

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7004741号

(P7004741)

(45)発行日 令和4年1月21日(2022.1.21)

(24)登録日 令和4年1月6日(2022.1.6)

(51)国際特許分類		F I		
G 0 6 N	3/06 (2006.01)	G 0 6 N	3/06	
G 0 6 F	12/04 (2006.01)	G 0 6 F	12/04	5 3 0
H 0 3 M	7/40 (2006.01)	H 0 3 M	7/40	

請求項の数 15 (全31頁)

(21)出願番号	特願2019-555659(P2019-555659)	(73)特許権者	314015767
(86)(22)出願日	平成30年4月16日(2018.4.16)		マイクロソフト テクノロジー ライセン
(65)公表番号	特表2020-517014(P2020-517014		シング,エルエルシー
	A)		アメリカ合衆国 ワシントン州 9 8 0 5
(43)公表日	令和2年6月11日(2020.6.11)		2 レッドモンド ワン マイクロソフト
(86)国際出願番号	PCT/US2018/027840		ウェイ
(87)国際公開番号	WO2018/194998	(74)代理人	100107766
(87)国際公開日	平成30年10月25日(2018.10.25)		弁理士 伊東 忠重
審査請求日	令和3年3月10日(2021.3.10)	(74)代理人	100070150
(31)優先権主張番号	62/486,432		弁理士 伊東 忠彦
(32)優先日	平成29年4月17日(2017.4.17)	(74)代理人	100091214
(33)優先権主張国・地域又は機関	米国(US)		弁理士 大貫 進介
(31)優先権主張番号	15/953,356	(72)発明者	コーカリー, ジョセフ レオン
(32)優先日	平成30年4月13日(2018.4.13)		アメリカ合衆国 ワシントン州 9 8 0 5
	最終頁に続く		2 レッドモンド ワン マイクロソフト
			最終頁に続く

(54)【発明の名称】 メモリ帯域幅利用を低減するために活性化データの圧縮及び復元を使用するニューラルネットワークプロセッサ

(57)【特許請求の範囲】

【請求項1】

ニューラルネットワークプロセッサであって、

1つ以上のニューロンと、圧縮ユニットとを有し、

前記圧縮ユニットは、

当該ニューラルネットワークプロセッサ内の前記ニューロンの少なくとも1つによって生成されるデータの非圧縮チャンクを受け取り、該データの非圧縮チャンクが一定数のバイトを含み、

圧縮された出力チャンクのマスク部分を生成し、該マスク部分が、前記データの非圧縮チャンク内の前記一定数のバイトに等しいビットの数を含み、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内のバイトに対応し、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内の対応するバイトがゼロである場合に論理0に設定され、前記データの非圧縮チャンク内の対応するバイトが非ゼロである場合に論理1に設定され、前記データの非圧縮チャンク内の非ゼロバイトの数を決定することと、前記データの非圧縮チャンク内の前記非ゼロバイトの数に基づき、前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記圧縮された出力チャンクのデータ部分内のビットの数を決定することと、該決定されたビットの数まで前記データの非圧縮チャンク内の前記非ゼロバイトを切り捨てることと、該切り捨てられた非ゼロバイトを前記圧縮された出力チャンクの前記データ部分に格納することとによって、前記圧縮された出力チャンクの前記データ部分を生成し、

前記マスク部分及び前記データ部分を含む前記圧縮された出力チャンクを出力するよう構成される、
ニューラルネットワークプロセッサ。

【請求項 2】

当該ニューラルネットワークプロセッサは、復元ユニットを更に有し、
前記復元ユニットは、
前記圧縮された出力チャンクを受け取り、
前記圧縮された出力チャンクの前記マスク部分に基づき、前記データの非圧縮チャンク内の非ゼロバイトの数を決定し、
前記圧縮された出力チャンクの前記マスク部分に基づき、前記データの非圧縮チャンク内の前記非ゼロバイトの位置を決定し、
前記圧縮された出力チャンクの前記データ部分に前記切り捨てられた非ゼロバイトを格納するために前記圧縮ユニットによって使用されたビットの数を決定し、
論理 0 である前記圧縮された出力チャンクの前記マスク部分内の各ビット位置について、ゼロバイトを、復元された出力チャンクの対応する位置に挿入し、
論理 1 である前記マスク部分内の各ビット位置について、前記圧縮された出力チャンクの前記データ部分内の対応する位置からの前記切り捨てられた非ゼロバイトを、前記復元された出力チャンクの対応する位置に、前記圧縮された出力チャンクの圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに挿入する
よう構成される、
請求項 1 に記載のニューラルネットワークプロセッサ。

10

20

【請求項 3】

前記圧縮ユニットは更に、
前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記圧縮された出力チャンクの前記データ部分内の追加ビットの数を決定し、
前記データの非圧縮チャンク内の前記非ゼロバイトのうちの 1 つ以上の非ゼロバイトを切り捨てる前に、該 1 つ以上の非ゼロバイトに前記追加ビットを割り当てる
よう構成される、
請求項 2 に記載のニューラルネットワークプロセッサ。

【請求項 4】

前記復元ユニットは更に、前記圧縮された出力チャンクの前記データ部分内に格納された前記 1 つ以上の非ゼロバイトに割り当てられている前記追加ビットの数を決定するよう構成される、
請求項 3 に記載のニューラルネットワークプロセッサ。

30

【請求項 5】

前記復元ユニットは更に、前記復元された出力チャンクに格納された前記切り捨てられた非ゼロバイトのうちの 1 つ以上の非ゼロバイトにオフセットを加えるよう構成される、
請求項 2 に記載のニューラルネットワークプロセッサ。

【請求項 6】

前記非ゼロバイトの 1 つ以上の最下位ビット (LSB) が切り捨てられる、
請求項 1 に記載のニューラルネットワークプロセッサ。

40

【請求項 7】

ニューラルネットワークプロセッサであって、
1 つ以上のニューロンと、復元ユニットとを有し、
前記復元ユニットは、
マスク部分及びデータ部分を含むデータの圧縮されたチャンクを受け取り、
前記マスク部分内のビットに基づき、データの復元されたチャンク内の非ゼロバイトの数を決定し、
前記非ゼロバイトの数に少なくとも部分的に基づき、前記データの圧縮されたチャンクの前記データ部分に切り捨てられた非ゼロバイトを格納するために使用されたビットの数を

50

決定し、

論理 0 である前記データの圧縮されたチャンクの前記マスク部分内の各ビット位置について、ゼロバイトを、前記データの復元されたチャンクの対応する部分に挿入し、

論理 1 である前記データの圧縮されたチャンクの前記マスク部分内の各ビット位置について、前記データの圧縮されたチャンクの前記データ部分内の対応する位置からの前記切り捨てられた非ゼロバイトを、前記データの復元されたチャンクの対応する位置に、前記データの圧縮されたチャンクの圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに挿入する

よう構成される、

ニューラルネットワークプロセッサ。

10

【請求項 8】

圧縮ユニットを更に有し、

前記圧縮ユニットは、

当該ニューラルネットワークプロセッサ内の前記ニューロンの少なくとも 1 つによって生成されるデータの非圧縮チャンクを受け取り、該データの非圧縮チャンクが一定数のバイトを含み、

前記データの圧縮されたチャンクの前記マスク部分を生成し、該マスク部分が、前記データの非圧縮チャンク内の前記一定数のバイトに等しいビットの数を含み、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内のバイトに対応し、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内の対応するバイトがゼロである場合に論理 0

20

を有し、前記データの非圧縮チャンク内の対応するバイトが非ゼロである場合に論理 1 を有し、
前記データの非圧縮チャンク内の非ゼロバイトの数を決定することと、前記データの非圧縮チャンク内の前記非ゼロバイトの数に基づき、前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記データの圧縮されたチャンクの前記データ部分内のビットの数を決定することと、該決定されたビットの数まで前記データの非圧縮チャンク内の前記非ゼロバイトを切り捨てることと、該切り捨てられた非ゼロバイトを前記データの圧縮されたチャンクの前記データ部分に格納することとによって、前記データの圧縮されたチャンクの前記データ部分を生成し、

前記マスク部分及び前記データ部分を含む前記データの圧縮されたチャンクを出力する

30

よう構成される、

請求項 7 に記載のニューラルネットワークプロセッサ。

【請求項 9】

前記圧縮ユニットは更に、前記データの非圧縮チャンク内の非ゼロバイトの数が前記データの圧縮されたチャンクの前記データ部分内のバイトの数以下である場合には切り捨てなしで、前記データの非圧縮チャンク内の前記非ゼロバイトを前記データの圧縮されたチャンクの前記データ部分に格納するよう構成される、

請求項 8 に記載のニューラルネットワークプロセッサ。

【請求項 10】

前記圧縮ユニットは更に、

前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記データの圧縮されたチャンクの前記データ部分内の追加ビットの数を決定し、

前記データの非圧縮チャンク内の前記非ゼロバイトのうちの 1 つ以上の非ゼロバイトを切り捨てる前に、該 1 つ以上の非ゼロバイトに前記追加ビットを割り当てる

よう構成される、

請求項 8 に記載のニューラルネットワークプロセッサ。

40

【請求項 11】

前記復元ユニットは更に、前記データの圧縮されたチャンクの前記データ部分内に格納された前記 1 つ以上の非ゼロバイトに割り当てられている前記追加ビットの数を決定するよう構成される、

50

請求項 10 に記載のニューラルネットワークプロセッサ。

【請求項 12】

ニューラルネットワークプロセッサの圧縮ユニットで、前記ニューラルネットワークプロセッサ内の少なくとも 1 つのニューロンによって生成されるデータの非圧縮チャンクを受け取り、該データの非圧縮チャンクが一定数のバイトを含む、ことと、

圧縮された出力チャンクのマスク部分を生成し、該マスク部分が、前記データの非圧縮チャンク内の前記一定数のバイトに等しいビットの数を含み、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内のバイトに対応し、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内の対応するバイトがゼロである場合に論理 0 を有し、前記データの非圧縮チャンク内の対応するバイトが非ゼロである場合に論理 1 を有する、ことと、

10

前記データの非圧縮チャンク内の非ゼロバイトの数を決定すること、前記データの非圧縮チャンク内の前記非ゼロバイトの数に基づき、前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記圧縮された出力チャンクのデータ部分内のビットの数を決定することと、該決定されたビットの数まで前記データの非圧縮チャンク内の前記非ゼロバイトを切り捨てることと、該切り捨てられた非ゼロバイトを前記圧縮された出力チャンクの前記データ部分に格納することとによって、前記圧縮された出力チャンクの前記データ部分を生成することと、

前記マスク部分及び前記データ部分を含む前記圧縮された出力チャンクを前記ニューラルネットワークプロセッサのメモリに記憶することと

20

を有する、コンピュータにより実施される方法。

【請求項 13】

前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記圧縮された出力チャンクの前記データ部分内の追加ビットの数を決定することと、前記データの非圧縮チャンク内の前記非ゼロバイトのうちの 1 つ以上の非ゼロバイトを切り捨てる前に、該 1 つ以上の非ゼロバイトに前記追加ビットを割り当てることとを更に有する、

請求項 12 に記載のコンピュータにより実施される方法。

【請求項 14】

前記データの非圧縮チャンク内の非ゼロバイトの数が前記圧縮された出力チャンクの前記データ部分内のバイトの数以下である場合には切り捨てなしで、前記データの非圧縮チャンク内の前記非ゼロバイトを前記圧縮された出力チャンクの前記データ部分に格納することを更に有する、

30

請求項 12 に記載のコンピュータにより実施される方法。

【請求項 15】

前記ニューラルネットワークプロセッサの復元ユニットで、前記圧縮された出力チャンクを受け取ることと、

前記圧縮された出力チャンクの前記マスク部分に基づき、前記データの非圧縮チャンク内の非ゼロバイトの数を決定することと、

前記圧縮された出力チャンクの前記マスク部分に基づき、前記データの非圧縮チャンク内の前記非ゼロバイトの位置を決定することと、

40

前記圧縮された出力チャンクの前記データ部分に前記切り捨てられた非ゼロバイトを格納するために前記圧縮ユニットによって使用されたビットの数を決定することと、

論理 0 である前記圧縮された出力チャンクの前記マスク部分内の各ビット位置について、ゼロバイトを、復元された出力チャンクの対応する位置に挿入することと、

論理 1 である前記マスク部分内の各ビット位置について、前記圧縮された出力チャンクの前記データ部分内の対応する位置からの前記切り捨てられた非ゼロバイトを、前記復元された出力チャンクの対応する位置に、前記圧縮された出力チャンクの圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに挿入することと

を更に有する、

50

請求項 1 2 に記載のコンピュータにより実施される方法。

【発明の詳細な説明】

【背景技術】

【0001】

ディープニューラルネットワーク (Deep Neural Network(s)) (“DNN”) は、人間の脳などの生体神経系における情報処理及び通信パターンを大まかにモデル化している。DNN は、制限なしに、物体検出、意味的標識付け (semantic labeling)、及び特徴抽出などの複雑な分類問題を解くために利用可能である。結果として、DNN は、コンピュータビジョン、音声認識、及び機械翻訳などの多くの人工知能 (Artificial Intelligence) (“AI”) アプリケーションの基礎を形成している。DNN は、それらの分野の多くで人間の精度以上であることができる。 10

【0002】

DNN の高い性能は、入力空間の効果的表現を得るために大規模データセットに対する統計学習を使用した後に入力データから高次特徴 (high-level features) を抽出するそれらの能力に由来する。しかし、DNN の優れた性能は、高い計算複雑性という代償の上に成り立っている。グラフィクスプロセッシングユニット (Graphics Processing Unit(s)) (“GPU”) などの高性能の汎用プロセッサが、多くの DNN アプリケーションによって必要とされる高度な計算性能を提供するために広く利用されている。

【0003】

GPU のような汎用プロセッサは、DNN を実装するための高度な計算性能を提供することができるが、他方で、そのようなプロセッサは、低電力消費が重要であるコンピュータ装置で DNN 動作を長期にわたって実行することにおける使用には、通常は適さない。例えば、GPU などの汎用プロセッサは、スマートフォン又は代替/仮想現実 (Alternate/Virtual Reality) (“AR/VR”) デバイスのような、バッテリー寿命を延ばすために電力消費の削減が求められるバッテリー駆動のポータブルデバイスで長期にわたる DNN タスクを実行することにおける使用には不適切であり得る。 20

【0004】

人間の動作の検出などの連続的な DNN タスクを実行しながら電力消費を低減することは、例えば、パワー・オーバー・イーサネット (Power-over-Ethernet) (“PoE”) などの非バッテリー駆動デバイスでも重要であり得る。この具体的な例において、PoE スイッチは、限られた量の電力しか供給することができず、防犯カメラなどの PoE デバイスの電力消費量を減らすことは、電力供給量が少ない PoE スイッチの使用を可能にする。 30

【0005】

高性能の DNN 処理を提供しながら同時に、汎用プロセッサと比較して電力消費量を削減することができる特定用途向け集積回路 (Application-Specific Integrated Circuit(s)) (“ASIC”) が開発されてきた。しかし、この分野の進歩にかかわらず、特に、低電力消費が重要であるコンピュータ装置での使用のために、DNN 処理を実行する ASIC の性能の改善及び電力消費量の削減に対するニーズが依然として存在する。

【0006】

本明細書でなされた開示が提示されるのは、これら及び他の技術的課題に関してである。 40

【発明の概要】

【0007】

メモリバス帯域幅の利用量を低減するために活性化データを圧縮及び復元することができる DNN モジュール又はプロセッサが、開示される。特に、DNN モジュールは、ニューロン出力とオンボード又はオフボードのメモリとの間のバス帯域幅の利用量を減らすために圧縮を利用することができる。DNN モジュールはまた、オンボード又はオフボードのメモリとニューロン入力との間のメモリバス帯域幅の利用量を減らすために復元を利用することもできる。帯域幅利用量の低減は、より高速な処理を可能にすることができ、その結果、電力量消費量も減らすことができる。本明細書で具体的に述べられない他の技術的利点も、開示されている対象の実施を通じて実現され得る。 50

【 0 0 0 8 】

先に簡単に述べられた技術的利点を実現するために、1つ以上のニューロン及び圧縮ユニットを含むDNNプロセッサが、開示される。圧縮ユニットは、1つ以上のニューロンによって生成されるデータの非圧縮チャンクを受け取ることができる。いくつかの実施形態において、データの非圧縮チャンクは、64バイトなどの、一定数のバイトを含む。

【 0 0 0 9 】

データの非圧縮チャンクを圧縮するために、圧縮ユニットは、圧縮された出力チャンクのマスク部分及びデータ部分を生成することができる。圧縮された出力チャンクのマスク部分は、データの非圧縮チャンク内の一定数のバイトに等しいビットの数を含む。例えば、データの非圧縮チャンクが64バイトのデータを含む場合には、マスク部分は64ビット（すなわち、8バイト）を含むことになる。

10

【 0 0 1 0 】

いくつかの実施形態において、圧縮された出力チャンクのマスク部分内の各ビットは、データの非圧縮チャンク内のバイトに対応する。例えば、マスク部分のビット1は、データの非圧縮チャンク内の第1バイトに対応することができ、マスク部分のビット2は、データの非圧縮チャンク内の第2バイトに対応することができる、など。他の実施形態では、圧縮された出力チャンクのマスク部分内の2つ以上のビットが、データの非圧縮チャンク内のあるバイトに対応する。かような実施形態で、圧縮された出力チャンクのマスク部分内のビットは、非圧縮チャンク内の対応するバイトだけでなくそのおおよその大きさも示すことができる。

20

【 0 0 1 1 】

マスク部分の個々のビットが非圧縮チャンク内のバイトに対応する場合に、圧縮ユニットは、圧縮された出力チャンクのマスク部分内の各ビットを、データの非圧縮チャンク内の対応するバイトが全てゼロを含むときに（すなわち、“ゼロバイト”）論理偽（false）（本明細書中「論理0」とも呼ばれ得る。）に設定する。圧縮ユニットはまた、データの非圧縮チャンク内の対応するバイトが少なくとも1つの非ゼロビットを含むときに（すなわち、“非ゼロバイト”）、圧縮された出力チャンクのマスク部分内の各ビットを論理真（true）（本明細書中「論理1」とも呼ばれ得る。）に設定する。このようにして、圧縮された出力チャンクのマスク部分は、データの非圧縮チャンクのゼロバイト及び非ゼロバイトの存在及び位置を符号化する。

30

【 0 0 1 2 】

圧縮ユニットは、データの非圧縮チャンク内の非ゼロバイトの数を決定することによって、圧縮された出力チャンクのデータ部分を生成する。圧縮ユニットは次いで、データの非圧縮チャンク内の非ゼロバイトの数と、圧縮された出力チャンクのデータ部分内で利用可能なバイトの数とに基づき、データの非圧縮チャンクの各非ゼロバイトを格納するために利用可能である圧縮された出力チャンクのデータ部分内のビットの数を決定する。例えば、データの圧縮されたチャンクのデータ部分が24バイト幅（すなわち、192ビット）であり、データの非圧縮チャンク内に47の非ゼロバイトがある場合に、データの非圧縮チャンクから各非ゼロバイトを格納するためには、4ビットがデータ部分内で利用可能である。

40

【 0 0 1 3 】

いくつかの実施形態において、圧縮ユニットはまた、データの非圧縮チャンクの非ゼロバイトを格納するために利用可能である圧縮された出力チャンクのデータ部分内で、もしあれば、追加ビットの数を決定することもできる。上記の例において、例えば、4つの追加ビットが、非ゼロバイトを格納するために利用可能である（すなわち、 $192 \bmod 47 = 4$ ビット）。圧縮ユニットは、それらの追加ビットを、データの非圧縮チャンク内の1つ以上の非ゼロバイトを切り捨てる前に、その1つ以上の非ゼロバイトに割り当てることができる。例えば、圧縮ユニットは、それらの追加ビットを、圧縮された出力チャンクのデータ部分内の最初の数バイトに割り当ててもよい。

【 0 0 1 4 】

50

圧縮ユニットは、次いで、データの非圧縮チャンク内の非ゼロバイトを、各非ゼロバイトを格納するためにデータ部分内で利用可能な決定されたビット数（すなわち、上記の例では、4）まで切り捨てる。一実施形態において、圧縮ユニットは、データ部分内の利用可能なビット数に収まるように非ゼロバイトの最下位ビット（Least Significant Bit(s)）（“LSB”）を切り捨てる。他の実施形態では、圧縮ユニットは、非ゼロバイトの最上位ビット（Most Significant Bit(s)）（“MSB”）を切り捨てる。圧縮ユニットは、次いで、切り捨てられた非ゼロバイトを、圧縮された出力チャンクのデータ部分に格納する。マスク部分及びデータ部分を含む圧縮された出力チャンクは、次いで、例えば、DNNプロセッサ内のオンボード・メモリ、又はDNNプロセッサのアプリケーションホストのオフボード・メモリへ、出力され得る。

10

【0015】

DNNモジュールはまた、上述されたように圧縮されたデータのチャンクを復元することができる復元ユニットを含むこともできる。例えば、復元ユニットは、DNNプロセッサ内のメモリ又はアプリケーションホストのメモリからデータの圧縮されたチャンクを受け取ることができる。復元ユニットは、次いで、圧縮された出力チャンクのマスク部分内の論理真ビットの数に基づき、データの非圧縮チャンク内の非ゼロバイトの数を決定することができる。復元ユニットはまた、圧縮された出力チャンクのマスク部分内の論理真ビットの位置に基づき、データの非圧縮チャンク内の非ゼロバイトの位置を決定することもできる。復元ユニットは、同様にして、データの非圧縮チャンク内のゼロバイトの位置を決定することができる。

20

【0016】

復元ユニットはまた、切り捨てられた非ゼロバイトを圧縮された出力チャンクのデータ部分に格納するために圧縮ユニットによって使用されたビットの数を決定することもできる。復元ユニットは、データの非圧縮チャンク内の非ゼロバイトの数と、圧縮された出力チャンクのデータ部分内で利用可能なバイトの数とに基づき、切り捨てられた各非ゼロバイトを格納するために使用されたビットの数を決定することができる。

【0017】

上記の例では、例えば、データの圧縮されたチャンクのデータ部分が24バイト幅（すなわち、192ビット）であり、データの非圧縮チャンク内に47の非ゼロバイトがある場合に、圧縮ユニットは、データの非圧縮チャンクの切り捨てられた各非ゼロバイトを格納するためにデータ部分において4ビットを利用した。復元ユニットはまた、もしあれば、圧縮ユニットが圧縮された出力チャンクのデータ部分に格納された切り捨てられた非ゼロバイトのうちの1つ以上に割り当てた追加ビットの数を決定することもできる。

30

【0018】

論理0である圧縮された出力チャンクのマスク部分内の各ビット位置について、復元ユニットは、ゼロバイトを、復元された出力チャンクの対応する部分に挿入する。論理1であるマスク部分内の各ビット位置について、復元ユニットは、圧縮された出力チャンクのデータ部分内の対応する位置からの切り捨てられた非ゼロバイトを、圧縮された出力チャンクの圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに、復元された出力チャンクの対応する位置に挿入する。ゼロビットは、どのビットが圧縮中に切り捨てられたかに応じて、切り捨てられた非ゼロバイトのLSB又はMSBに挿入され得る。

40

【0019】

いくつかの実施形態において、復元ユニットはまた、復元された出力チャンクに格納された切り捨てられた非ゼロバイトのうちの1つ以上にオフセット（例えば、00000001）を付加する。例えば、オフセットは、圧縮後にゼロバイトになるデータの非圧縮チャンクの非ゼロバイトに加えられ得る。このようにして、非ゼロバイトは、圧縮及び復元される場合にゼロバイトにならない。他の実施形態では、オフセットは、復元された出力チャンク内の全てのバイトに付加され得る。

【0020】

先に簡単に説明されたように、本明細書で開示される技術の実装は、DNNモジュールに

50

おけるメモリバス帯域幅利用量を低減し、DNNモジュールが処理動作をより速く完了することを可能にし、電力消費量を削減することができる。本明細書で具体的に特定されていない技術的利点も、開示されている技術の実施を通じて実現され得る。

【0021】

当然ながら、上記の対象は、コンピュータにより制御される装置、コンピュータにより実施される方法、コンピュータ装置として、又はコンピュータ可読媒体などの製品として実施され得る。それら及び様々な他の特徴は、続く詳細な説明を読み、添付の図面を参照することで、明らかになるだろう。

【0022】

この概要は、開示される技術のいくつかの態様について概要を簡単に紹介するために設けられており、それらの態様は、詳細な説明において更に後述される。この概要は、請求されている対象の重要な特徴又は必須の特徴を特定する意図はなく、この概要が、請求されている対象の適用範囲を制限するために使用されることも意図しない。更には、請求されている対象は、本開示のいずれかの部分で述べられているありとあらゆる欠点を解消する実施に限られない。

【図面の簡単な説明】

【0023】

【図1】一実施形態に従って、本明細書で開示される技術の態様を実装するDNNモジュールの構成及び動作の態様を示すコンピューティングアーキテクチャ図である。

【図2A】一実施形態に従って、活性化データを圧縮するDNNモジュールの構成及び動作の態様を示すコンピューティングシステムアーキテクチャ図である。

【図2B】一実施形態に従って、活性化データを圧縮するDNNモジュールの構成及び動作の態様を示すコンピューティングシステムアーキテクチャ図である。

【図3】一実施形態に従って、圧縮されていない活性化データのチャンクの例を参照して、活性化データを圧縮するためのDNNモジュールの動作の態様を説明するデータ構造図である。

【図4】本明細書で開示される一実施形態に従って、活性化データを圧縮するための開示されているDNNモジュールの動作の態様を説明するルーチンを示すフロー図である。

【図5A】一実施形態に従って、活性化データを復元するためのDNNモジュールの構成及び動作の態様を示すコンピューティングシステムアーキテクチャ図である。

【図5B】一実施形態に従って、活性化データを復元するためのDNNモジュールの構成及び動作の態様を示すコンピューティングシステムアーキテクチャ図である。

【図6】一実施形態に従って、圧縮された活性化データのチャンクの例を参照して、活性化データを復元するためのDNNモジュールの動作の態様を説明するデータ構造図である。

【図7】本明細書で開示される一実施形態に従って、活性化データを復元するための開示されているDNNモジュールの動作の態様を説明するルーチンを示すフロー図である。

【図8】一実施形態に従って、本明細書で提示されるDNNモジュールのためのアプリケーションホストとして動作することができるコンピュータ装置のためのコンピュータハードウェア及びソフトウェアアーキテクチャを例示するコンピュータアーキテクチャ図である。

【図9】本明細書で提示される様々な実施形態に従って、開示される技術の態様が実装され得る分散コンピューティング環境を説明するネットワーク図である。

【発明を実施するための形態】

【0024】

続く詳細な説明は、メモリバス帯域幅の利用量を低減するために活性化データを圧縮及び復元することができるDNNモジュールを対象とする。先に簡単に説明されたように、開示される技術の実装は、DNNモジュールにおけるメモリバス帯域幅利用量を低減し、DNNモジュールが処理動作をより速く完了することを可能にし、電力消費量を削減することができる。本明細書で具体的に述べられない他の技術的利点も、開示されている対象の実施を通じて実現され得る。

10

20

30

40

50

【 0 0 2 5 】

本明細書中に記載される対象は、ハードウェア DNN モジュールの一般的状況において提示されるが、当業者は、他のタイプのコンピューティングシステム及びモジュールと組み合わせて他の実施が実行され得ると認めるだろう。当業者はまた、本明細書中に記載される対象が、手持ち式デバイス、マルチプロセッサシステム、マイクロプロセッサに基づく又はプログラム可能な家電製品、デバイス（例えば、装用式コンピュータ装置、自動車、ホームオートメーション、など）に埋め込まれたコンピューティング又はプロセッシングシステム、ミニコンピュータ、メインフレームコンピュータ、などを含む他のコンピュータシステム構成により実施され得るであると認めるだろう。

【 0 0 2 6 】

以下で更に詳細に記載されるように、そのニューロンの出力を圧縮するよう構成される DNN モジュールが、開示される。圧縮された出力は、DNN モジュールのメモリに、又は DNN モジュールのためのアプリケーションホストによって提供されるメモリなどの、DNN モジュールの外にあるメモリに、記憶され得る。DNN モジュールは、前に圧縮されたデータを後に復元し、そして、復元されたデータをニューロンへ供給することができる。

【 0 0 2 7 】

一実施形態によれば、DNN プロセッサ内の圧縮ユニットは、圧縮されていない活性化データの固定長のチャンク（例えば、64 バイト）を一定の圧縮比（例えば、2 : 1）で圧縮する。圧縮ユニットによって生成される圧縮された活性化データは、固定長（例えば、32 バイト）を有するデータのチャンクを含むことができ、固定長のマスク部分（例えば、8 バイト）と、固定長のデータ部分（例えば、24 バイト）とを含む。

【 0 0 2 8 】

一実施形態において、圧縮された出力チャンクのマスク部分のビットは、圧縮されていない入力チャンク内のバイトに対応する。例えば、マスク部分の第 1 ビットは、圧縮されていない入力チャンク内の第 1 バイトに対応することができ、マスク部分の第 2 ビットは、圧縮されていない入力チャンク内の第 2 バイトに対応することができる、など。圧縮された活性化データのマスク部分内のビットは、圧縮されていない入力チャンク内の対応するバイトがゼロである場合には論理 0 にセットされ得、圧縮されていない入力チャンク内の対応するバイトが非ゼロである場合には論理 1 にセットされ得る。

【 0 0 2 9 】

先に簡単に説明されたように、いくつかの実施形態において、圧縮された出力チャンクのマスク部分内の 2 つ以上のビットが、データの非圧縮チャンク内のあるバイトに対応する。かような実施形態で、圧縮された出力チャンクのマスク部分内のビットは、非圧縮チャンク内の対応するバイトだけでなくそのおおよその大きさも示すことができる。

【 0 0 3 0 】

圧縮された出力チャンクのデータ部分は、圧縮されたデータ部分内の利用可能なビットの数をを用いて非圧縮入力チャンクの非ゼロバイトを表すよう、切り捨てられた非圧縮入力チャンクの非ゼロバイトを含む。いくつかの実施形態において、各非ゼロバイトごとの、圧縮された出力チャンクのデータ部分内の利用可能なビットの数は、データ部分内の利用可能なビットの総数（例えば、192 ビット）を非圧縮入力チャンク内の非ゼロバイトの数で割ることによって決定される。この計算の結果は、非圧縮入力チャンク内の非ゼロデータの各バイトを表すために利用可能である圧縮された出力チャンクのデータ部分内のビットの数を示す。如何なる残りのビットも、圧縮された出力チャンクのデータ部分において、非ゼロ値のうちのいくつかを表すための追加ビットを提供するために使用され得る。

【 0 0 3 1 】

非圧縮入力チャンク内の各非ゼロバイトを表すための圧縮された出力チャンクのデータ部分内で利用可能なビットの数が決定されると、非圧縮入力チャンク内の非ゼロ値の LSB が、利用可能なビットの数に収まるよう切り捨てられる。他の実施形態では、非ゼロ値の MSB が切り捨てられてもよい。次いで、切り捨てられた非ゼロ値が、圧縮された出力チャンクのデータ部分に格納され得る。このプロセスは、圧縮されていない入力活性化デー

10

20

30

40

50

タの各チャンクごとに繰り返され得る。次いで、圧縮された出力チャンクが、後の復元及びニューロンによる使用のためにオン又はオフモジュール・メモリに記憶され得る。

【0032】

開示されるDNNモジュールはまた、上述されたように圧縮ユニットによって圧縮された活性化値を復元する復元ユニットを含むこともできる。復元ユニットは、マスク部分及びデータ部分を含む圧縮された活性化データのチャンクを受け取る。復元ユニットは、復元された出力チャンクに存在することになる非ゼロバイトの数と、復元された出力チャンク内のそれらの位置とを特定するために、マスク部分のビットを利用することができる。マスクはまた、復元された出力チャンク内のゼロバイトの位置も示す。

【0033】

いくつかの実施形態において、復元ユニットは、圧縮されたチャンクのデータ部分内の利用可能なビットの総数（例えば、192ビット）を、マスクによって指定されている非圧縮入力チャンク内の非ゼロバイトの数で割ることによって、各非ゼロバイトを表すために圧縮ユニットによって使用されたビットの数を決定する。復元ユニットはまた、圧縮されたチャンクのデータ部分内の非ゼロ値のうちのいくつか（例えば、最初のN個の値）を表す追加ビットを供給するために圧縮ユニットがいずれかの残りのビットを使用したと仮定することもできる。

【0034】

論理0であるマスク内の各ビット位置について、復元ユニットは、ゼロバイトを、復元された出力チャンク内にその対応する位置で挿入することができる。論理1であるマスク内の各ビット位置について、復元ユニットは、圧縮された入力チャンクのデータ部分内の対応する位置からの切り捨てられた非ゼロバイトを、復元された出力チャンク内の対応する位置に挿入する。復元ユニットはまた、非ゼロ値の、必要に応じてLSB又はMSBにおいて、圧縮中に切り捨てられたそれらのビットを置き換えるようゼロを挿入する。

【0035】

いくつかの実施形態において、復元ユニットは、非ゼロの非圧縮値が復元されるときにゼロバイトにならないことを確かにするよう、切り捨てられた非ゼロ値にオフセット値を付加する。次いで、復元された出力チャンクが、ニューロンによる使用のためにオン又はオフモジュール・メモリに記憶され得る。DNNモジュール、圧縮ユニット、及び復元ユニットの動作に関する更なる詳細は、以下で与えられる。

【0036】

続く詳細な説明では、本願の部分を形成し、実例として具体的な構成又は例が示されている添付の図面が、参照される。これより図面を参照すると、図面では、いくつかの図を通して同じ番号が同じ要素を表しており、メモリバス帯域幅の利用量を低減するために活性化データを圧縮及び復元することができるDNNモジュールの態様が記載される。

【0037】

図1は、一実施形態に従って、本明細書で開示される技術の態様を実装するDNNモジュール105の構成及び動作の態様を示すコンピューティングアーキテクチャ図である。本明細書で開示されるDNNモジュール105は、例えば、制限なしに、物体検出、意味的標識付け、及び特徴抽出などの分類問題（及び関連する問題）を解くよういくつかの実施形態で構成される。

【0038】

この機能を提供するために、DNNモジュール105は、リコールオンリー（recall-only）ニューラルネットワークを実装し、多種多様なネットワーク構造をプログラムでサポートすることができる。DNNモジュール105によって実装されるネットワークの教育は、サーバーファーム、データセンター、又は他の適切なコンピューティング環境においてオンラインで実行され得る。DNNを教育する結果は、“重み”又は“カーネル”として知られ得るパラメータの組である。それらのパラメータは、入力に適用可能な変換関数を表し、その結果が分類又は意味的に標識付けされた出力である。

【0039】

10

20

30

40

50

本明細書で開示されるDNNモジュール105は、スーパースカラープロセッサと見なされ得る。DNNモジュール105は、ニューロン105Fと呼ばれる複数の実行ユニットへ1つ以上の命令をディパッチすることができる。実行ユニットは、“同時ディスパッチ同時完了”(simultaneous dispatch simultaneous complete)であることができ、各実行ユニットが他の実行ユニットの夫々と同期する。DNNモジュール105は、単一命令ストリーム多重データストリーム(single instruction stream, multiple data stream) (“SIMD”)アーキテクチャとして分類され得る。

【0040】

DNNモジュール105は、多数のニューロン105F(例えば、2の累乗)を含む。ニューロン105Fは、脳内の生体ニューロンをモデル化するために使用される人工ニューラルネットワーク内の基本単位である。ニューロン105Fのモデルは、活性化関数を用いて入力ベクトルと重みベクトルとの内積にバイアスを加えたものを含むことができる。本明細書中に記載されるDNNモジュール105内のニューロン105Fによって実行される処理は、人工ニューロンに密接にマッピングされる。

10

【0041】

DNNモジュール105内の各ニューロン105Fは、加重加算、最大プーリング(max pooling)、バイパス、及び場合によっては、他のタイプの演算を実行可能である。ニューロン105Fは、クロック周期ごとに入力及び重みデータを処理する。各ニューロン105Fは、DNNモジュール105内のカーネルデータのフローを最小限にするよう、カーネル内の進捗状況に関して、他の全てのニューロン105Fと同期する。

20

【0042】

各ニューロン105Fは、乗算器、加算器、比較器、及び多数の積算器(図1に図示せず。)を含むことができる。複数の積算器を有することによって、ニューロン105Fは、一度に複数の異なるアクティブなカーネルのコンテキストを保持することができる。各積算器は、BaSRAM150(後述される。)の読み出しからロード可能である。積算器は、それら自体を、他のニューロン105Fからの他の積算器の内容と加算することができる。

【0043】

DNNモジュール105は、画像データなどの平面(planar)データを入力として受け入れられる。DNNモジュール105への入力、しかしながら、画像データに限られない。むしろ、DNNモジュール105は、一様平面の形でDNNモジュール105へ与えられる如何なる入力データにも作用することができる。1つの特定の形態では、DNNモジュール105は、入力として、多断面(multi-planar)1バイト又は2バイトデータフレームを受け入れることができる。

30

【0044】

各入力フレームは、 $N \times K \times H \times W$ 組のカーネルにより畳み込まれ得る。ここで、 N は、カーネルの数であり、 K は、カーネルごとのチャネルの数であり、 H は高さであり、 W は幅である。畳み込みは、重なり合ったインターバルで入力データに対して実行され、インターバルは、 X 及び Y 方向におけるストライドによって定義される。これらの関数は、ニューロン105Fによって実行され、DNNモジュール105及びソフトウェア管理下の制御レジスタによって管理される。

40

【0045】

DNNモジュール105は、3つの主なデータタイプ、すなわち、重み、入力データ/特徴マップ、及び活性化データをサポートする。入力データ/特徴マップ及び活性化データは、ほとんどの場合に、レイヤの出力に言及するときに語「活性化データ」が使用されるとの区別により、同じデータについての2つの名称である。レイヤの入力を言及するときには、語「入力データ/特徴マップ」が使用される。

【0046】

DNNモジュール105内のニューロン105Fは、それらの入力の加重和を計算し、加重和を“活性化関数”又は“伝達関数”に通す。伝達関数は、一般にシグモイド関数を有する

50

が、区分線形関数 (piecewise linear function)、ステップ関数、又は他のタイプの関数でも構成され得る。活性化関数は、ニューロン 105 F が、分類境界が非線形である入力及び所望の出力のより大きい組を教育することを可能にする。

【0047】

DNNモジュール105は、ニューラルネットワークのレイヤに対応するレイヤ記述子のリストに作用する。レイヤ記述子のリストは、DNNモジュール105によって命令として教育され得る。それらの記述子は、メモリからDNNモジュール105にプリフェッチされ、順に実行され得る。記述子リストは、DNNモジュール105への命令の組として働く。ソフトウェアツール及び/又はコンパイラが、DNNモジュール105で実行される記述子リストを生成するよう、DNNモジュール105の外部にあるデバイスで実行され得る。

10

【0048】

一般に、記述子には2つの主な種類があり、メモリ間移動 (memory-to-memory move) (“M2M”) 記述子及びオペレーション記述子である。M2M記述子は、オペレーション記述子による消費のためにメインメモリからローカルバッファ (すなわち、後述されるラインバッファ125) へ又はローカルバッファからメインメモリへデータを移動させるために使用され得る。M2M記述子のターゲットパイプラインは、内部DMAエンジン105B又はコンフィグレーションレジスタ105Gであることができ、一方、オペレーション記述子のターゲットパイプラインは、ニューロン105Fであることができる。

【0049】

オペレーション記述子は、ローカルの静的ランダムアクセスメモリ (“SRAM”) に位置するデータ構造に対してニューロン105Fが実行すべき具体的な動作を指定する。オペレーション記述子は、順に処理され、多種多様なレイヤ動作が可能であり、そのうちの少なくともいくつかが本明細書中に記載されている。

20

【0050】

図1に表されているように、DNNモジュール105は、一意のL1及びL2バッファ構造を有するメモリサブシステムを具備する。図1に示されるL1及びL2バッファは、特にニューラルネットワーク処理のために設計される。例として、L2バッファ150は、選択された周波数で動作する高速プライベートインターフェイスとともに、選択された記憶容量を保持することができる。L1バッファ125は、カーネルと活性化データとの間で分けられ得る選択された記憶容量を保持することができる。L1バッファ125は、本明細書で“ラインバッファ125”と呼ばれてもよく、L2バッファ150は、本明細書でBaSRAM150と呼ばれてもよい。

30

【0051】

計算データ (すなわち、入力データ、重み及び活性化データ) は、いくつかの実施形態において、行優先 (row-major) でBaSRAM150に格納される。計算データは、2つのラインバッファとして体系化され得る。一方のラインバッファは入力データを含み、本明細書で“入力バッファ”と呼ばれてよく、他方のラインバッファは、本明細書で“重みバッファ”と呼ばれてよく、カーネル重みを含む。ラインバッファは、ロード/ストアユニット105CによってBaSRAM150から満たされる。データは、各ラインバッファにおいて、その所定の容量に達するまで蓄積される。ラインバッファデータは、次いで、いくつかの実施形態において、シャドウバッファにコピーされ、ニューロン105Fに渡される。

40

【0052】

DNNモジュール105は、プリフェッチユニット105A、セーブ/リストアユニット105E、レイヤコントローラ105D、及びレジスタインターフェイス105Gを含むがこれらに限られない多数の他のコンポーネントを有することもできる。DNNモジュール105は、いくつかの実施形態において、追加又は代替のコンポーネントを含むことができる。

【0053】

50

DNNモジュール105は、いくつかの構成において、他の外部コンピューティングコンポーネントと協働して動作する。例えば、DNNモジュール105は、いくつかの実施形態において、ホストアプリケーションプロセッサ・システムオンチップ(system on chip) (“ホストSOC”)130へ接続される。DNNモジュール105は、例えば、PCIeインターフェイスを通じてホストSOC130へ接続され得る。PCIeエンドポイント135などの適切なPCIeコンポーネントが、それらの接続を可能にするために利用され得る。

【0054】

ホストSOC130は、DNNモジュール105のためのアプリケーションプロセッサとなる。メインオペレーティングシステム、アプリケーション、及び補助センサ処理が、ホストSOC130によって実行される。ホストSOC130はまた、画像データなどの入力データをDNNモジュール105へ供給する外部カメラなどの入力データ源102へも接続され得る。

10

【0055】

DDR DRAM155も、ホストSOC130へ接続され得、メインシステムメモリとして使用され得る。このメモリは、メモリコントローラ145を用いて高帯域幅ファブリック120(例えば、PCIeバス)を越えてホストSOC130からアクセス可能である。高帯域幅ファブリック120は、双方向の直接メモリアクセス(“DMA”)スモールメッセージングトランザクション及びより大きいDMAトランザクションを提供する。ブリッジ115及び低帯域幅ファブリック110は、サブモジュール構成及び他の機能のためにDNNモジュール105をホストSOC130へ接続することができる。

20

【0056】

DNNモジュール105は、メインメモリ155へ及びそれからデータを移動させるよう構成されるDMAエンジン105Bを含むことができる。DMAエンジン105Bは、いくつかの実施形態において、2つのチャンネルを具備する。一方のチャンネルは、フェッチ動作記述子に専用であり、一方、他方のチャンネルは、M2M動作に専用である。本文脈中の記述子は、メモリの中身を移動させるために使用されるDMA記述子であり、上記のオペレーション記述子と混同されるべきではない。

【0057】

ローカルBaSRAM150をオフロードするよう、且つ、入力及び重みデータのための更なる空間を提供するよう、活性化出力は、任意に、DDMメモリ155に直接にストリーミングされ得る。データをDDRメモリ155にストリーミングするとき、DNNモジュール105は、高帯域幅ファブリック120上でのバーストトランザクションのために十分なデータを蓄積することになり、且つ、ニューロン105Fに対するバックプレッシャを最小限とするよう十分なトランザクションをバッファすることになる。DNNモジュール105の動作に関する更なる詳細は、以下で与えられる。

30

【0058】

図2A及び図2Bは、一実施形態に従って、活性化データを圧縮するDNNモジュール105の構成及び動作の態様を示すコンピューティングシステムアーキテクチャ図である。図2Aに示され且つ先に簡単に記載されたように、DNNモジュール105は、1つ以上のニューロン105Fと、圧縮ユニット200とを含む。圧縮ユニット200は、いくつかの実施形態において、ロード/ストアユニット105Cによって実装されるが、他の実施形態では他の方法で実装されてもよい。

40

【0059】

圧縮ユニット200は、ニューロン105Fの1つ以上によって生成される活性化データの非圧縮チャンク202を受け取ることができる。データの非圧縮チャンク202は、いくつかの実施形態において、64バイトなどの一定数のバイトを含む。

【0060】

圧縮ユニット200は、活性化データの圧縮されたチャンク204を生成するようデータの非圧縮チャンク202を圧縮することができる。活性化データの圧縮されたチャンク2

50

04は、次いで、メモリ206に格納され得る。例えば、活性化データの圧縮されたチャック204は、アプリケーションホストによって提供されるLPDDR4メモリ155に格納されても、あるいは、DNNモジュール105によって提供されるBaSRAM150に格納されてもよい。以下で更に詳細に開示されるように、本明細書で開示される技術は、LPDDR4メモリ155又はBaSRAM150から圧縮又は復元された活性化データを記憶し又は読み出すときにメモリバス利用の利用量を低減するために圧縮及び復元を利用することができる。これらの技術に関する更なる詳細は、図2A~9に関して以下で開示される。

【0061】

図2に表されているように、圧縮ユニット200は、データの圧縮された出力チャック204のマスク部分208及びデータ部分21を生成することができる。圧縮された出力チャック204のマスク部分208は、データの非圧縮チャック202内の一定数のバイトに等しいビットの数を含む。例えば、データの非圧縮チャック202が64バイトのデータを含む場合には、圧縮された出力チャック204のマスク部分208は64ビット（すなわち、8バイト）を含むことになる。

10

【0062】

圧縮された出力チャック204のマスク部分208内の各ビットは、いくつかの実施形態において、データの非圧縮チャック202内のバイトに対応する。例えば、マスク部分208のビット1は、データの非圧縮チャック202内の第1バイトに対応することができ、マスク部分208のビット2は、データの非圧縮チャック202内の第2バイトに対応することができる、など。

20

【0063】

圧縮ユニット200は、圧縮された出力チャック204のマスク部分208内の各ビットを、データの非圧縮チャック202内の対応するバイトがゼロバイトである場合に論理0にセットする。圧縮ユニット200はまた、圧縮された出力チャック204のマスク部分208内の各ビットを、データの非圧縮チャック202内の対応するバイトが非ゼロバイトである場合に論理1にセットする。このようにして、圧縮された出力チャック204のマスク部分208は、データの非圧縮チャック202内のゼロバイト及び非ゼロバイトの存在及び位置を符号化する。

【0064】

圧縮ユニット200は、データの非圧縮チャック202内の非ゼロバイトの数を決定することによって、圧縮された出力チャック204のデータ部分210を生成する。次いで、圧縮ユニット200は、データの非圧縮チャック202内の非ゼロバイトの数と、圧縮された出力チャック204のデータ部分210で利用可能なバイトの数とに基づき、データの非圧縮チャック202の各非ゼロバイトを格納するために利用可能である圧縮された出力チャック204のデータ部分210内のビットの数を決定する。例えば、データの圧縮されたチャック204のデータ部分210が24バイト幅（すなわち、192ビット）であり、データの非圧縮チャック202内に47の非ゼロバイトがある場合に、データの非圧縮チャック202から各非ゼロバイトを格納するためには、4ビットがデータ部分210内で利用可能である。

30

40

【0065】

いくつかの実施形態において、圧縮ユニット200はまた、データの非圧縮チャック202の非ゼロバイトを格納するために利用可能である圧縮された出力チャック204のデータ部分210内で、もしあれば、追加ビットの数を決定することもできる。上記の例において、例えば、4つの追加ビットが、非ゼロバイトを格納するために利用可能である（すなわち、 $192 \bmod 47 = 4$ ビット）。圧縮ユニット200は、それらの追加ビットを、データの非圧縮チャック202内の1つ以上の非ゼロバイトを切り捨てる前に、その1つ以上の非ゼロバイトに割り当てることができる。例えば、圧縮ユニット200は、それらの追加ビットを、圧縮された出力チャック204のデータ部分210内の最初のNバイトに割り当ててもよい。

50

【 0 0 6 6 】

圧縮ユニット 2 0 0 は、次いで、データの非圧縮チャンク 2 0 2 内の非ゼロバイトを、各非ゼロバイトを格納するためにデータ部分 2 1 0 内で利用可能な決定されたビット数（すなわち、上記の例では、4）まで切り捨てる。一実施形態において、圧縮ユニット 2 0 0 は、データ部分 2 1 0 内の利用可能なビット数に収まるように非ゼロバイトの L B S を切り捨てる。他の実施形態では、圧縮ユニット 2 0 0 は、非ゼロバイトの M S B を切り捨てる。圧縮ユニット 2 0 0 は、次いで、切り捨てられた非ゼロバイトを、圧縮された出力チャンク 2 0 4 のデータ部分 2 1 0 に格納する。マスク部分 2 0 8 及びデータ部分 2 1 0 を含む圧縮された出力チャンク 2 0 4 は、次いで、例えば、D N N モジュール 1 0 5 内のオンボード・メモリ、又は D N N モジュール 1 0 5 のアプリケーションホストのオフボード・メモリへ、出力され得る。

10

【 0 0 6 7 】

先に簡単に述べられたように、圧縮された出力チャンク 2 0 4 のマスク部分 2 0 8 内の 2 つ以上のビットは、いくつかの実施形態において、データの非圧縮チャンク 2 0 2 内のあるバイトに対応する。かような実施形態で、圧縮された出力チャンク 2 0 4 のマスク部分 2 0 8 内のビットは、非圧縮チャンク 2 0 2 内の対応するバイトだけでなくそのおおよその大きさも示すことができる。例えば、制限なしに、マスク部分 2 0 8 は、データの非圧縮チャンク 2 0 2 内のバイトごとに 2 ビットを含んでよい。この例では、0 0 は、データの非圧縮チャンク 2 0 2 内の対応する非ゼロバイトの M S B が 0 であることを示すことができ、0 1 は、M S B が < 6 4 であることを示すことができ、1 0 は、M S B が < 1 2 8 であることを示すことができ、1 1 は、M S B > 1 2 8 であることを示すことができる。これらの値は、データの非圧縮チャンク 2 0 2 内のバイトのどの M S B が切り捨てられ得るかを特定するために利用され得る。例えば、特定のバイトの M S B が < 6 4 である場合に、上から 2 つの M S B がデータの損失なしで切り捨てられ得る。

20

【 0 0 6 8 】

図 3 は、一実施形態に従って、圧縮されていない活性化データのチャンク 2 0 2 の例を参照して、圧縮されていない活性化データのチャンク 2 0 2 を圧縮するための D N N モジュール 1 0 5 の動作の態様を説明するデータ構造図である。図 3 に示される例では、活性化データの非圧縮チャンク 2 0 2 は、6 4 バイト長である。圧縮されていない活性化データのチャンク 2 0 2 のバイト 0、1 及び 6 3 はゼロバイトである。圧縮されていない活性化データのチャンク 2 0 2 のバイト 2、3 及び 6 2 は非ゼロバイトであり、夫々値 1 1 2、1 2 1 及び 2 を格納する。例となる圧縮されていない活性化データのチャンク 2 0 2 のバイト 4 乃至 6 1 は、ゼロ又は非ゼロバイトを格納することができる。

30

【 0 0 6 9 】

上述されたように、圧縮ユニット 2 0 0 は、活性化データの非圧縮チャンク 2 0 2 内のゼロバイト及び非ゼロバイトの存在及び位置を符号化するマスク部分 2 0 8 を生成することができる。この例では、例えば、マスク部分 2 0 8 のビット 0、1 及び 6 3 が、活性化データの非圧縮チャンク 2 0 2 内の対応する位置にあるゼロバイトの存在を示すよう、論理 0 にセットされている。同様に、マスク部分 2 0 8 のビット 2、3 及び 6 2 は、活性化データの非圧縮チャンク 2 0 2 のバイト 2、3 及び 6 2 が非ゼロバイトを格納していることを示すよう、論理 1 にセットされている。

40

【 0 0 7 0 】

上述されたように、圧縮ユニット 2 0 0 は、データの非圧縮チャンク 2 0 2 内の非ゼロバイトの数を決定することによって、圧縮された出力チャンク 2 0 4 のデータ部分 2 1 0 を生成する。図 3 に示される例では、例えば、データの非圧縮チャンク 2 0 2 は、4 7 個の非ゼロバイト（図 3 では、その全てが示されているわけではない。）を含む。圧縮ユニット 2 0 0 は、次いで、データの非圧縮チャンク 2 0 2 内の非ゼロバイトの数と、圧縮された出力チャンク 2 0 4 のデータ部分 2 1 0 で利用可能なバイトの数とに基づき、データの非圧縮チャンク 2 0 2 の各非ゼロバイトを格納するために利用可能である圧縮された出力チャンク 2 0 4 のデータ部分 2 1 0 内のビットの数を決定する。

50

【 0 0 7 1 】

図 3 に示される例では、例えば、データの圧縮されたチャンク 2 0 4 のデータ部分 2 1 0 は 2 4 バイト幅（すなわち、1 9 2 ビット）であり、データの非圧縮チャンク 2 0 2 内には 4 7 の非ゼロバイトがある。結果として、データの非圧縮チャンク 2 0 2 から各非ゼロバイトを格納するためには、4 ビットがデータ部分 2 1 0 内で利用可能である（すなわち、 $192 / 47 = 4$ 余り 4）。

【 0 0 7 2 】

やはり上述されたように、圧縮ユニット 2 0 0 はまた、データの非圧縮チャンク 2 0 2 の非ゼロバイトを格納するために利用可能である圧縮された出力チャンク 2 0 4 のデータ部分 2 1 0 内で、もしあれば、追加ビットの数を決定することもできる。図 3 に示される例では、例えば、4 つの追加ビットが、非ゼロバイトを格納するために利用可能である（すなわち、 $192 \bmod 47 = 4$ ビット）。圧縮ユニット 2 0 0 は、それらの追加ビットを、データの非圧縮チャンク 2 0 2 内の 1 つ以上の非ゼロバイトを切り捨てる前に、その 1 つ以上の非ゼロバイトに割り当てることができる。図 3 に示される例では、追加ビットのうちの 1 つが、圧縮されていない活性化データのチャンク 2 0 2 の最初の 4 つの非ゼロバイトの夫々に割り当てられている。結果として、活性化データの非圧縮チャンク 2 0 2 の最初の 4 バイトは、4 ではなく 5 ビットに切り捨てられる。

10

【 0 0 7 3 】

圧縮ユニット 2 0 0 は、次いで、データの非圧縮チャンク 2 0 2 内の非ゼロバイトを、各非ゼロバイトを格納するためにデータ部分 2 1 0 内で利用可能な決定されたビット数（すなわち、上記の例では、4 ビット。ただし、最初の 4 つの非ゼロバイトのためには 5 ビット。）まで切り捨てる。図 3 に示される例では、一実施形態において、圧縮ユニット 2 0 0 は、利用可能なビット数に収まるように非ゼロバイトの L B S を切り捨てる。他の実施形態では、圧縮ユニット 2 0 0 は、非ゼロバイトの M S B を切り捨てる。

20

【 0 0 7 4 】

図 3 に示されるように、活性化データの非圧縮チャンク 2 0 2 の第 2 のバイトは値 1 1 3（0 1 1 1 0 0 0 1）を格納する。活性化データの非圧縮チャンク 2 0 2 内の最初の 4 つの非ゼロ値には 5 ビットが割り当てられているので、この値の 3 つの L S B が切り捨てられ、その結果、値 0 1 1 1 0 が、活性化データの圧縮されたチャンク 2 1 0 内の最初の位置に格納される。活性化データの非圧縮チャンク 2 0 2 の第 3 のバイトは値 1 2 1（0 1 1 1 1 0 0 1）を格納する。活性化データの非圧縮チャンク 2 0 2 内の最初の 4 つの非ゼロ値には 5 ビットが割り当てられているので、この値の 3 つの L S B が切り捨てられ、その結果、値 0 1 1 1 1 が、活性化データの圧縮されたチャンク 2 1 0 内の第 2 の位置に格納される。

30

【 0 0 7 5 】

図 3 に示される例では、活性化データの非圧縮チャンク 2 0 2 の 6 2 番目のバイトは値 2（0 0 0 0 0 0 1 0）を格納する。活性化データの非圧縮チャンク 2 0 2 内の 5 番目から 6 3 番目の非ゼロ値には 4 ビットが割り当てられているので、この値の 4 つの L S B が切り捨てられ、その結果、値 0 0 0 0 が、活性化データの圧縮されたチャンク 2 1 0 内の 6 2 番目の位置に格納される。活性化データの圧縮されたチャンク 2 1 0 内の他の非ゼロバイトは、同様にして切り捨てられ、活性化データの圧縮されたチャンク 2 0 4 のデータ部分 2 1 0 に格納され得る。

40

【 0 0 7 6 】

圧縮されていない活性化データのチャンク 2 0 2 の全ての非ゼロバイトがデータ部分 2 1 0 に格納されると、圧縮ユニット 2 0 0 は、例えば、D N N モジュール 1 0 5 内のオンボード・メモリ又は D N N モジュール 1 0 5 のアプリケーションホストのオフボード・メモリに、マスク部分 2 0 8 及びデータ部分 2 1 0 を含む圧縮された出力チャンク 2 0 4 を格納する。圧縮プロセスに関する更なる詳細は、図 4 に関して以下で与えられる。

【 0 0 7 7 】

図 4 は、本明細書で開示される一実施形態に従って、圧縮されていない活性化データのチ

50

チャンク 202 を圧縮するための DNN モジュール 105 の動作の態様を説明するルーチン 400 を示すフロー図である。当然ながら、図 4 及び他の図に関して本明細書中に記載される論理動作は、(1) コンピュータにより実施される動作又はコンピュータ装置で実行されるプログラムモジュールのシーケンスとして、及び/又は(2) コンピュータ装置内の相互接続された機械論理回路又は回路モジュールとして、実装可能である。

【0078】

本明細書で開示される技術の特定の実施は、コンピュータ装置の性能及び他の要件に依存して選択できる問題である。従って、本明細書中に記載される論理動作は、状態、操作、構造的デバイス、動作、又はモジュールと様々に呼ばれる。これらの状態、操作、構造的デバイス、動作、及びモジュールは、ハードウェア、ソフトウェア、ファームウェア、特別目的のデジタルロジック、及びそれらの任意の組み合わせで実装され得る。当然ながら、図示及び本明細書中に記載されているよりも多くの又は少ない動作が実行されてよい。これらの動作は、本明細書中に記載されているのとは異なる順序で実行されてもよい。

10

【0079】

ルーチン 400 は動作 402 から開始する。402 で、圧縮ユニット 200 は、活性化データの圧縮されたチャンク 210 内の非ゼロバイトの数を決定する。ルーチン 400 は次いで、動作 404 へ進む。404 で、圧縮ユニット 200 は、活性化データの圧縮されたチャンク 210 内の非ゼロバイトの数が、活性化データの圧縮されたチャンク 204 のデータ部分 210 で利用可能なバイトの数以下であるかどうかを判定する。活性化データの非圧縮チャンク 202 の非ゼロバイトの数が、活性化データの圧縮されたチャンク 204 のデータ部分 210 で利用可能なバイトの数以下である場合には、非ゼロバイトは圧縮される必要がない。従って、この場合に、ルーチン 400 は動作 408 へ進む。408 で、非ゼロバイトは、切り捨てなしでデータ部分 210 に格納される。

20

【0080】

活性化データの非圧縮チャンク 202 内の非ゼロバイトの数が、活性化データの圧縮されたチャンク 204 のデータ部分 210 で利用可能なバイトの数よりも多い場合には、ルーチン 400 は動作 406 から動作 412 へ進む。動作 412 で、圧縮ユニット 200 は、上述されたようにして、活性化データの非圧縮チャンク 202 の切り捨てられた非ゼロバイトを格納するために利用可能な出力データの圧縮されたチャンク 204 のデータ部分 210 のビットの数を決定する。次いで、ルーチン 400 は動作 412 から動作 414 へ進む。

30

【0081】

動作 414 で、圧縮ユニット 200 は、データの非圧縮チャンク 202 の非ゼロバイトを格納するために利用可能である圧縮された出力チャンク 204 のデータ部分 210 内で、もしあれば、追加ビットの数を決定する。上述されたように、圧縮ユニット 200 は、それらの追加ビットを、データの非圧縮チャンク 202 内の非ゼロバイトの 1 つ以上を切り捨てる前に、その 1 つ以上の非ゼロバイトに割り当てることができる。これは動作 416 で行われる。

【0082】

動作 416 から、ルーチン 400 は動作 418 へ進む。418 で、圧縮ユニット 200 は、活性化データの圧縮されたチャンク 204 のマスク部分 208 内のビットを、活性化データの非圧縮チャンク 202 内の対応するバイトが非ゼロである場合に、論理 1 にセットする。圧縮ユニット 200 はまた、活性化データの圧縮されたチャンク 204 のマスク部分 208 内のビットを、活性化データの非圧縮チャンク 202 内の対応するバイトがゼロである場合に、論理 0 にセットする。

40

【0083】

動作 418 から、ルーチン 400 は次いで動作 420 へ進む。420 で、圧縮ユニット 200 は、データの非圧縮チャンク 202 内の非ゼロバイトの LSB 又は MSB を、各非ゼロバイトごとにデータ部分 210 内で利用可能な決定されたビット数まで切り捨てる。切り捨てられた非ゼロバイトは、次いで、活性化データの圧縮されたチャンク 204 のデー

50

タ部分 210 に格納される。圧縮ユニット 200 は次いで、マスク部分 208 及びデータ部分 210 を含む圧縮された出力チャンク 204 を、DNN モジュール 105 のオンボード・メモリ又は DNN モジュール 105 のアプリケーションホストのオフボード・メモリに格納する。動作 408 及び 420 から、メモリ 400 は動作 410 へ進み、終了する。

【0084】

図 5A 及び 5B は、一実施形態に従って、圧縮された活性化データを復元するための DNN モジュール 105 の構成及び動作の態様を示すコンピューティングシステムアーキテクチャ図である。先に簡単に説明されたように、且つ、図 5A 及び 5B に示されるように、DNN モジュール 105 はまた、復元ユニット 500 も含むことができ、復元ユニット 500 は、上述されたようにして圧縮された活性化データのチャンク 204 を復元することができる。

10

【0085】

例えば、復元ユニット 500 は、DNN プロセッサ内のメモリ又はアプリケーションホストのメモリなどのストレージ 206 から活性化データの圧縮されたチャンク 204 を受け取ることができる。復元ユニット 500 は、次いで、圧縮されたチャンク 204 のマスク部分 208 内の論理真ビットの数に基づき、データの圧縮されたチャンク 204 のデータ部分 210 内の非ゼロバイトの数を決定することができる。復元ユニット 500 はまた、圧縮された出力チャンク 204 のマスク部分 208 内の論理真ビットの位置に基づき、データの復元されたチャンク 502 内の非ゼロバイトの位置を決定することもできる。復元ユニット 500 は、同様にして、データの復元されたチャンク 502 内のゼロバイトの位置を決定することができる。

20

【0086】

復元ユニット 500 はまた、切り捨てられた非ゼロバイトの夫々を圧縮された出力チャンク 204 のデータ部分 210 に格納するために圧縮ユニット 200 によって使用されたビットの数を決定することもできる。復元ユニット 500 は、データの圧縮されたチャンク 204 内の非ゼロバイトの数（マスク部分 208 によって示される。）と、復元された出力チャンク 502 の目標サイズとに基づき、切り捨てられた各非ゼロバイトを格納するために使用されたビットの数を決定することができる。

【0087】

上記の例では、例えば、データの圧縮されたチャンク 204 のデータ部分が 24 バイト幅（すなわち、192 ビット）であり、データの非圧縮チャンク 202 内に 47 の非ゼロバイトがある場合に、これは、圧縮ユニット 200 が、データの非圧縮チャンク 202 の切り捨てられた各非ゼロバイトを格納するためにデータ部分 210 において 4 ビットを利用したことを意味する。復元ユニット 500 はまた、もしあれば、圧縮ユニット 200 が圧縮された出力チャンク 204 のデータ部分 210 に格納された切り捨てられた非ゼロバイトのうちの 1 つ以上に割り当てた追加ビットの数を決定することもできる。

30

【0088】

論理 0 である圧縮された出力チャンク 204 のマスク部分 208 内の各ビット位置について、復元ユニット 500 は、ゼロバイトを、復元された出力チャンク 502 の対応する部分に挿入する。論理 1 であるマスク部分 208 内の各ビット位置について、復元ユニット 500 は、圧縮された出力チャンク 204 のデータ部分 210 内の対応する位置からの切り捨てられた非ゼロバイトを、圧縮された出力チャンク 204 の圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに、復元された出力チャンク 502 の対応する位置に挿入する。ゼロビットは、どのビットが圧縮中に切り捨てられたかに応じて、切り捨てられた非ゼロバイトの LSB 又は MSB に挿入され得る。

40

【0089】

上述されたように、いくつかの実施形態において、復元ユニット 500 はまた、復元された出力チャンク 502 に格納された切り捨てられた非ゼロバイトのうちの 1 つ以上にオフセット（例えば、00000001）を付加する。例えば、オフセットは、圧縮後にゼロバイトになるデータの非圧縮チャンク 2 - 2 の非ゼロバイトに加えられ得る。このように

50

して、非ゼロバイトは、復元時にゼロバイトにならない。

【 0 0 9 0 】

図 6 は、一実施形態に従って、圧縮された活性化データのチャンクの例を参照して、活性化データを復元するための D N N モジュール 1 0 5 の動作の態様を説明するデータ構造図である。図 6 に示される例は、図 3 に関して上述された例で生成された圧縮された活性化データ 2 0 4 の復元について説明する。図 6 に示されるように、マスク部分 2 0 8 は、ビット 0、1 及び 6 3 に 0 を格納し、ビット 2、3 及び 6 2 に 1 を格納する。データ部分 2 1 0 は、図 6 に示されるように値 0 1 1 1 0、0 1 1 1 1 及び 0 0 0 0 を格納する。

【 0 0 9 1 】

復元ユニット 5 0 0 が上記の処理動作を実行する場合に、マスク部分 2 0 8 の最初のビット位置にある論理 0 は、復元ユニット 5 0 0 に、活性化データの復元されたチャンク 5 0 2 の最初のバイトとしてゼロバイトを格納させることになる。同様に、マスク部分 2 0 8 の第 2 のビット位置にある論理 0 は、復元ユニット 5 0 0 に、データの復元されたチャンク 5 0 2 の第 2 のバイトとしてゼロバイトを格納させることになる。

10

【 0 0 9 2 】

マスク部分 2 0 8 の第 3 のビット位置にある論理 1 は、復元ユニット 5 0 0 に、データ部分 2 1 0 の最初の 5 ビット（すなわち、0 1 1 1 0）を読み出させ且つ 3 つの L S B を挿入させて、その結果、値 0 1 1 1 0 0 0 0（1 1 2）が活性化データの復元されたチャンク 5 0 2 の第 3 のバイトとして格納されることになる。同様に、マスク部分 2 0 8 の第 4 のビット位置にある論理 1 は、復元ユニット 5 0 0 に、データ部分 2 1 0 の第 2 の 5 ビット（すなわち、0 1 1 1 1）を読み出させ且つ 3 つの L S B を挿入させて、その結果、値 0 1 1 1 1 0 0 0（1 2 0）が活性化データの復元されたチャンク 5 0 2 の第 4 のバイトとして格納されることになる。

20

【 0 0 9 3 】

マスク部分 2 0 8 の 6 2 番目のビット位置にある論理 1 は、復元ユニット 5 0 0 に、データ部分 2 1 0 の最後の 4 つのビット（すなわち、0 0 0 0）を読み出させ且つ 4 つのゼロ L S B を挿入させ、その結果、0 の値が活性化データの復元されたチャンク 5 0 4 の 6 2 番目のバイト部分に格納されることになる。マスク部分 2 0 8 の最後のビット位置にある論理 0 は、復元ユニット 5 0 0 に、データの復元されたチャンク 5 0 2 の最後のバイトとしてゼロバイトを格納させることになる。

30

【 0 0 9 4 】

上述されたように、復元ユニット 5 0 0 は、活性化データの復元されたチャンク 5 0 2 内の特定のバイトにオフセット値を加えることができる。例えば、復元ユニット 5 0 0 は、活性化データの非圧縮チャンク 2 0 2 内で非ゼロであったが、活性化データの圧縮されたチャンク 2 0 4 内でゼロバイトに圧縮されたバイトに、0 0 0 0 0 0 0 1 などのオフセット値を付加することができる。

【 0 0 9 5 】

図 6 に示される例では、データ部分 2 1 0 内の最後のバイトは、活性化データの非圧縮チャンク 2 0 2 では非ゼロ（すなわち、2）であったが、活性化データの圧縮されたチャンク 2 0 4 ではゼロになった。従って、復元ユニット 5 0 0 は、このバイトに 0 0 0 0 0 0 0 1 などのオフセット値を付加し、それによって、活性化データの非圧縮チャンク 2 0 2 内の非ゼロバイトがゼロバイトに圧縮されないことを確かに行うことができる。

40

【 0 0 9 6 】

図 7 は、本明細書で開示される一実施形態に従って、活性化データを復元するための D N N モジュール 1 0 5 の動作の態様を説明するルーチン 7 0 0 を示すフロー図である。ルーチン 7 0 0 は動作 7 0 2 から開始し、7 0 2 で、復元ユニット 5 0 0 は、活性化データの復元されたチャンク 5 0 2 内の非ゼロバイトの数及びそれらの位置を決定するために、活性化データの圧縮されたチャンク 2 0 4 のマスク部分 2 0 8 を利用する。

【 0 0 9 7 】

ルーチン 7 0 0 は、動作 7 0 2 から動作 7 0 4 へ進み、動作 7 0 4 で、復元ユニット 5 0

50

0 は、活性化データの圧縮されたチャンク内の非ゼロバイトの数が活性化データの復元されたチャンク 5 0 2 のバイトの数以下であるかどうかを判定する。上述されたように、活性化データの圧縮されたチャンク 2 0 4 の非ゼロバイトは、その非ゼロバイトの数が活性化データの復元されたチャンク 5 0 2 のバイトの数以下である場合には、復元される必要がない。従って、この場合に、ルーチン 7 0 0 は動作 7 0 8 へ進み、7 0 8 で、活性化データの圧縮されたチャンク 2 0 4 内の非ゼロバイトは、変更なしで、活性化データの復元されたチャンク 5 0 2 に格納される。

【 0 0 9 8 】

活性化データの圧縮されたチャンク内の非ゼロバイトの数が活性化データの復元されたチャンク 5 0 2 内のバイトの数よりも多い場合には、ルーチン 7 0 0 は、動作 7 0 6 から動作 7 1 2 へ進む。動作 7 1 2 で、復元ユニット 5 0 0 は、活性化データの非圧縮チャンク 2 0 2 の切り捨てられた各非ゼロバイトを格納するために圧縮ユニット 2 0 0 が使用した出力データの圧縮されたチャンク 2 0 4 のデータ部分 2 1 0 のビットの数を決定する。ルーチン 7 0 0 は次いで、上述されたように動作 7 1 2 から動作 7 1 4 へ進む。

【 0 0 9 9 】

動作 7 1 4 で、復元ユニット 5 0 0 は、もしあれば、データの非圧縮チャンク 2 0 2 の非ゼロバイトを格納するために使用された追加ビットの数を決定する。復元ユニット 5 0 0 は、上述されたように、それらの追加ビットを、データの復元されたチャンク 5 0 2 内の非ゼロバイトのうちの 1 つ以上に割り当てることができる。

【 0 1 0 0 】

動作 7 1 6 から、ルーチン 7 0 0 は動作 7 1 8 へ進み、動作 7 1 8 で、復元ユニット 5 0 0 は、論理 0 である圧縮された出力チャンク 2 0 4 のマスク部分 2 0 8 内の各ビット位置について、ゼロバイトを、復元された出力チャンク 5 0 2 の対応する位置に挿入する。論理 1 である圧縮された出力チャンク 2 0 4 のマスク部分 2 0 8 内の各ビット位置について、復元ユニット 5 0 0 は、圧縮された出力チャンク 2 0 4 の対応する位置からの切り捨てられた非ゼロバイトを、圧縮された出力チャンク 2 0 4 の圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに、復元された出力チャンク 5 0 2 の対応する位置に挿入する。ゼロビットは、どのビットが圧縮中に切り捨てられたかに応じて、切り捨てられた非ゼロバイトの L S B 又は M S B に挿入され得る。これは、動作 7 2 0 で行われる。

【 0 1 0 1 】

いくつかの実施形態において、復元ユニット 5 0 0 はまた、復元された出力チャンク 5 0 2 に格納された切り捨てられた非ゼロバイトのうちの 1 つ以上にオフセット値を付加することができる。例えば、オフセットは、圧縮後にゼロバイトになるデータの非圧縮チャンク 2 0 2 の非ゼロバイトに加えられ得る。このようにして、非ゼロバイトは、圧縮及び復元される場合にゼロバイトにならない。他の実施形態では、オフセットは、活性化データの復元されたチャンク 5 0 2 内の全てのバイトに付加され得る。

【 0 1 0 2 】

復元ユニット 5 0 0 は次いで、復元された出力チャンク 5 0 2 を、ニューロン 1 0 5 F による使用のために、D N N モジュール 1 0 5 のオンボード・メモリ又は D N N モジュール 1 0 5 のアプリケーションホストのオフボード・メモリに格納する。動作 7 0 8 及び 7 2 0 から、ルーチン 7 0 0 は動作 7 1 0 へ進み、終了する。

【 0 1 0 3 】

図 8 は、本明細書で提示される D N N モジュール 1 0 5 のためのアプリケーションホストとして動作することができるコンピュータ装置のためのコンピュータハードウェア及びソフトウェアアーキテクチャを例示するコンピュータアーキテクチャ図である。特に、図 8 に表されているアーキテクチャは、サーバーコンピュータ、携帯電話機、電子書籍リーダー、スマートフォン、デスクトップコンピュータ、AR / VR デバイス、タブレットコンピュータ、ラップトップコンピュータ、又は D N N モジュール 1 0 5 との使用に適した他のタイプのコンピュータ装置を実装するために利用され得る。

【 0 1 0 4 】

10

20

30

40

50

図 8 に表されているコンピュータ 800 は、中央演算処理装置 802 (“ CPU ”) と、ランダムアクセスメモリ (“ RAM ”) 及びリードオンリーメモリ (“ ROM ”) を含むシステムメモリ 804 と、メモリ 804 を CPU 802 へ結合するシステムバス 810 とを含む。起動中などのように、コンピュータ 800 内の要素間で情報を伝送するのを助ける基本ルーチンを含む基本入出力システム (“ BIOS ” 又は “ ファームウェア ”) が、ROM 808 に記憶され得る。コンピュータ 800 は、オペレーティングシステム 822、アプリケーションプログラム、及び他のタイプのプログラムを記憶する大容量記憶装置 812 を更に含む。大容量記憶装置 812 はまた、他のタイプのプログラム及びデータも記憶するよう構成され得る。

【 0105 】

大容量記憶装置 812 は、バス 810 へ接続されている大容量記憶コントローラ (図示せず。) を通じて CPU 802 へ接続されている。大容量記憶装置 812 及びその関連するコンピュータ可読媒体は、コンピュータ 800 のための不揮発性記憶を提供する。本明細書に含まれているコンピュータ可読媒体についての記載は、ハードディスク、CD-ROM ドライブ、又は USB ストレージキーなどの大容量記憶装置を指すが、当業者には当然ながら、コンピュータ可読媒体は、コンピュータ 800 によってアクセス可能な如何なる利用可能なコンピュータ記憶媒体又は通信媒体でもあることができる。

【 0106 】

通信媒体は、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、又は搬送波若しくは他の伝送メカニズムなどの変調データ信号内の他のデータを含み、如何なる配信媒体も含む。語「変調データ信号」は、信号内の情報を符号化するように変更又は設定されたその特徴の 1 つ以上が変更又は設定されている信号を意味する。例として、制限なしに、通信媒体は、有線ネットワーク又は直接配線接続などの有線媒体と、音響、無線周波数、赤外線及び他の無線媒体等の無線媒体とを含む。上記のうちのいずれかの組み合わせも、コンピュータ可読媒体の範囲内に含まれるべきである。

【 0107 】

例として、制限なしに、コンピュータ記憶媒体は、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、又は他のデータなどの情報の記憶のための如何なる方法及び技術でも実装される揮発性及び不揮発性のリムーバブル及び非リムーバブル媒体を含むことができる。例えば、コンピュータ記憶媒体は、RAM、ROM、EPROM、EEPROM、フラッシュメモリ若しくは他のソリッドステートメモリ技術、CD-ROM、デジタルバーサタイルディスク (“ DVD ”)、HD-DVD、BLUE-RAY、若しくは他の光学ストレージ、磁気カセット、磁気テープ、磁気ディスクストレージ若しくは他の磁気記憶デバイス、又は所望の情報を記憶するために使用可能であって、コンピュータ 800 によってアクセス可能であるあらゆる他の媒体を含むが、これらに限られない。特許請求の範囲のために、語句「コンピュータ記憶媒体」及びその変形は、波若しくは信号自体又は通信媒体を含まない。

【 0108 】

様々な構成によれば、コンピュータ 800 は、ネットワーク 820 などのネットワークを通じた遠隔のコンピュータへの論理接続を使用するネットワーク化された環境で動作することができる。コンピュータ 800 は、バス 810 へ接続されたネットワークインターフェイスユニット 816 を通じてネットワーク 820 へ接続することができる。当然ながら、ネットワークインターフェイスユニット 816 は、他のタイプのネットワーク及び遠隔コンピュータシステムへ接続するためにも利用可能である。コンピュータ 800 はまた、キーボード、マウス、タッチ入力、電子スタイラス (図 8 に図示せず。)、又はビデオカメラなどの物理センサを含む多数の他のデバイスから入力を受け取って処理する入出力 (“ I/O ”) コントローラ 818 も含むことができる。同様に、入出力コントローラ 818 は、表示スクリーン又は他のタイプの出力デバイス (図 8 に図示せず。) へ出力を供給することができる。

【 0109 】

当然ながら、本明細書で記載されるソフトウェアコンポーネントは、CPU 802にロードされ実行される場合に、CPU 802及びコンピュータ800の全体を、汎用のコンピュータ装置から、本明細書で提示される機能を容易にするようカスタマイズされた特別目的のコンピュータ装置へ変えることができる。CPU 802は、個々に又は集合的に任意数の状態を前提とすることができる任意数のトランジスタ又は他のディスクリート回路素子から構成され得る。より具体的には、CPU 802は、本明細書で開示されるソフトウェアモジュール内に含まれる実行可能命令に応答して有限状態機械として動作することができる。これらのコンピュータ実行可能命令は、如何にしてCPU 802が状態間を遷移するかを特定し、それによってCPU 802を構成するトランジスタ又は他のディスクリートハードウェア素子を変化させることによって、CPU 802を変化させることができる。

10

【0110】

本明細書で提示されるソフトウェアモジュールを符号化することは、本明細書で提示されるコンピュータ可読媒体の物理構造も変えることができる。物理構造の具体的な変形は、本明細書の種々の実施において、様々な因子に依存する。かような因子の例には、コンピュータ可読媒体を実装するために使用される技術、コンピュータ可読媒体が一次又は二次記憶と見なされるかどうか、などがあるが、これらに限られない。例えば、コンピュータ可読媒体が半導体に基づくメモリとして実装される場合には、本明細書で開示されるソフトウェアは、半導体メモリの物理状態を変えることによってコンピュータ可読媒体上で符号化され得る。例えば、ソフトウェアは、半導体メモリを構成するトランジスタ、キャパシタ、又は他のディスクリート回路素子の状態を変えることができる。ソフトウェアはまた、かような構成要素の物理状態を、それにデータを格納するために変えることもできる。

20

【0111】

他の例として、本明細書で開示されるコンピュータ可読媒体は、磁気又は光学技術を用いて実装可能である。かような実施において、本明細書で提示されるソフトウェアは、ソフトウェアが磁気又は光学媒体で符号化されるときに、磁気又は光学媒体の物理状態を変えることができる。これらの変形は、所与の磁気媒体内の特定の位置の磁気特性を変更することを含むことができる。これらの変形は、所与の光学媒体内の特定の位置の物理的な特徴又は特性を変更してそれらの位置の光学特性を変えることを含むこともできる。物理媒体の他の変形は、本明細書の適用範囲及び精神から外れることなしに可能であり、上記の例は、本議論を助けるためにのみ与えられている。

30

【0112】

上記を鑑みて、当然ながら、本明細書で提示されるソフトウェアコンポーネントを記憶及び実行するために、多数のタイプの物理変形がコンピュータ800において行われる。また、当然ながら、コンピュータ800に関して図8に示されるアーキテクチャ、又は同様のアーキテクチャは、手持ち式コンピュータ、ビデオゲーム機、埋込型コンピュータシステム、スマートフォン、タブレット、及びAR/VR装置などのモバイル機器、並びに当業者に知られている他のタイプのコンピュータ装置を含む他のタイプのコンピュータ装置を実装するために利用可能である。また、コンピュータ800は、図8に示されている全てのコンポーネントを含まなくてもよく、図8に明示的に示されていない他のコンポーネントを含むことができ、あるいは、図8に示されているのとは全く異なるアーキテクチャを利用することができる。

40

【0113】

図9は、本明細書で提示される様々な実施形態に従って、開示される技術の態様が実装され得る分散ネットワークコンピューティング環境900帯域幅を説明するネットワーク図である。図9に示されるように、1つ以上のサーバーコンピュータ900Aは、通信ネットワーク820(固定配線若しくは無線LAN、WAN、イントラネット、エクストラネット、ピア・ツー・ピアネットワーク、仮想プライベートネットワーク、インターネット、Bluetooth(登録商標)通信網、固有低電圧通信網、又は他の通信ネットワークのいずれか1つ又は組み合わせであってよい。)を介して、例えば、制限なしに、タブ

50

レットコンピュータ900B、ゲーム機900C、スマートウォッチ900D、スマートフォンなどの電話機900E、パーソナルコンピュータ900F、及びAR/VRデバイス900Gなどの多数のクライアントコンピュータ装置と相互接続され得る。

【0114】

例えば、通信ネットワーク820がインターネットであるネットワーク環境で、サーバーコンピュータ900Aは、ハイパーテキスト転送プロトコル(Hypertext Transfer Protocol) (“HTTP”)、ファイル転送プロトコル(File Transfer Protocol) (“FTP”)、又はシンプル・オブジェクト・アクセス・プロトコル(Simple Object Access Protocol) (“SOAP”)などの多数の既知のプロトコルのいずれかを介して、クライアントコンピュータ装置900B~900Gへのデータ及びそれらからのデータを処理及び通信するよう動作可能な専用サーバーコンピュータであることができる。更に、ネットワーク化されたコンピューティング環境900は、セキュアド・ソケット・レイヤ(Secured Socket Layer) (“SSL”)又はプリティ・グッド・プライバシー(Pretty Good Privacy) (“PGP”)などの様々なデータセキュリティプロトコルを利用することができる。クライアントコンピュータ装置900B~900Gの夫々は、サーバーコンピュータ900Aへのアクセスを得るようにウェブブラウザ(図9に図示せず。)などの1つ以上のコンピューティングアプリケーション若しくは端末セッション、又は他のグラフィカルユーザインターフェイス(図9に図示せず。)、又はモバイルデスクトップ環境(図9に図示せず。)をサポートするよう動作可能なオペレーティングシステムを装備され得る。

10

20

【0115】

サーバーコンピュータ900Aは、他のコンピューティング環境(図9に図示せず。)へ通信上結合され、参加ユーザのインタラクション/リソースネットワークに関するデータを受信することができる。実例となる動作では、ユーザ(図9に図示せず。)は、所望のデータを得よう及び/又は他のコンピューティングアプリケーションを実行するようクライアントコンピュータ装置900B~900Gで実行されるコンピューティングアプリケーションと相互作用してもよい。

【0116】

データ及び/又はコンピューティングアプリケーションは、サーバ900A又は複数のサーバ900Aに記憶され、実例となる通信ネットワーク820上でクライアントコンピュータ装置900B~900Gを通じて協調するユーザへ送られ得る。参加ユーザ(図9に図示せず。)は、完全に又は部分的にサーバーコンピュータ900Aに格納された特定のデータ及びアプリケーションへのアクセスを要求してもよい。これらのデータは、処理及び記憶のためにクライアントコンピュータ装置900B~900Gとサーバーコンピュータ900Aとの間で通信されてもよい。

30

【0117】

サーバーコンピュータ900Aは、データ及びアプリケーションの生成、認証、暗号化、及び通信のためのコンピューティングアプリケーション、プロセス及びアプレットをホストすることができ、他のサーバーコンピューティング環境(図9に図示せず。)、第三者サービスプロバイダ(図9に図示せず。)、ネットワーク・アタッチト・ストレージ(Network Attached Storage) (“NAS”)及びストレージ・エリア・ネットワーク(Storage Area Network(s)) (“SAN”)と協調してアプリケーション/データトランザクションを実現し得る。

40

【0118】

当然ながら、図8に示されているコンピューティングアーキテクチャ及び図9に示されている専用ネットワークコンピューティング環境は、議論を簡単にするために簡略化されている。また、当然ながら、コンピューティングアーキテクチャ及び専用コンピューティングネットワークは、本明細書で具体的に記載されていない多くの更なるコンピューティングコンポーネント、デバイス、ソフトウェアプログラム、ネットワークデバイス、及び他のコンポーネントを含み、利用することができる。

50

【 0 1 1 9 】

本明細書で提示される開示は、以下の付記に記載されている対象も包含する。

【 0 1 2 0 】

付記 1 . ニューラルネットワークプロセッサであって、1つ以上のニューロンと、圧縮ユニットとを有し、前記圧縮ユニットは、

当該ニューラルネットワークプロセッサ内の前記ニューロンの少なくとも1つによって生成されるデータの非圧縮チャンクを受け取り、該データの非圧縮チャンクが一定数のバイトを含み；

圧縮された出力チャンクのマスク部分を生成し、該マスク部分が、前記データの非圧縮チャンク内の前記一定数のバイトに等しいビットの数を含み、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内のバイトに対応し、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内の対応するバイトがゼロである場合に論理 0 に設定され、前記データの非圧縮チャンク内の対応するバイトが非ゼロである場合に論理 1 に設定され；前記データの非圧縮チャンク内の非ゼロバイトの数を決定することと、前記データの非圧縮チャンク内の前記非ゼロバイトの数に基づき、前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記圧縮された出力チャンクのデータ部分内のビットの数を決定することと、該決定されたビットの数まで前記データの非圧縮チャンク内の前記非ゼロバイトを切り捨てることと、該切り捨てられた非ゼロバイトを前記圧縮された出力チャンクの前記データ部分に格納することとによって、前記圧縮された出力チャンクの前記データ部分を生成し；

前記マスク部分及び前記データ部分を含む前記圧縮された出力チャンクを出力するよう構成される、ニューラルネットワークプロセッサ。

【 0 1 2 1 】

付記 2 . 当該ニューラルネットワークプロセッサは、復元ユニットを更に有し、該復元ユニットは、

前記圧縮された出力チャンクを受け取り；

前記圧縮された出力チャンクの前記マスク部分に基づき、前記データの非圧縮チャンク内の非ゼロバイトの数を決定し；

前記圧縮された出力チャンクの前記マスク部分に基づき、前記データの非圧縮チャンク内の前記非ゼロバイトの位置を決定し；

前記圧縮された出力チャンクの前記データ部分に前記切り捨てられた非ゼロバイトを格納するために前記圧縮ユニットによって使用されたビットの数を決定し；

論理 0 である前記圧縮された出力チャンクの前記マスク部分内の各ビット位置について、ゼロバイトを、復元された出力チャンクの対応する位置に挿入し；

論理 1 である前記マスク部分内の各ビット位置について、前記圧縮された出力チャンクの前記データ部分内の対応する位置からの前記切り捨てられた非ゼロバイトを、前記復元された出力チャンクの対応する位置に、前記圧縮された出力チャンクの圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに挿入する

よう構成される、付記 1 に記載のニューラルネットワークプロセッサ。

【 0 1 2 2 】

付記 3 . 前記圧縮ユニットは更に、

前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記圧縮された出力チャンクの前記データ部分内の追加ビットの数を決定し；

前記データの非圧縮チャンク内の前記非ゼロバイトのうちの1つ以上の非ゼロバイトを切り捨てる前に、該1つ以上の非ゼロバイトに前記追加ビットを割り当てる

よう構成される、付記 1 及び 2 のいずれかに記載のニューラルネットワークプロセッサ。

【 0 1 2 3 】

付記 4 . 前記復元ユニットは更に、前記圧縮された出力チャンクの前記データ部分内に格納された前記1つ以上の非ゼロバイトに割り当てられている前記追加ビットの数を決定するよう構成される、付記 1 乃至 3 のうちいずれかに記載のニューラルネットワークプロセ

10

20

30

40

50

ッサ。

【 0 1 2 4 】

付記 5 . 前記復元ユニットは更に、前記復元された出力チャンクに格納された前記切り捨てられた非ゼロバイトのうちの一つ以上の非ゼロバイトにオフセットを加えるよう構成される、付記 1 乃至 4 のうちいずれかに記載のニューラルネットワークプロセッサ。

【 0 1 2 5 】

付記 6 . 前記非ゼロバイトの一つ以上の最下位ビット (L S B) が切り捨てられる、付記 1 乃至 5 のうちいずれかに記載のニューラルネットワークプロセッサ。

【 0 1 2 6 】

付記 7 . 前記非ゼロバイトの一つ以上の最上位ビット (M S B) が切り捨てられる、付記 1 乃至 6 のうちいずれかに記載のニューラルネットワークプロセッサ。

10

【 0 1 2 7 】

付記 8 . ニューラルネットワークプロセッサであって、一つ以上のニューロンと、復元ユニットとを有し、前記復元ユニットは、

マスク部分及びデータ部分を含むデータの圧縮されたチャンクを受け取り；

前記マスク部分内のビットに基づき、データの復元されたチャンク内の非ゼロバイトの数を決定し；

前記非ゼロバイトの数に少なくとも部分的に基づき、前記データの圧縮されたチャンクの前記データ部分に切り捨てられた非ゼロバイトを格納するために使用されたビットの数を決定し；

20

論理 0 である前記データの圧縮されたチャンクの前記マスク部分内の各ビット位置について、ゼロバイトを、前記データの復元されたチャンクの対応する部分に挿入し；

論理 1 である前記データの圧縮されたチャンクの前記マスク部分内の各ビット位置について、前記データの圧縮されたチャンクの前記データ部分内の対応する位置からの前記切り捨てられた非ゼロバイトを、前記データの復元されたチャンクの対応する位置に、前記データの圧縮されたチャンクの圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに挿入する

よう構成される、ニューラルネットワークプロセッサ。

【 0 1 2 8 】

付記 9 . 圧縮ユニットを更に有し、該圧縮ユニットは、

当該ニューラルネットワークプロセッサ内の前記ニューロンの少なくとも一つによって生成されるデータの非圧縮チャンクを受け取り、該データの非圧縮チャンクが一定数のバイトを含み；

30

前記データの圧縮されたチャンクの前記マスク部分を生成し、該マスク部分が、前記データの非圧縮チャンク内の前記一定数のバイトに等しいビットの数を含み、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内のバイトに対応し、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内の対応するバイトがゼロである場合に論理 0 を有し、前記データの非圧縮チャンク内の対応するバイトが非ゼロである場合に論理 1 を有し；

前記データの非圧縮チャンク内の非ゼロバイトの数を決定することと、前記データの非圧縮チャンク内の前記非ゼロバイトの数に基づき、前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記データの圧縮されたチャンクの前記データ部分内のビットの数を決定することと、該決定されたビットの数まで前記データの非圧縮チャンク内の前記非ゼロバイトを切り捨てることと、該切り捨てられた非ゼロバイトを前記データの圧縮されたチャンクの前記データ部分に格納することとによって、前記データの圧縮されたチャンクの前記データ部分を生成し；

40

前記マスク部分及び前記データ部分を含む前記データの圧縮されたチャンクを出力するよう構成される、付記 8 に記載のニューラルネットワークプロセッサ。

【 0 1 2 9 】

付記 1 0 . 前記圧縮ユニットは更に、前記データの非圧縮チャンク内の非ゼロバイトの数

50

が前記データの圧縮されたチャンクの前記データ部分内のバイトの数以下である場合には切り捨てなしで、前記データの非圧縮チャンク内の前記非ゼロバイトを前記データの圧縮されたチャンクの前記データ部分に格納するよう構成される、請求項 8 及び 9 のいずれかに記載のニューラルネットワークプロセッサ。

【 0 1 3 0 】

付記 1 1 . 前記圧縮ユニットは更に、

前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記データの圧縮されたチャンクの前記データ部分内の追加ビットの数を決定し；

前記データの非圧縮チャンク内の前記非ゼロバイトのうちの 1 つ以上の非ゼロバイトを切り捨てる前に、該 1 つ以上の非ゼロバイトに前記追加ビットを割り当てる

よう構成される、付記 8 乃至 1 0 のうちいずれかに記載のニューラルネットワークプロセッサ。

【 0 1 3 1 】

付記 1 2 . 前記復元ユニットは更に、前記データの圧縮されたチャンクの前記データ部分内に格納された前記 1 つ以上の非ゼロバイトに割り当てられている前記追加ビットの数を決定するよう構成される、付記 8 乃至 1 1 のうちいずれかに記載のニューラルネットワークプロセッサ。

【 0 1 3 2 】

付記 1 3 . 前記非ゼロバイトの 1 つ以上の最下位ビット (L S B) が、前記データの圧縮されたチャンクの圧縮中に切り捨てられる、付記 8 乃至 1 2 のうちいずれかに記載のニューラルネットワークプロセッサ。

【 0 1 3 3 】

付記 1 4 . 前記非ゼロバイトの 1 つ以上の最上位ビット (M S B) が、前記データの圧縮されたチャンクの圧縮中に切り捨てられる、付記 8 乃至 1 3 のうちいずれかに記載のニューラルネットワークプロセッサ。

【 0 1 3 4 】

付記 1 5 . コンピュータにより実施される方法であって、

ニューラルネットワークプロセッサの圧縮ユニットで、前記ニューラルネットワークプロセッサ内の少なくとも 1 つのニューロンによって生成されるデータの非圧縮チャンクを受け取り、該データの非圧縮チャンクが一定数のバイトを含む、ことと、

圧縮された出力チャンクのマスク部分を生成し、該マスク部分が、前記データの非圧縮チャンク内の前記一定数のバイトに等しいビットの数を含み、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内のバイトに対応し、前記マスク部分内の各ビットが、前記データの非圧縮チャンク内の対応するバイトがゼロである場合に論理 0 を有し、前記データの非圧縮チャンク内の対応するバイトが非ゼロである場合に論理 1 を有する、ことと、

前記データの非圧縮チャンク内の非ゼロバイトの数を決定すること、前記データの非圧縮チャンク内の前記非ゼロバイトの数に基づき、前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記圧縮された出力チャンクのデータ部分内のビットの数を決定することと、該決定されたビットの数まで前記データの非圧縮チャンク内の前記非ゼロバイトを切り捨てることと、該切り捨てられた非ゼロバイトを前記圧縮された出力チャンクの前記データ部分に格納することとによって、前記圧縮された出力チャンクの前記データ部分を生成することと、

前記マスク部分及び前記データ部分を含む前記圧縮された出力チャンクを前記ニューラルネットワークプロセッサのメモリに記憶することと

を有する、方法。

【 0 1 3 5 】

付記 1 6 . 前記データの非圧縮チャンクの切り捨てられた非ゼロバイトを格納するために利用可能な前記圧縮された出力チャンクの前記データ部分内の追加ビットの数を決定することと、前記データの非圧縮チャンク内の前記非ゼロバイトのうちの 1 つ以上の非ゼロバ

10

20

30

40

50

イトを切り捨てる前に、該 1 つ以上の非ゼロバイトに前記追加ビットを割り当てることとを更に有する、付記 15 に記載のコンピュータにより実施される方法。

【0136】

付記 17 . 前記データの非圧縮チャンク内の非ゼロバイトの数が前記圧縮された出力チャンクの前記データ部分内のバイトの数以下である場合には切り捨てなしで、前記データの非圧縮チャンク内の前記非ゼロバイトを前記圧縮された出力チャンクの前記データ部分に格納することを更に有する、付記 15 及び 16 のうちいずれかに記載のコンピュータにより実施される方法。

【0137】

付記 18 . 前記ニューラルネットワークプロセッサの復元ユニットで、前記圧縮された出力チャンクを受け取ることと、

10

前記圧縮された出力チャンクの前記マスク部分に基づき、前記データの非圧縮チャンク内の非ゼロバイトの数を決定することと、

前記圧縮された出力チャンクの前記マスク部分に基づき、前記データの非圧縮チャンク内の前記非ゼロバイトの位置を決定することと、

前記圧縮された出力チャンクの前記データ部分に前記切り捨てられた非ゼロバイトを格納するために前記圧縮ユニットによって使用されたビットの数を決定することと、

論理 0 である前記圧縮された出力チャンクの前記マスク部分内の各ビット位置について、ゼロバイトを、復元された出力チャンクの対応する位置に挿入することと、

論理 1 である前記マスク部分内の各ビット位置について、前記圧縮された出力チャンクの
前記データ部分内の対応する位置からの前記切り捨てられた非ゼロバイトを、前記復元された出力チャンクの対応する位置に、前記圧縮された出力チャンクの圧縮中に切り捨てられたビットの数に等しい数のゼロビットとともに挿入することと

20

を更に有する、付記 15 乃至 17 のうちいずれかに記載のコンピュータにより実施される方法。

【0138】

付記 19 . 前記復元された出力チャンクに格納される前記切り捨てられた非ゼロバイトのうち
の 1 つ以上にオフセットを付加することを更に有する、付記 15 乃至 18 のうちいずれかに記載のコンピュータにより実施される方法。

【0139】

30

付記 20 . 前記オフセットは、前記復元された出力チャンクに格納される前記切り捨てられた非ゼロバイトの 1 つ以上の最下位ビット (L S B) に加えられる、付記 15 乃至 19 のうちいずれかに記載のコンピュータにより実施される方法。

【0140】

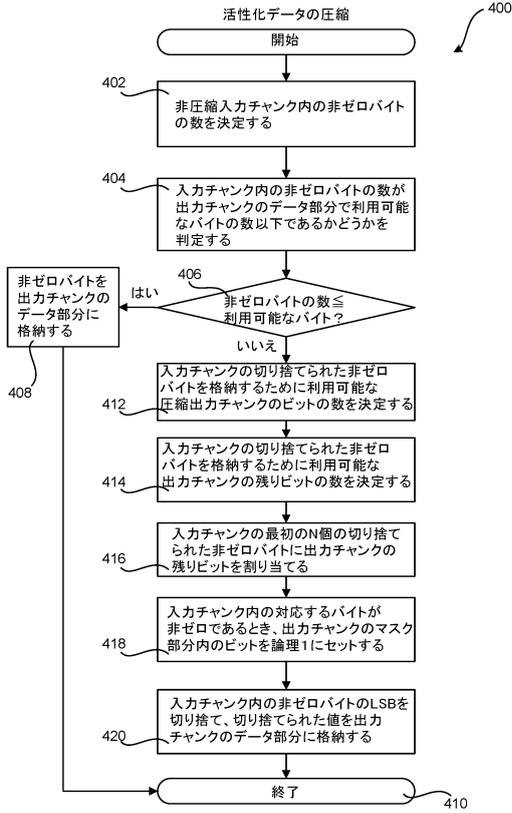
上記に基づき、メモリバス帯域幅の利用量を低減するために活性化データを圧縮及び復元することができる D N N モジュールが本明細書で開示されてきた、ことが認識されるべきである。本明細書で提示される対象は、コンピュータ構造特徴、方法論的及び変形可能な動作、具体的な計算器、並びにコンピュータ可読媒体に特有の言語で記載されてきたが、添付の特許請求の範囲に示されている対象は、必ずしも、本明細書で記載される具体的な特徴、動作、又は媒体に限られないことが理解されるべきである。むしろ、具体的な特徴、動作及び媒体は、請求される対象を実施する形態の例として開示されている。

40

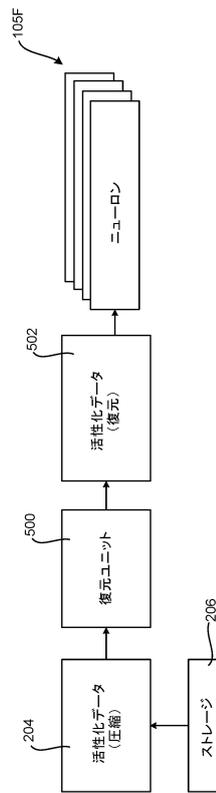
【0141】

上記の対象は、実例としてのみ与えられており、限定として解釈されるべきではない。様々な修正及び変更が、図示及び記載されている例となる構成及び適用に従うことなく、且つ、続く特許請求の範囲で示されている本開示の適用範囲から外れることなしに、本明細書で記載されている対象に対して行われ得る。

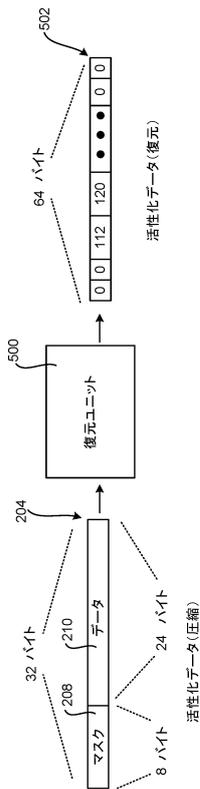
【図4】



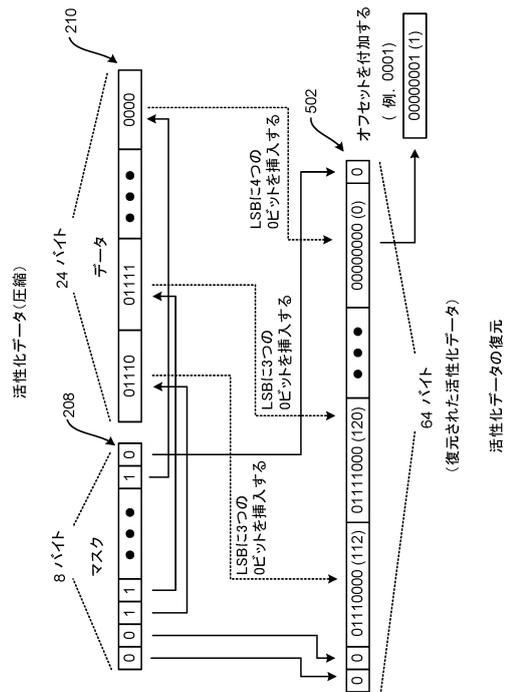
【図5A】



【図5B】



【図6】



10

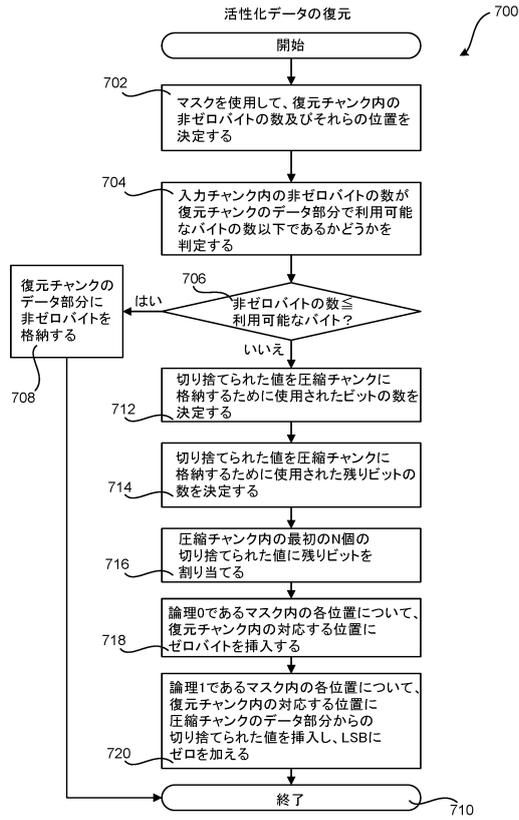
20

30

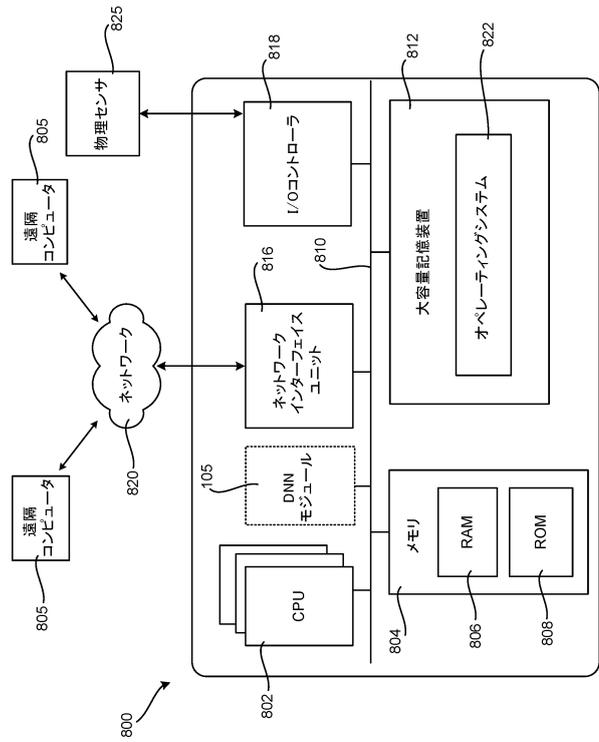
40

50

【図7】



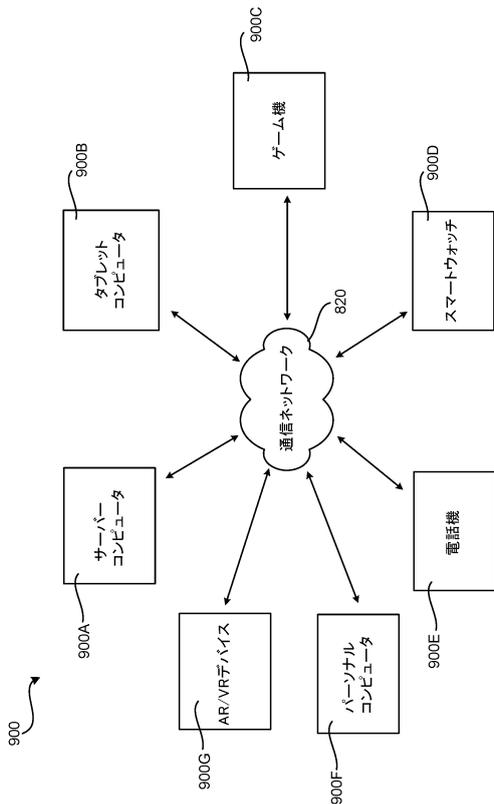
【図8】



10

20

【図9】



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)

ウェイ マイクロソフト テクノロジー ライセンシング, エルエルシー (番地無し)

(72)発明者 ルンデル, ベンジャミン エリオット

アメリカ合衆国 ワシントン州 98052 レッドモンド ワン マイクロソフト ウェイ マイクロソフト テクノロジー ライセンシング, エルエルシー (番地無し)

(72)発明者 ウォール, ラリー マーヴィン

アメリカ合衆国 ワシントン州 98052 レッドモンド ワン マイクロソフト ウェイ マイクロソフト テクノロジー ライセンシング, エルエルシー (番地無し)

(72)発明者 マクブライド, チャド ボーリング

アメリカ合衆国 ワシントン州 98052 レッドモンド ワン マイクロソフト ウェイ マイクロソフト テクノロジー ライセンシング, エルエルシー (番地無し)

(72)発明者 アンバルデカル, アモール アショク

アメリカ合衆国 ワシントン州 98052 レッドモンド ワン マイクロソフト ウェイ マイクロソフト テクノロジー ライセンシング, エルエルシー (番地無し)

(72)発明者 ピーター, ジョージ

アメリカ合衆国 ワシントン州 98052 レッドモンド ワン マイクロソフト ウェイ マイクロソフト テクノロジー ライセンシング, エルエルシー (番地無し)

(72)発明者 セドラ, ケント ディー

アメリカ合衆国 ワシントン州 98052 レッドモンド ワン マイクロソフト ウェイ マイクロソフト テクノロジー ライセンシング, エルエルシー (番地無し)

(72)発明者 ボブロフ, ボリス

アメリカ合衆国 ワシントン州 98052 レッドモンド ワン マイクロソフト ウェイ マイクロソフト テクノロジー ライセンシング, エルエルシー (番地無し)

審査官 北川 純次

(56)参考文献 米国特許出願公開第2010/312735 (US, A1)

米国特許第5369773 (US, A)

特開2000-59227 (JP, A)

(58)調査した分野 (Int.Cl., DB名)

G06N 3/06

G06F 12/04

H03M 7/40