

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2019-133662

(P2019-133662A)

(43) 公開日 令和1年8月8日(2019.8.8)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06	302E
G06F 13/10 (2006.01)	G06F 13/10	340A
G06F 13/12 (2006.01)	G06F 13/12	340B
	G06F 3/06	301G

審査請求 未請求 請求項の数 20 O L (全 15 頁)

(21) 出願番号 特願2019-13888 (P2019-13888)
 (22) 出願日 平成31年1月30日 (2019. 1. 30)
 (31) 優先権主張番号 62/625, 532
 (32) 優先日 平成30年2月2日 (2018. 2. 2)
 (33) 優先権主張国・地域又は機関
 米国 (US)
 (31) 優先権主張番号 15/942, 218
 (32) 優先日 平成30年3月30日 (2018. 3. 30)
 (33) 優先権主張国・地域又は機関
 米国 (US)

(71) 出願人 390019839
 三星電子株式会社
 Samsung Electronics
 Co., Ltd.
 大韓民国京畿道水原市靈通区三星路129
 129, Samsung-ro, Yeon
 gtong-gu, Suwon-si, G
 yeonggi-do, Republic
 of Korea
 (74) 代理人 110000051
 特許業務法人共生国際特許事務所
 (72) 発明者 李 周 桓
 アメリカ合衆国, 95134, カリフ
 オルニア州, サン ノゼ, スカイトッ
 プ ストリート #217, 52
 最終頁に続く

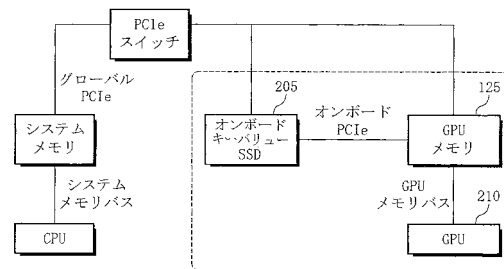
(54) 【発明の名称】 キーバリュアアクセスを含むマシンラーニングのためのシステム及び方法

(57) 【要約】

【課題】 シリアライゼーションされたキーバリュアアクセスを防止するマシンラーニングシステムのためのシステム及び方法を提供する。

【解決手段】 本発明のマシンラーニングのためのシステム及び方法において、システムは、GPUメモリを有するGPU及びGPUメモリに連結されたキーバリュアストレージデバイスを含む。方法は、GPUにより、キーを含むキーバリュア要請をGPUメモリの入出力領域内のキーバリュア要請キューに書き込むステップと、キーバリュアのストレージデバイスにより、キーバリュア要請をキーバリュア要請キューから読み取るステップと、キーバリュアストレージデバイスにより、キーバリュア要請にตอบสนองして、キーバリュア要請のキーに対応するバリュアをGPUメモリの入出力領域に書き込むステップと、を有する。

【選択図】 図2



【特許請求の範囲】

【請求項 1】

第 1 グラフィックプロセッシングユニットにより、キーを含む第 1 キーバリューストレージデバイスに接続された第 1 メモリの第 1 入出力領域内のキーバリューストレージデバイスに書き込むステップと、

前記第 1 メモリに接続された第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスの第 1 入出力領域内のキーバリューストレージデバイスから読み取るステップと、

前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスに書き込むステップと、を有することを特徴とするマシンラーニングのための方法。

10

【請求項 2】

前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスで、キーバリューストレージデバイスをリトリートするステップを更に含むことを特徴とする請求項 1 に記載のマシンラーニングのための方法。

【請求項 3】

前記第 1 キーバリューストレージデバイスによる書き込み領域は、前記第 1 キーバリューストレージデバイスに割り当てられた領域であるリターンバリューストレージ領域を含むことを特徴とする請求項 1 に記載のマシンラーニングのための方法。

20

【請求項 4】

前記第 1 キーバリューストレージデバイスを前記第 1 メモリの第 1 入出力領域に書き込むステップは、前記第 1 キーバリューストレージデバイスを前記リターンバリューストレージ領域に書き込むステップを含むことを特徴とする請求項 3 に記載のマシンラーニングのための方法。

【請求項 5】

前記第 1 キーバリューストレージデバイスを前記第 1 メモリの第 1 入出力領域に書き込むステップは、前記第 1 キーバリューストレージデバイスを前記第 1 メモリの第 1 入出力領域内のリターンバリューストレージ領域に書き込むステップを含むことを特徴とする請求項 1 に記載のマシンラーニングのための方法。

【請求項 6】

前記第 1 キーバリューストレージデバイス及び前記第 1 グラフィックプロセッシングユニットに接続されたホストによって、

前記第 1 キーバリューストレージデバイスが、前記第 1 メモリの第 1 入出力領域にアクセスしてキーバリューストレージデバイスを受信し、前記第 1 キーバリューストレージデバイスに書き込むステップと、

前記第 1 グラフィックプロセッシングユニットが、前記第 1 メモリの第 1 入出力領域内のキーバリューストレージデバイスを格納し、前記第 1 メモリの第 1 入出力領域からキーバリューストレージデバイスを読み取るように構成するステップと、を更に含むことを特徴とする請求項 1 に記載のマシンラーニングのための方法。

30

【請求項 7】

前記ホストに接続された第 2 グラフィックプロセッシングユニットにより、キーを含む第 2 キーバリューストレージデバイスを前記第 2 グラフィックプロセッシングユニットに接続された第 2 メモリの入出力領域内のキーバリューストレージデバイスに書き込むステップと、

前記ホスト及び前記第 2 メモリに接続された第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスの第 2 入出力領域内のキーバリューストレージデバイスから読み取るステップと、

前記第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスに書き込むステップと、を更に含むことを特徴とする請求項 6 に記載のマシンラーニングのための方法。

40

【請求項 8】

50

前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスで、キールックアップを遂行して前記第 1 バリューストレージデバイスをリトリブするステップと、前記第 1 キーバリューストレージデバイスにより、前記キールックアップを遂行するステップと同時に、前記第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスで、キールックアップを遂行して前記第 2 バリューストレージデバイスをリトリブするステップと、を更に含むことを特徴とする請求項 7 に記載のマシンラーニングのための方法。

【請求項 9】

前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスから読み取るステップは、P2P (peer-to-peer) DMA (direct memory access) を介して前記第 1 キーバリューストレージデバイスから読み取るステップを含むことを特徴とする請求項 1 に記載のマシンラーニングのための方法。

10

【請求項 10】

前記第 1 キーバリューストレージデバイスにより、前記第 1 バリューストレージデバイスに書き込むステップは、P2P (peer-to-peer) DMA (direct memory access) を介して前記第 1 バリューストレージデバイスに書き込むステップを含むことを特徴とする請求項 1 に記載のマシンラーニングのための方法。

【請求項 11】

前記第 1 キーバリューストレージデバイスは、PCI (peripheral component interconnect) 連結により、前記第 1 グラフィックプロセッシングユニットに連結されることを特徴とする請求項 10 に記載のマシンラーニングのための方法。

20

【請求項 12】

前記第 1 グラフィックプロセッシングユニットにより、前記第 1 キーバリューストレージデバイスから読み取るステップは、前記第 1 メモリの第 1 入出力領域内のキーバリューストレージデバイスに書き込むステップの後に、且つ前記第 1 キーバリューストレージデバイスにより、前記第 1 バリューストレージデバイスを前記第 1 メモリの第 1 入出力領域に書き込むステップの前に、前記第 1 グラフィックプロセッシングユニットにより、第 2 キーバリューストレージデバイスから読み取るステップを更に含むことを特徴とする請求項 1 に記載のマシンラーニングのための方法。

30

【請求項 13】

前記第 1 グラフィックプロセッシングユニットにより、キーを含む第 2 キーバリューストレージデバイスから読み取るステップは、前記第 1 メモリの第 2 入出力領域内のキーバリューストレージデバイスに書き込むステップと、

前記第 1 メモリに連結された第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスから読み取るステップと、

前記第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスから読み取るステップに、前記第 2 キーバリューストレージデバイスのキーに対応する第 2 バリューストレージデバイスを前記第 1 メモリの第 2 入出力領域に書き込むステップと、を更に含むことを特徴とする請求項 1 に記載のマシンラーニングのための方法。

40

【請求項 14】

前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスで、キールックアップを遂行して前記第 1 バリューストレージデバイスをリトリブするステップと、

前記第 1 キーバリューストレージデバイスにより、前記キールックアップを遂行するステップと同時に、前記第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスで、キールックアップを遂行して前記第 2 バリューストレージデバイスをリトリブするステップと、を更に含むことを特徴とする請求項 13 に記載のマシンラーニングのための方法。

【請求項 15】

50

グラフィックプロセッシングユニットと、
前記グラフィックプロセッシングユニットに連結されたメモリと、
キーバリューストレージデバイスと、を備え、
前記キーバリューストレージデバイスは、PCI (peripheral component interconnect) 連結により前記グラフィックプロセッシングユニットに連結され、

前記グラフィックプロセッシングユニットは、前記メモリの入出力領域内のメモリがマッピングされた入力及び出力動作を遂行し、1つ以上のキーバリューストレージデバイスを前記メモリの入出力領域内のキーバリューストレージデバイスに書き込み、

前記キーバリューストレージデバイスは、
前記メモリの入出力領域で、メモリがマッピングされた入力及び出力動作を遂行し、
前記1つ以上のキーバリューストレージデバイスを前記キーバリューストレージデバイスから読み取り、
前記1つ以上のキーバリューストレージデバイスの中のキーバリューストレージデバイスに
前記キーバリューストレージデバイスのキーに対応するバリューストレージデバイスを前記メモリの入出力領域に書き込むことを特徴とする
マシンラーニングのためのシステム。

【請求項16】

前記キーバリューストレージデバイスによる書き込み領域は、前記バリューストレージデバイスに割り当てられた領域である
リターンバリューストレージ領域を含むことを特徴とする請求項15に記載のマシンラーニングのためのシステム。

【請求項17】

前記バリューストレージデバイスを前記メモリの入出力領域に書き込むことは、前記バリューストレージデバイスを前記リターンバリューストレージ領域に書き込むことを含むことを特徴とする請求項16に記載のマシンラーニングのためのシステム。

【請求項18】

前記バリューストレージデバイスを前記メモリの入出力領域に書き込むことは、前記バリューストレージデバイスを前記メモリの入出力領域内のリターンバリューストレージキューに書き込むことを含むことを特徴とする請求項15に記載のマシンラーニングのためのシステム。

【請求項19】

グラフィックプロセッシングユニットと、
キーバリューストレージデバイスと、
前記グラフィックプロセッシングユニットと前記キーバリューストレージデバイスとの間の通信のために共有されたメモリ手段と、を備え、

前記グラフィックプロセッシングユニットは、前記通信のために共有されたメモリ手段を介して1つ以上のキーバリューストレージデバイスを前記キーバリューストレージデバイスに伝送し、

前記キーバリューストレージデバイスは、
前記1つ以上のキーバリューストレージデバイスを受信し、
前記1つ以上のキーバリューストレージデバイスの中のキーバリューストレージデバイスに
前記キーバリューストレージデバイスのキーに対応するバリューストレージデバイスを前記通信のために共有されたメモリ手段を介して、前記グラフィックプロセッシングユニットに伝送することを特徴とするマシンラーニングのためのシステム。

【請求項20】

前記通信のために共有されたメモリ手段は、前記グラフィックプロセッシングユニットに連結されたメモリを含み、PCI (peripheral component interconnect) 連結を通じたP2P (peer-to-peer) DMA (direct memory access) を介して、前記キーバリューストレージデバイスによってアクセスされることを特徴とする請求項19に記載のマシンラーニングのためのシステム。

【発明の詳細な説明】

【技術分野】

【0001】

10

20

30

40

50

本発明は、マシンラーニング (machine learning) に関し、より詳細には、シリアライゼーション (serialization) されたキーバリューストア (key value access) を防止するマシンラーニングシステムのためのシステム及び方法に関する。

【背景技術】

【0002】

ブロックインターフェース (block interface) を有するいくつかの従来技術の SSD (solid state drive) において、SSD に格納されたデータに対するキーバリューストア (key value access) は、CPU (central processing unit) を含めて全体のトレーニングデータのサブセットをランダムに (randomly) サンプルする確率的 (stochastic) マシンラーニング中にキーバリューストアを提供することを要求する。ホスト CPU は、ファイルインデックスルックアップ (file index lookup) 及びファイルシステムアクセス (file system access) を遂行してデータの位置を識別するが、これはシリアライゼーションされたキーバリューストアをもたず。このようなシリアライゼーションされたキーバリューストアは、パフォーマンスを制限する。

10

【0003】

従って、データに対するキーバリューストアを含むマシンラーニングを遂行するための改善されたシステム及び方法が必要である。

20

【先行技術文献】

【特許文献】

【0004】

【特許文献1】米国特許第8996781号明細書

【特許文献2】米国特許第9336217号明細書

【特許文献3】米国特許出願公開第2012/0310370号明細書

【特許文献4】米国特許出願公開第2016/0034809号明細書

【特許文献5】米国特許出願公開第2016/0283156号明細書

【特許文献6】米国特許出願公開第2017/0039269号明細書

【特許文献7】米国特許出願公開第2017/0148431号明細書

30

【特許文献8】米国特許出願公開第2017/0286284号明細書

【発明の概要】

【発明が解決しようとする課題】

【0005】

本発明は、上記従来の問題点に鑑みてなされたものであって、本発明の目的は、シリアライゼーションされたキーバリューストアを防止するマシンラーニングシステムのためのシステム及び方法を提供することにある。

【課題を解決するための手段】

【0006】

上記目的を達成するためになされた本発明の一態様によるマシンラーニングのための方法は、第1グラフィックプロセッシングユニット (GPU) により、キーを含む第1キーバリューストアを前記第1グラフィックプロセッシングユニットに連結された第1メモリの第1入出力領域内のキーバリューストアキューに書き込むステップと、前記第1メモリに連結された第1キーバリューストアデバイス (a first key value storage device (例えば、SSD)) により、前記第1キーバリューストアを前記第1メモリの第1入出力領域内のキーバリューストアキューから読み取るステップと、前記第1キーバリューストアデバイスにより、前記第1キーバリューストアに回答して、前記第1キーバリューストアのキーに対応する第1バリューストアを前記第1メモリの第1入出力領域に書き込むステップと、を有する。

40

【0007】

50

前記方法は、前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスで、キールックアップを遂行して前記第 1 バリューストレージデバイスをリトリブ (r e t r i e v e) するステップを更に含み得る。

前記第 1 キーバリューストレージデバイスによる書き込み領域は、前記第 1 バリューストレージデバイスに割り当てられた領域であるリターンバリューストレージ領域を含み得る。

前記第 1 バリューストレージデバイスを前記第 1 メモリの第 1 入出力領域に書き込むステップは、前記第 1 バリューストレージデバイスを前記リターンバリューストレージ領域に書き込むステップを含み得る。

前記第 1 バリューストレージデバイスを前記第 1 メモリの第 1 入出力領域に書き込むステップは、前記第 1 バリューストレージデバイスを前記第 1 メモリの第 1 入出力領域内のリターンバリューストレージキュー (r e t u r n v a l u e q u e u e) に書き込むステップを含み得る。

10

前記方法は、前記第 1 キーバリューストレージデバイス及び前記第 1 グラフィックプロセッシングユニットに連結されたホストによって、前記第 1 キーバリューストレージデバイスが、前記第 1 メモリの第 1 入出力領域にアクセスしてキーバリューストレージデバイスの要請を受信し、前記第 1 キーバリューストレージデバイスに回答してバリューストレージデバイスに書き込む (w r i t e) ように構成するステップと、前記第 1 グラフィックプロセッシングユニットが、前記第 1 メモリの第 1 入出力領域内のキーバリューストレージデバイスの要請を格納し、前記第 1 メモリの第 1 入出力領域からバリューストレージデバイスを読み取 (r e a d) るように構成するステップと、を更に含み得る。

前記方法は、前記ホストに連結された第 2 グラフィックプロセッシングユニットにより、キーを含む第 2 キーバリューストレージデバイスの要請を前記第 2 グラフィックプロセッシングユニットに連結された第 2 メモリの入出力領域内のキーバリューストレージデバイスに書き込むステップと、前記ホスト及び前記第 2 メモリに連結された第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスの要請を前記第 2 メモリの入出力領域内のキーバリューストレージデバイスから読み取るステップと、前記第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスの要請に回答して、前記第 2 キーバリューストレージデバイスのキーに対応する第 2 バリューストレージデバイスを前記第 2 メモリの入出力領域に書き込むステップと、を更に含み得る。

20

前記方法は、前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスで、キールックアップを遂行して前記第 1 バリューストレージデバイスをリトリブ (r e t r i e v e) するステップと、前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスを遂行するステップと同時に、前記第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスで、キールックアップを遂行して前記第 2 バリューストレージデバイスをリトリブするステップと、を更に含み得る。

30

前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスの要請を読み取るステップは、P2P (p e e r - t o - p e e r) DMA (d i r e c t m e m o r y a c c e s s) を介して前記第 1 キーバリューストレージデバイスの要請を読み取るステップを含み得る。

前記第 1 キーバリューストレージデバイスにより、前記第 1 バリューストレージデバイスを書き込む (w r i t e) ステップは、P2P (p e e r - t o - p e e r) DMA (d i r e c t m e m o r y a c c e s s) を介して前記第 1 バリューストレージデバイスを書き込む (w r i t e) ステップを含み得る。

前記第 1 キーバリューストレージデバイスは、PCI (p e r i p h e r a l c o m p o n e n t i n t e r c o n n e c t) 連結により、前記第 1 グラフィックプロセッシングユニットに連結され得る。

40

前記方法は、前記第 1 グラフィックプロセッシングユニットにより、前記第 1 キーバリューストレージデバイスの要請を前記第 1 メモリの第 1 入出力領域内のキーバリューストレージデバイスに書き込むステップの後に、且つ前記第 1 キーバリューストレージデバイスにより、前記第 1 バリューストレージデバイスを前記第 1 メモリの第 1 入出力領域に書き込むステップの前に、前記第 1 グラフィックプロセッシングユニットにより、第 2 キーバリューストレージデバイスの要請を前記第 1 メモリの第 1 入出力領域内のキーバリューストレージデバイスに書き込むステップを更に含み得る。

前記方法は、前記第 1 グラフィックプロセッシングユニットにより、キーを含む第 2 キーバリューストレージデバイスの要請を前記第 1 メモリの第 2 入出力領域内のキーバリューストレージデバイスに書き込むステップと、前記第 1 メモリに連結された第 2 キーバリューストレージデバイスにより

50

、前記第 2 キーバリューストレージデバイスを前記第 1 メモリの第 2 入出力領域内のキーバリューストレージキューから読み取るステップと、前記第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスに書き込むステップと、を更に含み得る。

前記方法は、前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスで、キールックアップを遂行して前記第 1 バリューストレージキューをリトリブ (retrieve) するステップと、前記第 1 キーバリューストレージデバイスにより、前記第 1 キーバリューストレージデバイスで、キールックアップを遂行するステップと同時に、前記第 2 キーバリューストレージデバイスにより、前記第 2 キーバリューストレージデバイスで、キールックアップを遂行して前記第 2 バリューストレージキューをリトリブするステップと、を更に含み得る。

10

【0008】

上記目的を達成するためになされた本発明の一態様によるマシンラーニングのためのシステムは、グラフィックプロセッシングユニットと、前記グラフィックプロセッシングユニットに連結されたメモリと、キーバリューストレージデバイスと、を備え、前記キーバリューストレージデバイスは、PCI (peripheral component interconnect) 連結により前記グラフィックプロセッシングユニットに連結され、前記グラフィックプロセッシングユニットは、前記メモリの入出力領域内のメモリがマッピング (mapping) された入力及び出力動作を遂行し、1 つ以上のキーバリューストレージキューを前記メモリの入出力領域内のキーバリューストレージキューに書き込み、前記キーバリューストレージデバイスは、前記メモリの入出力領域で、メモリがマッピングされた入力及び出力動作を遂行し、前記 1 つ以上のキーバリューストレージキューを前記キーバリューストレージキューから読み取り (read)、前記 1 つ以上のキーバリューストレージキューの中のキーバリューストレージキューに書き込む。

20

【0009】

前記キーバリューストレージデバイスによる書き込み領域は、前記バリューストレージキューに割り当てられた領域であるリターンバリューストレージキュー領域を含み得る。

前記バリューストレージキューを前記メモリの入出力領域に書き込むことは、前記バリューストレージキューを前記リターンバリューストレージキュー領域に書き込むことを含み得る。

前記バリューストレージキューを前記メモリの入出力領域に書き込むことは、前記バリューストレージキューを前記メモリの入出力領域内のリターンバリューストレージキュー (return value queue) に書き込むことを含み得る。

30

【0010】

上記目的を達成するためになされた本発明の他の態様によるマシンラーニングのためのシステムは、グラフィックプロセッシングユニットと、キーバリューストレージデバイスと、前記グラフィックプロセッシングユニットと前記キーバリューストレージデバイスとの間の通信のために共有されたメモリ手段と、を備え、前記グラフィックプロセッシングユニットは、前記通信のために共有されたメモリ手段を介して 1 つ以上のキーバリューストレージキューを前記キーバリューストレージデバイスに伝送し、前記キーバリューストレージデバイスは、前記 1 つ以上のキーバリューストレージキューを受信し、前記 1 つ以上のキーバリューストレージキューの中のキーバリューストレージキューに書き込む。

40

【0011】

前記通信のために共有されたメモリ手段は、前記グラフィックプロセッシングユニットに連結されたメモリを含み、PCI (peripheral component interconnect) 連結を通じた P2P (peer-to-peer) DMA (direct memory access) を介して、前記キーバリューストレージデバイスによってアクセスされ得る。

【発明の効果】

50

【0012】

本発明によると、ホストアプリケーションによって遂行されるタスク (task) が、GPUとSSDとの間の通信のための経路を設定することのみであるため、GPU計算のシリアライゼーションを防止することができる。従って、本発明は、規模拡張 (scale out) のための多数のGPUの使用を可能にして、マシンラーニングトレーニングを加速化することができる。

【図面の簡単な説明】

【0013】

【図1】従来技術によるマシンラーニングのためのシステムの機能ブロック図である。

【図2】本発明の一実施形態によるオンボードSSDを備えたグラフィックスカードのブロック図である。

10

【図3】本発明の一実施形態によるデータの流れ図である。

【図4】従来技術と本発明の一実施形態とを比較したタイミング図である。

【図5】従来技術と本発明の他の実施形態とを比較したタイミング図である。

【発明を実施するための形態】

【0014】

以下、本発明のキーバリューストトレージアクセスを含むマシンラーニングを遂行するためのシステム及び方法を実施するための形態の具体例を、図面を参照しながら詳細に説明する。本発明の実施形態は、本発明が構成されたり、活用されたりする唯一の形態を示すものとして意図しない。実施形態の説明は、図示した実施形態に関連する本発明の特徴を提供する。しかし、同一又は同等の機能及び構造は、本発明の範囲内に含まれるものとして意図する他の実施形態によって達成される。本明細書の他の箇所でも表示しているように、類似の構成要素の符号は類似の構成要素又は特徴を示すものとして意図する。

20

【0015】

従来技術のマシンラーニング (machine learning) プラットフォーム (platform) は、全体のトレーニング (training) データのサブセットをランダムに (randomly) サンプルする確率的 (stochastic) マシンラーニングトレーニングに利用される場合に欠点を有する。このようなマシンラーニングのプラットフォームは、確率的マシンラーニング中にキーバリューストアクセスによる低いGPU (graphics processing unit: グラフィックスプロセッシングユニット) の活用で問題になり、これはCPU (central processing unit) を含めてPCIe (peripheral component interconnect express) バスを横切って (traversing)、キーバリューストインターフェース及びデータ伝送を提供することを要求するからである。上述した通り、いくつかの従来技術のシステムにおいて、ホストCPUは、ファイルインデックスルックアップ (lookup) 及びファイルシステムアクセスを遂行してデータの位置を識別するが、これはシリアライゼーションされたキーバリューストアクセスにつながる。一方、いくつかの実施形態で、CPUがオンボード (onboard) SSD (solid state drive) に格納されたデータに対するキーバリューストアクセスでは、関与しないため、性能は向上する。GPUは、オンボードのキーバリューストストレージデバイス (例えば、オンボードのキーバリューストSSD)、例えば非同期 (asynchronous) のキーバリューストアクセスを可能にしてアクセス待機時間 (latency) の効果を減らす、GPU及びオンボードのキーバリューストを含むグラフィックスカード (on a graphics card) 上にキーバリューストのコマンドを直接伝送する。本明細書で使用するものとして、「キーバリューストストレージデバイス」は、このような各要請に回答して、要請に含まれるキーに対応するバリューストを返すことで、キーバリュースト要請 (それぞれは、キーを含む) に回答するように構成される (SSDのような) 永久ストレージデバイスである。

30

40

【0016】

図1は、従来技術によるマシンラーニングのためのシステムの機能ブロック図であり、

50

GPUマシンラーニングトレーニング中に、ソフトウェアキーバリューストアを利用して、ブロックインターフェースを有するSSD上に格納されたデータに対するキーバリューストアアクセスの全体的な流れを示している。先ず、ホストアプリケーション105は、「GET」要請（或いは取得要請）をソフトウェアキーバリューストア110（ソフトウェアキーバリューストア（S/W KV Store）とも称される）に伝送することにより、キーバリューストアアクセスを開始する。あるキーに対応するデータの位置を識別するために、ソフトウェアキーバリューストア110は、指定された（specified）キーに対応するデータのファイルオフセット（offset）を格納するインデックステーブルにアクセスする。そうすると、ソフトウェアキーバリューストア110は、ファイルオフセットを有するファイルシステム115にアクセスし、ファイルシステム115は、ブロックインターフェース120を有するSSDにアクセスし、指定されたキーに対応するデータをフェッチ（fetch）する。一旦バリューがホストアプリケーションで使用可能になると、ホストアプリケーションは、GPU計算のためのGPUメモリ125にバリューを伝送してGPUカーネル（kernel）を開始する（launch）。ブロックインターフェースを有する従来技術のSSDを有するソフトウェアキーバリューストアのため、これらに対する全ての動作は、順次的に（sequentially）遂行される。

10

20

30

40

50

【0017】

図1に示したように、ブロックインターフェースを有するSSDに対するキーバリューストアアクセスは、指定されたキーに対応するデータの位置を識別するためにホスト上で実行されるいくつかの計算ステップを含む。GPUはソフトウェアキーバリューストアの計算が終わった後にのみ計算を遂行することができる。GPUは、他のGPUからのキーバリューストアアクセスの完了を待たなければならず、これはGPU計算のシリアライゼーションにつながるため、より多くのGPUがシステムで使用されている場合、キーバリューストアアクセスの待ち時間は増加する。結果的に、ブロックインターフェースを有する従来技術のSSDに対するキーバリューストアアクセスは、システムで生産的に使用可能なGPUの数を制限する。

【0018】

本実施形態において、キーバリューストアインターフェース（又は「キーバリューストアSSD」）を有するオンボードSSDを搭載したグラフィックスカードが、従来技術のシステムの短所の一部を克服するために使用される。図2は、本発明の一実施形態によるオンボードSSDを備えたグラフィックスカードのブロック図であり、このようなデバイスを有する例示的なシステムを示す。ストレージとGPUとの間のデータの移動がグローバルPCIeバスを横切るデータ伝送を要請し、これはGPUからのデータアクセスの待機時間を増加させる従来技術のシステムとは異なり、オンボードのキーバリューストアSSD205を有するグラフィックスカードは、オンボードのキーバリューストアSSD205とGPU210との間のP2P（peer to peer）DMA（direct memory access）を利用することにより、オーバーヘッド（overhead）を減らすことができ、完全なP2P DMA制御をGPU210に与える。いくつかの実施形態で、オンボードのキーバリューストアSSD205は、非標準的な（non-standard）コマンドとしてキーバリューストアコマンドを提供する。例えば、キーバリューストア要請キュー（queue、以下でより詳細に説明する）は、NVMe（non volatile memory express）コマンドキューとして具現され、NVMeコマンドは、キーバリューストアSSD205のファームウェア（firmware）及びGPU上のドライバソフトウェアで定義されたベンダー固有の（vendor-specific）コマンドであり、キーバリューストアコマンド、即ちキーバリューストアSSD205からコマンドの一部として提供されるキーに対応するバリューを要請するために使用されるコマンドに対応する。

【0019】

本実施形態におけるシステムは、オンボードのキーバリューストアSSDにおける非同期のキーバリューストアアクセスを提供するために使用され、本実施形態は、トレーニングデータのランダムサンプリングのためにグラフィックスカードにおけるキーバリューストアSSDを活用す

る。図3は、本発明の一実施形態によるデータの流れ図であり、マシンラーニングトレーニング中のキーバリュアクセスの流れを示している。本実施形態と従来技術のシステムとの間の重要な違いは、本実施形態におけるGPUは、キーバリュコマンドをキーバリュSSD205に直接伝送することである。先ず、ホストアプリケーションの実行の初期ステップで、ホストアプリケーション105は、特定のGPUデバイスのメモリをPCI(peripheral component interconnect)BAR(base address registers)メモリ領域にマッピング(mapping)して、キーバリュSSD205とGPUとの間の直接的な通信を設定(establish)する。このプロセスにより、キーバリュSSD205とGPUとの間の通信(例えば、メモリがマッピング(mapping)された入出力によって)に割り当てられたGPUメモリの領域は、GPUメモリの「入出力領域」として本明細書で指称する。GPU及びキーバリュSSD205の両方によって直接的にアクセス可能なGPUメモリの入出力領域は、共有されたメモリとして機能的に動作する。GPUアプリケーション305は、キーバリュSSD205上でメモリがマッピングされたIO(入出力)を遂行し、露出されたGPUメモリのバスアドレスを供給することにより、キーバリュSSD205にGET要請を発行する。キーバリュSSD205内のファームウェアは、キールックアップを遂行して、キーに対応するバリュをリトリブした後に、ホストアプリケーション105の仲裁(intermediation)なしに、そのバリュをマッピングされたGPUデバイスのメモリに(即ち、GPUメモリの入出力領域に)書き込む(write)。

10

20

【0020】

本実施形態において、GPUが二番目の次の要請をする前に、一番目の要請に対する応答を待つ必要がないという意味で、キーバリュ要請キュー(key value request queue: KVRQ)310が使用され、キーバリュアクセスは、非遮断的(non-blocking)である。代わりに、GPUは、キーバリュ要請をキーバリュ要請キュー310に配置し、そして要請は、キーバリュSSD205によって順番に処理される。このように、GPUアプリケーションが要請をキーバリュ要請キュー310に入れることで、要請動作が完了する。キーバリュ要請キュー310のエントリの個数がキーバリュ要請の個数であるように、キーバリュ要請キュー310は、完了していない要請を有している(hold)。キーバリュSSD205内のファームウェアは、バリュがGPUメモリに伝送された場合、指定されたキーに対応するキーバリュ要請キューのエントリを開放する(release)。

30

【0021】

各GPUの個々のキーバリュアクセスは、多数のGPUからのキーバリュアクセスを重ねる(重畳させる)こと(又はGPUからのオーバーラッピング(overlapping)キーバリュアクセス)を可能にする。例えば、二つのGPUを搭載したシステムで、各GPUは、それぞれのキーバリュSSDに連結され、二つのGPUは、同時に要請を発行し、そしてそれらのそれぞれのキーバリュSSDは同時に応答することができる。図4は、従来技術と本発明の一実施形態とを比較したタイミング図であり、それぞれのキーバリュSSDにそれぞれ連結された二つのGPUを含む例の動作を示す。図4は、またGPU計算がシリアルライゼーションされた従来技術のアプローチと比較して、キーバリュアクセスを重ねることが二つのGPUにより遂行されることを示す、本実施形態で、システムにより節約された時間を示している。いくつかの実施形態で、3つ以上の(例えば、任意の個数)GPUは、それぞれのキーバリュSSDにそれぞれ連結され、キーバリュの動作を(例えば、同時に)重ねることを遂行する。

40

【0022】

本実施形態において、キーバリュアクセスのための要請及び応答の分離は、非同期のキーバリュアクセスを可能にし、例えばGPUから多数の要請のバッチング(batching)が可能である。図5は、従来技術と本発明の他の実施形態とを比較したタイミング図であり、二つのキーバリュコマンドをバッチングする場合の非同期のキーバリュ

50

ーアクセスの例を図示する。従来技術によるGPU計算及びSSDデバイスのアクセスがシリアライゼーションされる同期のキーバリュアクセスと比較して、本実施形態の非同期のキーバリュアクセスは、GPUの計算を有する多数のキーバリュコマンドの重なり(オーバーラッピング)を可能にする。この例で、GPUは、前の要請の完了のためにその度に待機する代わりに、GET要請を連続的に発行する。いくつかの実施形態で、3つ以上の(例えば、任意の個数)GPUは、それぞれのキーバリュSSDにそれぞれ連結され、キーバリュの動作を(例えば、同時に)重ねることを遂行する。

【0023】

本実施形態において、キーバリュSSDがキーバリュ要請に応答してバリュをリトリブ(retrieve)する場合、それはキーバリュ要請キューに、即ちこの目的のためにキーバリュ要請内の割り当てられたメモリの領域(又は「リターンバリュ(return-value)領域」)に再びリトリブ(retrieve)されたバリュを書き込む(write)。他の実施形態で、キーバリュSSDはGPUメモリの入出力領域に割り当てられた分離されたキュー(又は「リターンバリュキュー」)にリトリブ(retrieve)されたバリュを代わりに書き込む(write)。本実施形態で、各GPUは、それがキーバリュ要請を伝送する一つの専用のキーバリュSSDを有する代わりに、一つのGPUが複数のキーバリュSSDを有する。本実施形態で、複数のキーバリュ要請キューは、それぞれのキーバリュSSDに対してそれぞれGPUメモリ内に割り当てられる。他の実施形態で、いくつかのGPUは1つのキーバリュSSDに連結され、1つのキーバリュSSDは、例えばラウンドロビン(round-robin)方式で、GPUのそれぞれのキーバリュ要請キューで、キーバリュ要請を提供する。

10

20

【0024】

本実施形態において、ホストアプリケーションによって遂行されるタスク(task)は、GPUとSSDとの間の通信のための経路を設定することのみを必要とし、GPUの計算のシリアライゼーションを防止することにより本実施形態の拡張性(scalability)を向上させる。従来技術ではGPU計算のシリアライゼーションが防止されず、CPU上のホストアプリケーションによって遂行されるキーバリュアクセス動作が発生する。このように、本実施形態は、規模拡張(scale out)された多数のGPUの使用がマシンラーニングトレーニングを加速化させる。複雑なキーバリュソフトウェアを単純なデバイスのインターフェースに交替することにより、本実施形態は、またそれ以外のホスト上に課せられるCPUコアの個数の要請を含むリソースの要請を減らす。これらの要請を防止することは、より良いエネルギー効率を生じさせる。

30

【0025】

上述の実施形態は、一つ以上のプロセッシング回路を用いて構成される。本明細書で使用した「プロセッシング回路」という用語は、データ又はデジタル信号を処理するために利用されるハードウェア、ファームウェア、及びソフトウェアの任意の組み合わせを意味する。プロセッシング回路のハードウェアは、例えばASICs(application specific integrated circuits)、一般的な目的の又は特定の目的のCPUs(central processing units)、DSPs(digital signal processors)、GPUs(graphics processing units)、及びFPGAs(field programmable gate arrays)のようなプログラマブル(programmable)ロジックデバイスを含む。本明細書で使用したように、プロセッシング回路の各機能は、機能を遂行するために構成されたハードウェア(即ち、ハードウェアに内蔵(hard-wired))、又は非一時的(non-transitory)記憶媒体に格納されたコマンドを実行するように構成されたCPUのようなより一般的な目的のハードウェアのいずれか一つによって行われる。プロセッシング回路は、一つのPCB(printed circuit board)上で製造されるか、又は相互連結された複数のPCBに分散される。プロセッシング回路は、他のプロセッシング回路を含み、例え

40

50

ばプロセッシング回路は、PCB上に相互連結された二つのプロセッシング回路、FPGA、及びCPUを含む。

【0026】

多様な構成要素、成分、領域、階層、及び/又はセクションを説明するために、「第1」、「第2」、「第3」などを本明細書で使用したが、これらの構成要素、成分、階層、及び/又はセクションはこれらの用語によって制限されない。これらの用語は、単に一つの構成要素、成分、領域、階層、又はセクションを他の構成要素、成分、領域、階層、又はセクションから区別するために使用される。従って、本明細書で議論された第1の構成要素、成分、領域、階層、又はセクションは、本発明の思想及び範囲を逸脱せずに、第2の構成要素、成分、領域、階層、又はセクションと称される。

10

【0027】

本明細書で使用した用語は、特定の実施形態を説明するためのものであり、本発明を制限しようとするものではない。本明細書で使用したように、「実質的に」、「約」という用語及び類似の用語は、近似の用語として使用されて、程度の用語として使用されず、当業者によって識別されて測定された、又は計算されたバリューの固有の変動を考慮するためのものである。本明細書で使用したように、「主(major)成分」という用語は、組成物、ポリマー、又は組成物若しくは生成物の任意の他の一つの成分の量よりも大きい量の生成物に存在する成分を示す。一方、「主(primary)成分」という用語は、組成物、ポリマー、又は生成物の重量の少なくとも50%以上を構成する成分を示す。本明細書で使用したように、「主な部分」という用語は、複数の項目に適用される場合、項目の少なくとも半分を意味する。

20

【0028】

本明細書で使用したように、文脈上明らかに違うように示さない限り、単数表現は複数の形態も含むものとして意図する。「包含する(comprises)」及び/又は「包含する(comprising)」という用語は、本明細書で使用する場合、記述された特徴、数字、ステップ、動作、構成要素、及び/又は成分の存在を明示するが、一つ以上の他の特徴、数字、ステップ、動作、構成要素、成分、及び/又はこれらのグループの存在又は付加を排除しない。本明細書で使用したように、用語の「及び/又は」は、関連して列挙された項目の一つ以上の任意且つ全ての組み合わせを含む。「少なくとも一つの」のような表現は、要素のリストに先立つ場合、要素の全体リストを修正し、リストの個々の要素を修正しない。なお、「できる(may)」の使用は、本発明の実施形態を説明する場合、本発明の一つ以上の実施形態を示す。なお、「例示的な」という用語は、例示又は図示を示すものとして意図する。本明細書で使用したように、「利用する(use)」、「利用する(using)」、及び「使用された(used)」という用語は、「活用する(utilize)」、「活用する(utilizing)」、及び「活用された(utilized)」という用語と同じことを意味するものと見なされる。

30

【0029】

構成要素又は階層が、他の構成要素又は階層「上の」、「連結された」又は「隣接の」として称される場合、それは直接的に、他の構成要素又は階層上の、連結された又は隣接のものであり、一つ以上の中間構成要素又は階層が存在する。一方、構成要素又は階層が、他の構成要素又は階層「に直接的に」、「に直接的に連結された」、「のすぐ隣接して」として称される場合、中間要素又は階層は存在しない。

40

【0030】

以上、本発明の実施形態について図面を参照しながら詳細に説明したが、本発明は、上述の実施形態に限定されるものではなく、本発明の技術的範囲から逸脱しない範囲内で多様に変更実施することが可能である。

【産業上の利用可能性】

【0031】

本発明は、GPU計算におけるシリアライゼーションを防止することにより、多数のGPUの使用を可能にして、マシンラーニングトレーニングを加速化するマシンラーニング

50

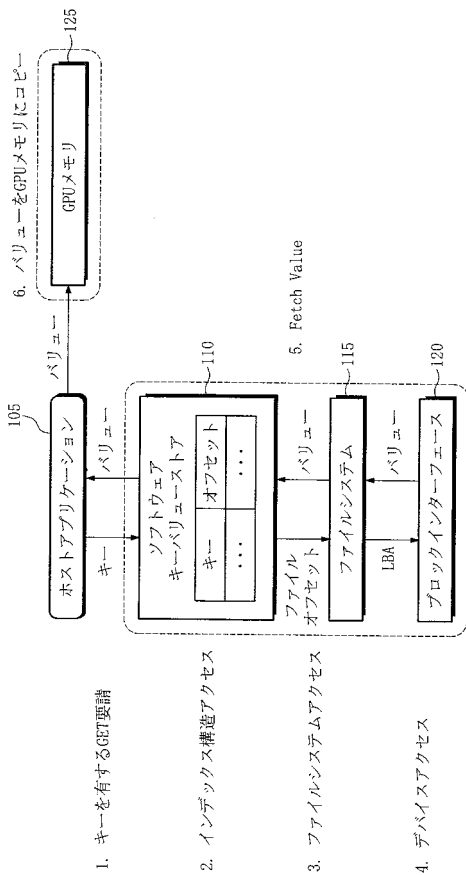
システムに有用である。

【符号の説明】

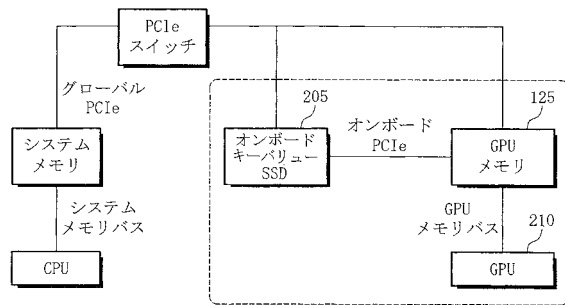
【0032】

- 105 ホストアプリケーション
- 110 ソフトウェアキーバリューストア
- 115 ファイルシステム
- 120 ブロックインターフェース
- 125 GPUメモリ
- 205 (オンボードの)キーバリュースSD
- 210 GPU
- 305 GPUアプリケーション
- 310 キーバリュース要請キュー

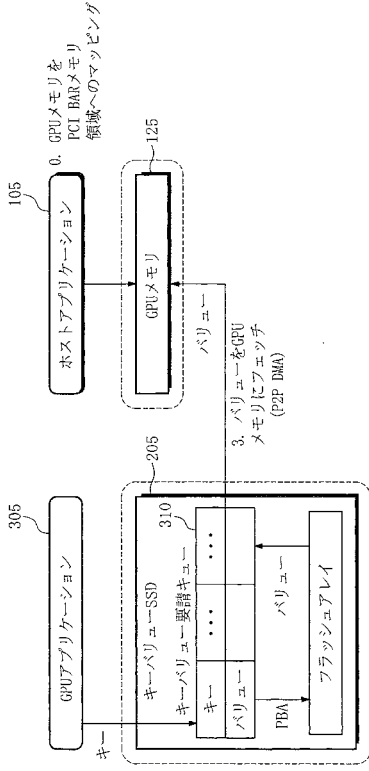
【図1】



【図2】



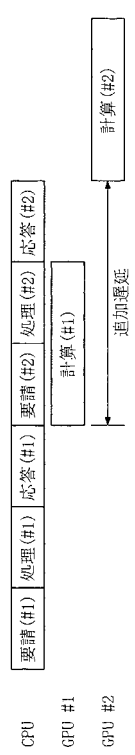
【 図 3 】



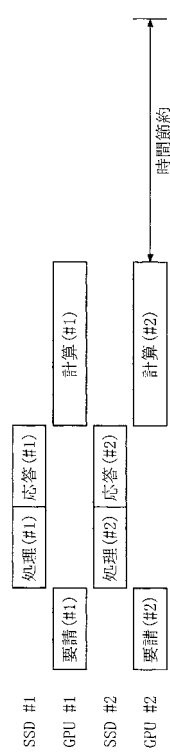
1. キーを有するGPU要求(MMIOを介してアバブルを鳴らす)
2. フラッシュメモリアレイアクセス

【 図 4 】

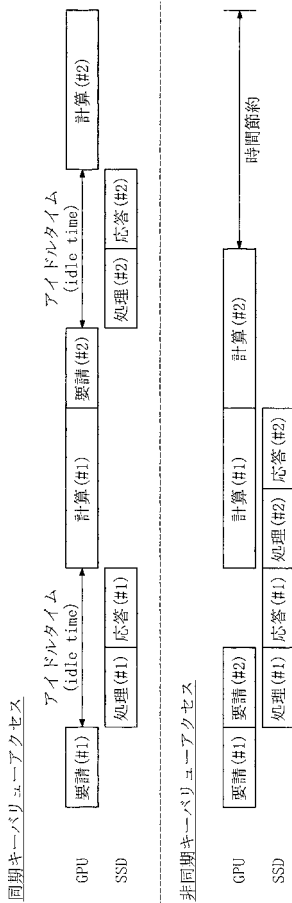
ソフトウェアバリューストア



キーバリュースSD



【 図 5 】



フロントページの続き

(72)発明者 奇 亮 ソク

アメリカ合衆国, 94303, カリフォルニア州, パロ アルト, アルテア ウォーク
873