

【公報種別】特許法第17条の2の規定による補正の掲載

【部門区分】第6部門第3区分

【発行日】令和4年8月8日(2022.8.8)

【公開番号】特開2022-70955(P2022-70955A)

【公開日】令和4年5月13日(2022.5.13)

【年通号数】公開公報(特許)2022-084

【出願番号】特願2022-19764(P2022-19764)

【国際特許分類】

G 0 6 N 3/063(2006.01)

10

【F I】

G 0 6 N 3/063

【手続補正書】

【提出日】令和4年7月28日(2022.7.28)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

20

【特許請求の範囲】

【請求項1】

ハードウェア集積回路上で実現されるニューラルネットワークを使用してニューラルネットワーク入力のバッチを処理するための方法であって、前記ニューラルネットワークは、有向グラフの状態で配置された複数のニューラルネットワークレイヤを備え、前記方法は、
前記ニューラルネットワーク入力のバッチを受信するステップと、
複数のスーパーレイヤを備えるシーケンスにパーティショニングされる前記ニューラルネットワークのレイヤを識別するステップとを備え、各スーパーレイヤは、2つ以上のニューラルネットワークレイヤを備え、前記有向グラフのパーティションであり、前記方法は、
さらに、

前記ハードウェア集積回路を使用して前記ニューラルネットワーク入力のバッチを処理するステップを備え、前記ハードウェア集積回路を使用して前記ニューラルネットワーク入力のバッチを処理するステップは、前記複数のスーパーレイヤの各スーパーレイヤについて、

前記バッチ内の前記ニューラルネットワーク入力の各々に対応するそれぞれのスーパーレイヤ入力を取得するステップと、

前記スーパーレイヤにおける前記ニューラルネットワークレイヤのいずれかを使用してスーパーレイヤ入力を処理する前に、前記スーパーレイヤにおける前記レイヤの各々のためのそれぞれのパラメータセットを前記ハードウェア集積回路のパラメータメモリにロードするステップと、

前記パラメータメモリから取得された前記ニューラルネットワークレイヤのための前記それぞれのパラメータセットを使用して、前記スーパーレイヤにおけるニューラルネットワークレイヤを介して前記スーパーレイヤ入力を処理するステップとを備える、方法。

【請求項2】

前記スーパーレイヤにおけるニューラルネットワークレイヤを介した前記スーパーレイヤ入力の前記処理に基づいてスーパーレイヤ出力を生成するステップをさらに備える、請求項1に記載の方法。

【請求項3】

前記スーパーレイヤ出力は、前記シーケンス内の第1のスーパーレイヤの出力であり、前

50

記方法はさらに、

前記シーケンス内の第2のスーパーレイヤにおけるニューラルネットワークレイヤへのスーパーレイヤ入力として前記スーパーレイヤ出力を受信するステップと、

前記第1のスーパーレイヤの前記スーパーレイヤ出力に対応する前記スーパーレイヤ入力について、前記シーケンス内の前記第2のスーパーレイヤにおけるニューラルネットワークレイヤを介して前記スーパーレイヤ入力を処理するステップとを備える、請求項2に記載の方法。

【請求項4】

前記スーパーレイヤにおける前記ニューラルネットワークレイヤの各々について前記それぞれのパラメータセットをロードするステップは、

前記ハードウェア集積回路および前記パラメータメモリの、外部にあるホストから受信されたデータ値に基づいて、前記それぞれのパラメータセットを前記パラメータメモリに予めロードするステップを備える、請求項1～3のいずれか1項に記載の方法。

10

【請求項5】

前記ハードウェア集積回路を使用して前記ニューラルネットワーク入力のバッチを処理するステップは、

スケジューリングプロセスに基づいて前記ニューラルネットワーク入力のバッチを処理するステップを備え、前記スケジューリングプロセスは、前記ハードウェア集積回路において実現されるニューラルネットワークモデルのバッチ次元およびレイヤ次元に対してニューラルネットワーク計算のグローバルスケジューリングを実行する、請求項4に記載の方法。

20

【請求項6】

前記ニューラルネットワーク計算の前記グローバルスケジューリングは、前記ホストを使用して実行される、請求項5に記載の方法。

【請求項7】

前記シーケンス内のスーパーレイヤの各ニューラルネットワークレイヤは、それぞれのワーキングセットに関連付けられ、

前記それぞれのワーキングセットは、部分的に、前記ワーキングセットにおけるスーパーレイヤ入力を処理するために使用される前記ニューラルネットワークレイヤのためのパラメータを格納するのに必要なメモリの量によって定義される、請求項1～6のいずれか1項に記載の方法。

30

【請求項8】

前記シーケンス内の第1のスーパーレイヤは、前記有向グラフの第1のパーティションを表し、

前記シーケンス内の第2のスーパーレイヤは、前記有向グラフの第2の異なるパーティションを表す、請求項1～7のいずれか1項に記載の方法。

【請求項9】

ハードウェア集積回路上で実現されるニューラルネットワークを使用してニューラルネットワーク入力のバッチを処理するためのシステムであって、前記ニューラルネットワークは、有向グラフの状態で配置された複数のニューラルネットワークレイヤを備え、前記システムは、

40

前記ハードウェア集積回路と、プロセッサと、命令を格納するための非一時的なコンピュータ読み取り可能記憶装置とを備え、前記命令は、動作の実行をさせるように前記プロセッサによって実行可能であり、前記動作は、

前記ニューラルネットワーク入力のバッチを受信するステップと、

複数のスーパーレイヤを備えるシーケンスにパーティショニングされる前記ニューラルネットワークのレイヤを識別するステップとを備え、各スーパーレイヤは、2つ以上のニューラルネットワークレイヤを備え、前記有向グラフのパーティションであり、前記動作はさらに、

前記ハードウェア集積回路を使用して前記ニューラルネットワーク入力のバッチを処理す

50

るステップを備え、前記ハードウェア集積回路を使用して前記ニューラルネットワーク入力のバッチを処理するステップは、前記複数のスーパーレイヤの各スーパーレイヤについて、

前記バッチ内の前記ニューラルネットワーク入力の各々に対応するそれぞれのスーパーレイヤ入力を取得するステップと、

前記スーパーレイヤにおける前記ニューラルネットワークレイヤのいずれかを使用してスーパーレイヤ入力を処理する前に、前記スーパーレイヤにおける前記レイヤの各々のためのそれぞれのパラメータセットを前記ハードウェア集積回路のパラメータメモリにロードするステップと、

前記パラメータメモリから取得された前記ニューラルネットワークレイヤのための前記それぞれのパラメータセットを使用して、前記スーパーレイヤにおけるニューラルネットワークレイヤを介して前記スーパーレイヤ入力を処理するステップとを備える、システム。
10

【請求項 10】

前記動作は、前記スーパーレイヤにおけるニューラルネットワークレイヤを介した前記スーパーレイヤ入力の前記処理に基づいてスーパーレイヤ出力を生成するステップを備える、請求項 9 に記載のシステム。

【請求項 11】

前記スーパーレイヤ出力は、前記シーケンス内の第 1 のスーパーレイヤの出力であり、前記動作はさらに、

前記シーケンス内の第 2 のスーパーレイヤにおけるニューラルネットワークレイヤへのスーパーレイヤ入力として前記スーパーレイヤ出力を受信するステップと、
20

前記第 1 のスーパーレイヤの前記スーパーレイヤ出力に対応する前記スーパーレイヤ入力について、前記シーケンス内の前記第 2 のスーパーレイヤにおけるニューラルネットワークレイヤを介して前記スーパーレイヤ入力を処理するステップとを備える、請求項 10 に記載のシステム。

【請求項 12】

前記スーパーレイヤにおける前記ニューラルネットワークレイヤの各々について前記それぞれのパラメータセットをロードするステップは、

前記ハードウェア集積回路および前記パラメータメモリの、外部にあるホストから受信されたデータ値に基づいて、前記それぞれのパラメータセットを前記パラメータメモリに予めロードするステップを備える、請求項 9 ~ 11 のいずれか 1 項に記載のシステム。
30

【請求項 13】

前記ハードウェア集積回路を使用して前記ニューラルネットワーク入力のバッチを処理するステップは、

スケジューリングプロセスに基づいて前記ニューラルネットワーク入力のバッチを処理するステップを備え、前記スケジューリングプロセスは、前記ハードウェア集積回路において実現されるニューラルネットワークモデルのバッチ次元およびレイヤ次元に対してニューラルネットワーク計算のグローバルスケジューリングを実行する、請求項 12 に記載のシステム。

【請求項 14】

前記ニューラルネットワーク計算の前記グローバルスケジューリングは、前記ホストを使用して実現される、請求項 13 に記載のシステム。

【請求項 15】

前記シーケンス内のスーパーレイヤの各ニューラルネットワークレイヤは、それぞれのワーキングセットに関連付けられ、

前記それぞれのワーキングセットは、部分的に、前記ワーキングセットにおけるスーパーレイヤ入力を処理するために使用される前記ニューラルネットワークレイヤのためのパラメータを格納するのに必要なメモリの量によって定義される、請求項 9 ~ 14 のいずれか 1 項に記載のシステム。
40

【請求項 16】

50

前記シーケンス内の第1のスーパーレイヤは、前記有向グラフの第1のパーティションを表し、

前記シーケンス内の第2のスーパーレイヤは、前記有向グラフの第2の異なるパーティションを表す、請求項9～15のいずれか1項に記載のシステム。

【請求項17】

ハードウェア集積回路上で実現されるニューラルネットワークを使用してニューラルネットワーク入力のバッチを処理するための命令を格納するように構成された非一時的なコンピュータ読取可能記憶装置であつて、

前記ニューラルネットワークは、有向グラフの状態で配置された複数のニューラルネットワークレイヤを備え、前記命令は、動作の実行をさせるようにプロセッサによって実行可能であり、前記動作は、

前記ニューラルネットワーク入力のバッチを受信するステップと、

複数のスーパーレイヤを備えるシーケンスにパーティショニングされる前記ニューラルネットワークのレイヤを識別するステップとを備え、各スーパーレイヤは、2つ以上のニューラルネットワークレイヤを備え、前記有向グラフのパーティションであり、前記動作はさらに、

前記ハードウェア集積回路を使用して前記ニューラルネットワーク入力のバッチを処理するステップを備え、前記ハードウェア集積回路を使用して前記ニューラルネットワーク入力のバッチを処理するステップは、前記複数のスーパーレイヤの各スーパーレイヤについて、

前記バッチ内の前記ニューラルネットワーク入力の各々に対応するそれぞれのスーパーレイヤ入力を取得するステップと、

前記スーパーレイヤにおける前記ニューラルネットワークレイヤのいずれかを使用してスーパーレイヤ入力を処理する前に、前記スーパーレイヤにおける前記レイヤの各々のためのそれぞれのパラメータセットを前記ハードウェア集積回路のパラメータメモリにロードするステップと、

前記パラメータメモリから取得された前記ニューラルネットワークレイヤのための前記それぞれのパラメータセットを使用して、前記スーパーレイヤにおけるニューラルネットワークレイヤを介して前記スーパーレイヤ入力を処理するステップとを備える、非一時的なコンピュータ読取可能記憶装置。

【請求項18】

前記動作は、前記スーパーレイヤにおけるニューラルネットワークレイヤを介した前記スーパーレイヤ入力の前記処理に基づいてスーパーレイヤ出力を生成するステップをさらに備える、請求項17に記載の非一時的なコンピュータ読取可能記憶装置。

【請求項19】

前記スーパーレイヤ出力は、前記シーケンス内の第1のスーパーレイヤの出力であり、前記動作はさらに、

前記シーケンス内の第2のスーパーレイヤにおけるニューラルネットワークレイヤへのスーパーレイヤ入力として前記スーパーレイヤ出力を受信するステップと、

前記第1のスーパーレイヤの前記スーパーレイヤ出力に対応する前記スーパーレイヤ入力について、前記シーケンス内の前記第2のスーパーレイヤにおけるニューラルネットワークレイヤを介して前記スーパーレイヤ入力を処理するステップとを備える、請求項18に記載の非一時的なコンピュータ読取可能記憶装置。

【請求項20】

前記スーパーレイヤにおける前記ニューラルネットワークレイヤの各々について前記それぞれのパラメータセットをロードするステップは、

前記ハードウェア集積回路および前記パラメータメモリの外部にあるホストから受信されたデータ値に基づいて、前記それぞれのパラメータセットを前記パラメータメモリに予めロードするステップを備える、請求項17～19のいずれか1項に記載の非一時的なコンピュータ読取可能記憶装置。

10

20

30

40

50