

【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7月13日(2022.7.13)

【国際公開番号】WO2020/190809

【公表番号】特表2022-523761(P2022-523761A)

【公表日】令和4年4月26日(2022.4.26)

【年通号数】公開公報(特許)2022-075

【出願番号】特願2021-547450(P2021-547450)

【国際特許分類】

G 06 F 17/16(2006.01)

G 06 T 1/20(2006.01)

G 06 F 9/38(2006.01)

G 06 F 15/80(2006.01)

G 06 F 17/10(2006.01)

10

【F I】

G 06 F 17/16 P

G 06 T 1/20 A

G 06 F 9/38 3 7 0 C

20

G 06 F 15/80

G 06 F 17/10 A

【手続補正書】

【提出日】令和4年7月5日(2022.7.5)

【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】

30

【請求項1】

キャッシュメモリと結合される複数の処理リソースを含み、少なくとも1つの処理リソースが行列アクセラレータを含み、該行列アクセラレータは、スペース内積命令に応答してスペース第1行列及び第2行列の複数の要素に対して内積演算を実行するよう構成され、前記スペース第1行列の要素は、要素の組を含む圧縮表現に圧縮され、前記要素の組は、少なくとも1つの非ゼロ値要素及び該少なくとも1つの非ゼロ値要素の指示を含む、計算クラスタを有し、

前記圧縮表現は、圧縮された形式で前記キャッシュメモリに格納され、

前記少なくとも1つの処理リソースは、

前記圧縮表現を前記キャッシュメモリから前記少なくとも1つの処理リソース内のメモリにロードし、

前記第2行列を前記キャッシュメモリから前記少なくとも1つの処理リソース内の前記メモリにロードし、

前記圧縮表現からの要素及び前記第2行列の選択された要素に対して前記内積演算を実行し、前記第2行列の前記選択された要素が、前記圧縮表現内に格納された前記スペース第1行列の非ゼロ値と対応し、前記少なくとも1つの非ゼロ値の前記指示に基づき選択され、

前記内積演算の出力を前記少なくとも1つの処理リソース内の前記メモリに書き込むよう構成される、

汎用グラフィックスプロセッサ。

50

【請求項 2】

前記キャッシュメモリは、レベル2(L2)キャッシュメモリである。
請求項1に記載の汎用グラフィクスプロセッサ。

【請求項 3】

前記少なくとも1つの処理リソース内の前記メモリは、レベル1(L1)キャッシュメモリを含む。

請求項1に記載の汎用グラフィクスプロセッサ。

【請求項 4】

前記少なくとも1つの処理リソース内の前記メモリは、共有メモリを含む。
請求項1に記載の汎用グラフィクスプロセッサ。

10

【請求項 5】

前記少なくとも1つの処理リソース内の前記メモリは、レジスタファイルを含む。
請求項1に記載の汎用グラフィクスプロセッサ。

【請求項 6】

前記少なくとも1つの処理リソース内の前記メモリは、前記行列アクセラレータ内のメモリを含む。

請求項1に記載の汎用グラフィクスプロセッサ。

【請求項 7】

前記スパース第1行列は、ニューラルネットワークに関連した重みデータを含む。
請求項1に記載の汎用グラフィクスプロセッサ。

20

【請求項 8】

前記第2行列は、前記ニューラルネットワークに関連した入力活性化データを含む。
請求項7に記載の汎用グラフィクスプロセッサ。

【請求項 9】

前記内積演算の前記出力は、前記ニューラルネットワークに関連した出力活性化データを含む。

請求項8に記載の汎用グラフィクスプロセッサ。

【請求項 10】

前記内積演算の前記出力は、密行列である。
請求項9に記載の汎用グラフィクスプロセッサ。

30

【請求項 11】

前記行列アクセラレータは、処理要素のストリックアレイを含む。
請求項1に記載の汎用グラフィクスプロセッサ。

【請求項 12】

前記スパース第1行列は、構造化されたスパース性を有する。
請求項1に記載の汎用グラフィクスプロセッサ。

【請求項 13】

前記スパース第1行列の要素は、前記構造化されたスパース性に基づき圧縮表現に圧縮される。

請求項12に記載の汎用グラフィクスプロセッサ。

40

【請求項 14】

前記内積演算は、8ビット整数内積演算である。
請求項1に記載の汎用グラフィクスプロセッサ。

【請求項 15】

前記スパース第1行列は、8ビット整数要素を含む。
請求項14に記載の汎用グラフィクスプロセッサ。

【請求項 16】

メモリデバイスと、

請求項1乃至15のうち何れか一項に記載の汎用グラフィクスプロセッサと
を有するデータ処理システム。

50

【請求項 17】

スペース内積命令に応答してスペース第1行列及び第2行列の複数の要素に対して内積演算を実行することであり、前記内積演算は、キャッシュメモリと結合される複数の処理リソースを含む計算クラスタにより実行され、少なくとも1つの処理リソースは行列アクセラレータを含み、前記スペース第1行列の要素は、要素の組を含む圧縮表現に圧縮され、前記要素の組は、少なくとも1つの非ゼロ値要素及び該少なくとも1つの非ゼロ値要素の指示を含む、前記実行することと、

前記圧縮表現を、圧縮された形式で前記キャッシュメモリに格納することと、

前記少なくとも1つの処理リソースにより、

前記圧縮表現を前記キャッシュメモリから前記少なくとも1つの処理リソース内のメモリにロードし、10

前記第2行列を前記キャッシュメモリから前記少なくとも1つの処理リソース内の前記メモリにロードし、

前記圧縮表現からの要素及び前記第2行列の選択された要素に対して前記内積演算を実行し、前記第2行列の前記選択された要素が、前記圧縮表現内に格納された前記スペース第1行列の非ゼロ値と対応し、前記少なくとも1つの非ゼロ値の前記指示に基づき選択され、

前記内積演算の出力を前記少なくとも1つの処理リソース内の前記メモリに書き込むことと

を有する方法。20

【請求項 18】

前記スペース第1行列の要素を、前記少なくとも1つの処理リソースのメモリ内で、前記圧縮表現に圧縮することを更に有する、

請求項17に記載の方法。

【請求項 19】

前記スペース第1行列は、構造化されたスペース性を有し、

前記スペース第1行列の要素は、前記構造化されたスペース性に基づき圧縮表現に圧縮される、

請求項17に記載の方法。

【請求項 20】

内積演算を実行することは、8ビット整数内積演算を実行することを含む、30

請求項17に記載の方法。

【請求項 21】

前記スペース第1行列は、8ビット整数要素を含む、

請求項20に記載の方法。

【請求項 22】

実行される場合にマシンに請求項17乃至21のうち何れか一項に記載の方法を実行させるコンピュータプログラム。

【請求項 23】

請求項22に記載のコンピュータプログラムを記憶しているマシン可読記憶媒体。40

【請求項 24】

請求項17乃至21のうち何れか一項に記載の方法を実行する手段を有する装置。