

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第6987860号
(P6987860)

(45) 発行日 令和4年1月5日 (2022. 1. 5)

(24) 登録日 令和3年12月3日 (2021. 12. 3)

(51) Int. Cl.

F I

GO 6 N 3/08 (2006. 01)

GO 6 N 3/063 (2006. 01)

GO 6 N 3/08

GO 6 N 3/063

請求項の数 13 (全 28 頁)

(21) 出願番号	特願2019-524156 (P2019-524156)	(73) 特許権者	502208397
(86) (22) 出願日	平成29年8月23日 (2017. 8. 23)		グーグル エルエルシー
(65) 公表番号	特表2019-537139 (P2019-537139A)		G o o g l e L L C
(43) 公表日	令和1年12月19日 (2019. 12. 19)		アメリカ合衆国 カリフォルニア州 94
(86) 国際出願番号	PCT/US2017/048123		043 マウンテン ビュー アンフィシ
(87) 国際公開番号	W02018/089079		アター パークウェイ 1600
(87) 国際公開日	平成30年5月17日 (2018. 5. 17)		1600 Amphitheatre P
審査請求日	令和2年5月7日 (2020. 5. 7)		arkway 94043 Mounta
(31) 優先権主張番号	15/348, 199	(74) 代理人	110001195
(32) 優先日	平成28年11月10日 (2016. 11. 10)		特許業務法人深見特許事務所
(33) 優先権主張国・地域又は機関	米国 (US)	(72) 発明者	ヤング, レジナルド・クリフォード
(31) 優先権主張番号	15/467, 382		アメリカ合衆国、94043 カリフォル
(32) 優先日	平成29年3月23日 (2017. 3. 23)		ニア州、マウンテン・ビュー、アンフィシ
(33) 優先権主張国・地域又は機関	米国 (US)		アター・パークウェイ、1600
			最終頁に続く

(54) 【発明の名称】 ハードウェアにおけるカーネルストライドの実行

(57) 【特許請求の範囲】

【請求項 1】

コンピュータにより実現される方法であって、
専用ハードウェア回路上に畳み込みニューラルネットワークを実装する要求を受け取り、前記畳み込みニューラルネットワークを用いて、前記専用ハードウェア回路に命令を実行させることによって、ニューラルネットワーク入力を受け取って処理することを備え、前記畳み込みニューラルネットワークは、1より大きいストライドを有する第1の畳み込みニューラルネットワーク層を含み、前記専用ハードウェア回路は、ニューラルネットワーク計算を実行するための集積回路であり、ベクトル-行列乗算を実行するようにされた行列計算ユニットと、前記行列計算ユニットの出力に対してプーリングを実行するようにされたプーリング回路を含むベクトル計算ユニットとを含み、前記方法はさらに、前記専用ハードウェア回路によって実行されると、前記専用ハードウェア回路に、前記畳み込みニューラルネットワークによる入力テンソルの処理中に、動作を実行することによって前記第1の畳み込みニューラルネットワーク層の出力と等価な層出力テンソルを生成させる命令を生成することを備え、前記動作は、
前記行列計算ユニットが、前記第1の畳み込みニューラルネットワーク層への前記入力テンソルを、1に等しいストライドを有するがそれ以外は前記第1の畳み込みニューラルネットワーク層と同等である第2の畳み込みニューラルネットワーク層を用いて処理することにより、第1のテンソルを生成することと、
前記ベクトル計算ユニットが、前記第2の畳み込みニューラルネットワーク層が前記第

1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう前記第1のテンソルの要素を零出力して、第2のテンソルを生成することと、

前記ベクトル計算ユニットの前記プーリング回路が、前記第2のテンソルに対して最大プーリングを実行して前記層出力テンソルを生成することとを含み、

前記ベクトル計算ユニットが前記第1のテンソルの要素を零出力することは、

前記ベクトル計算ユニットがマスキングテンソルと前記第1のテンソルとの要素ごとの乗算を実行して前記第2のテンソルを生成することを含み、前記マスキングテンソルは、
(i) 前記第2の畳み込みニューラルネットワーク層が前記第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう前記第1のテンソルの要素に対応する前記マスキングテンソルの各要素位置において0を含み、(ii) 前記マスキングテンソルの各他の要素位置において1を含む、コンピュータにより実現される方法。

10

【請求項2】

前記ベクトル計算ユニットが前記第1のテンソルの要素を零出力することは、

前記第1のテンソルの要素のサブセットに0を乗算することと、

前記サブセットに含まれていない前記第1のテンソルの要素に1を乗算することとを含み、

前記サブセットは、前記第2の畳み込みニューラルネットワーク層が前記第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう前記第1のテンソルの要素に対応する前記マスキングテンソルの各要素からなる、請求項1に記載の方法。

20

【請求項3】

前記マスキングテンソルは、前記専用ハードウェア回路によってアクセス可能なメモリに格納される、請求項1または2に記載の方法。

【請求項4】

コンピュータにより実現される方法であって、

専用ハードウェア回路上に畳み込みニューラルネットワークを実装する要求を受け取り、前記畳み込みニューラルネットワークを用いて、前記専用ハードウェア回路に命令を実行させることによって、ニューラルネットワーク入力を受け取って処理することを備え、前記畳み込みニューラルネットワークは、1より大きいストライドを有する第1の畳み込みニューラルネットワーク層を含み、前記専用ハードウェア回路は、ニューラルネットワーク計算を実行するための集積回路であり、ベクトル・行列乗算を実行するようにされた行列計算ユニットと、前記行列計算ユニットの出力に対してプーリングを実行するようにされたプーリング回路を含むベクトル計算ユニットとを含み、前記方法はさらに、

30

前記専用ハードウェア回路によって実行されると、前記専用ハードウェア回路に、前記畳み込みニューラルネットワークによる入力テンソルの処理中に、動作を実行することによって前記第1の畳み込みニューラルネットワーク層の出力と等価な層出力テンソルを生成させる命令を生成することを備え、前記動作は、

前記行列計算ユニットが、前記第1の畳み込みニューラルネットワーク層への前記入力テンソルを、1に等しいストライドを有するがそれ以外は前記第1の畳み込みニューラルネットワーク層と同等である第2の畳み込みニューラルネットワーク層を用いて処理することにより、第1のテンソルを生成することと、

40

前記ベクトル計算ユニットが、前記第2の畳み込みニューラルネットワーク層が前記第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう前記第1のテンソルの要素を零出力して、第2のテンソルを生成することと、

前記ベクトル計算ユニットの前記プーリング回路が、前記第2のテンソルに対して最大プーリングを実行して前記層出力テンソルを生成することとを含み、

前記ベクトル計算ユニットが前記第1のテンソルの要素を零出力することは、

前記ベクトル計算ユニットが、第1のマスキングテンソルと前記第1のテンソルとの要素ごとの乗算を実行して、修正された第1のテンソルを生成することを含み、前記第1の

50

マスキングテンソルは、(i) 前記第 2 の畳み込みニューラルネットワーク層が前記第 1 の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう前記第 1 のテンソルの要素に対応する前記第 1 のマスキングテンソルの各要素位置において 0 を含み、(i i) 前記第 2 の畳み込みニューラルネットワーク層が前記第 1 の畳み込みニューラルネットワーク層のストライドを有する場合には生成されたであろう前記第 1 のテンソルの要素に対応する前記第 1 のマスキングテンソルの各要素位置においてそれぞれの非 0 値を含み、前記ベクトル計算ユニットが前記第 1 のテンソルの要素を零出力することはさらに、

前記ベクトル計算ユニットが第 2 のマスキングテンソルと前記修正された第 1 のテンソルとの要素ごとの乗算を実行することを含み、前記第 2 のマスキングテンソルは、前記第 2 の畳み込みニューラルネットワーク層が前記第 1 の畳み込みニューラルネットワーク層のストライドを有する場合に生成されるであろう前記第 1 のテンソルの要素に対応する各要素位置において、前記第 1 のマスキングテンソルの前記それぞれの非 0 値の逆数を含む、コンピュータにより実現される方法。

【請求項 5】

前記第 1 のマスキングテンソルおよび前記第 2 のマスキングテンソルは、前記専用ハードウェア回路によってアクセス可能なメモリに格納される、請求項 4 に記載の方法。

【請求項 6】

1 より大きい複数のストライドにそれぞれ対応するように、複数のマスキングテンソルが前記メモリに格納され、

前記方法はさらに、前記複数のマスキングテンソルの中から、前記第 1 の畳み込みニューラルネットワーク層のストライドに対応するマスキングテンソルを選択することを備える、請求項 5 に記載の方法。

【請求項 7】

前記ベクトル計算ユニットの前記プーリング回路が最大プーリングを実行することは、前記第 1 の畳み込みニューラルネットワーク層のストライドによって定義される前記第 2 のテンソルの 1 つまたは複数のウィンドウの各々について、前記ウィンドウ内の要素の最大値要素を取得することを含む、請求項 1 ~ 6 のいずれか 1 項に記載の方法。

【請求項 8】

前記第 2 のテンソルの前記 1 つまたは複数のウィンドウの各々は、前記第 1 の畳み込みニューラルネットワーク層のストライドに対応する次元を有する矩形ウィンドウであり、前記第 2 のテンソルの異なる要素を含む、請求項 7 に記載の方法。

【請求項 9】

前記ベクトル計算ユニットの前記プーリング回路が最大プーリングを実行することは、前記第 2 のテンソルの要素の 1 つまたは複数のサブセットの各々について、前記サブセットの最大値要素を取得することを含む、請求項 1 ~ 8 のいずれか 1 項に記載の方法。

【請求項 10】

前記第 1 の畳み込みニューラルネットワーク層は、前記畳み込みニューラルネットワーク内の複数の畳み込みニューラルネットワーク層のうちの第 1 のニューラルネットワーク層であり、前記入力テンソルは、デジタル画像の、前記デジタル画像の画素に対応する要素を含む表現である、請求項 1 ~ 9 のいずれか 1 項に記載の方法。

【請求項 11】

前記入力テンソルは前記専用ハードウェア回路のユニファイドバッファに格納され、前記第 2 の畳み込みニューラルネットワーク層の重みは前記専用ハードウェア回路のダイナミックメモリに格納され、前記第 2 の畳み込みニューラルネットワーク層を用いて前記第 1 の畳み込みニューラルネットワーク層への前記入力テンソルを処理することは、

前記入力テンソルを前記ユニファイドバッファから前記行列計算ユニットに送ることと、

前記ダイナミックメモリから前記行列計算ユニットに前記第 2 の畳み込みニューラルネットワーク層の前記重みを送ることと、

10

20

30

40

50

前記行列計算ユニットによって、前記第 2 の畳み込みニューラルネットワーク層の前記重みを用いて前記入力テンソルを処理して、前記第 1 のテンソルを生成することを含む、請求項 1 ~ 1 0 のいずれか 1 項に記載の方法。

【請求項 1 2】

システムであって、

請求項 1 に記載の専用ハードウェア回路と、

命令を格納する 1 つまたは複数の記憶装置とを備え、前記命令は、前記専用ハードウェア回路によって実行されると、前記専用ハードウェア回路に請求項 1 ~ 1 1 のいずれか 1 項に記載の方法を実行させるよう動作可能である、システム。

【請求項 1 3】

1 つまたは複数のコンピュータによって実行されると、前記 1 つまたは複数のコンピュータに請求項 1 ~ 1 1 のいずれか 1 項に記載の方法を実行させる命令を含むコンピュータプログラム。

【発明の詳細な説明】

【背景技術】

【0 0 0 1】

背景

本明細書は、ハードウェアにおけるニューラルネットワーク推論の計算に関する。

【0 0 0 2】

ニューラルネットワークは、受け取られた入力に対する出力、たとえば分類を生成するために 1 つまたは複数の層を用いる機械学習モデルである。いくつかのニューラルネットワークは、出力層に加えて 1 つまたは複数の隠れ層を含む。各隠れ層の出力は、ネットワーク内の別の層、たとえばネットワークの次の隠れ層または出力層への入力として用いられる。ネットワークの各層は、それぞれのパラメータセットの現在の値に従って、受け取られた入力から出力を生成する。

【発明の概要】

【課題を解決するための手段】

【0 0 0 3】

概要

一般に、この明細書はニューラルネットワーク推論を計算する専用ハードウェア回路について記載する。

【0 0 0 4】

一般に、この明細書に記載される主題の 1 つの革新的な局面は、ハードウェア回路上でニューラルネットワークを処理するよう要求を受け取ることを備え、ニューラルネットワークは、1 より大きいストライドを有する第 1 の畳み込みニューラルネットワーク層を含み、さらに、これに応答して、ハードウェア回路によって実行されると、ハードウェア回路に、ニューラルネットワークによる入力テンソルの処理中に、以下の動作を実行することによって第 1 の畳み込みニューラルネットワーク層の出力と等価な層出力テンソルを生成させる命令を生成することを備える方法およびシステムにおいて実施され得る。この動作は、第 1 の畳み込みニューラルネットワーク層への入力テンソルを、1 に等しいストライドを有するがそれ以外は第 1 の畳み込みニューラルネットワーク層と同等である第 2 の畳み込みニューラルネットワーク層を用いて処理することにより、第 1 のテンソルを生成することと、第 2 の畳み込みニューラルネットワーク層が第 1 の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第 1 のテンソルの要素を零出力して、第 2 のテンソルを生成することと、第 2 のテンソルに対して最大プーリングを実行して層出力テンソルを生成することを含む。

【0 0 0 5】

実装形態は、以下の特徴の 1 つ以上を含み得る。いくつかの実装形態では、第 1 のテンソルの要素を零出力することは、第 1 のテンソルの要素のサブセットに 0 を乗算することと、サブセットに含まれていない第 1 のテンソルの要素に 1 を乗算することを含む。第

10

20

30

40

50

1のテンソルの要素を零出力することは、マスキングテンソルと第1のテンソルとの要素ごとの乗算を実行して第2のテンソルを生成することを含み、マスキングテンソルは、(i)第2の畳み込みニューラルネットワーク層が第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第1のテンソルの要素に対応するマスキングテンソルの各要素位置において0を含み、(ii)マスキングテンソルの各他の要素位置において1を含む。いくつかの実装形態では、マスキングテンソルは、ハードウェア回路によってアクセス可能なメモリに格納され、マスキングテンソルと第1のテンソルとの要素ごとの乗算は、ハードウェア回路に含まれる、ハードウェアで実装されるベクトル計算ユニットによって実行される。

【0006】

10

実装形態は、さらに、以下の特徴の1つ以上を含み得る。いくつかの実装形態では、第1のテンソルの要素を零出力することは、第1のマスキングテンソルと第1のテンソルとの要素ごとの乗算を実行して、修正された第1のテンソルを生成することを含み、第1のマスキングテンソルは、(i)第2の畳み込みニューラルネットワーク層が第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第1のテンソルの要素に対応するマスキングテンソルの各要素位置において0を含み、(ii)第2の畳み込みニューラルネットワーク層が第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第1のテンソルの要素に対応するマスキングテンソルの各要素位置においてそれぞれの非0値を含み、第1のテンソルの要素を零出力することはさらに、第2のマスキングテンソルと修正された第1のテンソルとの要素ごとの乗算を実行することを含み、第2のマスキングテンソルは、第2の畳み込みニューラルネットワーク層が第1の畳み込みニューラルネットワーク層のストライドを有する場合に生成されるであろう第1のテンソルの要素に対応する各要素位置において、第1のマスキングテンソルのそれぞれの非0値の逆数を含む。

20

【0007】

実装形態は、さらに、以下の特徴の1つ以上を含み得る。いくつかの実装形態では、最大プーリングを実行することは、第1の畳み込みニューラルネットワーク層のストライドによって定義される第2のテンソルの1つまたは複数のウィンドウの各々について、ウィンドウ内の要素の最大値要素を取得することを含む。第2のテンソルの1つまたは複数のウィンドウの各々は、畳み込みニューラルネットワーク層のストライドに対応する次元を有する矩形ウィンドウであり、第2のテンソルの異なる要素を含む。いくつかの実装形態では、最大プーリングを実行することは、第2のテンソルの要素の1つまたは複数のサブセットの各々について、サブセットの最大値要素を取得することを含む。第2のテンソル上で実行される最大プーリングは、ハードウェア回路のプーリング回路によって実行される。畳み込みニューラルネットワーク層は、ニューラルネットワーク内の第1のニューラルネットワーク層であり、入力テンソルは、デジタル画像の、当該デジタル画像の画素に対応する要素を含む表現である。

30

【0008】

実装形態は、さらに、以下の特徴の1つ以上を含み得る。いくつかの実装形態では、入力テンソルはハードウェア回路のユニファイドバッファに格納され、第2の畳み込みニューラルネットワーク層の重みはハードウェア回路のダイナミックメモリに格納され、第2の畳み込みニューラルネットワーク層を用いて第1の畳み込みニューラルネットワーク層への入力テンソルを処理することは、入力テンソルをユニファイドバッファからハードウェアで実装されるハードウェア回路の行列計算ユニットに送ることと、ダイナミックメモリからハードウェア回路の行列計算ユニットに第2の畳み込みニューラルネットワーク層の重みを送ることと、ハードウェア回路の行列計算ユニットによって、第2の畳み込みニューラルネットワーク層の重みを用いて入力テンソルを処理して、第1のテンソルを生成することを含む。

40

【0009】

この明細書において記載される主題の特定の実施形態は、以下の利点の1つ以上を実現

50

するように実現することができる。ハードウェア回路が1より大きいストライドを有する畳み込みニューラルネットワークを用いて入力テンソルを直接処理することができない場合でも、1より大きいストライドを有する畳み込みニューラルネットワーク層に対応する出力テンソルを、専用ハードウェア回路によってハードウェアで生成できる。専用ハードウェア回路を用いて適切な出力を生成することにより、1より大きいストライドを有するニューラルネットワーク層の処理は、たとえ専用ハードウェア回路がそのような処理を直接サポートしていなくても、データをホストコンピュータに返送することなく、すなわち計算の少なくとも一部をオフチップで実行することなく、実行できる。これにより、専用ハードウェア回路のハードウェアアーキテクチャを変更することなく、1より大きいストライドを有する畳み込み層を含むニューラルネットワークの推論を効率的に判断することが可能になる。すなわち、処理の一部をオフチップで、ソフトウェアで、またはその両方で実行することから生じる処理遅延が回避される。

10

【0010】

本明細書に記載される主題は、たとえば、ニューラルネットワーク推論を計算するときにカーネルストライドを実行するための、開示された技法およびハードウェアを用いる画像認識または分類方法およびシステムにも関する。

【0011】

この明細書において記載される主題の1つ以上の実施形態の詳細は、添付の図面および以下の詳細な説明において述べられる。主題の他の特徴、局面および利点は、詳細な説明、図面および特許請求の範囲から明らかになる。

20

【図面の簡単な説明】

【0012】

【図1】例示的なニューラルネットワーク処理システムを示す。

【図2】ニューラルネットワークの所与の層について計算を実行するための例示的な方法の流れ図である。

【図3】例示的なニューラルネットワーク処理システムを示す。

【図4】行列計算ユニットを含む例示的なアーキテクチャを示す。

【図5】シストリックアレイ内のセルの例示的なアーキテクチャを示す。

【図6】ベクトル計算ユニットの例示的なアーキテクチャを示す。

【図7】プーリング回路のための例示的なアーキテクチャを示す。

30

【図8】1より大きいストライドでニューラルネットワークの所与の層に対して計算を実行するようにニューラルネットワーク処理システムに命令するための例示的な方法の流れ図である。

【図9】1より大きいストライドを有するニューラルネットワークの所与の層に対して計算を実行するための例示的な方法の流れ図である。

【図10】1より大きいストライドでのニューラルネットワークの所与の層に対する計算の例である。

【発明を実施するための形態】

【0013】

さまざまな図面中の同様の参照番号および指定は、同様の要素を示す。

40

詳細な説明

複数の層を有するニューラルネットワークを用いて推論を計算することができる。たとえば、入力を与えられると、ニューラルネットワークはその入力に対する推論を計算することができる。ニューラルネットワークは、ニューラルネットワークの各層を通して入力を処理することによって、この推論を計算する。各層は入力を受け取り、その層に対する重みのセットに従って入力を処理して出力を生成する。

【0014】

したがって、受け取った入力から推論を計算するために、ニューラルネットワークは入力を受け取り、それを各ニューラルネットワーク層の各々を通して処理して推論を生成し、1つのニューラルネットワーク層からの出力は次のニューラルネットワーク層への入力

50

として与えられる。ニューラルネットワーク層へのデータ入力、たとえば、ニューラルネットワークへの入力、またはシーケンス内におけるその層の下層の、あるニューラルネットワーク層への出力は、その層への活性化入力と呼ぶことができる。

【0015】

いくつかの実装形態では、ニューラルネットワークの層はシーケンスで配置される。ある他の実装形態では、層は有向グラフとして配される。つまり、任意の特定の層が複数の入力、複数の出力、またはそれらの両方を受け取ることができる。ニューラルネットワークの層は、ある層の出力を前の層への入力として送り返すことができるように構成することもできる。

【0016】

いくつかのニューラルネットワークは、1つまたは複数のニューラルネットワーク層からの出力をプーリングして、後続のニューラルネットワーク層への入力として用いられるプーリングされた値を生成する。いくつかの実装形態では、ニューラルネットワークは、出力のグループの最大値、最小値、または平均値を判断し、最大値、最小値、または平均値をグループのプーリングされた出力として用いることによって、出力のグループをプーリングする。出力をプーリングすることは、空間的不変性のある程度維持することができるので、さまざまな構成で配置される出力を、同じ推論を有するように処理することができる。出力をプーリングすることはまた、プーリングする前の出力の所望の特性を維持しながら後続のニューラルネットワーク層で受け取られる入力の次元を低減することができる、それはニューラルネットワークによって生成される推論の品質を著しく損なうことなく効率を改善できる。

【0017】

いくつかのニューラルネットワークは、1より大きいストライドを有する1つまたは複数の畳み込みニューラルネットワーク層を含む。概念的には、1のストライドの場合、畳み込みニューラルネットワーク層は、重みのセットを活性化入力に順次適用することができる。すなわち、活性化入力アレイの場合、重みを活性化入力のサブセットに適用し、畳み込み計算が完了するまで、活性化入力の各他のサブセットに、1つの位置、たとえば行または列だけ、移動させることができる。1より大きい整数であるストライドを有する畳み込みニューラルネットワーク層の場合、重みを活性化入力のサブセットに適用し、畳み込み計算が完了するまで、活性化入力の各他のサブセットに、ストライドに等しい数の位置だけ、たとえば、ストライドによって示される行または列の数だけ、移動させることができる。

【0018】

本明細書は、ニューラルネットワーク層を処理し、任意選択で1つまたは複数のニューラルネットワーク層の出力に対してプーリングを実行する専用ハードウェア回路を記載する。専用ハードウェア回路は、1のストライドを有するニューラルネットワーク層を処理することができる回路を含む。専用ハードウェア回路は、1より大きいストライドを有するニューラルネットワーク層の処理を直接サポートしないが、専用ハードウェア回路は、1より大きいストライドを有するニューラルネットワーク層の出力と等価の出力を生成するように制御され得る。したがって、開示された技術の1つの技術的效果および利点は、1のストライドを有するニューラルネットワーク層を処理することができる回路を、より柔軟な態様で、1より大きいストライドを有するニューラルネットワーク層についてニューラルネットワーク推論を計算するために用いることができることである。

【0019】

図1は、例示的なニューラルネットワーク処理システム100を示す。ニューラルネットワーク処理システム100は、以下に記載されるシステム、コンポーネント、および技術が実装され得る1つまたは複数の位置に1つまたは複数コンピュータとして実装されるシステムの例である。

【0020】

ニューラルネットワーク処理システム100は、専用ハードウェア回路110を用いて

10

20

30

40

50

ニューラルネットワーク計算を実行するシステムである。ハードウェア回路 110 は、ニューラルネットワーク計算を実行するための集積回路であり、ハードウェアでベクトル・行列乗算を実行する行列計算ユニット 120 を含む。ハードウェア回路 110 はまた、行列計算ユニット 120 の出力に対してプーリングを実行するためのプーリング回路を含むベクトル計算ユニット 140 を含む。例示的な専用ハードウェア回路 120 が、図 3 を参照して以下により詳細に記載される。

【0021】

特に、ニューラルネットワーク処理システム 100 は、専用ハードウェア回路 110 上にニューラルネットワークを実装する要求を受信し、専用ハードウェア回路 110 上にニューラルネットワークを実装し、所与のニューラルネットワークが実装されると、ニューラルネットワーク推論を生成するために専用集積回路 110 を用いてニューラルネットワークへの入力を処理する。

10

【0022】

すなわち、ニューラルネットワーク処理システム 100 は、入力を処理するために用いられるべきニューラルネットワークのためのニューラルネットワークアーキテクチャを指定する要求を受け取ることができる。ニューラルネットワークアーキテクチャは、ニューラルネットワーク内の層の数および構成、ならびにパラメータを有する各層のパラメータの値を定義する。

【0023】

専用集積回路 110 上にニューラルネットワークを実装するために、ニューラルネットワーク処理システム 100 は、1 つまたは複数の物理的位置にある 1 つまたは複数のコンピュータ上の 1 つまたは複数のコンピュータプログラムとして実装されるニューラルネットワーク実装エンジン 150 を含む。

20

【0024】

ニューラルネットワーク実装エンジン 150 は命令を生成し、命令は、専用ハードウェア回路 110 によって実行されると、ハードウェア回路 110 に、ニューラルネットワークによって指定される動作を実行させて、受け取られたニューラルネットワーク入力からニューラルネットワーク出力を生成させる。

【0025】

命令がニューラルネットワーク実装エンジン 150 によって生成され、ハードウェア回路 110 に与えられると、ニューラルネットワーク処理システム 100 は、ニューラルネットワーク入力を受け取り、ニューラルネットワークを用いて、ハードウェア回路 110 に、生成された命令を実行させることによって、ニューラルネットワーク入力を処理することができる。

30

【0026】

しかしながら、いくつかのニューラルネットワークは、1 つまたは複数の互換性のないニューラルネットワーク層を含む。本明細書で用いられる互換性のないニューラルネットワーク層という用語は、専用ハードウェア回路 110 によってハードウェアで直接実行することができない操作を指定するニューラルネットワーク層を指す。ハードウェア回路 110 上にこれらのニューラルネットワークを実装するために、ニューラルネットワーク実装エンジン 150 は、ハードウェア回路 110 によって実行されると、ハードウェア回路 110 に、ハードウェアにおいて以下の操作を実行することによって互換性のないニューラルネットワーク層についての出力を生成させる命令を生成する。それらの操作は、ニューラルネットワーク層によって指定されるものとは異なる操作であるが、互換性のないニューラルネットワーク層の仕様を満たす層出力、たとえば層出力テンソル、つまり、層によって指定される操作を直接実行することによって生成されたであろう出力と同じ層出力が生成される結果となる。

40

【0027】

特に、いくつかのニューラルネットワークは、1 より大きいストライドを有する畳み込みニューラルネットワーク層を含む。そのようなニューラルネットワーク層は、入力テン

50

ソルを用いて非順次的に処理される1つまたは複数のカーネルを特徴とする。たとえば、1のストライドでカーネルストライドを実行するとき、カーネルは入力テンソルの要素に順次適用される。しかしながら、2のストライドでカーネルストライドを実行するとき、ニューラルネットワーク層のカーネルは、カーネルの特定の要素が入力テンソルの1つおきの要素に適用されて出力テンソルを生成するように、シフトされる。出力テンソルは、ニューラルネットワークの別の層により入力として用いることができる。

【0028】

ハードウェア回路110上で行列演算を実行する主ハードウェアユニットは行列計算ユニット120であるので、集積回路は1より大きいストライドを有するニューラルネットワーク層を直接計算することはできない。1より大きいストライドを有する層を含むニューラルネットワークを実現するために、ニューラルネットワーク実現エンジン150は命令を生成し、命令は、ニューラルネットワークによるニューラルネットワーク入力の処理中に専用ハードウェア回路110によって実行されると、ハードウェア回路110にハードウェアで他の操作を実行させて、行列乗算ユニット120およびブリング回路を特徴とするベクトル計算ユニット140を用いて、1より大きいストライドを有するニューラルネットワーク層の仕様を満たす出力テンソルを生成する。これらの命令および他の操作は、図7～図10を参照して以下により詳細に説明される。

【0029】

図2は、専用ハードウェア回路を用いてニューラルネットワークの所与の層について計算を実行するための例示的なプロセス200の流れ図である。便宜上、方法200は、方法200を実行する1つまたは複数の回路を有するシステムに関して説明される。方法200は、受け取られた入力から推論を計算するために、ニューラルネットワークの各層に対して実行することができる。

【0030】

システムは、所与の層について重み入力のセット(ステップ202)および活性化入力のセット(ステップ204)を受け取る。重み入力のセットおよび活性化入力のセットは、専用ハードウェア回路のダイナミックメモリおよびユニファイドバッファからそれぞれ受け取ることができる。いくつかの実装形態では、重み入力のセットと活性化入力のセットとの両方をユニファイドバッファから受け取ることができる。

【0031】

システムは、専用ハードウェア回路の行列乗算ユニットを用いて、重み入力および活性化入力から累積値を生成する(ステップ206)。いくつかの実装形態では、累積値は、重み入力のセットと活性化入力のセットとの内積である。すなわち、層内のすべての重みのサブセットである1セットの重みについて、システムは各重み入力を各活性化入力と乗算し、それらの積を合計して累積値を形成することができる。システムは、次いで、他のセットの重みと他のセットの活性化入力との内積を計算することができる。いくつかの実装形態では、専用ハードウェア回路は、特定のニューラルネットワーク層のストライド、すなわちニューラルネットワーク層が1のストライドまたは1より大きいストライドを有するかどうかにかかわらず、同様にそのような操作を実行し得る。行列乗算ユニットからの出力のその後の処理は、ニューラルネットワーク層が1より大きい指定されたストライドで処理された場合に生成されるであろう出力と等価な出力を生成するよう実行することができる。

【0032】

システムは、専用ハードウェア回路のベクトル計算ユニットを用いて累積値から層出力を生成することができる(ステップ208)。いくつかの実装形態では、ベクトル計算ユニットは、累積値に活性化関数を適用する。これについては、図5を参照して以下でさらに説明する。層の出力は、ニューラルネットワーク内の次の層への入力として用いるためにユニファイドバッファに格納することができ、または推論を決めるために用いることができる。いくつかの実装形態では、ニューラルネットワーク層は、1より大きいストライドを指定し得、システムは、1より大きいストライドを有するニューラルネットワーク層

10

20

30

40

50

の出力と等価である層出力を得るために、累積値に対して追加の処理を行い得る。受け取られた入力ニューラルネットワークの各層を介して処理されて、受け取られた入力に対する推論を生成すると、システムはニューラルネットワークの処理を終了する。

【0033】

図3は、ニューラルネットワーク計算を実行するための例示的な専用ハードウェア回路300を示す。システム300はホストインターフェース302を含む。ホストインターフェース302は、ニューラルネットワーク計算のためのパラメータを含む命令を受け取ることができる。パラメータは、以下のうちの1つまたは複数を含むことができる。処理すべき層の数、モデルの各層に対する対応する重み入力のセット、活性化入力の初期セット、すなわち推論が計算されるニューラルネットワークへの入力、各層の対応する入力および出力サイズ、ニューラルネットワーク計算のためのストライド値、および処理されるべき層のタイプ、たとえば畳み込み層または全結合層。

10

【0034】

ホストインターフェース302は、命令をシーケンサ306に送ることができる。シーケンサ306は、命令を、ニューラルネットワーク計算を実行するように回路を制御する低レベル制御信号に変換する。いくつかの実装形態では、制御信号は、回路内のデータフロー、たとえば重み入力のセットおよび活性化入力のセットが回路をどのように流れるか、を調整する。シーケンサ306は、制御信号をユニファイドバッファ308、行列計算ユニット312、およびベクトル計算ユニット314に送ることができる。いくつかの実装形態では、シーケンサ306はまた、ダイレクトメモリアクセスエンジン304およびダイナミックメモリ310に制御信号を送る。いくつかの実装形態では、シーケンサ306は制御信号を生成するプロセッサである。シーケンサ306は、適切なときに制御信号を回路300の各構成要素に送るために、制御信号のタイミングを用いることができる。いくつかの他の実装形態では、ホストインターフェース302は外部プロセッサから制御信号を渡す。

20

【0035】

ホストインターフェース302は、重み入力のセットおよび活性化入力の初期セットをダイレクトメモリアクセスエンジン304に送ることができる。ダイレクトメモリアクセスエンジン304は、ユニファイドバッファ308に活性化入力のセットを格納することができる。いくつかの実装形態では、ダイレクトメモリアクセスは、メモリユニットであり得るダイナミックメモリ310に重みのセットを格納する。いくつかの実装形態では、ダイナミックメモリ310は回路の外に配置されている。

30

【0036】

ユニファイドバッファ308はメモリバッファである。それは、ダイレクトメモリアクセスエンジン304からの活性化入力のセットおよびベクトル計算ユニット314の出力を格納するために用いることができる。ベクトル計算ユニット314は、図6を参照して以下により詳細に説明される。ダイレクトメモリアクセスエンジン304は、ユニファイドバッファ308からベクトル計算ユニット314の出力を読み出すこともできる。

【0037】

ダイナミックメモリ310およびユニファイドバッファ308は、重み入力のセットおよび活性化入力のセットをそれぞれ行列計算ユニット312に送ることができる。いくつかの実装形態では、行列計算ユニット312は二次元シストリックアレイである。行列計算ユニット312は、数学的演算、たとえば乗算および加算を実行することができる二次元シストリックアレイまたは他の回路とすることもできる。いくつかの実装形態では、行列計算ユニット312は汎用の行列プロセッサである。

40

【0038】

行列計算ユニット312は、重み入力および活性化入力を処理し、出力のベクトルをベクトル計算ユニット314に与えることができる。いくつかの実装形態では、行列計算ユニット312は、出力ベクトルをユニファイドバッファ308に送り、ユニファイドバッファ308は、出力ベクトルをベクトル計算ユニット314に送る。ベクトル計算ユニッ

50

ト 3 1 4 は、出力のベクトルを処理し、処理された出力のベクトルをユニファイドバッファ 3 0 8 に格納することができる。1 より大きいストライドを有するニューラルネットワーク層の場合、ベクトル計算ユニット 3 1 4 は、出力のベクトルを処理して、1 より大きいストライドを有するニューラルネットワーク層の出力と等価である層出力テンソルを生成し、層出力テンソルをユニファイドバッファ 3 0 8 に格納することができる。処理された出力のベクトルは、たとえばニューラルネットワーク内の後続の層で用いるために、行列計算ユニット 3 1 2 への活性化入力として用いることができる。行列計算ユニット 3 1 2 およびベクトル計算ユニット 3 1 4 は、図 4 および図 6 をそれぞれ参照して以下により詳細に説明される。

【 0 0 3 9 】

10

図 4 は、行列計算ユニットを含む例示的アーキテクチャ 4 0 0 を示す。行列計算ユニットは、二次元シストリックアレイ 4 0 6 である。アレイ 4 0 6 は複数のセル 4 0 4 を含む。いくつかの実装形態では、シストリックアレイ 4 0 6 の第 1 の次元 4 2 0 はセルの列に対応し、シストリックアレイ 4 0 6 の第 2 の次元 4 2 2 はセルの行に対応する。シストリックアレイは、列よりも多い行、行よりも多い列、または同数の列と行とを有することができる。

【 0 0 4 0 】

図示の例では、値ローダ 4 0 2 は活性化入力をアレイ 4 0 6 の行に送り、重みフェッチインターフェース (weight fetcher interface) 4 0 8 は重み入力をアレイ 4 0 6 の列に送る。しかしながら、いくつかの他の実装形態では、活性化入力は列に転送され、重み

20

【 0 0 4 1 】

値ローダ 4 0 2 は、ユニファイドバッファ、たとえば、図 3 のユニファイドバッファ 3 0 8 から、活性化入力を受け取ることができる。各値ローダは、対応する活性化入力をアレイ 4 0 6 の最も左側の異なるセルに送ることができる。たとえば、値ローダ 4 1 2 はセル 4 1 4 に活性化入力を送ることができる。

【 0 0 4 2 】

重みフェッチインターフェース 4 0 8 は、メモリユニット、たとえば図 3 のダイナミックメモリ 3 1 0 から重み入力を受け取ることができる。重みフェッチインターフェース 4 0 8 は、対応する重み入力をアレイ 4 0 6 の最も上の異なるセルに送ることができる。たとえば、重みフェッチインターフェース 4 0 8 は、重み入力をセル 4 1 4 および 4 1 6 に送ることができる。重みフェッチインターフェース 4 0 8 はさらに、メモリユニット、たとえばダイナミックメモリ 3 1 0 から複数の重みを受け取り、複数の重みをアレイ 4 0 6 の最も上の別個のセルに並列に送ることができる。たとえば、重みフェッチインターフェース 4 0 8 は、異なる重みをセル 4 1 4 および 4 1 6 に同時に送ることができる。

30

【 0 0 4 3 】

いくつかの実装形態では、ホストインターフェース、たとえば、図 3 のホストインターフェース 3 0 2 は、活性化入力をアレイ 4 0 6 全体にわたって 1 つの次元に沿って、たとえば右にシフトしながら、重み入力をアレイ 4 0 6 全体にわたって別の次元に沿って、たとえば下にシフトする。たとえば、1 クロックサイクルにわたって、セル 4 1 4 における活性化入力は、セル 4 1 4 の右にあるセル 4 1 6 の活性化レジスタにシフトすることができる。同様に、セル 4 1 6 における重み入力は、セル 4 1 4 の下にあるセル 4 1 8 における重みレジスタにシフトすることができる。

40

【 0 0 4 4 】

各クロックサイクルで、各セルは、所与の重み入力、所与の活性化入力、および隣接セルからの累積出力を処理して、累積出力を生成することができる。累積出力は、与えられた重み入力と同じ次元に沿って隣接セルに渡すこともできる。各セルは、隣接セルからの累積出力を処理することなく、所与の重み入力および所与の活性化入力を処理して出力を生成することなく、与えられた重み入力および出力と

50

同じ次元に沿って隣接セルに渡されることができる。個々のセルは、図 5 を参照して以下にさらに説明される。

【 0 0 4 5 】

いくつかの実装形態では、恒等行列、すなわち、主対角線上に 1 および他の場所に 0 を有する行列をアレイ 4 0 6 に渡すことができ、それによって値ローダ 4 0 2 で与えられる入力を修正なしでアキュムレータ 4 1 0 に渡すことができる。これは、2 つの入力の要素ごとの乗算を実行するために用いられ得、ここで、アキュムレータでの第 1 の出力は、`output = MatMul(input1, identity)` として表すことができ、`MatMul` は、行列計算ユニットが行列乗算を実行するための命令であり、要素ごとの乗算結果に対応する第 2 の出力は、`output *= MatMul(input2, identity)` として表される。`*=` 演算、すなわち演算 `output = output * MatMul(input2, identity)` を実行するために、アーキテクチャ 4 0 0 は、`+=` または `*=` 計算を実行するための構成要素を含み得る。`+=` または `*=` 演算を実行するための構成要素は、アキュムレータ 4 1 0 の前、すなわちセル 4 0 4 の最後の行の後に配置することができる。いくつかの実装形態では、図 3 のベクトル計算ユニット 3 1 4 が、`+=` または `*=` 演算を実行するための構成要素を含んでもよく、すなわち、その場合、ベクトル計算ユニット 3 1 4 が、要素ごとの乗算を実行するために、`output = output * MatMul(input2, identity)` 演算を実行する。

【 0 0 4 6 】

累積出力は、重み入力と同じ列に沿って、たとえばアレイ 4 0 6 内の列の一番下に向かって渡すことができる。いくつかの実装形態では、各列の一番下において、アレイ 4 0 6 は、行よりも多い活性化入力を有する層の計算を実行するときに各列から出力される各累積出力を格納および累積する、アキュムレータユニット 4 1 0 を含み得る。いくつかの実装形態では、各アキュムレータユニットは複数の並列累積値を格納する。アキュムレータユニット 4 1 0 は、各累積出力を累積して最終累積値を生成することができる。最終累積値はベクトル計算ユニット、たとえば図 6 のベクトル計算ユニットに転送することができる。いくつかの他の実装形態では、アキュムレータユニット 4 1 0 は、行よりも少ない活性化入力を有する層を処理するときに累積を実行せずに累積値をベクトル計算ユニットに渡す。

【 0 0 4 7 】

図 5 は、シストリックアレイ内のセル、たとえば図 4 のシストリックアレイ 4 0 6 のセル 4 1 4、4 1 6、または 4 1 8 のうちの 1 つ、の例示的アーキテクチャ 5 0 0 を示す。

【 0 0 4 8 】

セルは、活性化入力を格納する活性化レジスタ 5 0 6 を含み得る。活性化レジスタは、シストリックアレイ内のセルの位置に応じて、左側の隣接セル、すなわち所与のセルの左側に位置する隣接セルから、またはユニファイドバッファから、活性化入力を受け取ることができる。セルは、重み入力を格納する重みレジスタ 5 0 2 を含み得る。重み入力は、シストリックアレイ内のセルの位置に応じて、上の隣接セルまたは重みフェッチインターフェースから転送され得る。セルは総和レジスタ 5 0 4 を含むこともできる。総和レジスタ 5 0 4 は、上の隣接セルからの累積値を格納することができる。乗算回路 5 0 8 を用いて、重みレジスタ 5 0 2 からの重み入力を活性化レジスタ 5 0 6 からの活性化入力と乗算することができる。乗算回路 5 0 8 は積を合計回路 5 1 0 に出力することができる。

【 0 0 4 9 】

合計回路 5 1 0 は、積と総和レジスタ 5 0 4 からの累積値とを合計して新たな累積値を生成することができる。次いで、合計回路 5 1 0 は、新たな累積値を、下の隣接セルに位置する別の総和レジスタに送ることができる。新たな累積値は、下の隣接セルにおける合計のためのオペランドとして用いることができる。合計回路 5 1 0 はまた、総和レジスタ 5 0 4 からの値を受け入れ、総和レジスタ 5 0 4 からの値を、乗算回路 5 0 8 からの積と合計することなく、下の隣接セルに送ることもできる。

【 0 0 5 0 】

セルは、重み入力および活性化入力を、処理のために、隣接セルにシフトすることでも

きる。たとえば、重み経路レジスタ512は、重み入力を下の隣接セル内の別の重みレジスタに送ることができる。活性化レジスタ506は、活性化入力を右の隣接セル内の別の活性化レジスタに送ることができる。したがって、重み入力と活性化入力との両方を、後続のクロックサイクルでアレイ内の他のセルによって再利用することができる。

【0051】

いくつかの実装形態では、セルは制御レジスタも含む。制御レジスタは、セルが重み入力または活性化入力のいずれかを隣接セルにシフトすべきかを決定する制御信号を記憶することができる。いくつかの実装形態では、重み入力または活性化入力をシフトすることは、1つまたは複数のクロックサイクルを要する。制御信号は、活性化入力または重み入力が乗算回路508に転送されるかどうかにも決定し得るか、または乗算回路508が活性化入力および重み入力で行うかどうかにも決定し得る。制御信号は、たとえば配線を用いて、1つまたは複数の隣接セルに渡すこともできる。

【0052】

いくつかの実装形態では、重みは重み経路レジスタ512に事前にシフトされる。重み経路レジスタ512は、重み入力をたとえば上の隣接セルから受け取り、重み入力を制御信号に基づいて重みレジスタ502に転送することができる。重みレジスタ502は、活性化入力が複数のクロックサイクルにわたってたとえば活性化レジスタ506を介してセルに転送されるとき、重み入力がセル内に留まり、隣接セルに転送されないように、重み入力を静的に格納することができる。したがって、重み入力は、たとえば乗算回路508を用いて、複数の活性化入力に適用することができ、それぞれの累積値は隣接セルに転送することができる。

【0053】

図6は、ベクトル計算ユニット602の例示的アーキテクチャ600を示す。ベクトル計算ユニット602は、行列計算ユニット、たとえば図3を参照して説明した行列計算ユニット312または図4の行列計算ユニットのアキュムレータ410から、累積値のベクトルを受け取ることができる。

【0054】

ベクトル計算ユニット602は、活性化ユニット604で累積値のベクトルを処理することができる。いくつかの実装形態では、活性化ユニットは、活性化値を生成するために各累積値に非線形関数を適用する回路を含む。たとえば、非線形関数は $\tanh(x)$ とすることができ、ここで、 x は累積値である。

【0055】

任意選択で、ベクトル計算ユニット602は、プーリング回路608を用いて、値、たとえば活性化値をプーリングすることができる。プーリング回路608は、プーリングされた値を生成するために値の1つまたは複数に集約関数を適用することができる。いくつかの実装形態では、集約関数は、値、または値のサブセットの、最大値、最小値、もしくは平均値を返す関数である。

【0056】

制御信号610は、たとえば、図3のシーケンサ306によって転送することができ、ベクトル計算ユニット602がどのように累積値のベクトルを処理するかを調整することができる。すなわち、制御信号610は、活性化値がプーリングされるかどうかを調整することができ、その場合、活性化値はたとえばユニファイドバッファ308に格納され、またはそうでなければ、制御信号610は、活性化値の取り扱いを調整することができる。制御信号610は、活性化関数またはプーリング関数、および活性化値またはプーリング値、たとえばストライド値を処理するための他のパラメータを指定することもできる。

【0057】

ベクトル計算ユニット602は、値、たとえば活性化値またはプーリングされた値を、ユニファイドバッファ、たとえば図3のユニファイドバッファ308に送ることができる。いくつかの実装形態では、プーリング回路608は、活性化値またはプーリングされた値を受け取り、活性化値またはプーリングされた値をユニファイドバッファに格納する。

【 0 0 5 8 】

図 7 は、プーリング回路のための例示的アーキテクチャ 7 0 0 を示す。プーリング回路は、プーリングされた値を生成するために、1 つまたは複数の活性化された値に集約関数を適用することができる。例示として、アーキテクチャ 7 0 0 は、 4×4 セットの活性化された値のプーリングを実行することができる。図 7 に示されるプーリングは正方形の領域、すなわち 4×4 を有するが、長方形の領域も可能である。たとえば、領域が $n \times m$ のウィンドウを有する場合、アーキテクチャ 7 0 0 は $n * m$ 個のレジスタ、すなわち n 個の列および m 個の行を有することができる。

【 0 0 5 9 】

プーリング回路アーキテクチャ 7 0 0 は、値のベクトルから、たとえば図 6 の活性化回路 6 0 4 から、要素のシーケンスを受け取ることができる。たとえば、シーケンスは画像の 8×8 部分の画素を表すことができ、プーリング回路アーキテクチャ 7 0 0 は 8×8 部分の 4×4 サブセットからの値をプーリングすることができる。いくつかの実装形態では、プーリングされた値は、いったんプーリング回路アーキテクチャ 7 0 0 によって計算されるとシーケンスに追加される。いくつかの実装形態では、ニューラルネットワークプロセッサは、複数の並列プーリング回路を含む。各クロックサイクルにわたって、各プーリング回路は活性化回路 6 0 4 から、値のベクトルからのそれぞれの要素を受け取ることができる。各プーリング回路は、活性化回路 6 0 4 から受け取った要素を、ラスタ順に到着する二次元画像として解釈することができる。

【 0 0 6 0 】

プーリング回路は、一連のレジスタおよびメモリユニットを含み得る。各レジスタは、レジスタ内部に格納されている値にわたって集約関数を適用する集約回路 7 0 6 に出力を送ることができる。集約関数は、値のセットから最小値、最大値、または平均値を返すことができる。

【 0 0 6 1 】

第 1 の値は、レジスタ 7 0 2 に送られてその内部に格納され得る。後続のクロックサイクルで、第 1 の値は後続のレジスタ 7 0 8 にシフトしてメモリ 7 0 4 に格納されることができ、第 2 の値はレジスタ 7 0 2 に送られてレジスタ 7 0 2 内に格納されることができる。

【 0 0 6 2 】

4 クロックサイクル後、4 つの値が最初の 4 つのレジスタ 7 0 2、7 0 8 ~ 7 1 2 の内部に格納される。いくつかの実装形態では、メモリユニット 7 0 4 は先入れ先出し (FIFO) の下で動作する。各メモリユニットは最大 8 つの値を格納できる。メモリユニット 7 0 4 が画素の完全な行を含んだ後、メモリユニット 7 0 4 はレジスタ 7 1 4 に値を送ることができる。

【 0 0 6 3 】

任意の所与の時点で、集約回路 7 0 6 は各レジスタからの値にアクセスすることができる。レジスタ内の値は、画像の 4×4 部分の値を表す。

【 0 0 6 4 】

プーリング回路は、集約回路 7 0 6 を用いることによって、アクセスされた値から、プーリングされた値、たとえば最大値、最小値、または平均値を生成することができる。プーリングされた値は、ユニファイドバッファ、たとえば図 3 のユニファイドバッファ 3 0 8 に送ることができる。

【 0 0 6 5 】

第 1 のプーリングされた値を生成した後、プーリング回路は、新たな値がレジスタに格納され集約回路 7 0 6 によってプーリングされることできるように、各レジスタを通して値をシフトすることによって、プーリングされた値を生成し続けることができる。たとえば、アーキテクチャ 7 0 0 では、プーリング回路は値をさらに 4 クロックサイクルにわたってシフトすることができ、それによってメモリユニット内の値をレジスタにシフトする。いくつかの実装形態では、プーリング回路は、新たな値が最後の最上位レジスタ、た

例えばレジスタ 7 1 6 に格納されるまで、新たな値をシフトする。

【 0 0 6 6 】

次いで、集約回路 7 0 6 は、レジスタに格納されている新たな値をプーリングすることができる。新たな値をプーリングした結果は、ユニファイドバッファに格納できる。

【 0 0 6 7 】

図 8 は、1 より大きいストライドでニューラルネットワークの所与の畳み込み層について計算を実行するための例示的なプロセス 8 0 0 のフローチャートである。一般に、プロセス 7 0 0 は、専用ハードウェア回路を含む 1 つまたは複数のコンピュータのシステムによって実行される。いくつかの実装形態では、例示的なプロセス 8 0 0 は、図 1 のシステムによって実行され得る。

10

【 0 0 6 8 】

システムは、専用ハードウェア回路上にニューラルネットワークを実装するための要求を受信する（ステップ 8 0 2）。特に、ニューラルネットワークは、1 より大きいストライドを有する畳み込みニューラルネットワーク層を含む。要求は、さらに、ニューラルネットワークを用いて処理する入力、ニューラルネットワークによって生成された出力テンソルを格納するための位置、または他のパラメータなど、ニューラルネットワークを実装するための他のパラメータを指定し得る。

【 0 0 6 9 】

システムは、要求に基づいて、1 より大きいストライドを有するニューラルネットワーク層を処理する際に用いられるべきマスキングテンソルを生成する（ステップ 8 0 4）。たとえば、ニューラルネットワークを実施する要求およびニューラルネットワークへの入力を特定する情報を受け取ることに基づいて、システムは、1 より大きいストライドを有するニューラルネットワーク層を処理するためのマスキングテンソルを生成する。

20

【 0 0 7 0 】

マスキングテンソルのサイズは、特定される入力の次元、または1 より大きいストライドを有するニューラルネットワーク層への入力テンソルの予想サイズに基づいて、判断することができる。マスキングテンソルに含まれる値は、1 より大きいストライドを有するニューラルネットワーク層の指定されたストライドに基づいて判断されてもよい。たとえば、ニューラルネットワーク層の指定されたストライドが4である場合、マスキングテンソルの4つおきの要素は1に設定され得る一方、マスキングテンソルの他のすべてのエントリは0に設定され得る。いくつかの実装形態では、ニューラルネットワークは、1 より大きいストライドを有する複数の層を含み得、システムは、1 より大きいストライドを有する層の各々について対応するマスキングテンソルを生成し得る。さらに、いくつかの実装形態では、システムは、たとえばメモリに、マスキング行列またはマスキング行列成分のライブラリを格納し、そのライブラリを用いることに基づいてマスキング行列を選択または生成することができる。

30

【 0 0 7 1 】

システムは、専用ハードウェア回路 1 1 0 によって実行されると、専用ハードウェア回路 1 1 0 に、ニューラルネットワークによる入力テンソルの処理中に、マスキングテンソルを用いて1 より大きいストライドを有する畳み込みニューラルネットワーク層の出力と等価な層出力テンソルを生成させる命令を生成する（ステップ 8 0 6）。たとえば、要求に応答して、ニューラルネットワーク実装エンジン 1 5 0 は、専用ハードウェア回路 1 1 0 に指示またはそれを制御して、専用ハードウェア回路 1 1 0 が1 より大きいストライドを有する畳み込みニューラルネットワーク層を用いて入力テンソルを処理した場合に等価である出力テンソル、すなわち出力ベクトルを生成するよう、命令を生成することができる。

40

【 0 0 7 2 】

システムは命令およびマスキングテンソルを専用ハードウェア回路 1 1 0 に送る（ステップ 8 0 8）。たとえば、ニューラルネットワーク実装エンジン 1 5 0 は、命令を専用ハードウェア回路 1 1 0 に与えることができ、専用ハードウェア回路 1 1 0 は、たとえば図

50

3のホストインターフェース302において命令を受け取ることができる。ニューラルネットワーク実装エンジン150は、ホストインターフェース302によっても受け取られ得る、ニューラルネットワーク計算のための他の命令および/またはパラメータも与え得る。

【0073】

図9は、1より大きいストライドを有するニューラルネットワーク計算層を計算するための例示的なプロセス900のフローチャートである。たとえば、プロセス900は、ニューラルネットワーク実装エンジン150から受け取られる命令に基づいて、図1の専用ハードウェア回路110によって実行することができる。

【0074】

たとえば、1より大きいストライドを有するニューラルネットワーク層を実装するための命令を受け取ると、ホストインターフェース302は命令を図3のシーケンサ306に送ることができる。シーケンサ306は、ニューラルネットワーク計算を実行するよう、命令を、図3の専用ハードウェア回路300を制御する低レベル制御信号に変換することができる。

【0075】

受け取られた命令に基づいて、専用ハードウェア回路300は、畳み込みニューラルネットワーク層への入力テンソルを、1のストライドを有する第2の畳み込みニューラルネットワーク層を用いて処理する(ステップ902)。たとえば、受け取られた命令から生成される制御信号は、畳み込まれたテンソルを生成するべく、専用ハードウェア回路300を制御して、1に等しいがそれ以外は畳み込みニューラルネットワーク層に等しいストライドを有する第2の畳み込みニューラルネットワーク層を用いて、入力テンソル、たとえばユニファイドバッファ308に格納されるニューラルネットワークの先行する層の出力または指定されたもしくは専用ハードウェア回路300に与えられるニューラルネットワークへの入力を処理する。

【0076】

第2の畳み込みニューラルネットワーク層を用いて入力テンソルを処理するために、制御信号は、入力テンソル、すなわちニューラルネットワークへの入力または先行するニューラルネットワークの出力に対応し得る活性化入力を図3の行列計算ユニット312に供給するように、ユニファイドバッファ308を制御し得る。制御信号はまた、図3のダイレクトメモリアクセスエンジン304および/またはダイナミックメモリ310に命令して、1のストライド、すなわち単位ストライドを有するが、それ以外は1より大きいストライドを有するニューラルネットワーク層と同等である第2のニューラルネットワーク層に対応する重みを行列計算ユニット312に与えてもよい。

【0077】

シーケンサ306はさらに、重みを用いて、たとえば図3に関して説明したプロセスを用いて、入力テンソルを処理するように行列計算ユニット312を制御する命令を生成することができる。いくつかの実装形態では、行列計算ユニット312は、2015年9月3日に提出された米国特許出願第14/844,738号に記載されている技法を用いて畳み込みを実行し、その全体をここに引用により援用する。

【0078】

行列計算ユニット312は、制御信号に基づいて計算を行い、畳み込まれたテンソルをベクトル計算ユニット314に出力する。たとえば、行列計算ユニット312は、行列計算ユニット312が生成した出力のベクトルをベクトル計算ユニット314に送る。出力のベクトルは、1のストライドを有するがそれ以外の点では1より大きいストライドを有するニューラルネットワーク層と同等であるニューラルネットワーク層に対応する重みを用いて入力テンソルを処理することに基づいて判定され得る。ベクトル計算ユニット314は、畳み込まれたテンソルをユニファイドバッファ308に格納することができる。

【0079】

畳み込まれたテンソルを生成するために畳み込みニューラルネットワーク層を介して1

10

20

30

40

50

のストライドで活性化入力を処理した後、専用ハードウェア回路 300 は、第 2 の畳み込みニューラルネットワーク層が 1 より大きいストライドを有する畳み込みネットワーク層のストライドを有すると仮定した場合に生成されなかったであろう要素を零出力する（ステップ 904）。要素を零出力するとは、通常、その要素の現在の値を 0 に置き換えることを指す。値を取り消す、すなわち 0 にすることは、畳み込まれたテンソルとマスキングテンソル、すなわちニューラルネットワーク処理エンジン 150 によって生成され専用ニューラルネットワークに送られるマスキングテンソルとの要素ごとの乗算を実行することによって達成することができる。

【0080】

入力テンソルが指定されたストライドで畳み込みニューラルネットワーク層によって処理された場合に生成されなかったであろう畳み込まれたテンソルの値を取り消すために、シーケンサ 306 は、制御信号を送って、行列乗算ユニット 312 に、畳み込まれたテンソルとマスキングテンソルとの要素ごとの乗算を実行させることができる。畳み込まれたテンソルは、シーケンサ 306 からの他の制御信号に基づいてユニファイドバッファ 308 から行列乗算ユニット 312 に送られてもよく、マスキングテンソルは、シーケンサ 306 からダイレクトメモリアクセスエンジン 304 またはダイナミックメモリ 310 への制御信号に基づいて、すなわち、マスキングテンソルが専用ハードウェア回路 300 によって受け取られ、ダイナミックメモリ 310 に格納された後、行列計算ユニット 312 に送られてもよい。

【0081】

一般に、図 8 に関して説明したように、マスキングテンソルは、1 より大きいストライドを有する畳み込みニューラルネットワーク層を用いて入力テンソルを処理することによって生成されるであろう要素に対応する要素位置に単位値要素、すなわち 1 の値を含み、他のすべての位置、すなわち 1 より大きいストライドを有する畳み込みニューラルネットワーク層を用いて活性化値を処理することによって生成されないであろう要素に対応する位置に 0 値要素を含むベクトルである。

【0082】

マスキングテンソルは、たとえば、ダイナミックメモリ 310 に格納されてもよく、シーケンサ 306 は、マスキングテンソルをダイナミックメモリ 310 から行列計算ユニット 312 に送るよう制御信号を送信してもよい。たとえば、専用ハードウェア回路 300 に与えられる命令は、マスキングテンソルを識別してもよく、たとえばマスキングテンソルのダイナミックメモリ 310 内の位置を与えてもよく、またはその場合にダイナミックメモリ 310 に格納されるマスキングテンソルを定義するデータを含んでもよく、シーケンサ 306 は、ダイナミックメモリ 310 内のその位置に格納されるマスキングテンソルを行列計算ユニット 312 に送らせる制御信号を送信してもよい。さらに、シーケンサ 306 は、ユニファイドバッファ 308 に格納される畳み込まれたテンソルを行列計算ユニット 312 に与えさせるための制御信号を与えることができる。そして、行列計算ユニット 312 は、畳み込まれたテンソルとマスキングテンソルとの要素ごとの乗算を行い、修正された畳み込まれたテンソルを生成する。修正された畳み込まれたテンソルは、行列計算ユニット 312 からベクトル計算ユニット 314 によって受け取られることができる。ベクトル計算ユニット 314 は、任意選択で、修正された畳み込まれたテンソルをユニファイドバッファ 308 に格納することができる。

【0083】

マスキングテンソルでの要素ごとの乗算のため、修正された畳み込まれたテンソルは、入力テンソルが 1 より大きい指定されたストライドを有するニューラルネットワーク層を用いて処理された場合に出力されるであろう値を含む。修正された畳み込まれたテンソルは、入力テンソルが指定されたストライドを有する畳み込みニューラルネットワークで処理された場合に出力されなかったであろう、1 のストライドでの畳み込みニューラルネットワーク層を用いた入力テンソルの計算において出力される値に対応する位置に 0 を含む。他の実装形態では、畳み込まれたテンソルの要素を 0 にする他の方法を利用することが

10

20

30

40

50

できる。たとえば、畳み込まれた行列は、修正された形でユニファイドバッファ 308 または他のメモリにおいて書き直され得、指定されたストライドを有する畳み込みニューラルネットワークを用いての入力テンソルの計算において出力される値に対応する要素は変更されず、他の要素は 0 として書き込まれる。

【0084】

ベクトル計算ユニット 314 は、修正された畳み込まれたテンソルを受け取り、修正された畳み込まれたテンソルに対して最大プーリングを実行して、1 より大きいストライドを有する畳み込みニューラルネットワーク層の層出力テンソルを生成する（ステップ 906）。たとえば、ベクトル計算ユニット 314 は、行列計算ユニット 312 から修正された畳み込まれたテンソルを受け取り、プーリング回路 608 を用いて、修正された畳み込まれたテンソルに対して最大プーリングを実行することができる。最大プーリングは、データのセットを受け取り、そのデータの 1 つまたは複数のサブセットの各々について、サブセット内の要素の最大値を出力する操作である。修正された畳み込まれたテンソルに対して最大プーリングを実行すると、修正された畳み込まれたテンソルの要素の複数のサブセットの各々について、サブセットの最大値を含むテンソルが得られる結果となる。ベクトル計算ユニット 314 は、畳み込みニューラルネットワーク層の指定されたストライドに基づいて決定される修正された畳み込まれたテンソルのウィンドウについて最大プーリングを実行することができる。たとえば、ストライドが 2 の場合、プーリング回路 608 は、 2×2 ウィンドウを用いて最大プーリングを実行し、各 2×2 ウィンドウからの最大値要素を含む層出力テンソルを生成する。4 のストライドを有するニューラルネットワーク層の場合、プーリング回路 608 は、 4×4 ウィンドウを用いて最大プーリングを実行し、各 4×4 ウィンドウからの最大値要素を含む層出力テンソルを生成する。最大プーリング操作の結果は、ベクトル計算ユニット 314 によってユニファイドバッファ 308 に格納され、その結果とは、専用ハードウェア回路 300 が 1 より大きいストライドを有するニューラルネットワーク層を用いて入力テンソルを処理した場合に生成されるであろう出力と等価の出力テンソルである。ニューラルネットワークの後続の層の処理は、最終的にニューラルネットワークの推論を取得するよう、層出力テンソルを用いて実行されてもよい。

【0085】

図 10 は、1 より大きいストライドでのニューラルネットワークの所与の層に対する計算の例を示す。図 10 の例は、図 7 のプロセスおよび図 2 の専用ハードウェア回路 300 を用いて実行することができる。例として、図 10 の例は、4 のストライド有する畳み込みニューラルネットワーク層を活性化値の 8×8 アレイに適用する。畳み込みニューラルネットワーク層は、活性化値の 8×8 アレイに適用される重みの 4×4 カーネルを有してもよい。活性化値は、ニューラルネットワークに入力される画像の 8×8 部分、すなわち画像の 8×8 部分に対応する値のシーケンスを表すことができる。代替的に、活性化値の 8×8 アレイは、別の入力テンソル、たとえばニューラルネットワークの先行する層の出力に対応する入力テンソルの 8×8 部分を表すことができる。

【0086】

図 10 の部分 (a) において、 8×8 入力テンソルは、1 のストライドを有するが他の点では 1 より大きいストライドを有する畳み込みニューラルネットワーク層と同等である畳み込みニューラルネットワーク層を用いて処理される。したがって、部分 (a) に示される重みの 4×4 のカーネルは、最初に、入力テンソルの最初の 4 行および最初の 4 列に対応する入力テンソルの要素に適用され得る（値は示されてはいない）。処理の結果は、結果として得られる畳み込まれたテンソルの第 1 の要素、すなわち、図 10 の部分 (a) に示される、結果として得られる畳み込まれたテンソルの要素「a」であり得る。

【0087】

入力テンソルの処理は、指定された 4 のストライドではなく、1 のストライドで畳み込みニューラルネットワーク層を用いて実行されるので、部分 (a) に示された重みの 4×4 セットは、活性化値アレイの最初の 4 行および入力テンソルの第 2 列から第 5 列に対応

する入力テンソルの要素に適用されてもよい（値は示されず）。処理結果は、畳み込まれたテンソルの第2の要素、すなわち、図10の部分（a）に示す畳み込み結果の要素「b」である。重みの4×4セットを活性化値アレイに1のストライドを用いて適用することによって、すなわち重みの4×4セットを活性化値アレイに増分的に列方向および行方向の両方に適用することによって、プロセスを繰り返してもよい。この処理の結果、図10の部分（a）に示す8×8の畳み込まれたテンソルが得られる。

【0088】

次に、図9の部分（b）に示すように、畳み込まれたテンソルとマスキングテンソルとの間で要素ごとの乗算が行われ、修正された畳み込まれたテンソルが得られる。マスキングテンソルのサイズは、入力テンソルのサイズまたは畳み込まれたテンソルのサイズに基づいて判断され、それは、1のストライドを有する畳み込みニューラルネットワーク層を用いる図10の部分（a）での処理のため、一般的に等しい。マスキングテンソルは、入力テンソルが指定されたストライドを有する畳み込みニューラルネットワーク層を用いて処理された場合に生成されるであろう値に対応する位置に、単位値、すなわち1を含む。その場合、一般に、マスキングテンソルにおける単位値エントリの位置は、畳み込みニューラルネットワーク層の指定されたストライドに依存する。図10の例では、畳み込みニューラルネットワーク層は4のストライドを有するので、マスキングテンソルは列方向および行方向の両方において4つおきの位置に単位値を含むであろう。マスキングテンソルの他のエントリには0値が割り当てられ、畳み込まれたテンソルとマスキングテンソルとの要素ごとの乗算は、畳み込みニューラルネットワークが指定されたストライドを有する状態で入力テンソルが処理された場合に生成されないであろうすべての値を0にする結果となることになる。

【0089】

修正された畳み込まれたテンソルを生成するために、畳み込まれたテンソルとマスキングテンソルとの要素ごとの乗算が実行される。図10に示すように、要素ごとの乗算の後、畳み込まれたテンソルの4つおきの要素は維持され、畳み込まれたテンソルの要素の残りはマスキング行列の対応する0値要素との乗算により0になる。したがって、8×8の畳み込まれたテンソルの要素のうち、4つの要素だけが非0のままである。

【0090】

いくつかの実装形態では、最初に、畳み込まれたテンソルの要素に非単位係数を乗算し、続いて、それらの要素に第2の非単位係数を乗算することによって、同様の結果を得ることができる。たとえば、マスキングテンソルは、入力テンソルが指定されたストライドを有する畳み込みニューラルネットワーク層を用いて処理された場合に生成されるであろう値に対応する位置に、2（または他の値）を含んでもよい。したがって、上記の例に従って、畳み込まれたテンソルとマスキングテンソルとの要素ごとの乗算は、畳み込まれたテンソルの4つおきの要素が2倍になり、要素の残りが0である、修正された畳み込まれたテンソルを生成する。その後、修正された畳み込まれたテンソルの半分（または他の値の逆数）によるスカラー乗算が実行されてもよい。代替的に、修正された畳み込まれたテンソルと第2のマスキングテンソルとの要素ごとの乗算が実行されてもよく、第2のマスキングテンソルは、入力テンソルが指定されたストライドを有する畳み込みニューラルネットワーク層を用いて処理された場合に生成されるであろう値に対応する位置に2分の1の値を含む。

【0091】

続いて、図10の部分（c）において修正された畳み込み結果アレイに対して最大プーリングが行われる。最大プーリングの結果は、入力テンソルがストライドが4の畳み込みニューラルネットワーク層によって処理された場合に得られるであろう結果と同等である。図6のプロセスを用いて、修正された畳み込まれたテンソルに対して最大プーリングが実行され、修正された畳み込まれたテンソルの各4×4ウィンドウの最大値が識別される。次いで、最大プーリングの結果は、4のストライドを有する畳み込みニューラルネットワーク層の出力テンソルとして格納される。入力テンソルは8×8アレイであったため、

4のストライドを有するニューラルネットワーク層による処理は2×2出力アレイをもたらす。2×2出力アレイは、図2のユニファイドバッファ308に、たとえばラスト順で格納することができる。2×2出力アレイの値は、ニューラルネットワークの次の層への入力として与えられてもよい。

【0092】

本明細書において記載される主題および機能的動作の実施形態は、本明細書に開示される構造およびそれらの構造的等価物を含む、デジタル電子回路系において、有形で実施されるコンピュータソフトウェアもしくはファームウェアにおいて、コンピュータハードウェアにおいて、またはそれらの1つ以上の組合せにおいて実現され得る。本明細書に記載される主題の実施形態は、1つ以上のコンピュータプログラムとして、すなわち、データ処理装置による実行のために、または、データ処理装置の動作を制御するために有形の非一時的なプログラム担体上でエンコードされたコンピュータプログラム命令の1つ以上のモジュールとして実現され得る。代替的に、または加えて、プログラム命令は、データ処理装置による実行に対して好適な受信側装置への送信のために情報をエンコードするように生成される、たとえばマシンにより生成された電気信号、光信号、または電磁気信号などの、人為的に生成された伝搬される信号上でエンコードすることができる。コンピュータ記憶媒体は、コンピュータ可読記憶装置、コンピュータ可読記憶基板、ランダムもしくはシリアルアクセスメモリデバイス、または、それらの1つ以上の組合せであり得る。

【0093】

「データ処理装置」という用語は、例としてプログラマブルプロセッサ、コンピュータ、または複数のプロセッサもしくはコンピュータを含む、データを処理するためのすべての種類の装置、デバイスおよびマシンを包含する。当該装置は、たとえばFPGA（フィールドプログラマブルゲートアレイ）またはASIC（特定用途向け集積回路）といった特定目的論理回路を含み得る。当該装置は、ハードウェアに加えて、たとえばプロセッサファームウェア、プロトコルスタック、データベース管理システム、オペレーティングシステム、または、それらの1つ以上の組合せを構成するコードといった、当該コンピュータプログラムについて実行環境を作成するコードをさらに含み得る。

【0094】

（プログラム、ソフトウェア、ソフトウェアアプリケーション、モジュール、ソフトウェアモジュール、スクリプトまたはコードとも称され、または記載され得る）コンピュータプログラムは、コンパイル型もしくはインタープリタ型言語、または宣言型もしくは手続き型言語を含む任意の形態のプログラミング言語で記述され得、スタンドアロンプログラムとして、または、モジュール、コンポーネント、サブルーチン、もしくは、コンピューティング環境で使用するのに好適な他のユニットとして任意の形態で展開され得る。コンピュータプログラムは、ファイルシステムにおけるファイルに対応し得るが、対応する必要があるわけではない。プログラムは、当該プログラムに専用である単一のファイルにおいて、または、複数の連携ファイル（*coordinated files*）（たとえばコードの1つ以上のモジュール、サブプログラムまたは部分を格納するファイル）において、他のプログラムまたはデータ（たとえばマークアップ言語ドキュメントに格納される1つ以上のスクリプト）を保持するファイルの一部に格納され得る。コンピュータプログラムは、1つの場所に位置するかもしくは複数の場所にわたって分散され通信ネットワークによって相互接続される1つのコンピュータまたは複数のコンピュータ上で実行されるように展開され得る。

【0095】

本明細書に記載されるプロセスおよび論理フローは、入力データ上で動作し出力を生成することにより機能を実行するよう1つ以上のプログラマブルコンピュータが1つ以上のコンピュータプログラムを実行することによって実行され得る。本プロセスおよび論理フローの実行、ならびに本装置の実施は、さらに、たとえばFPGA（フィールドプログラマブルゲートアレイ）またはASIC（特定用途向け集積回路）といった特殊目的論理回路系によってもなされ得る。

【 0 0 9 6 】

コンピュータプログラムの実行に好適であるコンピュータは、例として、汎用マイクロプロセッサもしくは特殊目的マイクロプロセッサもしくはその両方または任意の種類の中央処理ユニットに基づき得る。一般に、中央処理ユニットは、リードオンリメモリもしくはランダムアクセスメモリまたはその両方から命令およびデータを受け取る。コンピュータの必須の要素は、命令を実行するための中央処理ユニットと、命令およびデータを格納するための1つ以上のメモリデバイスとである。一般に、コンピュータはさらに、たとえば磁気ディスク、光磁気ディスクまたは光ディスクといった、データを格納するための1つ以上の大容量記憶装置を含むか、当該1つ以上の大容量記憶装置からデータを受け取るかもしくは当該1つ以上の大容量記憶装置にデータを転送するよう動作可能に結合されるか、またはその両方を行う。しかしながら、コンピュータはそのような装置を有する必要はない。さらに、コンピュータはたとえば、携帯電話、携帯情報端末（PDA）、モバイルオーディオまたはビデオプレーヤ、ゲームコンソール、全地球測位システム（GPS）受信機、またはポータブル記憶装置（たとえばユニバーサルシリアルバス（USB）フラッシュドライブ）といった別のデバイスに埋め込まれ得る。

10

【 0 0 9 7 】

コンピュータプログラム命令およびデータを格納するのに好適であるコンピュータ可読媒体は、例として、たとえばEPROM、EEPROMおよびフラッシュメモリデバイスといった半導体メモリデバイスを含むすべての形態の不揮発性メモリ、媒体およびメモリデバイス；たとえば内部ハードディスクまたはリムーバブルディスクといった磁気ディスク；光磁気ディスク；ならびにCD-ROMおよびDVD-ROMディスクを含む。プロセッサおよびメモリは、特殊目的論理回路によって補足され得るか、または特殊目的論理回路に組み込まれ得る。

20

【 0 0 9 8 】

ユーザとの対話を求めて、本明細書に記載される主題の実施形態は、たとえばCRT（陰極線管）またはLCD（液晶ディスプレイ）モニタといったユーザに対して情報を表示するための表示デバイスと、たとえばマウス、トラックボールといったユーザがコンピュータに入力を送ることができるキーボードおよびポインティングデバイスとを有するコンピュータ上で実現され得る。他の種類のデバイスが、同様に、ユーザとの対話を求めて用いられ得；たとえば、ユーザに提供されるフィードバックは、たとえば視覚フィードバック、聴覚フィードバックまたは触覚フィードバックといった任意の形態の感覚フィードバックであり得；ユーザからの入力、音響入力、音声入力、または触覚入力を含む任意の形態で受け取られ得る。加えて、コンピュータは、ユーザが使用するデバイスにドキュメントを送信しユーザが使用するデバイスからドキュメントを受信することによって、たとえば、ウェブブラウザから受信された要求に応答してユーザのクライアントデバイス上のウェブブラウザにウェブページを送信することによって、ユーザと対話し得る。

30

【 0 0 9 9 】

本明細書に記載される主題の実施形態は、たとえばデータサーバとしてバックエンドコンポーネントを含む計算システムにおいて実現され得るか、たとえばアプリケーションサーバといったミドルウェアコンポーネントを含む計算システムにおいて実現され得るか、たとえば本明細書に記載される主題の実現例とユーザが対話することが可能であるグラフィカルユーザインターフェイスもしくはウェブブラウザを有するクライアントコンピュータといったフロントエンドコンポーネントを含む計算システムにおいて実現され得るか、または1つ以上のそのようなバックエンドコンポーネント、ミドルウェアコンポーネントもしくはフロントエンドコンポーネントの任意の組合せの計算システムにおいて実現され得る。システムのコンポーネントは、たとえば通信ネットワークといったデジタルデータ通信の任意の形態または媒体によって相互接続され得る。通信ネットワークの例は、ローカルエリアネットワーク（「LAN」）およびワイドエリアネットワーク（「WAN」）、たとえばインターネットを含む。

40

【 0 1 0 0 】

50

計算システムはクライアントおよびサーバを含むことができる。クライアントとサーバとは一般に互いから遠隔にあり、典型的には通信ネットワークを通じて対話する。クライアントとサーバとの関係は、それぞれのコンピュータ上で実行されるとともに互いに対してクライアント - サーバ関係を有するコンピュータプログラムによって発生する。

【0101】

本明細書は多くの特定の實現例の詳細を含んでいるが、これらは如何なる発明の範囲または請求され得るものの範囲に対する限定としても解釈されるべきではなく、特定の発明の特定の実施形態に特有の特徴であり得る記載として解釈されるべきである。別個の実施形態の文脈で本明細書において記載されるある特徴は、単一の実施形態において組合せでも實現され得る。反対に、単一の実施形態の文脈において記載されるさまざまな特徴は、複数の実施形態において別々に、または任意の好適な部分的組合せでも實現され得る。さらに、特徴は、ある組合せにおいて作用すると上で記載され、最初はそのように請求されていさえする場合もあるが、請求される組合せからの1つ以上の特徴はいくつかの場合には当該組合せから削除され得、請求される組合せは、部分的組合せまたは部分的組合せの変形例に向けられ得る。

【0102】

同様に、動作が図においては特定の順に示されているが、そのような動作は、望ましい結果を達成するために、示された当該特定の順もしくは連続した順で実行される必要があると理解されるべきではなく、または、すべての示された動作が実行される必要があると理解されるべきではない。ある状況においては、マルチタスキングおよび並列処理が有利であり得る。さらに、上述の実施形態におけるさまざまなシステムモジュールおよびコンポーネントの分離は、すべての実施形態においてそのような分離を必要とすると理解されるべきではなく、記載されるプログラムコンポーネントおよびシステムは一般に単一のソフトウェア製品に統合され得るかまたは複数のソフトウェア製品にパッケージ化され得ることが理解されるべきである。

【0103】

その他の實現例は、以下の例にまとめられる。

例1：ハードウェア回路上でニューラルネットワークを処理するよう要求を受け取るとを備え、ニューラルネットワークは、1より大きいストライドを有する第1の畳み込みニューラルネットワーク層を含み、さらに、これに応答して、ハードウェア回路によって実行されると、ハードウェア回路に、ニューラルネットワークによる入力テンソルの処理中に、動作を実行することによって第1の畳み込みニューラルネットワーク層の出力と等価な層出力テンソルを生成させる命令を生成することを備え、動作は、第1の畳み込みニューラルネットワーク層への入力テンソルを、1に等しいストライドを有するがそれ以外は第1の畳み込みニューラルネットワーク層と同等である第2の畳み込みニューラルネットワーク層を用いて処理することにより、第1のテンソルを生成することと、第2の畳み込みニューラルネットワーク層が第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第1のテンソルの要素を零出力して、第2のテンソルを生成することと、第2のテンソルに対して最大プーリングを実行して層出力テンソルを生成することを含む、方法。

【0104】

例2：第1のテンソルの要素を零出力することは、第1のテンソルの要素のサブセットに0を乗算することと、サブセットに含まれていない第1のテンソルの要素に1を乗算することを含む、例1の方法。

【0105】

例3：第1のテンソルの要素を零出力することは、マスキングテンソルと第1のテンソルとの要素ごとの乗算を実行して第2のテンソルを生成することを含み、マスキングテンソルは、(i)第2の畳み込みニューラルネットワーク層が第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第1のテンソルの要素に対応するマスキングテンソルの各要素位置において0を含み、(ii)マスキング

テンソルの各他の要素位置において 1 を含む、例 1 の方法。

【 0 1 0 6 】

例 4 : マスキングテンソルは、ハードウェア回路によってアクセス可能なメモリに格納され、マスキングテンソルと第 1 のテンソルとの要素ごとの乗算は、ハードウェア回路に含まれる、ハードウェアにおいて実装されるベクトル計算ユニットによって実行される、例 3 の方法。

【 0 1 0 7 】

例 5 : 第 1 のテンソルの要素を零出力することは、第 1 のマスキングテンソルと第 1 のテンソルとの要素ごとの乗算を実行して、修正された第 1 のテンソルを生成することを含み、第 1 のマスキングテンソルは、(i) 第 2 の畳み込みニューラルネットワーク層が第 1 の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第 1 のテンソルの要素に対応するマスキングテンソルの各要素位置において 0 を含み、(i i) 第 2 の畳み込みニューラルネットワーク層が第 1 の畳み込みニューラルネットワーク層のストライドを有する場合には生成されたであろう第 1 のテンソルの要素に対応するマスキングテンソルの各要素位置においてそれぞれの非 0 値を含み、第 1 のテンソルの要素を零出力することはさらに、第 2 のマスキングテンソルと修正された第 1 のテンソルとの要素ごとの乗算を実行することを含み、第 2 のマスキングテンソルは、第 2 の畳み込みニューラルネットワーク層が第 1 の畳み込みニューラルネットワーク層のストライドを有する場合に生成されるであろう第 1 のテンソルの要素に対応する各要素位置において、第 1 のマスキングテンソルのそれぞれの非 0 値の逆数を含む、例 1 の方法。

【 0 1 0 8 】

例 6 : 最大プーリングを実行することは、第 1 の畳み込みニューラルネットワーク層のストライドによって定義される第 2 のテンソルの 1 つまたは複数のウィンドウの各々について、ウィンドウ内の要素の最大値要素を取得することを含む、例 1 ~ 例 5 の 1 つの方法。

【 0 1 0 9 】

例 7 : 第 2 のテンソルの 1 つまたは複数のウィンドウの各々は、畳み込みニューラルネットワーク層のストライドに対応する次元を有する矩形ウィンドウであり、第 2 のテンソルの異なる要素を含む、例 6 の方法。

【 0 1 1 0 】

例 8 : 最大プーリングを実行することは、第 2 のテンソルの要素の 1 つまたは複数のサブセットの各々について、サブセットの最大値要素を取得することを含む、例 1 ~ 例 7 の 1 つの方法。

【 0 1 1 1 】

例 9 : 第 2 のテンソル上で実行される最大プーリングは、ハードウェア回路のプーリング回路によって実行される、例 1 ~ 例 8 の 1 つの方法。

【 0 1 1 2 】

例 10 : 畳み込みニューラルネットワーク層は、ニューラルネットワーク内の第 1 のニューラルネットワーク層であり、入力テンソルは、デジタル画像の、デジタル画像の画素に対応する要素を含む表現である、例 1 ~ 例 9 の 1 つの方法。

【 0 1 1 3 】

例 11 : 入力テンソルはハードウェア回路のユニファイドバッファに格納され、第 2 の畳み込みニューラルネットワーク層の重みはハードウェア回路のダイナミックメモリに格納され、第 2 の畳み込みニューラルネットワーク層を用いて第 1 の畳み込みニューラルネットワーク層への入力テンソルを処理することは、入力テンソルをユニファイドバッファからハードウェアで実装されるハードウェア回路の行列計算ユニットに送ることと、ダイナミックメモリからハードウェア回路の行列計算ユニットに第 2 の畳み込みニューラルネットワーク層の重みを送ることと、ハードウェア回路の行列計算ユニットによって、第 2 の畳み込みニューラルネットワーク層の重みを用いて入力テンソルを処理して、第 1 のテンソルを生成することとを含む、例 1 ~ 例 10 の 1 つの方法。

【 0 1 1 4 】

例 1 2 : システムであって、ハードウェア回路と、命令を格納する 1 つまたは複数の記憶装置とを備え、命令は、ハードウェア回路によって実行されると、ハードウェア回路に動作を実行させるよう動作可能であり、動作は、1 より大きいストライドを有する畳み込みニューラルネットワーク層への入力テンソルを、1 に等しいストライドを有するがそれ以外は畳み込みニューラルネットワーク層と同等である第 2 の畳み込みニューラルネットワーク層を用いて処理することにより、第 1 のテンソルを生成することと、第 2 の畳み込みニューラルネットワーク層が畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第 1 のテンソルの要素を零出力して、第 2 のテンソルを生成することと、第 2 のテンソルに対して最大プーリングを実行して層出力テンソルを生成することとを含む、システム。

10

【 0 1 1 5 】

例 1 3 : 第 1 のテンソルの要素を零出力することは、マスキングテンソルと第 1 のテンソルとの要素ごとの乗算を実行して第 2 のテンソルを生成することを含み、マスキングテンソルは、(i) 第 2 の畳み込みニューラルネットワーク層が第 1 の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第 1 のテンソルの要素に対応するマスキングテンソルの各要素位置において 0 を含み、(i i) マスキングテンソルの各他の要素位置において 1 を含む、例 1 2 のシステム。

【 0 1 1 6 】

例 1 4 : マスキングテンソルは、ハードウェア回路によってアクセス可能なメモリに格納され、マスキングテンソルと第 1 のテンソルとの要素ごとの乗算は、ハードウェア回路に含まれる、ハードウェアで実装されるベクトル計算ユニットによって実行される、例 1 3 のシステム。

20

【 0 1 1 7 】

例 1 5 : 最大プーリングを実行することは、第 1 の畳み込みニューラルネットワーク層のストライドによって定義される第 2 のテンソルの 1 つまたは複数のウィンドウの各々について、ウィンドウ内の要素の最大値要素を取得することを含む、例 1 2 ~ 例 1 4 の 1 つのシステム。

【 0 1 1 8 】

例 1 6 : 第 2 のテンソルの 1 つまたは複数のウィンドウの各々は、畳み込みニューラルネットワーク層のストライドに対応する次元を有する矩形ウィンドウであり、第 2 のテンソルの異なる要素を含む、例 1 5 のシステム。

30

【 0 1 1 9 】

例 1 7 : 第 2 のテンソル上で実行される最大プーリングは、ハードウェア回路のプーリング回路によって実行される、例 1 2 ~ 例 1 6 の 1 つのシステム。

【 0 1 2 0 】

例 1 8 : 畳み込みニューラルネットワーク層は、ニューラルネットワーク内の第 1 のニューラルネットワーク層であり、入力テンソルは、デジタル画像の、デジタル画像の画素に対応する要素を含む表現である、例 1 2 ~ 例 1 7 の 1 つのシステム。

【 0 1 2 1 】

例 1 9 : 入力テンソルはハードウェア回路のユニファイドバッファに格納され、第 2 の畳み込みニューラルネットワーク層の重みはハードウェア回路のダイナミックメモリに格納され、第 2 の畳み込みニューラルネットワーク層を用いて第 1 の畳み込みニューラルネットワーク層への入力テンソルを処理することは、入力テンソルをユニファイドバッファからハードウェアで実装されるハードウェア回路の行列計算ユニットに送ることと、ダイナミックメモリからハードウェア回路の行列計算ユニットに第 2 の畳み込みニューラルネットワーク層の重みを送ることと、ハードウェア回路の行列計算ユニットによって、第 2 の畳み込みニューラルネットワーク層の重みを用いて入力テンソルを処理して、第 1 のテンソルを生成することとを含む、例 1 2 ~ 例 1 8 の 1 つのシステム。

40

【 0 1 2 2 】

50

例 20：コンピュータプログラムでエンコードされたコンピュータ可読記憶装置であって、コンピュータプログラムは、1つまたは複数のコンピュータによって実行されると、1つまたは複数のコンピュータに動作を実行させる命令を含み、動作は、ハードウェア回路上でニューラルネットワークを処理するよう要求を受け取ることを含み、ニューラルネットワークは、1より大きいストライドを有する第1の畳み込みニューラルネットワーク層を含み、動作はさらに、これにตอบสนองして、ハードウェア回路によって実行されると、ハードウェア回路に、ニューラルネットワークによる入力テンソルの処理中に、動作を実行することによって第1の畳み込みニューラルネットワーク層の出力と等価な層出力テンソルを生成させる命令を生成することを含み、動作は、第1の畳み込みニューラルネットワーク層への入力テンソルを、1に等しいストライドを有するがそれ以外は第1の畳み込みニューラルネットワーク層と同等である第2の畳み込みニューラルネットワーク層を用いて処理することにより、第1のテンソルを生成することと、第2の畳み込みニューラルネットワーク層が第1の畳み込みニューラルネットワーク層のストライドを有する場合には生成されなかったであろう第1のテンソルの要素を零出力して、第2のテンソルを生成することと、第2のテンソルに対して最大プーリングを実行して層出力テンソルを生成することを含む、コンピュータ可読記憶装置。

【0123】

主題の特定の実施形態が記載された。他の実施形態は以下の請求の範囲内にある。たとえば、請求項において記載されるアクションは、異なる順で実行され得、それでも望ましい結果を達成し得る。一例として、添付の図において示されるプロセスは、望ましい結果を達成するために、示された特定の順または連続する順であることを必ずしも必要としない。ある実現例においては、マルチタスキングおよび並列処理が有利であり得る。

【図1】

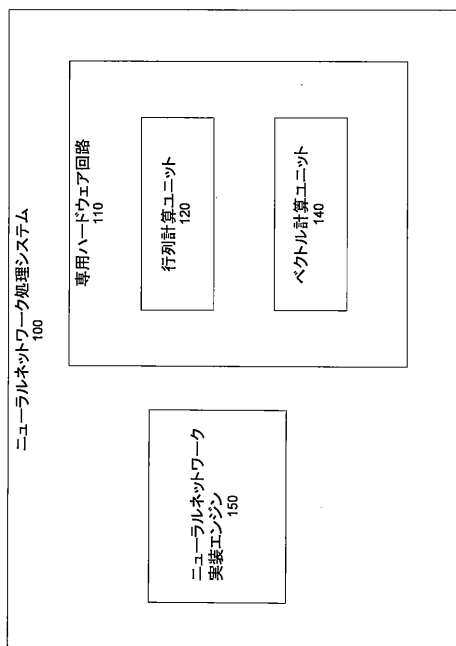
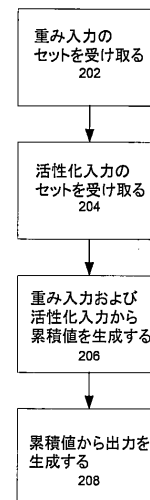


FIG. 1

【図2】



200

FIG. 2

【 図 3 】

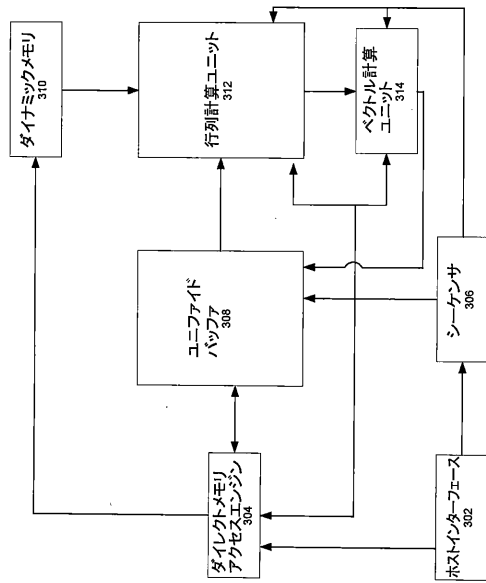


FIG. 3

【 図 4 】

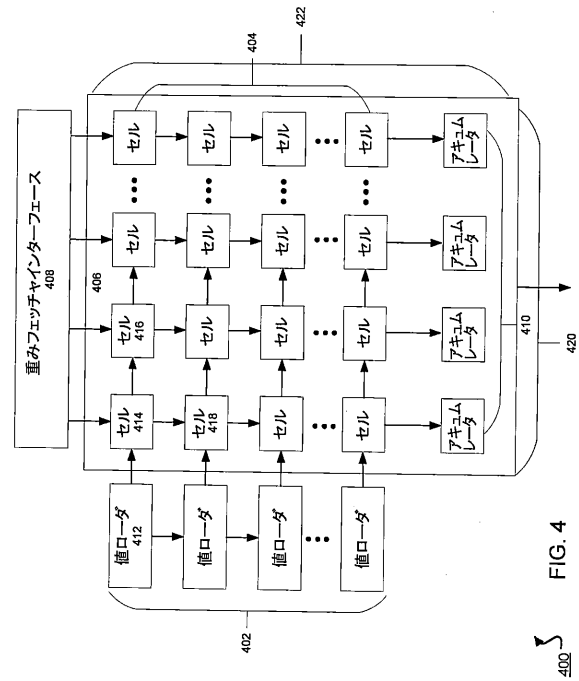


FIG. 4

【 図 5 】

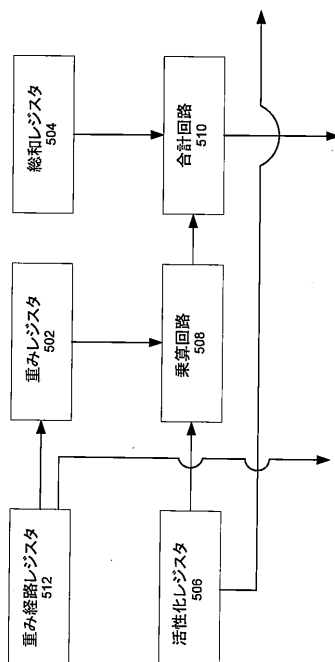


FIG. 5

【 図 6 】

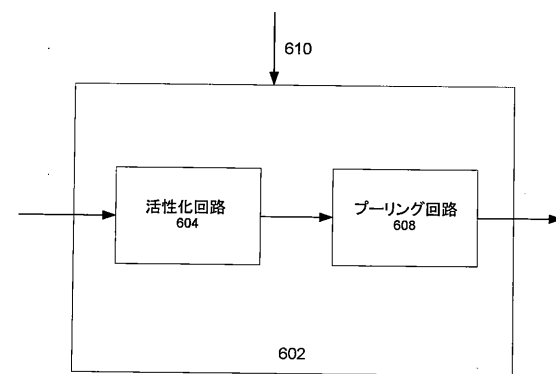
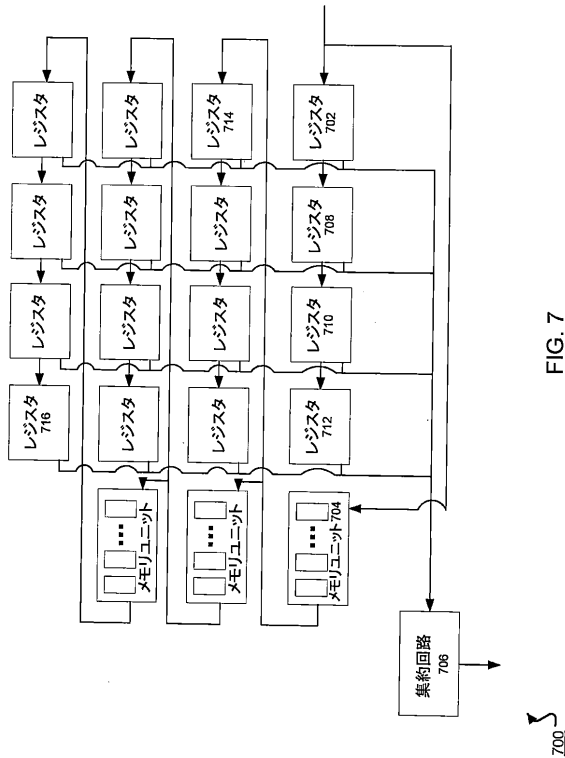


FIG. 6

【図 7】



【図 8】

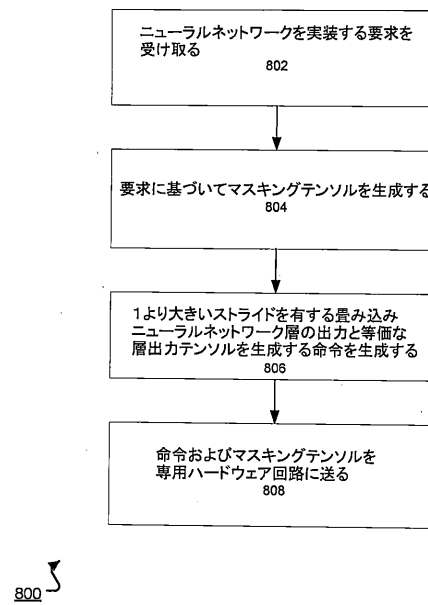


FIG. 8

【図 9】

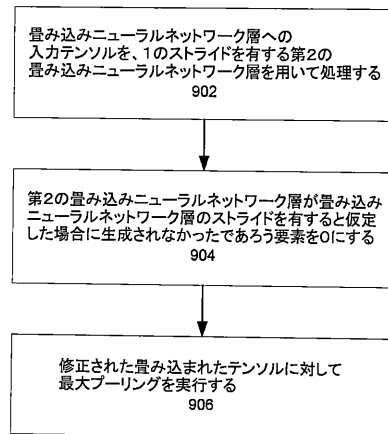
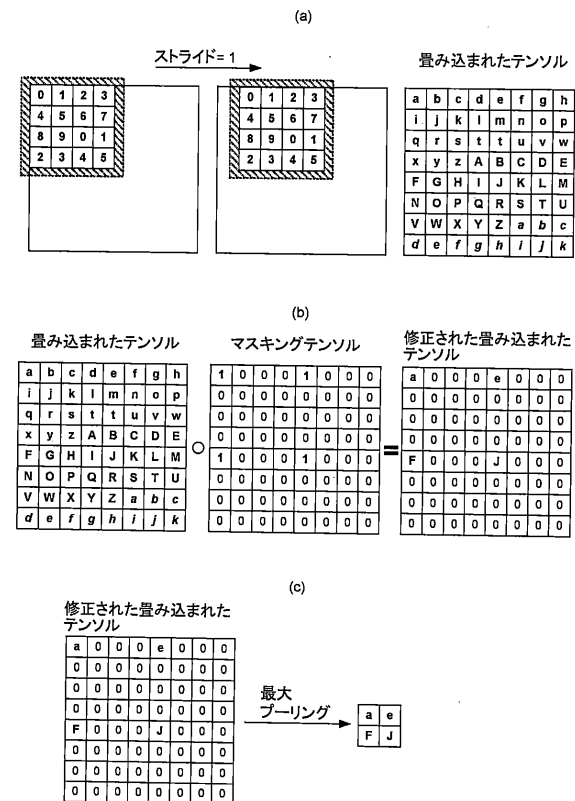


FIG. 9

【図 10】



フロントページの続き

(72)発明者 ガランド, ウィリアム・ジョン
アメリカ合衆国、9 4 0 4 3 カリフォルニア州、マウンテン・ビュー、アンフィシアター・パークウェイ、1 6 0 0

審査官 金田 孝之

(56)参考文献 特開2 0 1 6 - 1 5 3 9 8 4 (J P , A)
米国特許出願公開第2 0 1 5 / 0 2 7 8 6 4 2 (U S , A 1)
Chen Zhang et al. , Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks , Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays , 2015年02月22日 , pp.161-170
Clement Farabet et al. , CNP: An FPGA-based processor for Convolutional Networks , 2009 International Conference on Field Programmable Logic and Applications , 2009年08月31日 , pp.32-37

(58)調査した分野(Int.Cl. , D B 名)
G 0 6 N 3 / 0 0 - 9 9 / 0 0