

【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2年8月27日(2020.8.27)

【公表番号】特表2020-521195(P2020-521195A)

【公表日】令和2年7月16日(2020.7.16)

【年通号数】公開・登録公報2020-028

【出願番号】特願2019-552217(P2019-552217)

【国際特許分類】

G 06 N 3/063 (2006.01)

【F I】

G 06 N 3/063

【手続補正書】

【提出日】令和2年6月23日(2020.6.23)

【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】

【特許請求の範囲】

【請求項1】

方法であって、

ハードウェア回路上のニューラルネットワークを使用して処理される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を備え、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レイヤを有し、各レイヤは、それぞれのパラメータセットを有し、前記方法はさらに、

前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへのパーティショニングを決定するステップを備え、各スーパーレイヤは、1つ以上のレイヤを含む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、前記シーケンスにおける各スーパーレイヤについて、

前記スーパーレイヤにおける前記レイヤのための前記それぞれのパラメータセットを前記ハードウェア回路のメモリにロードするステップと、

前記バッチ内の各ニューラルネットワーク入力について、前記ハードウェア回路の前記メモリ内の前記パラメータセットを使用して、前記スーパーレイヤの各レイヤを介して前記ニューラルネットワーク入力に対応するスーパーレイヤ入力を処理して、前記ニューラルネットワーク入力のためのスーパーレイヤ出力を生成するステップとを備える、方法。

【請求項2】

前記シーケンスにおける第1のスーパーレイヤでは、前記ニューラルネットワーク入力に対応する前記スーパーレイヤ入力が前記ニューラルネットワーク入力である、請求項1に記載の方法。

【請求項3】

前記第1のスーパーレイヤ出力の後の各スーパーレイヤへの前記スーパーレイヤ入力は、前記シーケンスにおける先行するスーパーレイヤによって生成されたスーパーレイヤ出力である、請求項2に記載の方法。

【請求項4】

前記ハードウェア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、各スーパーレイヤについて、

前記バッチ内の第2のニューラルネットワーク入力に対応するスーパーレイヤ入力が前記スーパーレイヤの各レイヤを介して後に処理される前に前記バッチ内の第1のニューラルネットワーク入力のための前記スーパーレイヤ入力が前記スーパーレイヤの各レイヤを介して処理されるように、前記スーパーレイヤの各レイヤを介して前記ニューラルネットワーク入力のバッチに対応する前記スーパーレイヤ入力をシーケンシャルに処理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。

【請求項5】

スーパーレイヤのそれぞれのレイヤは、ワーキングセットに関連付けられ、各ワーキングセットは、少なくとも

i) 前記ハードウェア回路上の前記ニューラルネットワークを使用して処理される前記ニューラルネットワーク入力のバッチの1つ以上の入力または前記スーパーレイヤの先行するレイヤの1つ以上の出力、および

ii) 前記スーパーレイヤの各レイヤを介して前記1つ以上の入力を処理するのに必要なメモリの量を示すサイズパラメータ

によって定義される、請求項1～4のいずれか1項に記載の方法。

【請求項6】

前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへの前記パーティショニングを決定するステップは、

i) 少なくとも1つのワーキングセットのための特定のサイズパラメータを決定するステップと、

ii) 前記ハードウェア回路の前記メモリの特定の集約パラメータ容量を決定するステップと、

iii) 前記少なくとも1つのワーキングセットのための前記特定のサイズパラメータまたは前記ハードウェア回路の前記メモリの特定の集約パラメータ容量のうちの少なくとも1つに基づいて、前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへの前記パーティショニングを決定するステップとを備える、請求項5に記載の方法。

【請求項7】

前記ハードウェア回路の前記メモリは、閾値記憶容量を有し、前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへの前記パーティショニングを決定するステップは、

前記ハードウェア回路の前記メモリの前記閾値記憶容量に基づいて、前記ニューラルネットワークのレイヤをスーパーレイヤのシーケンスにパーティショニングするステップを備える、請求項1～6のいずれか1項に記載の方法。

【請求項8】

前記ニューラルネットワークのレイヤは、前記ハードウェア回路が前記ニューラルネットワーク入力のバッチを処理する際に前記メモリの前記閾値記憶容量を超えないようにスーパーレイヤのシーケンスにパーティショニングされる、請求項7に記載の方法。

【請求項9】

前記ニューラルネットワーク入力のバッチおよび前記それぞれのパラメータセットは、前記ハードウェア回路の外部のソースから受信され、前記スーパーレイヤの各レイヤを介して前記ニューラルネットワーク入力に対応する前記スーパーレイヤ入力を処理するステップは、前記外部のソースから追加のパラメータを受信することなく前記スーパーレイヤ入力を処理するステップを備える、請求項1～8のいずれか1項に記載の方法。

【請求項10】

コンピューティングシステムであって、

前記コンピューティングシステムに配設されたハードウェア回路を備え、前記ハードウェア回路は、1つ以上の処理装置を含み、前記コンピューティングシステムはさらに、

動作を実行するように前記1つ以上の処理装置によって実行可能な命令を格納するための1つ以上の機械読み取り可能記憶装置を備え、前記動作は、

ハードウェア回路上のニューラルネットワークを使用して処理される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を備え、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レイヤを有し、各レイヤは、それぞれのパラメータセットを有し、前記動作はさらに、

前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへのパーティショニングを決定するステップを備え、各スーパーレイヤは、1つ以上のレイヤを含む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、前記シーケンスにおける各スーパーレイヤについて、

前記スーパーレイヤにおける前記レイヤのための前記それぞれのパラメータセットを前記ハードウェア回路のメモリにロードするステップと、

前記バッチ内の各ニューラルネットワーク入力について、前記ハードウェア回路の前記メモリ内の前記パラメータを使用して、前記スーパーレイヤの各レイヤを介して前記ニューラルネットワーク入力に対応するスーパーレイヤ入力を処理して、前記ニューラルネットワーク入力のためのスーパーレイヤ出力を生成するステップとを備える、コンピューティングシステム。

【請求項 1 1】

前記シーケンスにおける第1のスーパーレイヤでは、前記ニューラルネットワーク入力に対応する前記スーパーレイヤ入力が前記ニューラルネットワーク入力である、請求項10に記載のコンピューティングシステム。

【請求項 1 2】

前記第1のスーパーレイヤ出力の後の各スーパーレイヤへの前記スーパーレイヤ入力は、前記シーケンスにおける先行するスーパーレイヤによって生成されたスーパーレイヤ出力である、請求項11に記載のコンピューティングシステム。

【請求項 1 3】

前記ハードウェア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、各スーパーレイヤについて、

前記バッチ内の第2のニューラルネットワーク入力に対応するスーパーレイヤ入力が前記スーパーレイヤの各レイヤを介して後に処理される前に前記バッチ内の第1のニューラルネットワーク入力のための前記スーパーレイヤ入力が前記スーパーレイヤの各レイヤを介して処理されるように、前記スーパーレイヤの各レイヤを介して前記ニューラルネットワーク入力のバッチに対応する前記スーパーレイヤ入力をシーケンシャルに処理するステップを備える、請求項10～12のいずれか1項に記載のコンピューティングシステム。

【請求項 1 4】

スーパーレイヤのそれぞれのレイヤは、ワーキングセットに関連付けられ、各ワーキングセットは、少なくとも

i) 前記ハードウェア回路上の前記ニューラルネットワークを使用して処理される前記ニューラルネットワーク入力のバッチの1つ以上の入力または前記スーパーレイヤの先行するレイヤの1つ以上の出力、および

ii) 前記スーパーレイヤの各レイヤを介して前記1つ以上の入力を処理するのに必要なメモリの量を示すサイズパラメータ

によって定義される、請求項10～13のいずれか1項に記載のコンピューティングシステム。

【請求項 1 5】

前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへの前記パーティショニングを決定するステップは、

i) 少なくとも1つのワーキングセットのための特定のサイズパラメータを決定するステップと、

ii) 前記ハードウェア回路の前記メモリの特定の集約パラメータ容量を決定するステ

ップと、

i i i) 前記少なくとも 1 つのワーキングセットのための前記特定のサイズパラメータまたは前記ハードウェア回路の前記メモリの特定の集約パラメータ容量のうちの少なくとも 1 つに基づいて、前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへの前記パーティショニングを決定するステップとを備える、請求項 14 に記載のコンピューティングシステム。

【請求項 16】

前記ハードウェア回路の前記メモリは、閾値記憶容量を有し、前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへの前記パーティショニングを決定するステップは、

前記ハードウェア回路の前記メモリの前記閾値記憶容量に基づいて、前記ニューラルネットワークのレイヤをスーパーレイヤのシーケンスにパーティショニングするステップを備える、請求項 10 ~ 15 のいずれか 1 項に記載のコンピューティングシステム。

【請求項 17】

前記ニューラルネットワークのレイヤは、前記ハードウェア回路が前記ニューラルネットワーク入力のバッチを処理する際に前記メモリの前記閾値記憶容量を超えないようにスーパーレイヤのシーケンスにパーティショニングされる、請求項 16 に記載のコンピューティングシステム。

【請求項 18】

前記ニューラルネットワーク入力のバッチおよび前記それぞれのパラメータセットは、前記ハードウェア回路の外部のソースから受信され、前記スーパーレイヤの各レイヤを介して前記ニューラルネットワーク入力に対応する前記スーパーレイヤ入力を処理するステップは、前記外部のソースから追加のパラメータを受信することなく前記スーパーレイヤ入力を処理するステップを備える、請求項 10 ~ 17 のいずれか 1 項に記載のコンピューティングシステム。

【請求項 19】

動作を実行するように 1 つ以上の処理装置に命令を実行させるプログラムであって、前記動作は、

ハードウェア回路上のニューラルネットワークを使用して処理される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を備え、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レイヤを有し、各レイヤは、それぞれのパラメータセットを有し、前記動作はさらに、

前記ニューラルネットワークのレイヤの、スーパーレイヤのシーケンスへのパーティショニングを決定するステップを備え、各スーパーレイヤは、1 つ以上のレイヤを含む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、前記シーケンスにおける各スーパーレイヤについて、

前記スーパーレイヤにおける前記レイヤのための前記それぞれのパラメータセットを前記ハードウェア回路のメモリにロードするステップと、

前記バッチ内の各ニューラルネットワーク入力について、前記ハードウェア回路の前記メモリ内の前記パラメータセットを使用して、前記スーパーレイヤの各レイヤを介して前記ニューラルネットワーク入力に対応するスーパーレイヤ入力を処理して、前記ニューラルネットワーク入力のためのスーパーレイヤ出力を生成するステップとを備える、プログラム。

【請求項 20】

前記ハードウェア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、各スーパーレイヤについて、

前記バッチ内の第 2 のニューラルネットワーク入力に対応するスーパーレイヤ入力が前記スーパーレイヤの各レイヤを介して後に処理される前に前記バッチ内の第 1 のニューラ

ルネットワーク入力のための前記スーパーレイヤ入力が前記スーパーレイヤの各レイヤを介して処理されるように、前記スーパーレイヤの各レイヤを介して前記ニューラルネットワーク入力のバッチに対応する前記スーパーレイヤ入力をシーケンシャルに処理するステップを備える、請求項 1 9 に記載のプログラム。