPU 11 の集約処理用のスレッドは、ニューロン層 1 よりもニューロン層 1 が先に転送されるように、ニューロン層 1 の集約された変化量のメモリ転送のキューを登録する。そのようなキュー登録の結果、メモリ転送用スレッドはニューロン層 2 よりもニューロン層 1 の重みの変化量を先にメモリ転送する。

【0094】

図 16 は、学習処理後のメモリ転送において、層 1、2 が層 3 よりも優先される処理例のタイムチャートである。このタイムチャートでは、バックワード方向の処理において、ニューロン層 4 のメモリ転送中に、ニューロン層 3 とニューロン層 2 の学習が完了している。このような場合、階層が入力データに近いニューロン層 2 がニューロン層 3 よりも優先されてメモリ転送が開始される。

10

【0095】

さらに、ニューロン層 2 のメモリ転送中に、ニューロン層 1 の学習処理が完了する。すると、階層が入力データに近いニューロン層 1 がニューロン層 3 よりも優先されてメモリ転送が開始される。その後、ニューロン層 3 のメモリ転送が開始される。

【0096】

入力データが入力されるニューロン層 1 を最も優先し、ニューロン層 1 に近い層の順に優先してメモリ転送を実行することで、その後のノード間通信、集約、反映処理は、ニューロン層 1 を最も優先し、ニューロン層 1 に近い層の順に優先する結果となる。したがって、現在のバッチの学習終了後、次のバッチでは、現在のバッチで学習結果がニューロン層 1 から順に優先して重み w に反映される。したがって、現在のバッチのすべてのニューロン層の処理が完了する前であっても、GPU 13 は次のバッチでニューロン層 1 から学習を開始でき、次のバッチ全体の開始時期が早まる。

20

【0097】

図 15、図 16 のように、階層が低いニューロン層に対する処理の優先順位を高くするため、処理順序の変更は、MPI ALLReduce 階層ループの単位、もしくは実施例 3 における重みのパラメータ細分化後の部分列単位で実行される。各処理スレッドは、次のスレッドへのキューの登録時、通常 First In First Out (FIFO) 方式でキューを登録する。一方、実施形態 4 では、各処理スレッドは、処理順序の変更条件 (キューが優先順でない状態) が検知された場合には、優先順の位置にキューを登録する。

30

【0098】

1 つのノードの処理順序の変更により、処理順序が変更されたノードの処理順序が他ノードの処理順序とずれるとノード間転送がロックするため、計算ノード 10 同士が同期をとる。同期をとる手法としては、処理順序の変更を検知した計算ノード 10 がすべてのノードに処理順序の変更を配信し、各ノードは他ノードでの処理順序の変更に対して、同様に処理の順番を組み直す。

【0099】

図 17 は、実施形態 4 の学習処理を例示するフローチャートである。この処理では、GPU 13 は、ニューロン層 1 から N について、フォワード方向の処理を実行する (S11C)。ただし、S11C の処理は、前のバッチにおける全層についての学習処理が終了していても開始される点で実施形態 1 から 3 と相違する。そして、全層についてのフォワード方向の処理が終了すると、GPU 13 は、バックワード方向にニューロン層 N から 1 のループ (LAYER ループ (L) 開始 = N 、終了 = 1) の中で、S12、S13 の処理を実行する。S12 の処理は、実施形態 1 から 3 と同様である。

40

【0100】

S13 の処理では、GPU 13 は、ニューロン層のうち、入力側に近いニューロン層を優先して、CPU 11 にメモリ転送するとともに、集約処理を実行する CPU 11 のスレッドにキューを登録する (S13C)。S13C の処理は、複数階層のうち、演算処理の実行の順序が早い階層の係数の変化量を優先して処理部に転送することの一例である。

50

【 0 1 0 1 】

したがって、実施形態 1 では、GPU 13 は、各ニューロン層 (L) でバックワード方向の処理が終了するごとに、優先順の制御を実行する。すなわち、GPU 13 は、バックワード方向の処理が終了したニューロン層 (L) より上位のニューロン層 (L + k) で、メモリ転送および集約処理が未実行のニューロン層がキューに残っていないか否かを判定する。そして、バックワード方向の処理が終了したニューロン層 (L) より上位のニューロン層 (L + k) がキューに残っている場合には、GPU 13 は、入力側に近い下位のニューロン層 (L) を優先してキューを登録する。なお、このように、下位のニューロン層を優先するキューの登録は、CPU 11 がノード間通信およびメモリ転送 (CPU 11 から GPU 13) のキューを登録する場合も同様である。

10

【 0 1 0 2 】

そして、GPU 13 は、CPU 11 からの重み (w) の変化量 (Δw) の集約処理の完了を待つ。ただし、実施形態 4 では、GPU 13 は、ニューロン層 1 層ずつ、集約処理の完了を待つ (S 1 4 C)。

【 0 1 0 3 】

その後、CPU 11 で集約処理された各ニューロン層 (L) の重みの変化量 (Δw) が、CPU 11 から GPU 13 にメモリ転送される。あるニューロン層 (L) の集約処理が完了すると、GPU 13 は、当該ニューロン層において集約処理された変化量 (Δw) を重み (w) に反映する (S 1 5 C)。すなわち、GPU 13 は、次のバッチのフォワード処理およびバックワード処理で使用されるニューロン層 (L) の重み (w) を更新する。

20

【 0 1 0 4 】

そして、GPU 13 は、全層の集約処理が完了したか否かを判定する (S 1 6)。全層の集約処理が完了していない場合、GPU 13 は、次のバッチのニューロン層 L のフォワード処理の開始が可能か否かを判定する (S 1 7)。次のバッチのニューロン層 L のフォワード処理の開始が可能でない場合、GPU 13 は、制御を S 1 4 C に戻し、次のニューロン層の集約処理の完了を待つ。

【 0 1 0 5 】

一方、次のバッチのニューロン層 L のフォワード処理の開始が可能である場合、GPU 13 は、次のバッチのニューロン層 L のフォワード処理を開始させる (S 1 8)。S 1 7 の判定で、フォワード処理の開始が可能との判定される場合は、複数階層のうち、実行の順序が先の階層で使用される係数に対して積算された変化量を基に次回以降の演算処理で使用される係数が更新された場合の一例である。S 1 6 から S 1 8 の処理を実行することは、実行の順序が後の階層で使用される係数に対する積算された変化量の反映を待たないで、次の演算処理における実行順が先の階層の層別処理を開始することの一例である。

30

【 0 1 0 6 】

次のバッチのニューロン層 L のフォワード処理の開始が可能である場合とは、例えば、次のバッチのニューロン層 1 について、重みの変化量 (Δw) が集約処理され、重み (w) への反映が完了している場合をいう。また、例えば、次のバッチのニューロン層 1 から L - 1 のフォワード方向の処理が終了し、ニューロン層 L について、重みの変化量 (Δw) が集約処理され、重み (w) への反映が完了している場合をいう。このような場合には、GPU 13 は、現在処理中のバッチについて、全層の処理が終了していても、次のバッチのフォワード方向の処理を開始させる。そして、GPU 13 は、処理を S 1 4 C に戻す。

40

【 0 1 0 7 】

一方、全レイヤの集約処理が完了すると、GPU 13 は、学習の終わりが否かを判定する (S 1 9)。計算ノード 10 について用意された未学習のバッチが残っている場合には、GPU 13 は、処理を S 1 1 C に戻し、次のバッチを実行する。ただし、次のバッチにおけるニューロン層については、フォワード処理は、S 1 8 による処理開始によって、

50

すでに開始されているか、実行が完了しているものがあり得る。したがって、次のバッチでのS11Cの処理は、前のバッチについての全層についての学習処理が終了していなくても開始され、当該バッチでは、未実行のニューロン層から開始される。

【0108】

なお、図17では、反映処理は、S15CでGPU13が実施したが、実施形態2のようにCPU11が反映処理を実行してもよい。また、図17の処理は、ニューロン層ごとに実行されたが、実施形態3のように、ニューロン層の重み w のパラメータ列を部分列に分割し、部分列ごとに実行されるようにしてもよい。

【0109】

図18は、実施形態4の起動処理を例示するフローチャートである。この処理は、学習処理後のメモリ転送(GPU13からCPU11)、CPU11での集約処理、ノード間通信処理、反映処理、集約処理後のメモリ転送(CPU11からGPU13)を起動するときのキューの登録において適用可能である。なお、反映処理自体は、実施形態1のようにGPU13によって実行されてもよいし、実施形態2のようにCPU11によって集約処理とともに実行されてもよい。図18の処理の主体はGPU13またはCPU11である。また、この処理は、図14で説明した前処理(キュー発行スレッド)の処理である。そこで、以下の説明はキュー発行スレッドを主体として説明する。

10

【0110】

キュー発行スレッドは、キュー発行対象のニューロン層と処理対象データを取得する(S41)。例えば、キュー発行スレッドは、キュー発行スレッドの処理が完了したときに、キュー発行対象のニューロン層と処理対象データを取得することになる。

20

【0111】

次に、キュー発行スレッドは、現在登録済みのキューを読む(S42)。そして、キュー発行スレッドは、優先順位変更の要否を判定する(S43)。例えば、現在登録済みのキューのニューロン層がいずれも、キュー発行対象のニューロン層よりも、入力側に近い層(下位の層)があれば(S43でN)、キュー発行スレッドは、最後尾の位置にキュー発行対象のニューロン層のキューを登録する(S44)。

【0112】

一方、例えば、現在登録済みのキューに、キュー発行対象のニューロン層よりも、入力側から遠い層(上位の層)があれば(S43でY)、キュー発行スレッドは、当該上位の層より優先してキュー発行対象のニューロン層のキューを登録する(S45)。S43からS45の処理は、演算部が複数階層のうち、演算処理の実行の順序が早い階層の係数の変化量を優先して処理部に転送することの一例である。S43からS45の処理は、授受する処理の実行を要求することの一例でもある。S43からS45の処理は、処理部が、複数階層のうち演算処理の実行の順序が早い階層の係数を優先して演算部に次回以降の演算処理で使用される係数を更新させることの一例でもある。

30

そして、キュー発行スレッドは、処理順序の変更をMPI AllReduceのアルゴリズムで他の計算ノード10に通知する(S46)。

【0113】

以上述べたように、実施形態4によれば、入力側に近いニューロン層が優先して処理されるように、処理順序が変更される。1つのニューロン層 L の重みのパラメータ列(w_L)が複数の部分列(w_{Lk})に分割されて処理される実施形態3場合も同様である。このような処理順序の変更により、処理順序の変更が実施されたバッチの次のバッチにおいて、入力側に近い、階層の低いニューロン層が優先されて前のバッチの学習結果が重みに反映されることになる。すなわち、次のバッチでの入力データに近いニューロン層で使用される重みの更新を早めることができる。

40

【0114】

そして、S16からS18のように、全層の集約処理が完了していない場合であっても、次のバッチにおいて下位のニューロン層のフォワード処理の開始が可能である場合、GPU13は、次のバッチのニューロン層 L のフォワード処理を開始させる。したがっ

50

て、一部のニューロン層の重みに対して学習結果が反映されていなくても、次のバッチでの入力データに近いニューロン層での学習が早期に開始可能となる。

<実施形態 5 >

【0115】

図19および図20を参照して、実施形態5を説明する。実施形態1から4では、1つのバッチにおいて、学習、集約、ノード間通信、反映処理が完了した後に、次のバッチが開始された。実施形態5においては、現在のバッチ(N番目のバッチ)の学習処理が完了すると、集約、ノード間通信、反映処理が実行される前に、次のバッチ(N+1番目のバッチ)の学習処理が起動される。そして、現在のバッチ(N番目のバッチ)の学習処理の結果は、次のさらに次のバッチ(N+2番目のバッチ)の前に重みに反映される。実施形態5におけるこのような手順以外の手順および構成要素は、実施形態1から4と同様である。そこで、実施形態5の構成要素のうち、実施形態1から4と同一の構成要素については、同一の符号を付してその説明を省略する。

10

【0116】

図19に、実施形態5の処理のタイムチャートを実施形態4と対比して例示する。図19では、上側は実施形態4のタイムチャートであり、下側が実施形態5のタイムチャートである。実施形態5では、ニューロン層1から4までが想定されている。また、フォワード方向のニューロン層1から4の学習処理はF1からF4のラベルで示されている。一方、バックワード方向のニューロン層4から1の学習処理はB4からB1のラベルで示されている。

20

【0117】

図19のように、実施形態5では、N番目の学習処理(バッチ処理(N番目))が終了すると、N-1番目のバッチの学習処理の結果(集約済みの重みの変化量 w)が重み w に反映される。そして、N+1番目のバッチに対する学習処理(バッチ処理(N+1番目))が開始する。図19のように、バッチ処理(N番目)に続いてバッチ処理(N+1番目)の学習処理が実行されることは、演算処理と積算された変化量を基に次回以降の演算処理で使用される係数を更新する処理とが複数回繰り返して実行される場合の一例である。

【0118】

なお、実施形態2で説明したように、N+1番目の学習処理が開始するまでに、N-1番目のバッチによる学習処理の結果が重み w に反映されるようにすればさらに時間は短縮できる。また、N+1番目の各ニューロン層の学習処理の開始までに、N-1番目のバッチによる各層(k)の学習処理の結果(集約済みの $w(Lk)$)が各層の重みに反映されるようにすればさらに時間は短縮できる。なお、実施形態6とは異なり、実施形態5では、重み(w)を格納するバッファが1面だけ使用されるので、GPU13はバッチ処理(N番目)の学習処理後、直ちにバッチ処理(N+1番目)を開始できない。すなわち、GPU13は、バッチ処理(N+1番目)を開始する前に、学習処理の結果(集約済みの $w(Lk)$)が各層の重みに反映する時間を要する。また、実施形態2のように、CPU11が学習処理の結果が各層の重みに反映する場合には、GPU13は、バッチ処理(N+1番目)を開始する前に、学習処理の結果が反映された重みをメモリ14に保持する時間を要する。

30

【0119】

以上の処理の結果、実施形態5では、実施形態4と比較して、学習処理の結果の反映がバッチ1つ分遅れることになる。しかし、学習処理終了時に学習処理の結果を重みに反映しないため、実施形態4と比較して早期に次のバッチを開始できる。すなわち、実施形態4に対して、概ね少なくとも学習処理の結果を集約する時間が節約される。

【0120】

なお、図19の処理は、例えば、図7において、S14、S15の処理を実行しないで、S16において、未処理のバッチの有無を判定し、次のバッチの学習処理を実行することで実行される。図19でGPU13がN+1番目の学習処理が終了すると、N+2番目

40

50

のバッチに対する学習処理を開始することは、演算部が現在の演算処理による変化量を基に次回以降の演算処理で使用される係数が更新される前に次の演算処理を開始することの一例である。

【0121】

図20は、実施形態5におけるCPU11の学習処理結果の集約処理のフローチャートを例示する。図20の集約処理は、例えば、N番目のバッチでの学習処理が終了後に、N+1番目の学習処理と並列に実行される。この処理では、まず、CPU11は、バッチが2番目より後のバッチか否かを判定する(S51)。バッチが1番目または2番目のバッチの場合、CPU11は、処理を終了する。

【0122】

一方、バッチが2番目より後のバッチの場合、CPU11は、メモリ転送を実行し、N番目のバッチでの学習処理の結果を取得する(S52)。そして、メモリ転送したバッチの学習結果である(w)を集約する(S53)。そして、CPU11は、集約された(w)のGPU13へのメモリ転送を起動する(S54)。S54のメモリ転送を受け、N+2番目のバッチに対する学習処理の開始前に、GPU13は、集約された(w)を重み(w)に反映する。S52からS54の処理は、現在の演算処理による変化量を基に次々回の演算処理で使用される係数が更新されることの一例である。

【0123】

なお、変化量(w)の集約と重み(w)への反映は、実施形態2のように、CPU11において行ってもよい。つまり、GPU13は、集約された重み(w)が反映済みの重み(w)をメモリ転送で受け取るようにしてもよい。この場合には、反映処理は、単に、変化量(w)が反映済みの重み(w)をGPU13のメモリ14に保存する処理ということができる。

【0124】

また、メモリ転送(GPU13からCPU11)、変化量(w)の集約処理、ノード間通信、重み(w)への反映処理、および、メモリ転送(CPU11からGPU13)は、実施形態2のように、ニューロン層単位で行ってもよい。また、これらの処理は、実施形態3のように、ニューロン層単位よりも細かく分割したパラメータの部分列の単位で行ってもよい。

【0125】

以上述べたように、実施形態5では、N番目バッチの学習処理が終了すると、N+1番目のバッチに対する学習処理と並行してN番目バッチの学習処理結果の集約処理が実行される。したがって、図19のように、実施形態1から4の場合と比較して、集約処理の時間が短縮される。

【0126】

また、上記集約処理ともに、実施形態2と同様に、CPU11が反映処理を行った場合には、GPU13は、N+1番目のバッチの学習処理が開始するまでに、集約されたwを反映済みの重みをメモリ14に保存する処理を実行すればよい。この場合には、実施形態1から4の場合と比較して、集約処理および反映処理の時間が短縮される。

<実施形態6>

【0127】

図21および図22を参照して実施形態6を説明する。実施形態5では、計算ノード10は、N番目の学習処理の結果をN+2番目のバッチの学習開始までに集約し、重み(w)に反映した。このような処理によって、計算ノード10は、N番目の学習処理の終了後、直ちにN+1番目の学習処理を開始できた。実施形態6では、計算ノード10には、重み(w)を格納するバッファが複数、例えば、2面設けられる。すなわち、計算ノード10は、学習結果である重みの変化量(w)が反映された重み(w)を格納するバッファを2面有することで、実施形態5と同様に、第N番目のバッチが終了した後、直ちに、第N+1番目のバッチの学習処理を開始できる。

【0128】

10

20

30

40

50

図 2 1 に、実施形態 6 のタイムチャートを実施形態 4 と対比して例示する。図 2 1 のように、実施形態 4 では、バッファ $w a$ に格納した重みを用いた学習処理とバッファ $w b$ に格納した重みを用いた学習処理が交互に実行される。例えば、奇数番目のバッチの学習終了後に、次の偶数番目のバッチの学習処理と並行して集約処理と反映処理が実行される。そして、奇数番目のバッチの学習処理の結果である重みの変化量 (w) が反映された重み (w) がバッファ $w a$ に格納される。このとき、偶数番目のバッチの学習処理では、バッファ $w b$ に格納された重みを使用される。

【 0 1 2 9 】

一方、偶数番目のバッチの学習終了後に、次の奇数番目のバッチの学習処理と並行して集約処理と反映処理が実行される。そして、偶数番目のバッチの学習処理の結果である重みの変化量 (w) が反映された重み (w) がバッファ $w b$ に格納される。このとき、奇数番目のバッチ学習処理では、バッファ $w a$ に格納された重みを使用される。

10

【 0 1 3 0 】

したがって、図 2 1 のように、バッファ $w a$ に格納した重みによる第 N 番目のバッチの学習処理終了後、直ちに、バッファ $w b$ に格納した重みによる第 $N + 1$ 番目のバッチの学習処理が開始される。したがって、実施形態 4 の場合と比較して、実施形態 6 では、学習処理終了後の学習処理の結果である重みの変化量 (w) の集約処理と反映処理が次のバッチの学習処理と並行して実施できる。実施形態 6 の場合も、実施形態 5 と同様、第 N 番目のバッチの学習処理の結果を反映した重みは第 $N + 2$ 番目のバッチの学習に使用される。図 2 1 のバッファ $w a$ 、 $w b$ は、係数を格納するための 2 組以上の記憶部の一例である。

20

【 0 1 3 1 】

図 2 2 に、実施形態 6 における集約処理および反映処理のフローチャートを例示する。図 2 2 では、学習処理、集約反映処理、および格納処理の 3 つの処理が連携して実行される。GPU 1 3 が学習処理と格納処理を実行し、CPU 1 1 が集約反映処理を実行する。ここでは、 N 番目のバッチの学習処理が実行されるものとして説明する。

【 0 1 3 2 】

まず、GPU 1 3 は、 N 番目のバッチが奇数番目のバッチか否かを判定する (S 6 0)。 N 番目のバッチが奇数番目のバッチの場合、GPU 1 3 は、バッファ $w a$ に格納した重みによる学習処理を実行する (S 6 1)。一方、 N 番目のバッチが偶数番目のバッチの場合、GPU 1 3 は、バッファ $w b$ に格納した重みによる学習処理を実行する (S 6 2)。S 6 1、S 6 2 の処理は、第 1 の記憶部に格納した係数を用いて演算処理を実行することの一例である。そして、GPU 1 3 は、メモリ転送を CPU 1 1 に要求するとともに、集約反映処理のキューを登録する。そして、GPU 1 3 は当該バッチの学習処理を終了する。そして、GPU 1 3 は $N + 1$ 番目のバッチの学習処理を実行する。

30

【 0 1 3 3 】

CPU 1 1 は、 N 番目のバッチの学習結果である重みの変化量 (w) に対する集約処理と反映処理 (以下、単に集約反映処理) のキューを受け付け、集約反映処理を実行する。CPU 1 1 による集約反映処理は、GPU 1 3 による $N + 1$ 番目のバッチの学習処理と並行して実行される。

40

【 0 1 3 4 】

まず、CPU 1 1 は、GPU 1 3 による学習結果である重みの変化量 (w) をメモリ転送で取得する (S 6 3)。そして、CPU 1 1 は、重みの変化量 (w) を集約し、重み (w) に反映する (S 6 5)。S 6 5 の処理は、実施形態 2 (図 1 2) の S 2 2 から S 2 6 と同様である。そして、CPU 1 1 は、集約した重みの変化量 (w) が反映された重み (w) を GPU 1 3 にメモリ転送する (S 6 6)。

【 0 1 3 5 】

GPU 1 3 は、メモリ転送を受けると、バッチが奇数番目のバッチか否かを判定する (S 6 7)。バッチが奇数番目のバッチの場合、GPU 1 3 は、バッファ $w b$ に重みを格納する (S 6 8)。一方、バッチが偶数番目のバッチの場合、GPU 1 3 は、バッファ $w b$

50

に重みを格納する (S 6 9)。 S 6 8、 S 6 9 の処理は、演算処理による変化量を基に更新した係数を第 2 の記憶部に格納することの一例である。なお、 S 6 7 から S 6 9 の処理は、次のさらに次のバッチ (N + 2 番目のバッチ) の学習処理が開始されるまでに実行される。

【 0 1 3 6 】

以上述べたように、実施形態 6 においては、図 2 1 のように、バッファ w a に格納した重みによる第 N 番目のバッチの学習処理終了後、直ちに、バッファ w b に格納した重みによる第 N + 1 番目のバッチの学習処理が開始できる。

< 記録媒体 >

【 0 1 3 7 】

コンピュータその他の機械、装置 (以下、コンピュータ等) に上記いずれかの機能を実現させるプログラムをコンピュータ等が読み取り可能な記録媒体に記録することができる。そして、コンピュータ等に、この記録媒体のプログラムを読み込ませて実行させることにより、その機能を提供させることができる。

【 0 1 3 8 】

ここで、コンピュータ等が読み取り可能な記録媒体とは、データやプログラム等の情報を電氣的、磁氣的、光学的、機械的、または化学的作用によって蓄積し、コンピュータ等から読み取ることができる記録媒体をいう。このような記録媒体のうちコンピュータ等から取り外し可能なものとしては、例えばフレキシブルディスク、光磁気ディスク、Compact Disc (C D) - Read Only Memory (R O M)、C D - Recordable (R)、Digital Versatile Disk (D V D)、ブルーレイディスク、Digital Audio Tape (D A T)、8 m m テープ、フラッシュメモリなどのメモ리카ード等がある。また、コンピュータ等に固定された記録媒体としてハードディスク、R O M (リードオンリーメモリ) 等がある。さらに、Solid State Drive (S S D) は、コンピュータ等から取り外し可能な記録媒体としても、コンピュータ等に固定された記録媒体としても利用可能である。

【 符号の説明 】

【 0 1 3 9 】

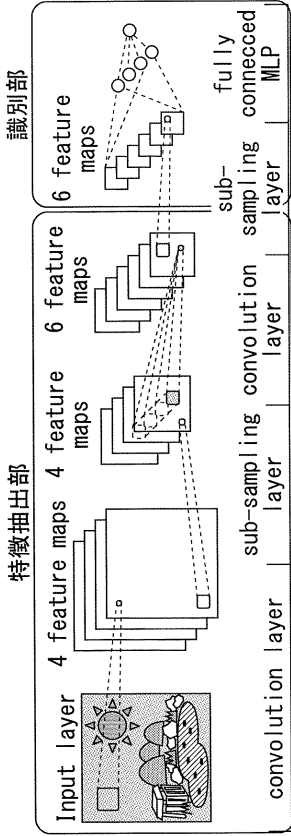
- 1 並列情報処理装置
- 1 0 計算ノード
- 1 1 C P U
- 1 2、1 4 メモリ
- 1 3 G P U
- 1 4 バス
- 1 5 ノード間インターフェース

10

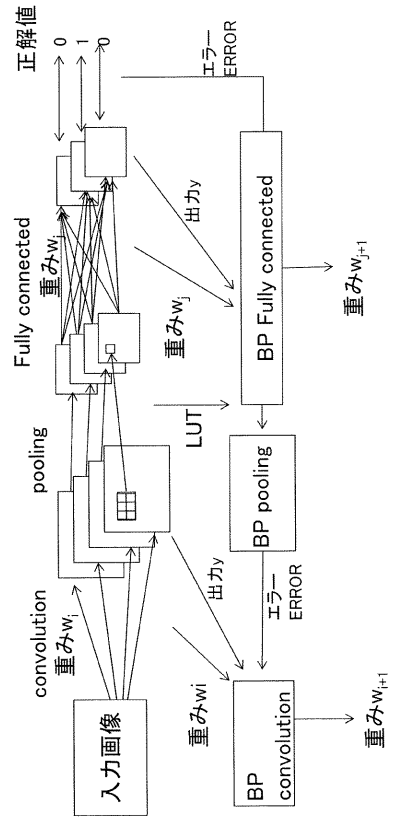
20

30

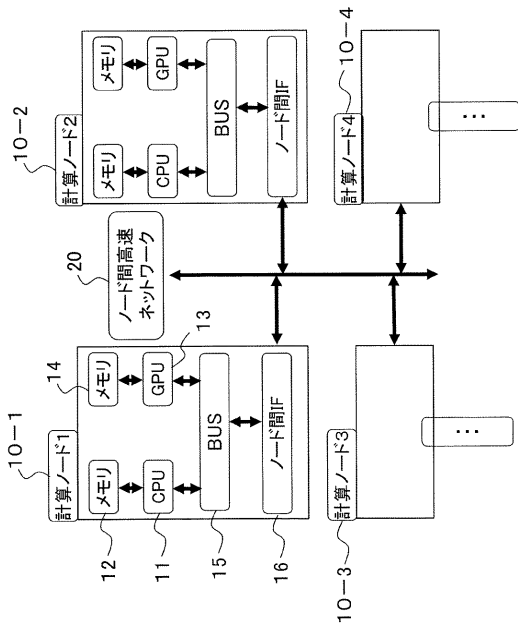
【 図 1 】



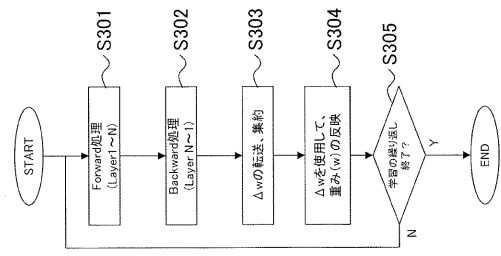
【 図 2 】



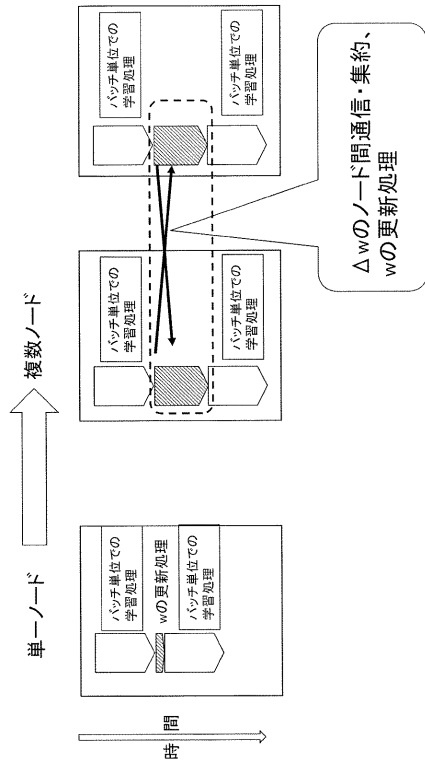
【 図 3 】



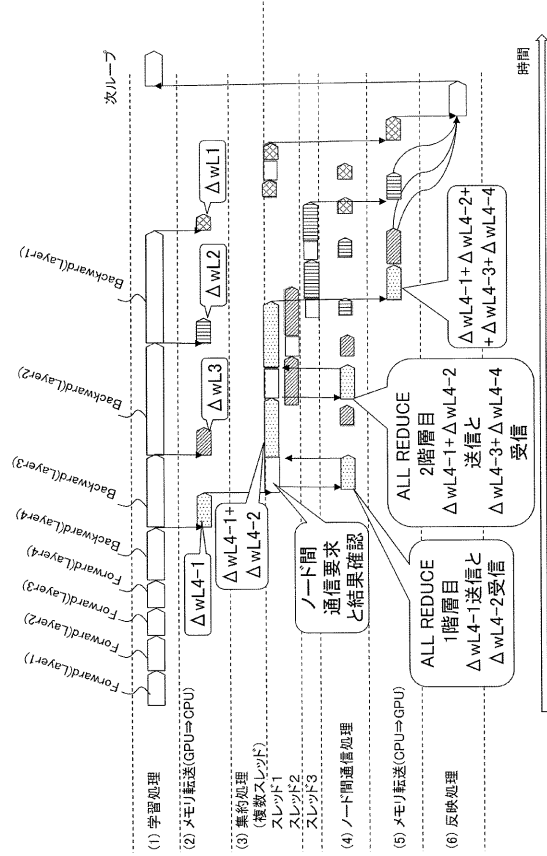
【 図 4 】



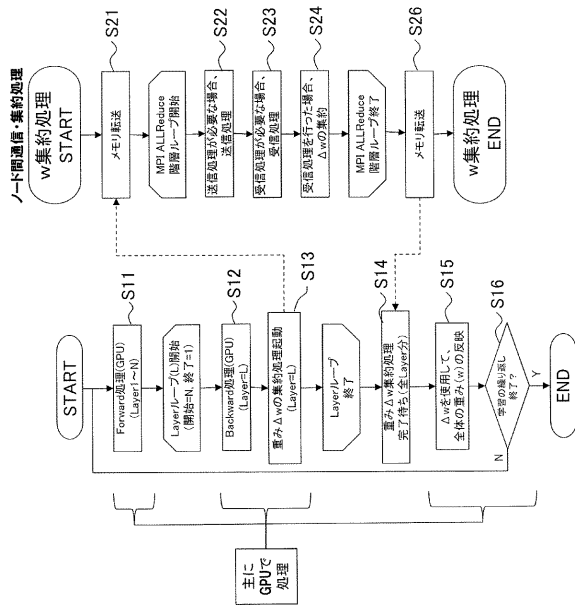
【図5】



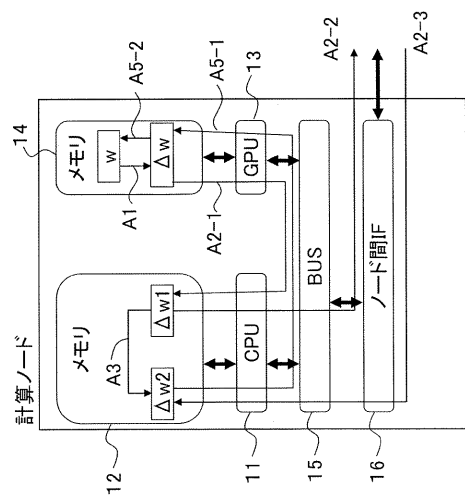
【図6】



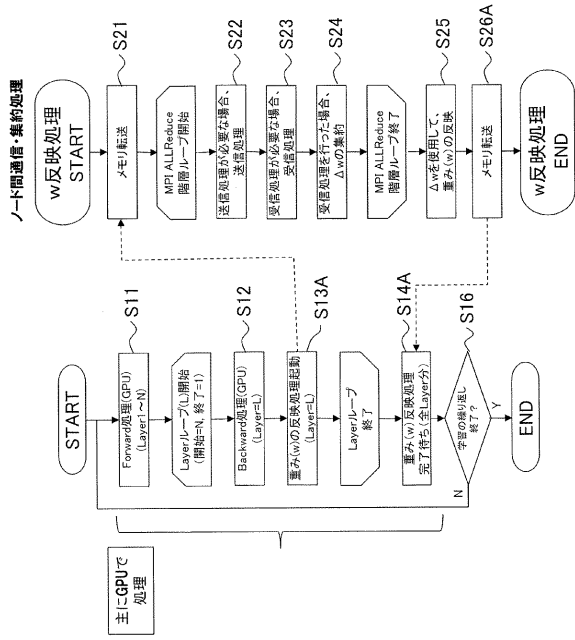
【図7】



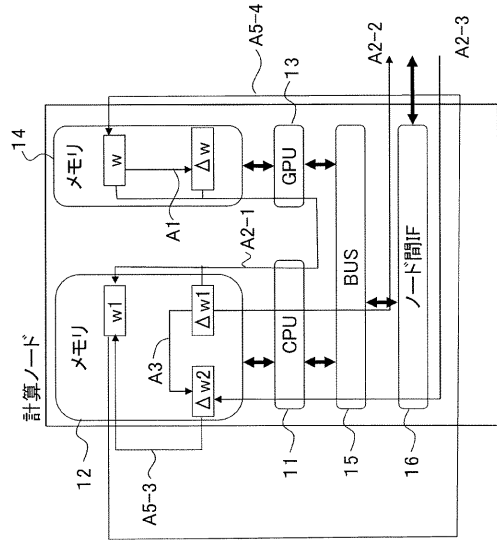
【図8】



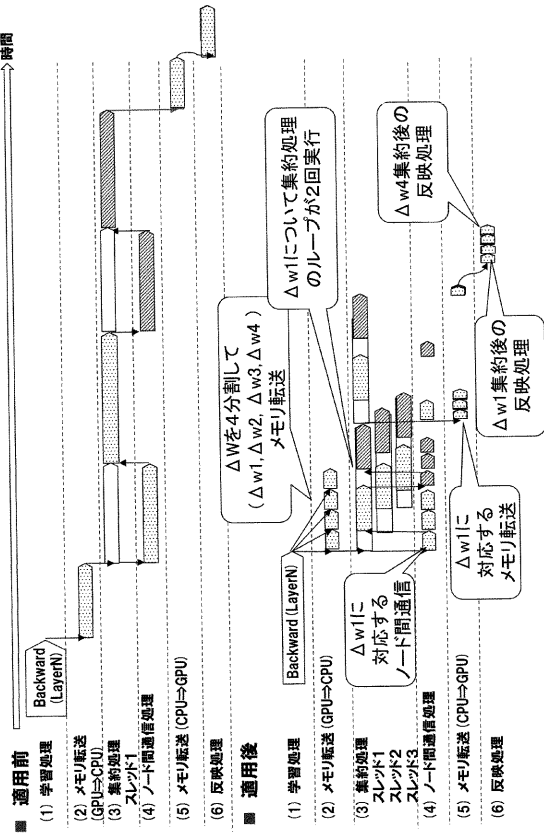
【図9】



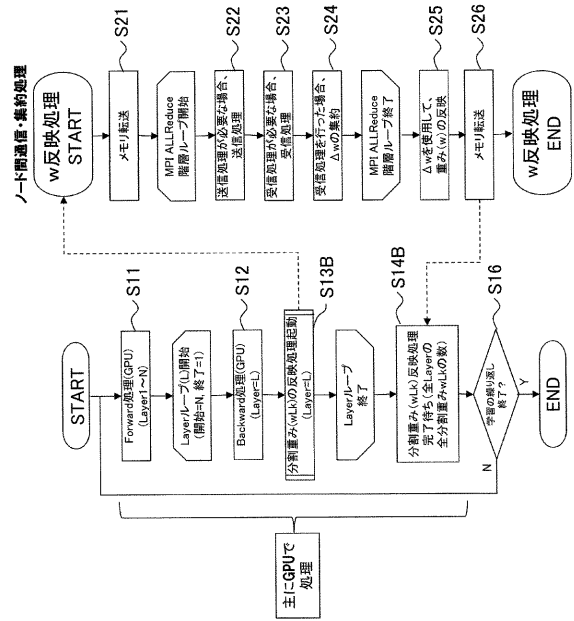
【図10】



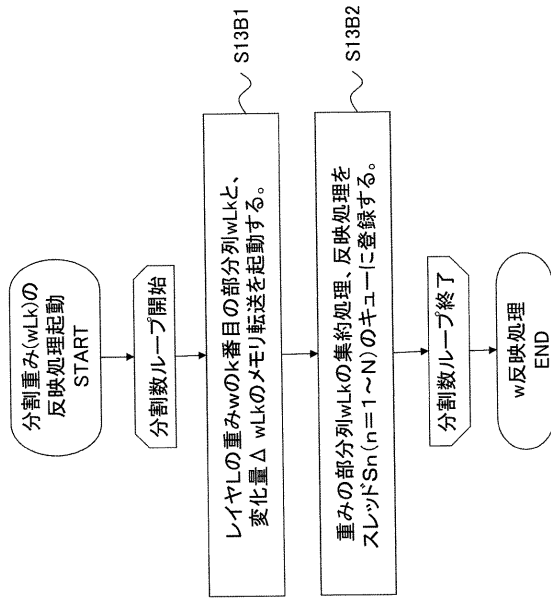
【図11】



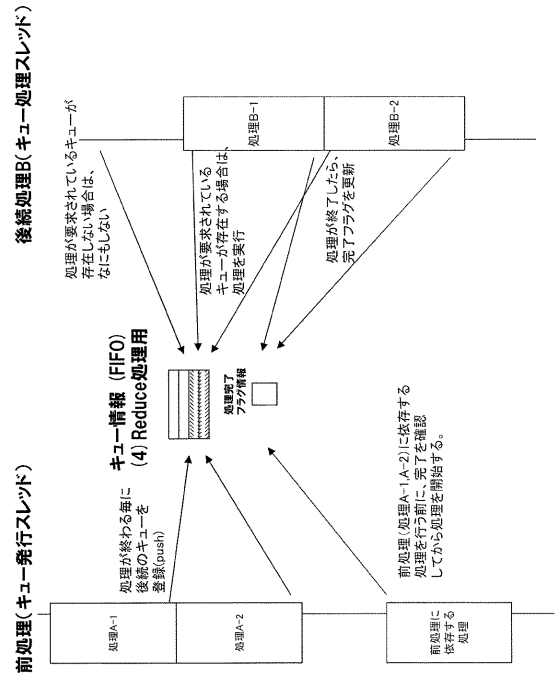
【図12】



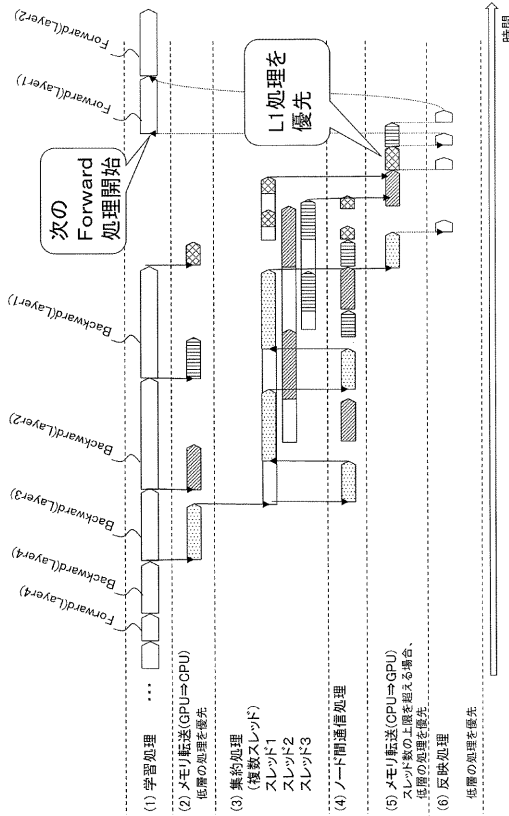
【図13】



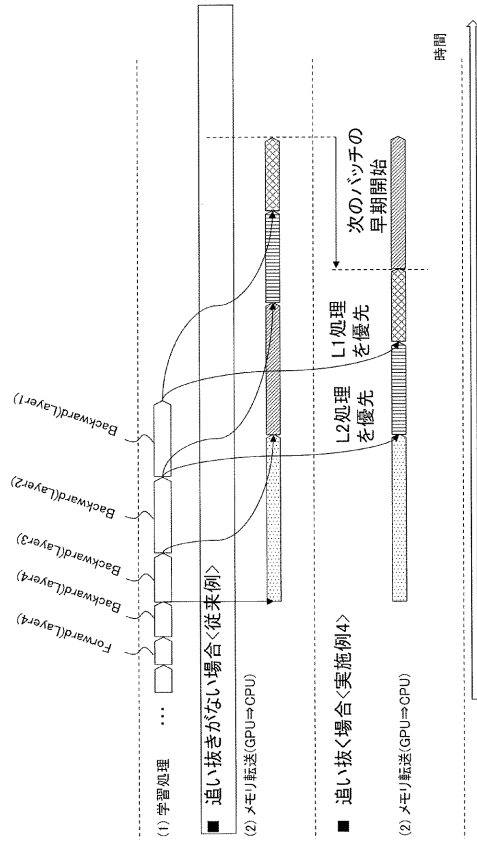
【図14】



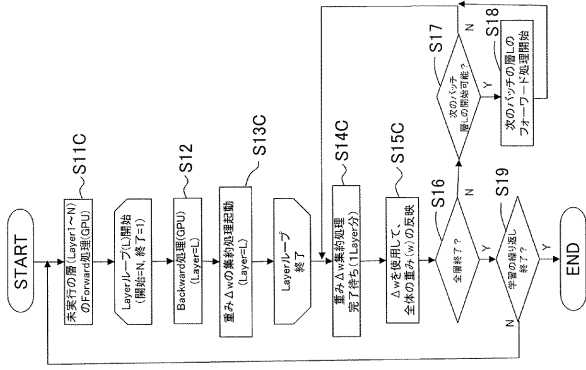
【図15】



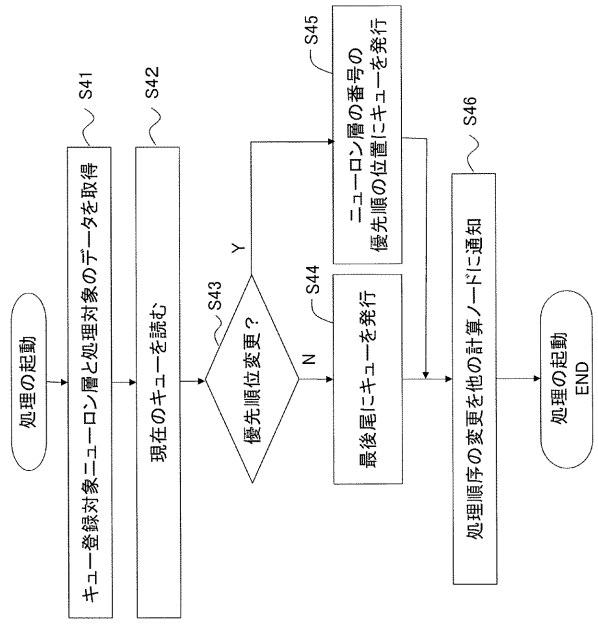
【図16】



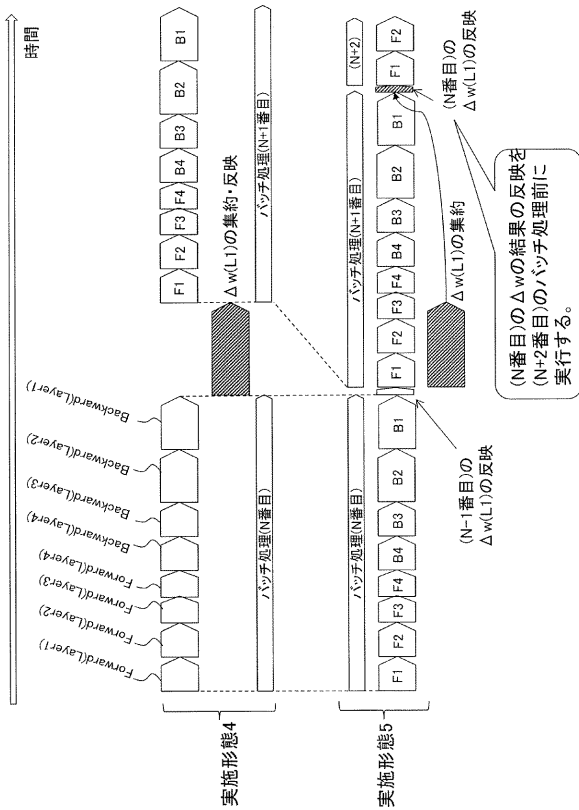
【 図 1 7 】



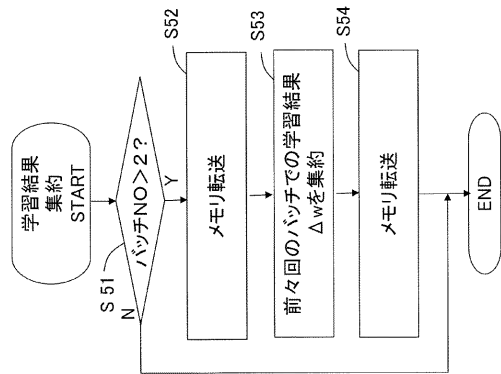
【 図 1 8 】



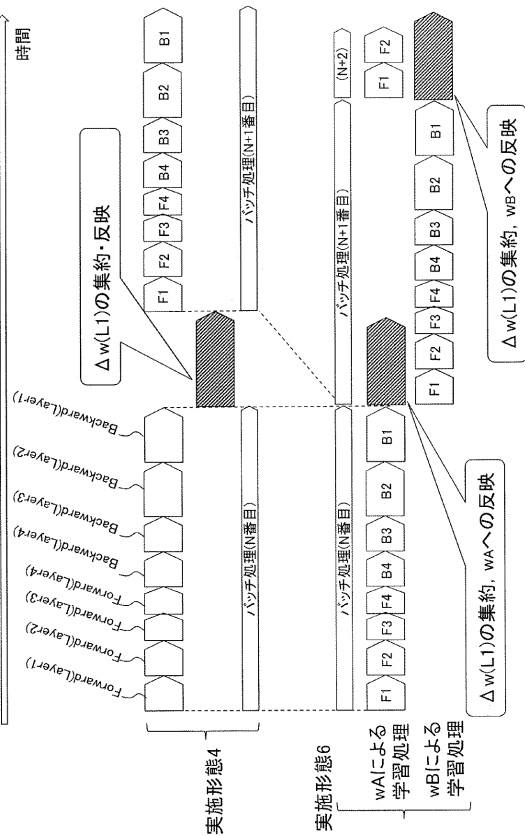
【 図 1 9 】



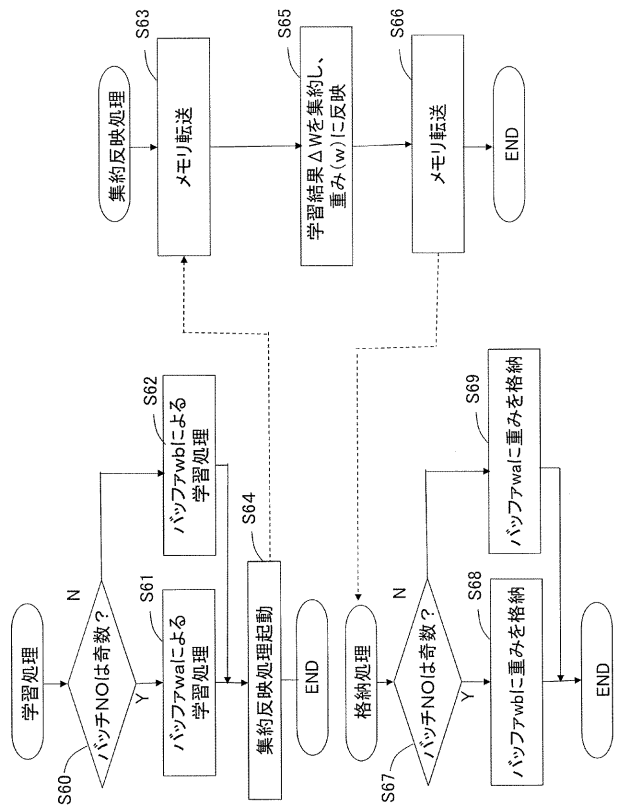
【 図 2 0 】



【図 2 1】



【図 2 2】



【手続補正書】

【提出日】平成29年6月22日(2017.6.22)

【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】0015

【補正方法】変更

【補正の内容】

【0015】

フォワード方向の処理は、入力画像に対して、畳み込み層の処理と、サブサンプリング層の処理を繰り返し実行する特徴抽出部の処理と、識別結果を出力する識別部の処理を含む。特徴抽出部は、入力画像に対して、畳み込み層の処理と、サブサンプリング層の処理を繰り返し実行することで、間引かれた画像を抽出する。畳み込み層の処理は、畳み込み演算ともいう。畳み込み演算は、例えば、N個×N個の画素を有する画像の情報(第N-1層)に対して、例えば、 $m \times m$ 個の重み $w_{a,b}$ ($a, b = 0, \dots, m-1$)のフィルタによる畳み込み演算を実行することで、次の層(第N層)の画像の情報を作る。サブサンプリング層の処理は、画像間引き処理であり、プーリング演算ともいう。

【手続補正2】

【補正対象書類名】明細書

【補正対象項目名】0019

【補正方法】変更

【補正の内容】

【0019】

勾配降下法によるニューラルネットワークの学習処理においては、エラーの評価関数(ERROR)の勾配と、学習係数イータの積が重みwの変化量(例えば、現在の重み w_t と次の重み w_{t+1} の差分値)となる。すなわち、深層学習においては、フォワード方

向に各ニューロン層の処理が実行され、バックワード方向に、各ニューロン層でのエラーの評価関数 (ERROR) が伝搬される。そして、各ニューロン層は、バックワード方向に伝搬するエラーの評価関数 (ERROR) から、エラーの評価関数 (ERROR) の勾配を求める。そして、各ニューロン層は、エラーの評価関数 (ERROR) が小さくする方向でのエラーの評価関数 (ERROR) の勾配と、学習係数イータの積から重み w_t の変化量 (勾配情報ともいう) を算出し、次回の重み w_{t+1} を求める。ここで、現在の重みを w_t で表し、次回の演算で使用される重みを w_{t+1} で表した。また、図 1 で説明したように、学習処理において、重み w は 1 以上の成分を有する係数列 (ベクトル) である。

【手続補正 3】

【補正対象書類名】明細書

【補正対象項目名】0048

【補正方法】変更

【補正の内容】

【0048】

図 7 は、実施形態 1 の計算ノード 10 の処理を例示するフローチャートである。図で左側のフローチャートは、主として GPU 13 が実行する学習処理と反映処理を例示する。また、右側のフローチャートは、主として CPU 11 が実行するノード間通信・集約処理を例示する。図 7 の処理では、まず、GPU 13 がニューロン層 (例えば、ニューロン層 1 から N) について、フォワード処理を実行する (S11)。

【手続補正 4】

【補正対象書類名】明細書

【補正対象項目名】0049

【補正方法】変更

【補正の内容】

【0049】

フォワード処理は、図 1 に例示したように、入力データと重み (w) とによる演算処理である。演算処理は、例えば、入力データの要素 $x(i, j)$ と $m \times m$ 個の重み w_{ab} ($a, b = 0, \dots, m-1$) のフィルタによる畳み込み演算、サブサンプリング層のプーリング演算、全結合層の演算等である。S11 の処理は、処理対象のデータに対する係数による演算処理の一例である。

【手続補正 5】

【補正対象書類名】明細書

【補正対象項目名】0062

【補正方法】変更

【補正の内容】

【0062】

そして、CPU 11 は、メモリ転送で GPU 13 に集約されたニューロン層 L の重み (w) の変化量 (Δw) を転送する (矢印 A5-1)。転送先の GPU 13 は、転送された重みの変化量を変化量 (Δw) に保存する。そして、GPU 13 は、集約された層 L の重み (w) の変化量 (Δw) を使用して、重み (w) を更新する (A5-2)。

【手続補正 6】

【補正対象書類名】明細書

【補正対象項目名】0076

【補正方法】変更

【補正の内容】

【0076】

例えば、あるニューロン層 L の重み (w) が $w = (p_1, p_2, \dots, p_X)$ のようなパラメータ列であるとする。パラメータ列は、係数列の一例である。つまり、ニューロン層 L の重み (w) は、複数使用され、係数列を形成する。そして、学習処理の結果、重みの変化量は $\Delta w = (p_1, p_2, \dots, p_X)$ のような多数のパラメータ

10

20

30

40

50

列として計算されるとする。このような場合に、GPU13は w を部分列に区切り、 $w1 = (p1, p2, \dots, pX1)$ 、 $w2 = (pX1 + 1, \dots, pX2)$ 、 $w3 = (pX2 + 1, \dots, pX3)$ 、 \dots 、 $wx = (pX3 + 1, \dots, pX)$ のように分割する。

【手続補正7】

【補正対象書類名】明細書

【補正対象項目名】0082

【補正方法】変更

【補正の内容】

【0082】

また、スレッド1は、分割された変化量 ($w1$) についてメモリ転送スレッドのキューを発行すると、分割された次の変化量 ($w4$) について、分割された変化量 ($w1$) と同様の処理を実行する。このようにして、例えば、CPU11が複数、例えば、5つのコアを有している場合には、CPU11はスレッド1から3およびメモリ転送スレッド、およびノード間通信スレッドを並行して実行できる。したがって、例えば、ある分割された変化量 (wk) についてのノード間通信の処理が、別の分割された変化量 (wj) についての集約処理の時間に実行できる。また、仮に、あるニューロン層Lの重み (wL) のパラメータ数が他の層よりも多いものであっても、GPU13およびCPU11は重み (wL) に含まれるパラメータを複数部分に分割し、複数スレッドで並行して処理できる。

【手続補正8】

【補正対象書類名】明細書

【補正対象項目名】0093

【補正方法】変更

【補正の内容】

【0093】

ところで、図15の例では、ニューロン層1の集約された重みの変化量がCPU11からGPU13にメモリ転送可能となったときに、まだ、ニューロン層2についても、集約された変化量のメモリ転送が起動されていない。例えば、ニューロン層2のメモリ転送処理 (CPU11からGPU13) は、キューが登録された状態で未実行の状態となっている。実施形態4では、このような場合にニューロン層1の集約処理が終了すると、集約処理用のスレッドは、ニューロン層2よりもニューロン層1のメモリ転送を優先する。すなわち、CPU11の集約処理用のスレッドは、ニューロン層2よりもニューロン層1が先に転送されるように、ニューロン層1の集約された変化量のメモリ転送のキューを登録する。そのようなキュー登録の結果、メモリ転送用スレッドはニューロン層2よりもニューロン層1の重みの変化量を先にメモリ転送する。

【手続補正9】

【補正対象書類名】明細書

【補正対象項目名】0099

【補正方法】変更

【補正の内容】

【0099】

図17は、実施形態4の学習処理を例示するフローチャートである。この処理では、GPU13は、ニューロン層1からNについて、フォワード方向の処理を実行する (S11C)。ただし、S11Cの処理は、前のバッチにおける全層についての学習処理が終了していなくても開始される点で実施形態1から3と相違する。そして、全層についてのフォワード方向の処理が終了すると、GPU13は、バックワード方向にニューロン層Nから1のループ (LAYER ループ (L) 開始 = N、終了 = 1) の中で、S12、S13Cの処理を実行する。S12の処理は、実施形態1から3と同様である。

【手続補正10】

10

20

30

40

50

【補正対象書類名】明細書

【補正対象項目名】0100

【補正方法】変更

【補正の内容】

【0100】

S13Cの処理では、GPU13は、ニューロン層のうち、入力側に近いニューロン層を優先して、CPU11にメモリ転送するとともに、集約処理を実行するCPU11のスレッドにキューを登録する(S13C)。S13Cの処理は、複数階層のうち、演算処理の実行の順序が早い階層の係数の変化量を優先して処理部に転送することの一例である。

【手続補正11】

【補正対象書類名】明細書

【補正対象項目名】0101

【補正方法】変更

【補正の内容】

【0101】

したがって、実施形態4では、GPU13は、各ニューロン層(L)でバックワード方向の処理が終了するごとに、優先順の制御を実行する。すなわち、GPU13は、バックワード方向の処理が終了したニューロン層(L)より上位のニューロン層(L+k)で、メモリ転送および集約処理が未実行のニューロン層がキューに残っていないか否かを判定する。そして、バックワード方向の処理が終了したニューロン層(L)より上位のニューロン層(L+k)がキューに残っている場合には、GPU13は、入力側に近い下位のニューロン層(L)を優先してキューを登録する。なお、このように、下位のニューロン層を優先するキューの登録は、CPU11がノード間通信およびメモリ転送(CPU11からGPU13)のキューを登録する場合も同様である。

【手続補正12】

【補正対象書類名】明細書

【補正対象項目名】0123

【補正方法】変更

【補正の内容】

【0123】

なお、変化量(w)の集約と重み(w)への反映は、実施形態2のように、CPU11において行ってもよい。つまり、GPU13は、集約された変化量(w)が反映済みの重み(w)をメモリ転送で受け取るようにしてもよい。この場合には、反映処理は、単に、変化量(w)が反映済みの重み(w)をGPU13のメモリ14に保存する処理とすることができる。

【手続補正13】

【補正対象書類名】明細書

【補正対象項目名】0128

【補正方法】変更

【補正の内容】

【0128】

図21に、実施形態6のタイムチャートを実施形態4と対比して例示する。図21のように、実施形態6では、バッファwaに格納した重みを用いた学習処理とバッファwbに格納した重みを用いた学習処理が交互に実行される。例えば、奇数番目のバッチの学習終了後に、次の偶数番目のバッチの学習処理と並行して集約処理と反映処理が実行される。そして、奇数番目のバッチの学習処理の結果である重みの変化量(w)が反映された重み(w)がバッファwaに格納される。このとき、偶数番目のバッチの学習処理では、バッファwbに格納された重みを使用される。

【手続補正14】

【補正対象書類名】明細書

10

20

30

40

50

【補正対象項目名】0132

【補正方法】変更

【補正の内容】

【0132】

まず、GPU13は、N番目のバッチが奇数番目のバッチか否かを判定する(S60)。N番目のバッチが奇数番目のバッチの場合、GPU13は、バッファwaに格納した重みによる学習処理を実行する(S61)。一方、N番目のバッチが偶数番目のバッチの場合、GPU13は、バッファwbに格納した重みによる学習処理を実行する(S62)。S61、S62の処理は、第1の記憶部に格納した係数を用いて演算処理を実行することの一例である。そして、GPU13は、メモリ転送をCPU11に要求するとともに、集約反映処理のキューを登録する(S64)。そして、GPU13は当該バッチの学習処理を終了する。そして、GPU13はN+1番目のバッチの学習処理を実行する。

10

【手続補正15】

【補正対象書類名】明細書

【補正対象項目名】0135

【補正方法】変更

【補正の内容】

【0135】

GPU13は、メモリ転送を受けると、バッチが奇数番目のバッチか否かを判定する(S67)。バッチが奇数番目のバッチの場合、GPU13は、バッファwbに重みを格納する(S68)。一方、バッチが偶数番目のバッチの場合、GPU13は、バッファwaに重みを格納する(S69)。S68、S69の処理は、演算処理による変化量を基に更新した係数を第2の記憶部に格納することの一例である。なお、S67からS69の処理は、次のさらに次のバッチ(N+2番目のバッチ)の学習処理が開始されるまでに実行される。

20

フロントページの続き

(72)発明者 笠置 明彦

神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内