

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4901097号  
(P4901097)

(45) 発行日 平成24年3月21日(2012.3.21)

(24) 登録日 平成24年1月13日(2012.1.13)

(51) Int.Cl. F I  
**G06N 3/00 (2006.01)**  
 G06N 3/00 550G  
 G06N 3/00 560E

請求項の数 29 (全 31 頁)

(21) 出願番号	特願2004-357352 (P2004-357352)	(73) 特許権者	500046438 マイクロソフト コーポレーション アメリカ合衆国 ワシントン州 9805 2-6399 レッドモンド ワン マイ クロソフト ウェイ
(22) 出願日	平成16年12月9日(2004.12.9)	(74) 代理人	100077481 弁理士 谷 義一
(65) 公開番号	特開2005-182785 (P2005-182785A)	(74) 代理人	100088915 弁理士 阿部 和夫
(43) 公開日	平成17年7月7日(2005.7.7)	(72) 発明者	デビッド ダブリュ. スタインクラウス アメリカ合衆国 98052 ワシントン 州 レッドモンド ワン マイクロソフト ウェイ マイクロソフト コーポレーシ ョン内
審査請求日	平成19年12月7日(2007.12.7)		
(31) 優先権主張番号	60/528, 578		
(32) 優先日	平成15年12月9日(2003.12.9)		
(33) 優先権主張国	米国 (US)		
(31) 優先権主張番号	10/837, 382		
(32) 優先日	平成16年4月30日(2004.4.30)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 グラフィックス処理ユニットを使用して機械学習技術の処理を速め、最適化するシステムおよび方法

(57) 【特許請求の範囲】

【請求項1】

コンピュータアプリケーションを処理するコンピュータ実施方法であって、  
 グラフィックス処理ユニットが、機械学習技術のトレーニング可能パラメータ、トレーニングデータ、及びターゲットデータを入力するステップと、

前記グラフィックス処理ユニットが、該入力したトレーニング可能パラメータ及び該入力したトレーニングデータに基づいた関数を計算することにより、第1の結果を取得するステップと、

前記第1の結果と前記ターゲットデータとの誤差が小さくなるように、前記トレーニング可能パラメータの値を更新するステップと、

前記グラフィックス処理ユニットが、クエリデータを入力するステップと、

前記グラフィックス処理ユニットが、該入力したクエリデータ及び該更新されたトレーニング可能パラメータに基づいた関数を計算することにより、第2の結果を取得するステップと、

前記第2の結果を使用して前記コンピュータアプリケーションによって使用できるように解を提供するステップとを備え、

前記計算するステップ及び前記取得するステップは、行列と列ベクトルとの内積を計算するステップをさらに含み、

前記内積を計算するステップは、

前記行列のうちの12列の行列と前記列ベクトルのうちの12列の列ベクトルとの内積

を計算するためのピクセルシェーダを使用することにより、複数の列ベクトル部分積を計算するステップと、

該計算された複数の列ベクトル部分積を合計することにより、単一の列ベクトルを取得するステップと

をさらに含むことを特徴とする方法。

【請求項 2】

前記コンピュータアプリケーションは ( a ) 音声認識アプリケーション、( b ) 手書き文字認識アプリケーションのうちの一方向の対話型使用コンピュータアプリケーションであることを特徴とする請求項 1 に記載のコンピュータ実施方法。

【請求項 3】

前記機械学習技術はニューラルネットワークであることを特徴とする請求項 1 に記載のコンピュータ実施方法。

【請求項 4】

前記グラフィックス処理ユニットに前記更新されたトレーニング可能パラメータの少なくとも一部を格納するステップをさらに含むことを特徴とする請求項 1 に記載のコンピュータ実施方法。

【請求項 5】

前記計算するステップ及び前記取得するステップは、前記ピクセルシェーダを使用して外積を計算するステップをさらに含むことを特徴とする請求項 1 に記載のコンピュータ実施方法。

【請求項 6】

テクスチャマッピングを使用して前記外積を計算するステップをさらに含むことを特徴とする請求項 5 に記載のコンピュータ実施方法。

【請求項 7】

前記計算するステップ及び前記取得するステップは、前記ピクセルシェーダを使用して行列の転置を行うステップをさらに含むことを特徴とする請求項 1 に記載のコンピュータ実施方法。

【請求項 8】

テクスチャマッピングを使用して前記行列の転置を行うステップをさらに含むことを特徴とする請求項 7 に記載のコンピュータ実施方法。

【請求項 9】

請求項 1 に記載の前記コンピュータ実施方法を実行するためのコンピュータ実行可能命令を格納したことを特徴とするコンピュータ可読記録媒体。

【請求項 10】

グラフィックス処理ユニットを使用して機械学習技術を処理するプロセスであって、グラフィックス処理ユニットが、前記機械学習技術のトレーニング可能パラメータ、トレーニングデータ、及びターゲットデータを入力するステップと、

前記グラフィックス処理ユニットが、該入力したトレーニング可能パラメータ及び該入力したトレーニングデータに基づいた関数を計算することにより、第 1 の結果を取得するステップと、

前記第 1 の結果と前記ターゲットデータとの誤差が小さくなるように、前記トレーニング可能パラメータの値を更新するステップと、

前記グラフィックス処理ユニットが、クエリデータを入力するステップと、

前記グラフィックス処理ユニットが、該入力したクエリデータ及び該更新されたトレーニング可能パラメータに基づいた関数を計算することにより、第 2 の結果を取得するステップと、

アプリケーションによって使用できるように前記第 2 の結果を出力するステップとを含み、

前記計算するステップ及び前記取得するステップは、行列と列ベクトルとの内積を計算するステップをさらに含み、

10

20

30

40

50

前記内積を計算するステップは、

前記行列のうちの12列の行列と前記列ベクトルのうちの12列の列ベクトルとの内積を計算するためのピクセルシェーダを使用することにより、複数の列ベクトル部分積を計算するステップと、

該計算された複数の列ベクトル部分積を合計することにより、単一の列ベクトルを取得するステップと

をさらに含むことを特徴とするプロセス。

【請求項11】

前記ピクセルシェーダはグラフィックカードに存在することを特徴とする請求項10に記載のプロセス。

【請求項12】

前記ピクセルシェーダを使用して、外積、行列の転置のうちの少なくとも1つを計算するステップをさらに含むことを特徴とする請求項10に記載のプロセス。

【請求項13】

前記機械学習技術はニューラルネットワークであることを特徴とする請求項10に記載のプロセス。

【請求項14】

前記機械学習技術は期待値最大化(EM)技術を使用することを特徴とする請求項10に記載のプロセス。

【請求項15】

前記機械学習技術は(a)K-means技術、(b)LVQ(Learning Vector Quantization)技術のうちの少なくとも一方であることを特徴とする請求項10に記載のプロセス。

【請求項16】

前記グラフィックス処理ユニットを使用して前記トレーニングデータを前処理するステップをさらに含むことを特徴とする請求項10に記載のプロセス。

【請求項17】

前記前処理するステップは、(a)前記トレーニングデータを標準化してそれを入力のためにより適した形式にするステップと、(b)前記トレーニングデータから情報および特徴を抽出するステップと、(c)入力データを取得して前記入力データを前記グラフィックス処理ユニットにロードするステップとのうちの少なくとも1つを含むことを特徴とする請求項16に記載のプロセス。

【請求項18】

1つまたは複数のプロセッサによって実行されると、前記1つまたは複数のプロセッサに請求項10の前記プロセスを実施させるためのコンピュータ可読命令を格納したことを特徴とする1つまたは複数のコンピュータ可読記録媒体。

【請求項19】

グラフィックス処理ユニットを使用して機械学習技術を処理するためのコンピュータ実行可能命令を格納したコンピュータ可読記録媒体であって、

グラフィックス処理ユニットが、前記機械学習技術のトレーニング可能パラメータ及びピクセルシェーダを入力するステップと、

前記グラフィックス処理ユニットが、前記機械学習技術のトレーニングデータ及びターゲットデータを入力するステップと、

前記グラフィックス処理ユニットが、該入力したトレーニング可能パラメータ及び該入力したトレーニングデータに基づいた関数を計算することにより、第1の結果を取得するステップと、

前記第1の結果と前記ターゲットデータとの誤差が小さくなるように、前記トレーニング可能パラメータの値を更新するステップと、

前記グラフィックス処理ユニットが、クエリデータを入力するステップと、

前記グラフィックス処理ユニットが、該入力したクエリデータ及び該更新されたトレー

10

20

30

40

50

ニング可能パラメータに基づいた関数を計算することにより、対話型使用アプリケーションによって使用するよう第2の結果を取得するステップとを含み、

前記計算するステップ及び前記取得するステップは、行列と列ベクトルとの内積を計算するステップをさらに含み、

前記内積を計算するステップは、

前記行列のうちの12列の行列と前記列ベクトルのうちの12列の列ベクトルとの内積を計算するための前記ピクセルシェーダを使用することにより、複数の列ベクトル部分積を計算するステップと、

該計算された複数の列ベクトル部分積を合計することにより、単一の列ベクトルを取得するステップと

10

をさらに含むことを特徴とするコンピュータ可読記録媒体。

【請求項20】

少なくとも一部のトレーニング可能パラメータを前記グラフィックス処理ユニットに配置するステップをさらに含むことを特徴とする請求項19に記載のコンピュータ可読記録媒体。

【請求項21】

前記クエリデータを前処理して入力データを取得するステップと、

前記グラフィックス処理ユニットに前記入力データをロードするステップと

をさらに含むことを特徴とする請求項19に記載のコンピュータ可読記録媒体。

【請求項22】

前記ピクセルシェーダを使用して、外積、行列の転置のうちの少なくとも1つを計算するステップをさらに含むことを特徴とする請求項19に記載のコンピュータ可読記録媒体。

20

【請求項23】

前記機械学習技術はニューラルネットワークであることを特徴とする請求項19に記載のコンピュータ可読記録媒体。

【請求項24】

前記機械学習技術は勾配降下技術であることを特徴とする請求項19に記載のコンピュータ可読記録媒体。

【請求項25】

前記勾配降下技術の調整可能パラメータの少なくとも一部を前記グラフィックス処理ユニットに配置するステップをさらに含むことを特徴とする請求項24に記載のコンピュータ可読記録媒体。

30

【請求項26】

前記対話型使用アプリケーションは手書き文字認識アプリケーションであることを特徴とする請求項19に記載のコンピュータ可読記録媒体。

【請求項27】

前記機械学習技術はニューラルネットワークであり、前記ニューラルネットワークを使用して分類を行うステップをさらに含むことを特徴とする請求項26に記載のコンピュータ可読記録媒体。

40

【請求項28】

前記対話型使用アプリケーションは光学文字認識アプリケーションであることを特徴とする請求項19に記載のコンピュータ可読記録媒体。

【請求項29】

前記機械学習技術はニューラルネットワークであり、前記ニューラルネットワークを使用して分類を行うステップをさらに含むことを特徴とする請求項28に記載のコンピュータ可読記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

50

本発明は、一般にコンピュータ処理技術に関し、より詳細にはグラフィックス処理ユニット（GPU）を使用して機械学習技術（ニューラルネットワークなど）および他の非グラフィックスのアプリケーションを処理し、その処理を速め、最適化するシステムおよび方法に関する。

【背景技術】

【0002】

グラフィックス処理ユニット（GPU）は、現代のパーソナルコンピュータ（PC）の不可欠な一部分である。GPUは、リアルタイムの3次元（3D）グラフィックスのユーザへの表示を速めるように設計された単一チッププロセッサである。当初ハイエンドのグラフィックスワークステーションの特徴であったGPUは、従来の中央処理装置（CPU）では適していなかった、または単に遅すぎたグラフィック機能のアクセラレータとしてパーソナルコンピュータバスに進出した。

10

【0003】

コンピュータグラフィックスは、一般に変更されたオシロスコープであるカリグラフィック表示装置上の線画として始まった。こうした表示の計算は、一般の座標変換、表示装置の境界へのクリッピング、および3D表示の透視変換を含むベクトル演算を必要とした。安価な商品である半導体メモリの出現がきっかけとなって、線画システムの代わりに、フレームバッファメモリを介してテレビのような表示を再生するラスタグラフィックスプロセッサが使用されるようになった。ユーザは一般に、ほとんどのアプリケーションの場合の線画ではなく、陰影付きのベタ塗りの表面を見ることを好むため、ラスタグラフィックスは、線画を迅速に表示した。線画で使用されていたような直線セグメントではなく、ラスタグラフィックスシステムの形状の構成単位（または基本要素）は、複数の三角形の配列から構成される多面体表面であった。表示の基本要素は、フレームバッファメモリに格納されたピクセルの矩形配列であった。配列の行は、ラスタ陰極線管（CRT）ディスプレイ上の個別の走査線に対応する。

20

【0004】

グラフィックスが線画からラスタグラフィックスに発展していくにつれて、より大きい処理能力が必要となった結果、数学コプロセッサがPCに含まれるようになった。数学コプロセッサとは、ホストCPU命令ストリームを共有し、CPUメモリにアクセスできる一体型の浮動小数点コプロセッサである。他のタイプの一体型コプロセッサは、パラレルデータバスを有し、非同期実行を行い、またCPUメモリにアクセスすることができるMultimedia Extension（MMX）、Streaming SIMD Extension（SSE）などのCPUの拡張機能である。

30

【0005】

3Dグラフィックスでの品質の高さおよび現実感の要求が絶えず高まっていった結果、より大きなグラフィックス処理能力が必要となった。この必要性を満たすために、計算集約的なグラフィックタスクを行うGPUが導入された。これによってCPUを楽にし、開放して、CPUが別の処理タスクを行えるようにした。現在の具体化では、GPUは、メインCPUメモリへのそれ自体の専用バス、およびそれ自体の専用のグラフィックスメモリを備えるPCの重要な構成要素である。数学コプロセッサとは対照的に、GPUは、それ自体の命令ストリーム、データバス、および専用メモリを備える自律した専用プロセッサである。

40

【0006】

GPUの設計および構成の現在の傾向によって、専用メモリがより大きくなり、グラフィックスメモリに対する帯域幅がより高くなり、また内部並列が向上している。さらに現在のGPUは、絶えず高まるプログラム可能性の度合いを考慮して設計されている。プログラム可能性の導入によって、GPUは、非グラフィックスのアプリケーションでの使い道を見つけるほど十分な柔軟性を得た。さらに、GPUのデータ並列アーキテクチャは、計算集約型アプリケーションでは、CPUに比べて劇的なパフォーマンスの向上をもたらす。代替のグラフィックスアルゴリズムへの拡張および科学的なコンピューティング問題

50

がいくつかの実例で探求されている。

【0007】

しかし、対話型使用（音声認識や手書き文字認識など）を対象としたアプリケーションには、相対的にあまり関心が集まっていない。この1つの理由は、GPUによって処理するためのこうしたアルゴリズムの実装が困難であり、いくつかの制限があったからである。汎用コンピューティングでは、GPUは本質的に制限を有するストリームプロセッサである。こうした制限を扱い、回避するには、明白でも直観的でもないプログラミングおよび処理のスタイルが必要である。

【0008】

こうした対話型使用アプリケーション（interactive use application）は、一般に自明ではない解を有し、大量のデータを扱う。こうした状況では、機械学習技術が好ましいソリューション技術である。機械学習技術は、トレーニング後、入力が正しく分類されるように、アルゴリズムのパラメータを自動的に調整することによって動作する。例えば、タスクが「A」のピクセル画像に正しいASCIIラベルを割り当てることであると仮定する。残念ながら、トレーニングは一般に、それら自体何十万もの演算を有するアルゴリズムに何十万もの対（入力、ターゲット）を提示することを伴う。その結果、最も早い使用可能なマシンでさえ、トレーニングには多大な時間がかかる可能性がある。また、実社会の様々な状態でアルゴリズムをテストし、または使用することは、法外に費用がかかる可能性がある。

【発明の開示】

【発明が解決しようとする課題】

【0009】

したがって、GPUを使用して機械学習技術処理するシステムおよび方法が必要である。さらに、GPUの制限を回避することによって機械学習技術の処理を速め、最適化するシステムおよび方法が必要である。このことによって、機械学習技術を使用した対話型使用アプリケーション（音声認識、手書き文字認識など）の解決にGPUの十分な処理能力が確実に利用されるようになる。

【課題を解決するための手段】

【0010】

本明細書で開示した本発明は、グラフィックス処理ユニット（GPU）を使用して機械学習技術処理するシステムおよび方法に関する。本発明は、CPU処理の一部をGPUに移植することによって中央処理装置（CPU）の計算の制限を軽減する。より具体的には、本発明は、従来CPUによって処理されていた機械学習アーキテクチャをGPUに移植する。機械学習技術は、GPUへの移植に特に適している。というのは、GPUは一般にCPUより強力であり、機械学習技術は、自明ではない解を有し、音声認識、手書き文字認識など、データアクセスよりかなり多くの計算を必要とする問題の解決を伴うからである。これは、計算に比べてデータアクセスが多いメモリベースの分類または検索とは異なる。GPUへのデータ転送にかなりの時間がかかる場合、GPUで計算を行うメリットは低下する。言い換えれば、計算に必要とされるのと同じ量のデータがGPUに転送される場合、データの転送がネックとなり、どんなネットの改良も無駄になる。

【0011】

本発明の方法は、GPUでのその後のテスト無しにGPUを使用して機械学習技術をトレーニングするステップと、GPUでの事前のトレーニング無しにGPUを使用して機械学習技術をテストするステップと、GPUでトレーニングおよびテストを行うステップとを含む。特に、本明細書で開示した本発明は、対話型使用のコンピュータアプリケーションを処理するコンピュータ実施方法を含む。この方法は、グラフィックス処理ユニットを使用して機械学習技術処理して関数の解を得るステップと、対話型使用のコンピュータアプリケーションによって使用できるようにその解を出力するステップとを含む。この方法は、グラフィックス処理ユニットを使用して機械学習技術のトレーニング用学習可能パラメータを得てトレーニング済み学習可能パラメータを得るステップと、トレーニング済

10

20

30

40

50

み学習可能パラメータを使用して解を得るステップとをさらに含む。対話型使用のコンピュータアプリケーションは、音声認識アプリケーションおよび手書き文字認識アプリケーションのうちの1つとすることができる。さらに機械学習技術は、ニューラルネットワークとすることができる。

【0012】

この方法は、学習可能パラメータを中央処理装置に格納するステップと、学習可能パラメータの少なくとも一部をグラフィックス処理ユニットに格納するステップとをさらに含むことができる。関数の解は、ピクセルシェーダを使用して、ベクトル内積および行列内積のうちの少なくとも一方である内積を計算することによって取得することができる。この方法は、内積をサブ問題に分解するステップと、ピクセルシェーダを使用してサブ問題にわたって複数回のパスを実行するステップとをさらに含むことができる。関数の解の取得は、ピクセルシェーダを使用して外積を計算することによって達成することができる。この方法は、テクスチャマッピングを使用して外積を計算するステップを含むこともできる。関数の解は、ピクセルシェーダを使用して行列の転置を行うことによって取得することができる。この方法は、テクスチャマッピングを使用して行列の転置を行うステップを含むこともできる。

10

【0013】

本明細書で開示した本発明は、グラフィックス処理ユニットを使用して機械学習技術を速め、最適化するプロセスも含む。この方法は、グラフィックス処理ユニット上でピクセルシェーダ (pixel shader) を使用して機械学習技術の学習可能パラメータをトレーニングするステップと、グラフィックス処理ユニット上でピクセルシェーダを使用し、またトレーニング済み学習可能パラメータを使用して機械学習技術から結果を取得するステップと、アプリケーションによって使用できるようにその結果を出力するステップとを含む。

20

【0014】

このプロセスは、グラフィックス処理ユニット上で学習可能パラメータのうちの少なくとも一部を探し出すステップと、ピクセルシェーダを使用して (a) ベクトル内積、(b) 行列内積、(c) 外積、(d) 行列の転置のうちの少なくとも1つを計算するステップとを含む。機械学習技術はニューラルネットワークとすることができ、ニューラルネットワークは、(a) マルチプレーヤの完全接続ニューラルネットワーク、(b) 重畳ニューラルネットワーク (convolutional neural network) のうちの少なくとも一方とすることができる。機械学習技術は、期待値最大化 (EM) アルゴリズムおよび K-means 技術および LVQ (Learning Vector Quantization) 技術を使用することもできる。このプロセスは、グラフィックス処理ユニットを使用してデータを前処理するステップも含む。

30

【0015】

本明細書で開示した本発明は、グラフィックス処理ユニットを使用して機械学習技術を処理するコンピュータ実行可能命令を有するコンピュータ可読媒体も含む。媒体は、シェーダをグラフィックス処理ユニットにロードするステップと、クエリデータを中央処理装置にロードするステップと、シェーダを呼び出して、機械学習技術を使用してグラフィックス処理ユニット上でトレーニング済み関数を処理するステップと、対話型使用アプリケーションによって使用できるように結果を取得するステップとを含む。

40

【0016】

コンピュータ可読媒体は、少なくとも一部の学習可能パラメータをグラフィックス処理ユニットに配置するステップと、入力データを取得するためにクエリデータを前処理し、入力データをグラフィックス処理ユニットにロードするステップとをさらに含む。さらに媒体は、シェーダを使用してグラフィックス処理ユニット上で基本的な演算 (primitive operation) を計算するステップを含む。基本的な演算は、(a) ベクトル内積、(b) 行列内積、(c) 外積、(d) 行列の転置のうちの少なくとも1つを含む。

50

## 【 0 0 1 7 】

機械学習技術は、ニューラルネットワーク、および勾配降下技術 (gradient descent technique) とすることができる。コンピュータ可読媒体は、勾配降下技術の調整可能パラメータのうち少なくとも一部をグラフィックス処理ユニットに配置するステップをさらに含む。対話型使用アプリケーションは、手書き文字認識アプリケーションとすることができる。さらに、機械学習技術は、ニューラルネットワークとすることができ、ニューラルネットワークを使用して分類を行うステップをさらに含む。また、対話型使用アプリケーションは、光学文字認識アプリケーションとすることもでき、機械学習技術は、ニューラルネットワークとし、ニューラルネットワークを使用して分類を行うステップをさらに含む。機械学習技術は、勾配降下技術とすることができ、グラフィックス処理ユニットに勾配降下パラメータを格納するステップと、勾配降下パラメータを2倍にすることによって勾配降下パラメータの仮数を拡張するステップとを含むことができる。

10

## 【 発明を実施するための最良の形態 】

## 【 0 0 1 8 】

本発明は、以下の説明および本発明の態様を示した添付図面を参照することによってさらに理解することができる。他の特徴および利点は、本発明の以下の詳細な説明を添付図面と併せ読めば明らかになる。添付図面は、一例として本発明の原理を示している。

## 【 0 0 1 9 】

図面を参照すると、図中、同様の参照番号は、図面を通じて対応する部分を表している。

20

## 【 0 0 2 0 】

本発明の以下の説明では、その一部を構成し、本発明を実施できる特定の例を実例として示した添付図面を参照する。他の実施形態を使用してもよく、また本発明の範囲から逸脱することなく構造的な変更を加えてもよいことを理解されたい。

## 【 0 0 2 1 】

## I . 導入

グラフィックス処理ユニット (GPU) は、従来、リアルタイム3Dグラフィックス表示を速めるために使用されてきた。しかし、処理能力およびプログラム可能性の向上の結果、GPUは、他の非グラフィックス関連のプロセスを効率的に処理することもできる。こうしたプロセスは一般に、大量のデータおよび処理 (計算流体力学など) を伴う代替のグラフィックスアルゴリズムおよび特定のコンピュータ問題に限定されている。しかし、音声認識および手書き文字認識などの対話型使用アプリケーションには、GPUアーキテクチャの制限および特徴のために、相対的にあまり関心が集まっていない。

30

## 【 0 0 2 2 】

本明細書に記載した機械学習GPU実行システムおよび方法は、CPU処理の一部またはすべてをGPUに移植することによって、CPUの計算の制限を軽減する。より具体的には、このシステムおよび方法は、様々な機械学習技術に使用できるアーキテクチャをCPUからGPUに転送する。GPUへの処理の転送は、制限を克服し、GPUアーキテクチャのフレームワーク内でよく働くいくつかの新しい技術を使用して達成される。こうした制限が克服された状態では、機械学習技術は、GPUでの処理に特に適している。というのは、一般にGPUは、一般のCPUよりかなり強力であるからである。さらに、グラフィックス処理と同様に、機械学習技術の処理は、自明ではない解および大量のデータの解決に関する問題を伴う。

40

## 【 0 0 2 3 】

## II . 動作環境の例

本明細書で開示した機械学習GPU実行エンジンおよび方法は、コンピューティング環境で動作するように設計されている。次の説明は、機械学習GPU実行エンジンおよび方法を実施できる適したコンピューティング環境の簡単な概要を提供するためのものである。

50

## 【 0 0 2 4 】

図 1 は、機械学習 GPU 実行エンジンおよび方法を実施できる適したコンピューティングシステム環境の例を示している。コンピューティングシステム環境 100 は、適したコンピューティング環境の一例にすぎず、本発明の使用または機能の範囲に関する限定を示唆するものではない。また、コンピューティング環境 100 を、動作環境 100 の例に示した構成要素のいずれか 1 つ、またはその組合せに関連する依存性または必要条件を有しているものと解釈すべきではない。

## 【 0 0 2 5 】

機械学習 GPU 実行エンジンおよび方法は、他の多くの汎用または専用コンピューティングシステム環境または構成で動作可能である。機械学習 GPU 実行エンジンおよび方法との使用に適したよく知られているコンピューティングシステム、環境、および/または構成の例には、それだけには限定されないが、パーソナルコンピュータ、サーバコンピュータ、例えばセルラー式電話や PDA などのハンドヘルド、ラップトップ、またはモバイルコンピュータまたは通信装置、マルチプロセッサシステム、マイクロプロセッサベースのシステム、セットトップボックス、プログラム可能家庭用電化製品、ネットワーク PC、ミニコンピュータ、メインフレームコンピュータ、上記の任意のシステムまたは装置を含む分散コンピューティング環境などがある。

## 【 0 0 2 6 】

機械学習 GPU 実行エンジンおよび方法は、コンピュータによって実行されるプログラムモジュールなどのコンピュータ実行可能命令の一般的な文脈で説明することができる。一般にプログラムモジュールは、特定のタスクを実行する、または特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、構成要素、データ構造などを含む。また、機械学習 GPU 実行エンジンおよび方法は、通信ネットワークによってリンクされているリモート処理装置によってタスクが実行される分散コンピューティング環境でも実施することができる。分散コンピューティング環境では、プログラムモジュールを、メモリ記憶装置を含むローカルおよびリモートのコンピュータ記憶媒体に配置することができる。図 1 を参照すると、機械学習 GPU 実行エンジンおよび方法を実施するシステムの例は、汎用コンピューティング装置をコンピュータ 110 の形で含んでいる。

## 【 0 0 2 7 】

コンピュータ 110 の構成要素は、それだけには限定されないが、処理ユニット 120 (中央処理装置、すなわち CPU など)、システムメモリ 130、およびシステムメモリを含む様々なシステム構成要素を処理ユニット 120 に結合するシステムバス 121 を含み得る。システムバス 121 は、様々なバスアーキテクチャのうちの任意のものを使用するメモリバスまたはメモリコントローラ、周辺バス、およびローカルバスを含むいくつかのタイプのバス構造のうちどんなものでもよい。こうしたアーキテクチャには、それだけには限定されないが一例として、ISA (Industry Standard Architecture) バス、MCA (Micro Channel Architecture) バス、EISA (Enhanced ISA) バス、VESA (Video Electronics Standards Association) ローカルバス、およびメザンバスとしても知られている PCI (Peripheral Component Interconnect) バスなどがある。

## 【 0 0 2 8 】

コンピュータ 110 は、一般に様々なコンピュータ可読媒体を含む。コンピュータ可読媒体は、コンピュータ 110 からアクセスできる使用可能な任意の媒体とすることができ、揮発性および不揮発性媒体、リムーバブルおよび非リムーバブル媒体を含む。コンピュータ可読媒体は、それだけには限定されないが一例として、コンピュータ記憶媒体および通信媒体を含み得る。コンピュータ記憶媒体には、コンピュータ可読命令、データ構造、プログラムモジュール、他のデータなど、情報を記憶するための任意の方法または技術で実施される揮発性および不揮発性のリムーバブルおよび非リムーバブル媒体がある。

## 【 0 0 2 9 】

10

20

30

40

50

コンピュータ記憶媒体には、それだけには限定されないが、RAM、ROM、EEPROM、フラッシュメモリまたは他のメモリ技術、CD-ROM、デジタル多用途ディスク(DVD)または他の光ディスク記憶装置、磁気カセット、磁気テープ、磁気ディスク記憶装置または他の磁気記憶装置、または所望の情報の格納に使用でき、コンピュータ110からアクセスできる他の任意の媒体などがある。通信媒体は一般に、コンピュータ可読命令、データ構造、プログラムモジュール、または他のデータを搬送波または他の移送機構などの変調されたデータ信号に組み込む。これには任意の情報配送媒体がある。

#### 【0030】

「変調されたデータ信号」という用語は、信号に情報を符号化するように1つまたは複数のその特性が設定または変更された信号を意味することに留意されたい。通信媒体には、それだけには限定されないが一例として、有線ネットワーク、直接配線された接続などの有線媒体、および音響、RF、赤外線、その他の無線媒体などの無線媒体がある。また、上記のどんな組合せでもコンピュータ可読媒体の範囲内に含まれるものとする。

#### 【0031】

システムメモリ130は、読取り専用メモリ(ROM)131やランダムアクセスメモリ(RAM)132など、揮発性および/または不揮発性メモリの形のコンピュータ記憶媒体を含む。基本入出力システム133(BIOS)は、例えば起動中など、コンピュータ110内の要素間での情報の転送を助ける基本ルーチンを含み、一般にROM131に格納されている。RAM132は一般に、処理ユニット120から直接アクセス可能な、かつ/または処理ユニット120が現在処理しているデータおよび/またはプログラムモジュールを含む。図1は、それだけには限定されないが一例として、オペレーティングシステム134、アプリケーションプログラム135、他のプログラムモジュール136、およびプログラムデータ137を示している。

#### 【0032】

コンピュータ110は、他のリムーバブル/非リムーバブル、揮発性/不揮発性コンピュータ記憶媒体を含むこともできる。一例にすぎないが、図1は、非リムーバブル不揮発性磁気媒体から読み取り、あるいはそこに書き込むハードディスクドライブ141、リムーバブル不揮発性磁気ディスク152から読み取り、あるいはそこに書き込む磁気ディスクドライブ151、およびCD-ROMや他の光媒体など、リムーバブル不揮発性光ディスク156から読み取り、あるいはそこに書き込む光ディスクドライブ155を示している。

#### 【0033】

動作環境の例で使用できる他のリムーバブル/非リムーバブル、揮発性/不揮発性コンピュータ記憶媒体には、それだけには限定されないが、磁気テープカセット、フラッシュメモリカード、デジタル多用途ディスク、デジタルビデオテープ、半導体RAM、半導体ROMなどがある。ハードディスクドライブ141は一般に、インターフェース140などの非リムーバブルメモリインターフェースを介してシステムバス121に接続され、磁気ディスクドライブ151および光ディスクドライブ155は一般に、インターフェース150などのリムーバブルメモリインターフェースによってシステムバス121に接続される。

#### 【0034】

上述し、図1に示したドライブおよびその関連のコンピュータ記憶媒体は、コンピュータ可読命令、データ構造、プログラムモジュール、およびコンピュータ110の他のデータの記憶域を提供する。図1では例えば、ハードディスクドライブ141は、オペレーティングシステム144、アプリケーションプログラム145、他のプログラムモジュール146、およびプログラムデータ147を格納するものとして示されている。これらの構成要素は、オペレーティングシステム134、アプリケーションプログラム135、他のプログラムモジュール136、およびプログラムデータ137と同じであっても、異なってもよいことに留意されたい。オペレーティングシステム144、アプリケーションプログラム145、他のプログラムモジュール146、およびプログラムデータ147は

10

20

30

40

50

少なくとも異なるコピーであることを示すために、ここではそれらに異なる番号を付している。ユーザは、キーボード 162、および一般にマウス、トラックボール、またはタッチパッドと呼ばれるポインティング装置 161 などの入力装置を介してコマンドおよび情報をコンピュータ 110 に入力することができる。

#### 【0035】

他の入力装置（図示せず）には、マイクロフォン、ジョイスティック、ゲームパッド、衛星パラボラアンテナ、スキャナ、無線受信機、またはテレビまたはブロードキャストビデオ受信機などがある。これらおよび他の入力装置は、しばしばシステムバス 121 に結合されているユーザ入力インターフェース 160 を介して処理ユニット 120 に接続されるが、パラレルポート、ゲームポート、ユニバーサルシリアルバス（USB）など他のインターフェースおよびバス構造で接続してもよい。モニタ 191 または他のタイプの表示装置もまた、ビデオインターフェース 190 などのインターフェースを介してシステムバス 121 に接続される。モニタに加えて、コンピュータは、出力周辺インターフェース 195 を介して接続できるスピーカ 197、プリンタ 196 などの他の周辺出力装置を含むこともできる。

#### 【0036】

コンピュータ 110 は、リモートコンピュータ 180 など 1 つまたは複数のリモートコンピュータへの論理接続を使用してネットワーク環境で動作することができる。リモートコンピュータ 180 は、パーソナルコンピュータ、サーバ、ルータ、ネットワーク PC、ピア装置、または他の一般のネットワークノードでよく、一般にコンピュータ 110 に関連して上述した多くまたはすべての要素を含むが、図 1 にはメモリ記憶装置 181 のみを示している。図 1 に示した論理接続は、ローカルエリアネットワーク（LAN）171 および広域ネットワーク（WAN）173 を含むが、他のネットワークを含んでいてもよい。こうしたネットワーキング環境は、オフィス、全社規模のコンピュータネットワーク、イントラネット、およびインターネットではごく一般的である。

#### 【0037】

LAN ネットワーキング環境で使用する場合、コンピュータ 110 は、ネットワークインターフェースまたはアダプタ 170 を介して LAN 171 に接続される。WAN ネットワーキング環境で使用する場合、コンピュータ 110 は一般に、モデム 172、またはインターネットなど WAN 173 を介して通信を確立する他の手段を含む。モデム 172 は、内蔵のものでも外付けのものでもよく、ユーザ入力インターフェース 160 または他の適切な機構を介してシステムバス 121 に接続することができる。ネットワーク環境では、コンピュータ 110 に関連して示したプログラムモジュール、またはその一部をリモートメモリ記憶装置に格納することができる。図 1 は、それだけには限定されないが一例として、リモートアプリケーションプログラム 185 をメモリ装置 181 上に存在するものとして示している。図示したネットワーク接続は例であり、コンピュータ間の通信リンクを確立する他の手段を使用してもよいことは理解されよう。

#### 【0038】

##### III. 概要

GPU は、リアルタイム 3D グラフィックス表示を速めるように設計されている。より優れたグラフィックスに対する要求が増えるにつれて、GPU は急速により強力なプログラム可能なものになりつつある。プログラム可能性の向上の結果、GPU は、他の多くのタイプの非グラフィックス関連プロセスの効率的な処理を行うこともできる。本明細書で開示した機械学習 GPU 実行エンジンおよび方法は、CPU 処理の一部を GPU に移植することによって CPU の計算の制約を軽減する。より具体的には、本発明は、従来 CPU によって処理されていた機械学習アーキテクチャを GPU に移植する。以下で詳述するように、これは、GPU のいくつかの制限を克服し、機械学習技術の GPU 処理を速め、最適化するためのいくつかの実装技術の使用を必要とする。汎用コンピューティングでは、機械学習は、GPU への移植に特に適している。というのは、GPU は一般の CPU より強力だからである。

## 【 0 0 3 9 】

図2は、本明細書で開示した機械学習GPU実行エンジンおよび方法の実装例を示すブロック図である。図2は、機械学習GPU実行エンジンおよび方法を実施し、使用できるいくつかの方法の1つにすぎないことに留意されたい。機械学習技術は、トレーニング段階およびテスト（または使用）段階中にGPUによって処理することができる。トレーニング段階とは、アルゴリズムのパラメータがトレーニングデータを使用して調整される計算を指す。テスト段階とは、システムを使用して、入力データおよびトレーニング済みパラメータに応じて有用な情報を計算する計算を指す。トレーニングは一般に時間がかかるが、行うのは一度だけでよい。テストは、配置したシステムを使用することを指し得るが、例えば文字認識または音声認識などの場合、時として非常に迅速な応答時間を必要とすることがある。トレーニングが長時間かかり、しかし最大限のハードウェア独立性が望まれる分野では、トレーニングはGPUで行われ、テストはCPUで行われる。また、テストアルゴリズムは簡単であるが、トレーニングアルゴリズムは非常に複雑でGPUでは稼働しない可能性もある（例えばニューラルネットワークにおけるトレーニング可能な重畳レイヤ（convolution layer）など）。この場合、学習アルゴリズムをCPUでトレーニングし、しかしテスト段階中はGPUを使用することができる。一部の場合には、当然、GPUでトレーニングおよびテストの両方を稼働させることが好ましいこともある。GPUの実施は、トレーニングおよびテストの両方について、機械学習処理速度を、CPUのみの実施に比べて最高で1桁スピードアップする。GPU技術は新しく、古いCPU技術より進化が早いため、この比率は増加している。

10

20

## 【 0 0 4 0 】

より具体的には、図2に示すように、コンピューティング装置110は、機械学習GPU実行エンジン200を含む。コンピューティング装置は、CPU120およびCPUメモリ130をさらに含む。CPUは、高速バス210（accelerated bus）を介してビデオインターフェース190と通信する。このバス210は、3Dグラフィックスのスループットに対する要求のために特に設計されたAccelerated Graphics Port（AGP）またはより新しいPCI Expressであることが好ましい。

## 【 0 0 4 1 】

ビデオインターフェース190は、GPU220およびGPUメモリ230を含む。GPU220は、バス210を介してCPU120にデータを転送することができる。また、機械学習GPU実行エンジン200もGPU220およびGPUメモリ230と通信する。機械学習GPU実行エンジン200は、トレーニングモジュール240およびテストモジュール250を含む。トレーニングモジュール240は、機械学習技術のトレーニング段階中に、GPUを使用して技術のパラメータを訓練するために使用される。テスト（または使用）モジュール250は、トレーニング済みパラメータおよび入力の所与の関数を計算するために使用される。計算の結果は、機械学習GPU実行エンジン200からCPU120に転送されて、音声認識など、対話型使用アプリケーションで使用される。

30

## 【 0 0 4 2 】

## IV. システムの構成要素

図3は、一例として勾配降下アルゴリズムを使用する3層学習機械のトレーニングモジュールを示すブロック/フロー図である。図3は、一例として示したにすぎず、本発明を勾配降下アルゴリズムに限定するためのものではない。この例では、勾配降下を使用して学習機械の最適化が行われる。他の機械学習技術では、ベイジアンネットワーク、グラフィカルモデルなどのように、勾配の代わりに確率を伝えることもある。層の数または構成もまた制限的なものではなく、本発明はより多くの層、より少ない層、または異なる層の構成（不連続）を有することができる。パラメータを $W = (W_1, W_2, W_3)$ とする。モジュール全体の入力は $X$ 、出力は $Y = Y_3$ である。

40

## 【 0 0 4 3 】

各層は、その入力 $Y$ 、 $X$ 、 $W$ の関数 $G$ を計算する。層の関数の観点から見ると、 $W$ また

50

は $Y$ は対称的であり、区別できない。これは、 $W_1$ 、 $W_2$ 、 $W_3$ は、それ自体他の層によって計算することができることを意味する。 $G$ は、 $Y$ および $W$ より多くの入力を有することもできる。各層は、その入力の間関数を計算し（前方パス）、 $dE/dY$ で示した出力勾配にヤコビの転置を掛けることによって得られるその出力の導関数を伝えることによって隣接した層と通信する。これは後方パスである。定義上、 $G$ のヤコビ $J$ は、行列 $J_{k,i} = dY_k/dX_i$ となる。式中、 $k$ は関数 $G$ の $k$ 番目の出力、および $i$ は関数 $G$ の $i$ 番目の入力の指数を表す。誤差関数 $E$ は、費用関数を使用して最後の層の出力をターゲット $T$ と比較する。異なる費用関数の例には、平均二乗誤差(MSE)、クロスエントロピー(CE)などがある。費用関数は、勾配、すなわち各変数が費用にどれだけ影響を与えるか、またどの方向に影響を与えるかを変数ごとに提供している。この勾配は、システムにおける変数およびパラメータごとに計算される。パラメータ $W$ は、費用関数を低減するように更新される。図4は、図3に示した3層学習機械のテストモジュールを示すブロック/フロー図である。

10

#### 【0044】

##### V. 動作の概要

本明細書で開示した機械学習GPU実行エンジン200は、機械学習GPU実行方法を使用して、GPUによる機械学習技術の処理を可能にする。一般にこの方法は、GPUを使用して機械学習技術をトレーニングするトレーニング段階、およびトレーニング済み関数を使用して適用された問題を解決するテスト段階に分けることができる。次にこれらの各段階について説明する。

20

#### 【0045】

##### トレーニング段階

トレーニング段階では、機械学習GPU実行方法を使用して、機械学習技術の学習可能パラメータをトレーニングする。図5は、機械学習GPU実行方法のトレーニング段階の動作の概要を示すフロー図である。図5は、GPUで実行されるトレーニングセッションのデータフローおよび制御構造を表す。図は、任意のタイプの機械学習技術に固有のものではないことに留意されたい。多くのタイプの機械学習技術(SVM、K-means、LVQ(Learning Vector Quantization)、期待値最大化(EM)など)がこのアーキテクチャを使用することができる。上記の機械学習技術は当分野ではよく知られており、詳しい説明は行わない。図5に示すように、CPUは、命令をGPUに与え、GPUがこれらの命令を完了するまでブロックする。あるいは、CPUおよびGPUに同時に計算させることもできる。

30

#### 【0046】

図5を参照すると、機械学習GPU実行方法のトレーニング段階は、シェードプログラム(P)、重みなど最初のトレーニング可能パラメータ(W)、および他の学習パラメータ(L)をロードする(ボックス500)ことによって開始する。このデータは次いでCPUからGPUに転送される(矢印505)。CPUからGPUへのデータ転送は、比較的費用がかかるため、CPUは、トレーニンググループに入る前に、GPUにできる限りプリロードしておく。トレーニング可能パラメータ(W)は、GPUメモリ230に格納される(ボックス510)。好ましい実装では、学習パラメータ(L)は、ニューラルネットワーク層のそれぞれの重み、および各ユニットの閾値である。学習パラメータ(L)は、学習率と呼ばれる単一のスカラである。シェードプログラム(P)は、順伝搬(forward propagation)および逆伝搬(backward propagation)および重みの更新に使用する様々なタイプのシェードのためのものである。

40

#### 【0047】

CPUは次いでトレーニングデータに対してループを開始し、トレーニングデータ(X, T)のバッチをロードすることによってトレーニングデータの組を蓄積する(ボックス515)。(X)はピクセル画像、(T)はターゲットラベルを表す。グループ(またはバッチ)を使用する理由は、CPUとGPUとの間のデータ転送を開始する費用があるからである。いくつかのパターンのグループごとに同時にデータの転送を行うことがより効

50

率的である。

【0048】

次いでピクセル画像(X)は、(X)が(X')に変換されるように前処理される(ボックス520)。次に(X')および(T)はGPUに送信される(矢印525)。次いでトレーニングデータ(X')は、トレーニングターゲットデータ(T)(ボックス535)とともにGPUメモリ230に格納される(ボックス530)。前処理には、いくつかあげると、データをよりよい形式にするためのデータの標準化、インテリジェントなまたは複雑な特徴の抽出、およびデータセットの質を高めるための歪みの生成など、多くの様々な機能があり得る。理論上、前処理は、GPUまたはCPUのいずれかで行うことができる。しかし実際には、GPUよりCPUでプログラムする方がかなり容易である。これは、前処理が計算的に高価でない場合、CPUで前処理を稼働させる方がかなり容易であることを意味する。一部の場

10

【0049】

前処理にトレーニング可能パラメータを含めることができることに留意されたい。これは、トレーニング可能パラメータが大域的最適化の一部であり、GPUに常駐する学習パラメータとともにトレーニングされる場合でさえも当てはまる。しかし、次のパターンが更新された前処理から利益を得るように、場合によっては、各パターンの後に、CPUに常駐するパラメータを更新するために情報(誤差勾配、負のフィードバックなど)をGPUからCPUに戻す必要があるという難問が生じる。現在、GPUからCPUへのデータフローは最適化されない。というのは、一般にグラフィックカードは、データを画面には送信するが、CPUには戻さないように設計されているからである。その結果、現在のアーキテクチャでは、すべてのトレーニング可能パラメータをGPUで保持することが好ましい。他のすべてのデータをCPUで保持し、処理はGPUで行うことが好ましい。あるいは、トレーニング可能パラメータをCPUおよびGPUの両方に配置することができる。

20

【0050】

トレーニングデータ(X')がGPUにロードされると、CPUは、機械学習技術の処理に必要な様々なシェーダを稼働させるようGPUに命令する(ボックス540)。一般の機械学習技術は、GPU220にロードされた学習モジュール545によって表される。学習モジュール545は、前処理された入力(X')およびトレーニング可能パラメータ(W)に応じて関数 $G(X', W)$ を計算する学習機械550を含む。目的は、この出力をターゲット値(T)にできるだけ近づけることである。 $G(X', W)$ と(T)との間の誤差555が計算され、誤差信号(Wに対する勾配など)が学習機械550に送り返される。次いで $G(X', W)$ と(T)との間の誤差を減らすために重み(W)が更新される。

30

【0051】

一例として、2層ニューラルネットワークをトレーニングするとき、前方および後方の伝達は、約20の異なるシェーダ(その一部は複数呼び出される)に相当する。シェーダの数および複雑度は、当然、使用するアルゴリズムによって異なり得る。シェーダは、1つのグループ内のパターンごとに呼び出される(例えば一部の場

40

50

## 【0052】

CPUがトレーニング中のGPUからフィードバックを受け取ることができるように、次の2つのプロセスを普遍性のために追加することができる。例えば、図5に示すように、CPUは、トレーニング統計値を収集し、トレーニング進捗データ (training progress data) を取得することができる (ボックス560)。これは、GPUが更新済みのトレーニング可能パラメータ (W) および誤差をCPUに送信することによって達成される (矢印565)。このプロセスは、図5の点線のボックスで示すように、オプションである。トレーニング進捗データを使用して、学習パラメータ、またはいくつかの種類のパターンの提示の頻度さえも調整することができる (ボックス570)。例えば、ニューラルネットワークでは、時として学習が進行するにつれて学習率を低下させるのが望ましいことがある。「ブースティング (boosting)」と呼ばれる別のクラスのアルゴリズムでは、いくつかのパターンの頻度、またはその学習の影響は、システムによる誤差に応じて変更することができる。

10

## 【0053】

前処理でのトレーニング可能パラメータも、GPUからの誤差信号に応じて変更することができる。これは、誤差信号がCPUに戻り、CPUで学習パラメータを更新することができるようにすることによって達成される。グループのサイズは、それに応じて変更することができる。極端には、グループのサイズは1であり、これは、前処理のセクションで上述したように、誤差信号がGPUから戻るとすぐに、CPUのトレーニングパラメータがそれぞれGPUに提示された後で更新されることを意味する。

20

## 【0054】

次いで、トレーニングが完了したかどうかに関する決定が行われる (ボックス575)。これは、一定の回数の反復の後、トレーニングデータのバッチのすべてが処理されたとき、または所望の誤差閾値が達成されたときに決定することができる。トレーニングが完了していない場合、トレーニンググループがデータのロードから再開する (ボックス515)。そうでない場合、最後にトレーニングされたパラメータ (W) が取得される (ボックス580)。これは、GPUがパラメータ (W) をCPUに送信することによって達成される (矢印585)。次いでトレーニングが完了する (ボックス590)。

## 【0055】

## テスト段階

テスト段階では、機械学習GPU実行方法を使用して入力およびトレーニング可能パラメータの関数を計算する。図6は、機械学習GPU実行方法のテスト段階の動作の概要を示すフロー図である。図6に示すように、機械学習GPU実行方法は、一部のトレーニング可能パラメータWおよび入力Xの所与の関数Gを計算する。トレーニング可能パラメータは、GPU、CPU、または他の任意の手段によって計算されている可能性があることに留意されたい。さらに、トレーニング可能パラメータは、トレーニングの結果であっても、またはそうでなくてもよい。この点で、図6は、図5に示したトレーニング可能パラメータのトレーニングには依存していない。

30

## 【0056】

図6では、図5のアーキテクチャと類似のアーキテクチャが使用されている。特に、図5および図6のアーキテクチャは、データ構造およびシェードプログラムPを共有する。しかし1つの重要な違いは、トレーニングがすでに完了しているため、逆伝搬シェードはもはや必要ないということである。機械学習GPU実行方法のテスト段階は、シェードプログラムPおよびトレーニング可能パラメータWをGPUにロードする (ボックス600) ことによって開始する。トレーニング段階と同様、テスト段階は、認識/使用ループ外にできる限りダウンロードし、事前に計算しようとする。次いでシェードプログラムPおよびトレーニング可能パラメータWは、GPUに送信される (矢印605)。特に、シェードプログラムPは、GPU220に送信されて処理を行い、トレーニング可能パラメータW610はGPUメモリ230に格納される。

40

## 【0057】

50

次に、関数  $G$  が計算される 1 組のパターン  $X$  が収集され、ロードされる (ボックス 6 1 5)。CPU と GPU との間の通信を開始するための費用があるため、グループ分けを行うことによって、この費用をいくつかのパターンにわたって消却することができる。次いでパターン  $X$  は、CPU から GPU に送信される前に前処理されて  $X'$  になる (ボックス 6 2 0)。次いで  $X'$  は GPU に送信される (矢印 6 2 5)。入力データ  $X' 6 3 0$  は、GPU メモリ 2 3 0 に格納される。

【 0 0 5 8 】

前処理は、CPU または GPU のいずれかで行うことができる。しかし、前処理は、計算的に費用がかかりすぎない限り CPU で行われることが好ましい。前処理は、標準化、および重要な情報および特徴をデータから抽出するなどの機能に有用である。前処理後、パターン  $X'$  は、グループとして GPU に送信される (矢印 6 2 5)。

10

【 0 0 5 9 】

次に、CPU は、シェードプログラム  $P$  を使用するよう GPU に命令する (ボックス 6 3 5)。関数モジュール 6 4 0 は、GPU に存在し、トレーニング済み関数 6 4 5 を含む。トレーニング済み関数 6 4 5 をシェードプログラムとともに使用して、関数  $Y = G(X', W)$  を計算する。GPU メモリ 2 3 0 内の一時計算スペース 6 5 0 を使用して、この計算を助けることができる。この計算から、グループの各パターンの結果  $Y 6 5 5$  が GPU メモリ 2 3 0 に蓄積され、関数モジュールの出力として送信される (ボックス 6 6 0)。次いで結果  $Y$  は、CPU に送り返され (矢印 6 6 5)、結果  $Y$  は CPU で取得される (ボックス 6 7 0)。繰り返しになるが、GPU から CPU への転送はかなり費用がかかるため、グループ分けを行い、 $Y$  のサイズを最低限に抑えることが有利である。手書き文字の分類の例では、 $Y$  は画像  $X$  のクラスであるにすぎず、したがって非常に小さい。クエリデータ  $X$  の追加のバッチがあるかどうかに関する決定が行われる (ボックス 6 7 5)。バッチがある場合は、認識 / 使用ループが再開する。そうでない場合、結果  $Y$  が出力として送信される (ボックス 6 8 0)。

20

【 0 0 6 0 】

VI . 動作の詳細

機械学習の問題の説明

いくつかの機械学習技術は、非常に大きい行列演算として実行することができる計算を核として含んでいる。特に興味深いのは、大きい行列 / ベクトル乗算である。一例として、手書き文字認識に使用されるニューラルネットワークは、一般に 1 対の層を含んでいる。図 7 は、ニューラルネットワークにおける 1 対の層の詳細図である。

30

【 0 0 6 1 】

図 7 を参照すると、次のように、隠れ変数 (hidden variable) の計算が行列乗算として実行され、次いでマッピングステップが行われる (1 層の場合の計算)

$$\begin{aligned} [w][i] &= [o'] \\ [o] &= f([o']) \end{aligned}$$

【 0 0 6 2 】

式中  $f$  はシグモイド関数

40

【 0 0 6 3 】

【数 1】

$$f(x) = \frac{1}{e^{-x} + 1}$$

【 0 0 6 4 】

または  $\tanh$  関数

【 0 0 6 5 】

## 【数2】

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

## 【0066】

である。こうした2つの層はカスケードされて、結果が生成される。入力の特徴ベクトルサイズ、隠れ変数、および出力は何千にも及ぶ。図7および上記はニューラルネットワークに基づいているが、他のいくつかの機械学習問題は、上述した大きい浮動小数点ベクトルに対する反復演算の一般的なモデルに適合することに留意されたい。

## 【0067】

GPUでの基本的な演算

幸いにも、ニューラルネットワーク、期待値最大化、および他の多くの勾配降下ベースのアルゴリズムなどの多くの機械学習技術は、簡単な基本要素（または構成単位）から成る。基本要素は、

- ・ 内積（ベクトルまたは行列とベクトルとの間）
- ・ 外積（ベクトル間）
- ・ 線形代数（ベクトルまたは行列に対するスカラによる加算、減算、乗算）
- ・ ベクトルまたは行列に適用される非線形（ $\tanh$ 、シグモイド、閾値処理（ $\text{thresholding}$ ））
- ・ 行列の転置

## 【0068】

機械学習GPU実行方法は、ピクセルシェーダを使用してこれらの各演算を実施する。さらに、この方法によって、これらの各演算を共にトレーニングのために使用したり、実際の設定に使用したりできるようになる。この方法は、ニューラルネットワーク、または後述したものと同様の基本要素またはその単なる拡張から成る他の任意の学習アルゴリズムとともに使用することができることに留意されたい。

## 【0069】

ピクセルシェーダを使用したGPU計算

機械学習GPU実行方法は、1つまたは複数のピクセルシェーダを使用して上記のすべての演算を実施する。ピクセルシェーダは、グラフィックス処理パイプラインで使用される。ピクセルシェーダは、ピクセルレンダリングハードウェア上で実行される個々のプログラムである。特に、三角形をピクセルに変えるプロセスは、「ラスタ化」と呼ばれる。ハードウェアは、三角形を入力としてみなし、三角形がレンダリングされる前にロードすることができるプログラムで指定する各ピクセルをレンダリングするよう命令することができる。これらのプログラム可能な三角形レンダラを「ピクセルシェーダ」と呼ぶ。シェーダでのプログラムの命令は、それぞれ直接的なハードウェアの実施を有するため、アセンブリ言語に近い。例えばHigh Level Shader Language (HLSL)、Nvidia（登録商標）社のC Graphics (Cg)、DirectXなど、シェーダの競合する言語（およびハードウェア）がある。ピクセルシェーダによって導入された新しい柔軟性によって、表面の写実的なレンダリングだけではなく、GPUを汎用並列プロセッサにより近づけることができるようになる。

## 【0070】

機械学習GPU実行方法は、GPUを汎用並列プロセッサとして使用する。さらにこの方法は、ピクセルシェーダを使用して、機械学習技術で使用される様々な基本的な演算を実施する。並列専用GPUハードウェアの結果、これらのアルゴリズムのトレーニングおよび使用は、GPUでは、CPUに比べて1桁速く稼働する。

## 【0071】

DirectXのDirect3D構成要素では、頂点シェーダおよびピクセルシェーダと呼ばれる2つの要素がある。DirectXは、GPUのハードウェア加速機能をアプリケーションに利用させるMicrosoft（登録商標）社によって開発された1組

10

20

30

40

50

のアプリケーションプログラムインターフェース (API) である。現世代のピクセルシェーダは、プログラム可能性が高い。両方のタイプのシェーダは、三角形 (グラフィックオブジェクトの基本単位) の出力装置へのレンダリングと関係がある。頂点シェーダは、三角形の頂点の空間変換および動画化などのタスクに使用することができる (したがってこのように名付けられている)。単一の三角形がラスタ化されると、ピクセルシェーダを使用して個々のピクセルに陰影を付ける (または個々のピクセルのカラー値を計算する)。

#### 【0072】

ピクセルシェーダは、DirectXで定義された限られたハードウェア依存型言語であるDirectXシェーダアセンブリ言語で一連の命令として表現される。シェーダでのコードは、レンダリングされる三角形のピクセルごとに一度実行され、その唯一の影響は、そのピクセルの4ベクトルの値を設定することである。シェーダ言語の制限、および副作用の欠如は、GPUが任意の順序で、そのハードウェアがサポートできるだけの並列処理を使用して自由にピクセルをレンダリングでき、その結果非常に高いパフォーマンスが得られることを意味する。ピクセルが4ベクトルであることから、さらに別の種類の並列処理を行うことができ、ピクセルシェーダの各実行によって、4つの要素 (ベクトルの4つの隣接する要素など) を同時に計算することができる。

10

#### 【0073】

アセンブリ言語プログラミングで使用する機能の多くは、ピクセルシェーダ内で使用することができる。これらには、定数、レジスタ、加算、減算、乗算、逆数、小さい1組の超越関数などがある。しかし、他のよく知られているループング、ブランチングなどの構成体は、一般には使用できない。これは、シェーダ当たりの命令の数が制限されており (一般的なGPUで96個の命令)、シェーダはレンダリングされるピクセル以外のメモリを変更することができないからである。これらの制限は、一部のアルゴリズムはGPUによる処理に適しておらず、仮に実施できたとしてもCPUに比べてGPUでの稼働の方が遅いことを意味する。特定の機械学習アルゴリズムをGPUで実行することができるかどうかを評価するために、アルゴリズムを構成する個々の演算のそれぞれを検査する必要がある。

20

#### 【0074】

単位演算 (unit operation) の実施

30

機械学習GPU実行方法は、1つまたは複数のピクセルシェーダを使用して黒丸付きにした (bulleted) 基本的な演算のそれぞれを実施する。第1の問題は、シェーダに結果を計算させることである。結果とはGPUメモリ内の矩形である。言い換えれば、結果は浮動小数点値 (4ベクトル) の2次元行列でなければならない。しかし、ピクセルシェーダは、所与の三角形のすべてのピクセルをレンダリングする。この問題を克服するために、機械学習GPU実行方法は、三角形で矩形のビューポートを指定する。図8は、機械学習GPU実行方法によって使用される矩形メモリ技術を示している。図8は特に、レンダリングされるピクセルを含んでいる三角形800、および三角形800内の矩形ビューポイント810を示している。機械学習GPU実行方法によれば、GPUは、ビューポイント810と三角形800との共通部分のピクセルを計算するだけでよい。したがって所望の矩形領域を計算するために、領域がビューポイント810として指定され、三角形800で囲まれる。これがレンダリングターゲットとなる。

40

#### 【0075】

行列の代わりにベクトルをレンダリングする必要がある場合、同じ技術が適用される。唯一の違いは、わずか1ピクセルの高さのビューポイントの矩形が指定されることである。同様に、スカラ値をレンダリングするために、1x1ピクセルのビューポイントを使用することができる。

#### 【0076】

この技術およびピクセルシェーダ内で使用可能な汎用プログラミング機能では、ベクトルおよび配列についての  $x = F(x)$  のような単項演算を実施することができる。他のオ

50

ペラントから読み取る必要がある演算では、GPUメモリ構成を検査する必要がある。特に、DirectXでは、メモリの矩形領域を作業領域として割り当てることができる。図9は、機械学習GPU実行方法のGPUメモリ構成を示している。図9に示すように、作業領域900は、GPUメモリ内に生成されている。現在のハードウェアでは、作業領域は、2048×2048ピクセル（各4値）の正方形領域の最小割り振りから生成することができる。このサブ矩形910から、シェーダ演算は、 $t_1$  920、 $t_2$  930、 $t_3$  940など他のサブ矩形からオペランドをレンダリングし、フェッチすることができる。これらのフェッチは、オペランドリード（operand read）950を介して達成される。

#### 【0077】

このメモリ構成を容易に行えるようにするDirectXの機構がテクスチャマッピングである。テクスチャマッピングは、レンダリングする面上に画像を置く必要性から生じる。例えば、木星、シマウマ、レンガ壁などのテクスチャの表面をレンダリングするとき、その面に描かれたパターン、光線の当たり方、角度、反射などを考慮する必要がある。この演算では、矩形ルックアップテーブルへのインデックスがラスタ化された三角形のターゲットエリアにわたって双線形に挿入される。（レンダリングされるエリア以外の）作業領域の矩形エリアをテクスチャとみなすことによって、こうした領域を宛先矩形にマッピングすることができる。したがって、宛先矩形内の任意の $x$ 、 $y$ 位置のピクセルを計算すると、テクスチャ矩形の形状的に対応するピクセルの値へのアクセスがある。例えば、サイズ $a \times b$ ピクセルの宛先矩形がレンダリングされている場合、同じく $a \times b$ ピクセルの別の領域をテクスチャマッピングすることができる。これによって、宛先の $i$ 番目のピクセルに対応する $a_i$ および $b_i$ のソース矩形値へのシェーダコード内での直接アクセスが得られる。この技術の簡単な適用によって、任意の行列、ベクトル、またはスカラー値を、作業領域内の同じサイズの行列、ベクトル、またはスカラーなどにコピーできるようになる。あるいは、 $x$ ピクセル×1ピクセルのテクスチャ領域を宛先矩形にマッピングすることができる。これは、レンダリングされる宛先ピクセルの $y$ 座標ではなく $x$ 座標に応じてその値が決まるルックアップテーブルへのアクセスを提供する。

#### 【0078】

テクスチャの有用性は、シェーダ内のレジスタ値に対する演算を使用することによって拡張することができる。レジスタは、シェーダが所与のピクセルをレンダリングするために使用できるローカル変数である。それらの値は、ピクセル間では共有することはできない（これは並列処理の前提を破ることになる）が、中間結果として（ローカルに）使用することはできる。例えば、その値を新しい位置にコピーしながら配列またはベクトルを転置することができる。その左、右、上、および下の座標が $l$ 、 $r$ 、 $t$ 、および $b$ のソース矩形を仮定する。次いで、その座標が $t$ 、 $b$ 、 $l$ 、および $r$ のテクスチャ矩形を指定する。ピクセルシェーダ内で $x$ および $y$ のテクスチャ座標が交換された後で、それらを使用して値をソースからフェッチし、宛先にコピーする。レンダリングの最後で、宛先はソースの転置を含む。

#### 【0079】

図10は、テクスチャ三角形を使用したテクスチャマッピングの使用を示している。図10に示すように、宛先1000の現在のピクセルのレンダリングは、指定されたソースベクトル1020のテクスチャ矩形1010からテクスチャ座標（15, 7）を取り出す。テクスチャ値をフェッチする前に、転置される実際の真のソースベクトル1040内に配置されている位置（7, 15）を有するテクスチャ矩形1030からの値が実際に読み取られるように、行および列の座標が逆にされる。

#### 【0080】

DirectXによって、テクスチャマッピングを使用して、複数のソース矩形を現在の宛先にマップすることができる。現在のハードウェアでは、各パスで少なくとも8つのこうしたマッピングを使用することができる。複数のソースでは、（ベクトルA - ベクトルB ->ベクトルC）などの演算を実施することができる。Cでの各ピクセルで、Aおよ

10

20

30

40

50

びBからテクスチャマッピングされた値がフェッチされ、レジスタ値に対して基本的な計算が行われ、結果が格納される。

#### 【0081】

2つのベクトルから行列へのテクスチャマッピングも、外積(ベクトルA×ベクトルB→C)を実施する方法を提供する。ベクトルAを1ピクセル幅の行ベクトル、およびBを1行高さの列ベクトルであると仮定する。これらの縮退矩形(degenerate rectangle)を行列Cの矩形にテクスチャマッピングする。次いでCのピクセルx、yをレンダリングすると、テクスチャのサンプリングによって、Aのy番目の要素およびBのx番目の要素が得られる。これらは単に、乗算し、格納するのに必要な値である。

10

#### 【0082】

##### 4つの成分の使用

内積の実施を説明する前に、GPU作業領域が4つの成分を有しているという点から各単位演算について説明する。各ピクセルがx、y、z、およびwの値から成るためにこれらの成分が生じる。ラベルxおよびyは、ここでは、上記の説明でピクセル座標を参照するxおよびyと混同しないものとする。4つの成分は従来、4次元物体空間に座標を格納するために使用されている。機械学習GPU実行方法は、GPUプログラミングモデルおよびハードウェアのこの特徴を活用して、より速く計算を行う。成分のうちの3つを無視して、すべての計算を例えばx平面で行うこともできるが、その結果得られるプログラミングの単純化は、パフォーマンスの面で代償が高くなる。

20

#### 【0083】

成分を利用するために、ピクセル平面の数学的配列またはベクトル内の位置へのマッピングが定義される。ベクトルの場合、最も簡単なマッピングは、

ピクセル0：x→要素0

ピクセル0：y→要素1

ピクセル0：z→要素2

ピクセル0：w→要素3

ピクセル1：x→要素4

などである。

#### 【0084】

これを行列に拡張するため、行列の各行(または列)がベクトルであることが観察される。上記のマッピングは、各行(または列)に適用される。4つの成分を行数(row dimension)にまとめる(collapse)か、列数(column dimension)にまとめるかは選択であり、この選択は、プログラミングを単純化する方法で行列ごとに個々に行うことができる。

30

#### 【0085】

行列またはベクトルの要素からピクセルおよび成分へのマッピングが与えられている場合、コピー演算はまったく影響を受けないことがわかる。シェーダ命令texldおよびmovは、多くの他のものと同様、一度に1つのピクセルに対して作用するため、それぞれ4つの値をテクスチャピクセルからレジスタに、またレジスタから宛先ピクセルに移動させる。

40

#### 【0086】

ベクトルに対する転置演算も変わらない。成分は常にベクトルの寸法にまとめられる。行列の場合、コードは変わらないが、まとめる方向がデータとともに転置されることに留意されたい。

#### 【0087】

数学的演算を適用することもできる。多くのシェーダ命令を適切な構文とともに4ベクトルの値、または単一の値に使用することができる。例えば、指数r1.x、r0.xは、レジスタ0のx平面をフェッチし、それをべき乗し、その結果をr1のx平面に格納する。

50

【0088】

外積

2つのベクトルの外積を、4つすべての成分の使用で達成することができるが、別の新しい技術を導入する必要がある。この技術は、機械学習GPU実行方法によって使用され、インデクサテクスチャ技術(indexer texture technique)と呼ばれる。一般にこれは、上述したように要素にマップされたベクトルの4つの成分のうちただ1つの成分の値を選択する方法である。

【0089】

一例として、ベクトルAがサイズaの列ベクトルであるとする。ベクトルBはサイズbの行ベクトルである。外積C、すなわちa行高さおよびb列幅の行列を計算することが望ましい。Cは、その4つの成分を列数すなわちyにまとめることである。言い換えれば、ピクセルに関しては、Cのメモリ矩形はb列幅であり、しかしa/4行高さである(各ピクセル行は行列の4つの行を格納するため)。aは4の倍数である必要がある。

10

【0090】

Cの各ピクセルを計算するシェーダルーチンが必要である。テクスチャマッピングは、Aについては単純であるが、Bについては単純ではない。Cのi番目のピクセルを計算するとき、Bのi番目の値(ピクセルではない)へのアクセスが必要である。例えば、Cのピクセル0, 0の4つの成分では、以下の値を計算する必要がある。

$$C_{0,0} = A_0 * B_0$$

$$C_{1,0} = A_1 * B_0$$

$$C_{2,0} = A_2 * B_0$$

$$C_{3,0} = A_3 * B_0$$

20

【0091】

この計算は、1つのGPU mul(乗算)命令で行うことができる。これには、(r2にA0~A3を、r3のw成分にB0を格納するなど)何らかの方法で値B0をレジスタのw平面に置くことが必要である。したがって次のようになる。

```
mul r1, r2, r3.wwww
```

【0092】

もう1ステップ戻って、そのピクセル内の隣接するB1~B3からB0を選び出す「ビットマスク」があった場合、dp4(ドット積)命令を使用してr3のw成分に値を抽出することができることになる。r5が(1, 0, 0, 0)を含み、r4が(Bのテクスチャマッピングによってロードされた)B0~B3を含むと仮定する。次いでこの命令は、(r4.x\*1+r4.y\*0+r4.z\*0+r4.w\*0)を計算する。これはr4.xに等しいが、B0:

30

```
dp4 r3.w, r4, r5
```

となる。

【0093】

内積は、行列Cの列0(0, 1, 0, 0)がレンダリングされる時、列1(0, 0, 1, 0)がレンダリングされる時、列2がレンダリングされる時など、値(1, 0, 0, 0)が使用可能な場合に計算することができる。これがインデクサテクスチャ技術の目的である。インデクサテクスチャ技術は、4ピクセル幅および1ピクセル高さで、次の値に初期設定される非常に小さいテクスチャを使用する。

40

【0094】

【表1】

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

## 【0095】

これら4つのピクセル値は、上述した技術を使用して個々の値をベクトルBから抽出するのに必要なすべての「ビットマスク」から成る。残るのは、インデクサテクスチャ技術の正しいピクセルをシェーダの各呼出で使用可能にする方法を見つけることである。このタスクを達成するために、機械学習GPU実行方法は、テクスチャマッピングのさらに別のアプリケーションを使用する。

## 【0096】

この時点まで、テクスチャピクセルと宛先ピクセルとの間の1対1マッピングを保つ方法でテクスチャマッピングが使用されてきた。これは、同じ数のピクセルを含むテクスチャ座標矩形および宛先矩形を定義することによって行われた。インデクサテクスチャ技術では、次のようにインデクサテクスチャ自体のサイズ(定義上1.0×1.0テクスチャ座標単位)の倍数であるテクスチャ座標矩形が定義される。行列Cの幅がwピクセルの場合、その幅がw/4(必ず整数)であるインデクサテクスチャのテクスチャ座標矩形が指定される。Cの幅にわたる各ピクセルでシェーダが呼び出されるので、インデクサテクスチャ座標は0からw/4に及ぶ。言い換えれば、Cのピクセルを4つ通過するごとに、インデクサ座標がある整数値から次の整数値に変わる。Cの4ピクセルごとにテクスチャ座標の端数部分を考える場合、これは0から1に及ぶ。Cの各ピクセルで正しい「ビットマスク」を得るために、インデクサテクスチャをサンプリングするのに必要な値がまさにこれらの値である。

`frc r3, t3` //インデクサテクスチャ座標の端数部分をレジスタ3にロードする

`texld r4, r3, s1` //レジスタ3のテクスチャ座標を使用してインデクサテクスチャをサンプリングする

上記の技術を使用して、2つのベクトルの外積を計算するピクセルシェーダが作成される。

## 【0097】

## 内積

内積とは、(行列A\*ベクトルB->ベクトルC)と表すことができる演算である。内積は問題を提起する。というのは、いくつかの値にわたってループしている間、合計を累積することが必要となるからである。これは、隣接するピクセル間に通信がないことを前提とするシェーダアーキテクチャに反する。このため、内積は、単一のピクセルシェーダとして実施することができない。しかし、以下に示すように、一連のシェーダおよび一部の追加メモリを使用することによって内積を実施することができる。

## 【0098】

サイズa行×b列の行列Aに、サイズaの行ベクトルであるベクトルBを掛ける。Bは通常、列ベクトルとして表されるが、行ベクトルの使用によってGPUの計算が単純化することに留意されたい。ベクトルの転置済みのコピーは、機械学習GPU実行方法を使用して作成することができるため、これは障害ではないことを思い出されたい。その結果のCは、高さbの列ベクトルである。行列Aは、上記のように、ピクセル当たりその4つの成分がy(行)数にまとめられた状態で構成される。

## 【0099】

図11は、内積を示すブロック図である。図11の対応する式は

$$C_0 = A_{0,0} * B_0 + A_{0,1} * B_1 + A_{0,2} * B_2 + A_{0,3} * B_3$$

$$C_1 = A_{1,0} * B_0 + A_{1,1} * B_1 + A_{1,2} * B_2 + A_{1,3} * B_3$$

などである。

## 【0100】

レンダーターゲットCは1次元であることに留意されたい。これは、前の手法を無効にする。というのは、定義された任意のテクスチャマップはそれぞれ、Cの任意の所与のピクセルの1つの値しか有することができないからである。しかしCの各値は、B、およびAの1行内のすべての値に応じて決まる。したがって別の手法を使用して、シェーダのAお

10

20

30

40

50

よびBから複数の値にアクセスする必要がある。この手法は、まず、複数のテクスチャマップ（現在のハードウェアでは8個が妥当な数である）を使用して、複数のオペランドを効率的に読み取る。次に、シェーダ内でアドレス計算を使用して、追加のオペランドを読み取る。

#### 【0101】

これらの技術でさえ、単一のレンダリングパスで大きい内積を計算することはできない。これは、ピクセルシェーダがブランチングやルーピングを使用することができず、また、限られた数（現在の1回の実施では96）のアセンブリ命令しか含めないからである。従って各シェーダは、命令制限に到達するまでにある量の計算しか行うことができない。幸いなことに、内積は、一度に考えられるAの列数を制限することによって1組のサブ問題に分解することができる。これは、1組の列ベクトル部分積を生成する。次いで内積結果を含む単一の列ベクトルを取得するまでこれらの部分積を繰り返し減らすことができる。

10

#### 【0102】

分解の技術は次のとおりである。まず、Aの12列のサブ行列にBの12列のサブベクトルを掛けるシェーダが作成される。8つのテクスチャマップが使用可能であり、6つはAに、2つはBに割り振られる。Aの6つのマップはすべて同じサイズのソース矩形を、しかし0ピクセルから5ピクセルまでの6つの異なる横方向のオフセットで有している。これによって行列の最初の6列への直接アクセスが得られる（横方向では1ピクセルは1列に等しいことを思い出されたい）。Bの2つのマップも1ピクセルだけオフセットして

20

#### 【0103】

シェーダ手順は、GPUの一時（スクラッチ）メモリを使用する宛先矩形を用い、12列の部分積について実行される。Cの4行（1ピクセル）を一度に計算するシェーダ手順は、次のとおりである。

#### 【0104】

a) 4つのテクスチャマップが互いに水平方向に1ピクセルだけオフセットしている状態で `texld`（テクスチャのロード）命令を4回使用して、Aの最初の4つのピクセルをレジスタにロードする。各ピクセルは4つの行値を含む。別の `texld` を使用してBの1ピクセル（4列）をロードする。

30

#### 【0105】

b) `mul` を使用して最初の4つの積をレジスタ0に格納する。これは、レンダリングされた第1のピクセルに関して  $A_{0,0} * B_0$ 、 $A_{1,0} * B_0$ 、 $A_{2,0} * B_0$ 、および  $A_{3,0} * B_0$  を計算する。表記 `xxxx` は、4つすべての積についてレジスタ1のx成分（つまり  $B_0$ ）を使用することを意味する。

`mul r0, r1, xxxx, r2`

#### 【0106】

c) `mad`（乗算/加算）を使用して次の4つの積をレジスタ0に蓄積する。これは、（`r3`を介して）Aの第2の水平方向のピクセルにアクセスし、前の結果に積  $A_{0,1} * B_1$ 、 $A_{1,1} * B_1$ 、 $A_{2,1} * B_1$ 、および  $A_{3,1} * B_1$  を追加する。B<sub>1</sub>へのアクセスに `yyyy` を使用する。

40

`mad r0, r1, yyyy, r3, r0`

#### 【0107】

d) `mad` を同じように2回以上使用して次の8つの積を累積する。

`mad r0, r1, zzzz, r4, r0`

`mad r0, r1, wwww, r5, r0`

#### 【0108】

e) ここで、5番目の列から8番目の列（列番号4～7）の準備をする。列4および5はすでに、Aに割り振られた6つのうちの残りの2つのテクスチャマップによってアクセ

50

ス可能である。列6および7については、列6の座標をとり、定数 $c_0$ を2回追加する。これは、1ピクセル（または作業領域幅の $1/2048$ ）に等しくなるように設定されている。列6および7のこれらの座標は、追加のレジスタに格納される。次いで4つの`texld`命令を使用して、これらの値、つまり $A_{0,4}$ から $A_{3,7}$ をレジスタに格納する。

【0109】

f) 1ピクセルだけオフセットされているBに割り振られた第2のテクスチャマップを使用して、 $B_4$ から $B_7$ （1ピクセルの幅）の値をロードする。

【0110】

g) 4つの`mad`命令を使用して、ステップ(b)から(d)で行われたように、 $r_0$ に格納されている現在の4つの合計に16以上の積を累積する。

【0111】

h) 次に、部分積の最後の4つの列について準備がなされる。この時点で、すべてのテクスチャマップが使用されている。したがってまだレジスタ内にある列7のアドレスが取得され、 $C_0$ が連続的にこの値に4回追加される。この結果がレジスタに格納される。同じ手順が行われて、 $C_0$ がBの第2のピクセルのアドレスを含むレジスタに追加される。これは、 $B_8$ から $B_{11}$ にアクセスするように1ピクセルだけそれを進めるステップを含む。正しい座標がレジスタに存在すると、`texld`を再度使用して $A_{0,8}$ から $A_{3,11}$ および $B_8$ から $B_{11}$ の値を他のレジスタにロードする。

【0112】

i) 4つ以上の`mad`命令を使用して、最後の16個の積を $r_0$ に累積する。ここで $r_0$ は、12列の部分内積の要素 $C_0$ から $C_3$ の値を含む。次いでシェーダルーチンが完了する。より多くの列がこのシェーダで処理されない理由は、シェーダごとの命令の制限に到達しているからである。

【0113】

上記は、12列部分内積を計算して一時メモリ内の列ベクトルにする方法である。この方法は、次のように、12列より広い行列に拡張することができる。まず、残りの幅が12以上である限り、12列部分積が処理される。次に、第1の部分結果に直接隣接するこれらのパスの結果が一時メモリの連続ブロックに配置される。最後に、残りの8つまたは4つ（行列幅は4の倍数でなければならないことを思い出されたい）の列がある場合、8つまたは4つの列を収容するように書き込まれた変更されたシェーダが使用される。これらの技術は、12列シェーダの単純化である。

【0114】

この結果、部分結果を表す1つまたは複数の隣接する列ベクトルのブロックが得られる。これらを「低減する」、つまり合計して、最後の内積を含む（また「スクラッチ」メモリの代わりに、他のルーチンによってアクセスできる明確な場所に配置される）単一の列ベクトルにする必要がある。効率上、（レジスタにおけるアドレス計算とは対照的に）テクスチャマップを使用して、低減する必要のある部分結果にアクセスする。8つのテクスチャマップをソースとして使用することによって、8つもの部分結果を一度に低減することができる。このシェーダコードは、非常に単純であり、`texld`命令および`add`命令から成る。プロセスを単純化するための1つの方法は、最終結果ではなく、低減された結果の位置の選択によるものである。つまり、低減された結果を次の列の既存の部分結果の右に入れることができる。このように、まさに最後の低減ステップまで、低減される1組の部分結果は常に一続きの矩形である。これらの技術および十分な一時メモリを使用して、任意のサイズの内積を計算することができる。

【0115】

上記では、機械学習GPU実行方法が、ニューラルネット（および同じ演算で構成される他の任意の機械学習アルゴリズム）のトレーニングに必要なすべての演算を含んでいることを示してきた。さらに、これらの演算は、GPUでピクセルシェーダとして実施することができ、いずれの場合にもアルゴリズムは高度の並列処理に役立つ。

10

20

30

40

50

【 0 1 1 6 】

V I I . 実施例

本明細書に開示した機械学習GPU実行システムおよび方法を完全に理解するために、実施例の動作上の詳細を提示する。この実施例は、機械学習GPU実行システムおよび方法を実施できる方法を1つだけ示していることに留意されたい。

【 0 1 1 7 】

この実施例では、機械学習GPU実行システムおよび方法は、単一のプロセッサ2.8GHz Pentium(登録商標)4、ATI(登録商標)Radeon9800グラフィックスカード付きで稼働させた。9800グラフィックスカードは、94箇所の命令メモリ、および256MBの作業領域を有しており、その16MBは、読み取り/書き込みの作業領域に使用し、残りを読み取り専用データに使用した。CPUのみの参照実装は、Pentium(登録商標)4のSSE SIMD拡張を利用するようにすでに最適化されているため、SIMD対SIMDでの比較を行う。さらにGPUは、テストされた例についてほぼ1桁高速である。

10

【 0 1 1 8 】

この実施例は、手書き文字認識を行うために使用される機械学習GPU実行システムおよび方法を示す。タスクは、手書きの数字のピクセル画像のクラス(「0」から「9」)を見つけるためのものであった。トレーニングデータベースは、正しいラベルを有する60,000のトレーニング画像で構成されていた。このトレーニングセットは、時としてMNISTと呼ばれることもあり、Web上で入手可能であり、機械学習アルゴリズムの資料ではよく知られている。このタスクを解決するために、ニューラルネットワーク手法が使用された。特に、ニューラルネットワークは2層完全接続ニューラルネットワークであった。

20

【 0 1 1 9 】

逆伝搬による2層完全接続ニューラルネットワークのトレーニングが次のパスに分解された。

【 0 1 2 0 】

【表2】

順伝搬	$W_1 I = H \xrightarrow{F_1(x)} H$
順伝搬	$W_2 H = O \xrightarrow{F_2(x)} O$
誤差計算	$\Delta_2 = \alpha(\text{target} - O) * D_2(O)$
重みの更新	$W_{2+} = (\Delta_2 \otimes H)$
誤算の伝搬	$\Delta_1 = (W_2^T \Delta_2) * D_1(H)$
重みの更新	$W_{1+} = (\Delta_1 \otimes I)$

30

【 0 1 2 1 】

ここで順伝搬は、行列 $W_1 \times$ ベクトル $I$ 、次いで関数マップ( $F_1(x)$ )。この場合 $F_1(x) = \tanh$ によりベクトル $H$ を算出する。このニューラルネットワークには2層あるため、 $W_2$ 、 $H$ 、および $F_2(H)$ (この場合 $F_2(x) = \text{シグモイド}$ )について同じ演算が繰り返されて、 $O$ が得られる。同じようにしてより多くの層を実装することができる。 $W_1$ および $W_2$ は、ニューラルネットのレベルごとの重みの行列、 $I$ は入力ベクトル、 $H$ は隠れ層ベクトル、および $O$ は出力ベクトルである。また、閾値も一般に隠れユニット(unit)および出力ユニットのそれぞれに追加される。閾値は、その値が常に1であるユニットを入力および隠れ層に追加することによってエミュレートすることができる。一定ユニットを他のユニットに接続している重みは、実質上こうしたユニットの閾値である。層 $I$ および $H$ が一定ユニットによって増大した場合、上記の式は、隠れ層および出力層ごとに閾値を正しく実施する。閾値パラメータは、 $W_1$ および $W_2$ に埋め込まれる。誤差計算は、ターゲットベクトルと、 $F(x)$ 関数の導関数 $D(x)$ を掛けたニューラルネットワークからの出力ベクトルとの差である。 $F(x)$ および対応する $D(x)$ の

40

50

式は次のとおりである。

【 0 1 2 2 】

【表 3】

	関数	導関数
シグモイド	$F(x) = \frac{1}{1 + e^{-x}}$	$D(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = F(x) * (1 - F(x))$
Tanh	$F(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$D(x) = \frac{4}{(e^x + e^{-x})^2} = (1 + F(x)) * (1 - F(x))$

10

【 0 1 2 3 】

シェーダ言語は、成分ごとにベクトルに適用でき、またはシェーダの観点からいえば、ピクセルごとに適用することができる指数関数を有することに留意されたい。

【 0 1 2 4 】

重みの更新は、その層の誤差ベクトルおよび入力ベクトルの外積を実行することによって実行された。結果として得られた行列は、学習係数 ( learning coefficient ) によってスケールされ、既存の重み行列に追加された。誤差をネットワークに伝搬するために、別の行列ベクトル乗算が計算された。つまり、誤差ベクトル付きの重み行列の転置が計算され、再度  $D(x)$  関数によってスケールされた。

20

【 0 1 2 5 】

まとめると、機械学習 GPU 実行方法を使用して GPU で次の演算が実施された。

- 1 . 行列 \* ベクトル - > ベクトル ( 内積 )
- 2 . ベクトルの要素ごとに  $x = f(x)$  ( この場合  $f$  は  $\tanh$  またはロジスティック関数、またはそれらの導関数 )
- 3 . ベクトル A - ベクトル B - > ベクトル C
- 4 . ベクトル \* ベクトル - > 行列 ( 外積 )
- 5 . 行列 A + 定数 \* 行列 B - > 行列 A
- 6 . 行列の置換
- 7 . GPU メモリ内の新しい場所へのコピー ( 行列またはベクトル )
- 8 . CPU メモリから GPU メモリへの、またはその逆のコピー ( 行列またはベクトル ) ( これは Direct X によって直接サポートされ、完全を期すためのみに言及される )

30

【 0 1 2 6 】

これらの演算のそれぞれは 1 回のパスで行われ、時として上述した制限によって単一のパスで計算できないときには、複数のパスで行われる。GPU での実際の実行は、パスのリストが前処理で作成されることを必要とする。このリストは、手動または自動で作成することができる。さらに、個々の関数がコンパイルされ、GPU にダウンロードされ ( これもまた前処理で )、入力値およびトレーニング値をグラフィックスメモリの作業領域にダウンロードし、三角形およびビューポイントをダウンロードして、関数をパスごとに指定する必要がある。こうした各ステップは、Direct 3D ( D3D ) グラフィックス API への呼出を介して実施された。三角形およびビューポイントのダウンロードは、パスの実行を開始することに留意されたい。上述したように、GPU でのクリップされた三角形のレンダリングは、黙示的なピクセルごとの DO ループを構成する。パスごとの特定のフラグメントシェーディングプロセス ( fragment shading process ) が各多角形の特性としてコード化された。反復型のトレーニングでは、反復ごとにこのシーケンスが繰り返される。プログラマが進捗を監視したい場合を除いて、反復を通じてグラフィックスメモリから CPU にデータを読み戻す必要はない。グラフィックスメモリからホストメモリへの転送は、現在のハードウェアでは遅く、プログラマは一般にこうした操作を避ける。

40

50

## 【 0 1 2 7 】

本発明の上記の説明は、例示および説明のために提示したものである。網羅的ではなく、また本発明を開示した正確な形状に限定するためのものでもない。上記の教示に鑑みて、多くの変更および変形が可能である。本発明の範囲は、本発明のこの詳細な説明によって限定されるものではなく、本明細書に添付した特許請求の範囲によって限定されるものである。

## 【 図面の簡単な説明 】

## 【 0 1 2 8 】

【 図 1 】 機械学習 GPU 実行エンジンおよび方法を実装できる適したコンピューティングシステム環境の例を示す図である。

10

【 図 2 】 本発明の一実施形態に係る機械学習 GPU 実行エンジンおよび方法の実装例を示すブロック図である。

【 図 3 】 一例として勾配降下アルゴリズムを使用した 3 層学習機械のトレーニングモジュールを示すブロック / フロー図である。

【 図 4 】 図 3 で示した 3 層学習機械のテストモジュールを示すブロック / フロー図である。

【 図 5 】 機械学習 GPU 実行方法のトレーニング段階の動作の概要を示すフロー図である。

【 図 6 】 機械学習 GPU 実行方法のテスト段階の動作の概要を示すフロー図である。

【 図 7 】 ニューラルネットワークでの 1 対の層を示す詳細図である。

20

【 図 8 】 機械学習 GPU 実行方法によって使用される矩形メモリ技術を示す図である。

【 図 9 】 機械学習 GPU 実行方法の GPU メモリ構成を示す図である。

【 図 1 0 】 テクスチャ三角形 ( texture triangle ) を使用したテクスチャマッピングの使用を示す図である。

【 図 1 1 】 内積を示すブロック図である。

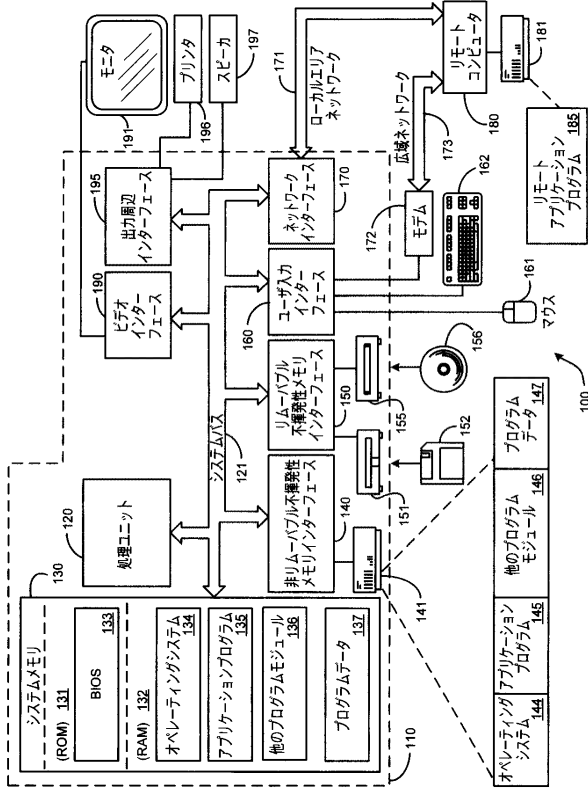
## 【 符号の説明 】

## 【 0 1 2 9 】

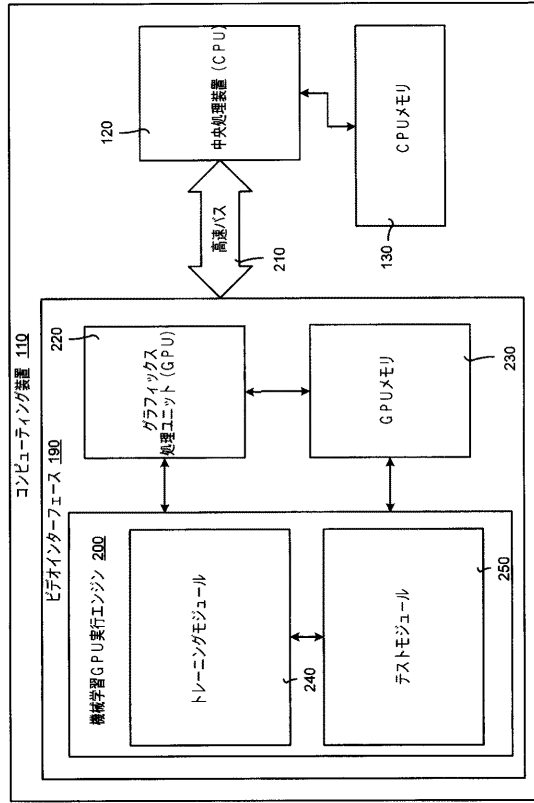
- 1 1 0 コンピューティング装置
- 1 2 0 中央処理装置 ( CPU )
- 1 3 0 CPU メモリ
- 1 9 0 ビデオインターフェース
- 2 0 0 機械学習 GPU 実行エンジン
- 2 1 0 高速バス
- 2 2 0 グラフィックス処理ユニット ( GPU )
- 2 3 0 GPU メモリ
- 2 4 0 トレーニングモジュール
- 2 5 0 テストモジュール

30

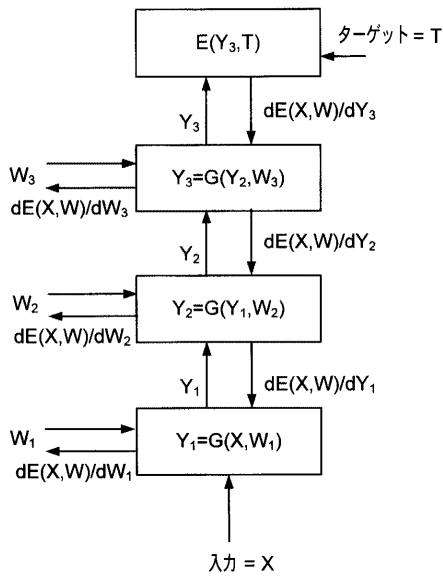
【図1】



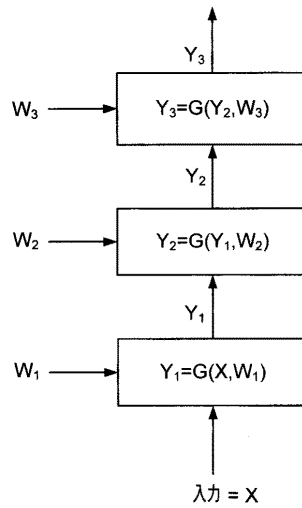
【図2】



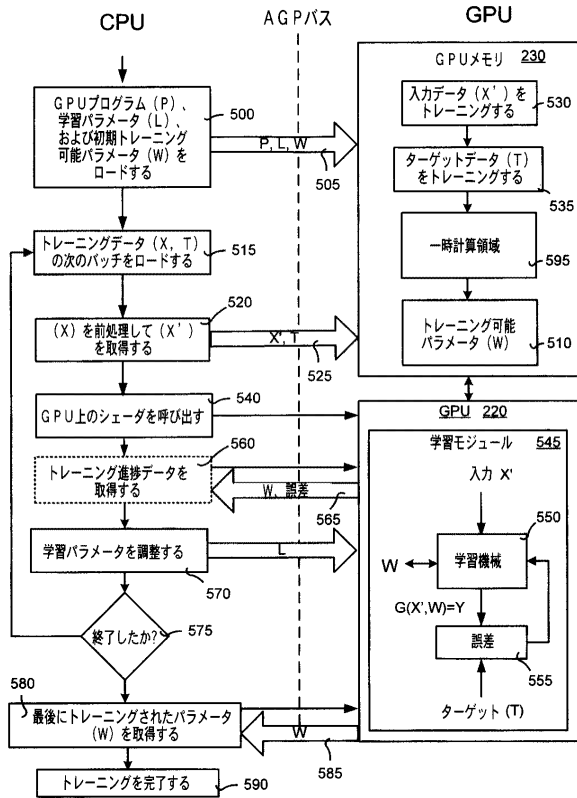
【図3】



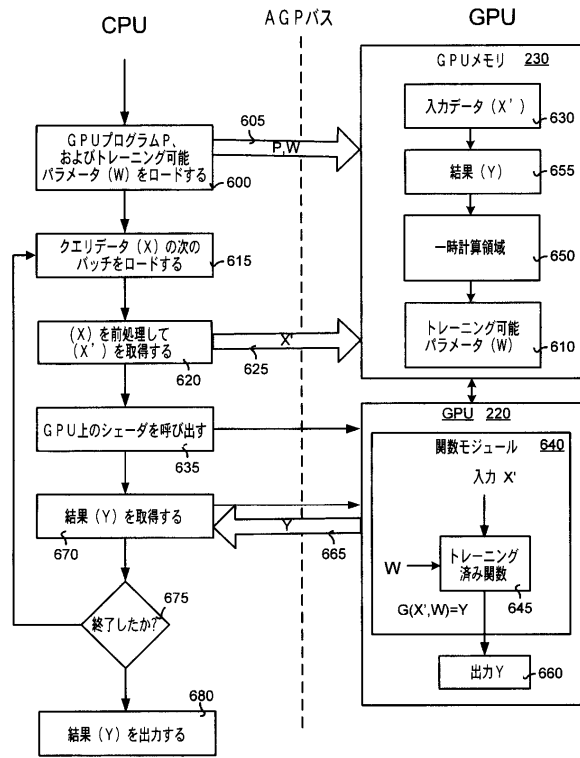
【図4】



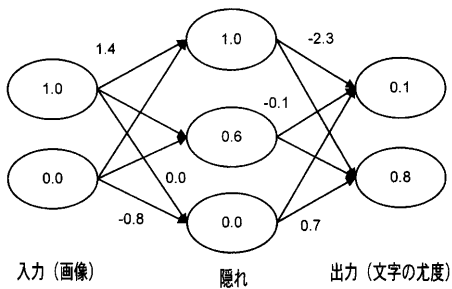
【図5】



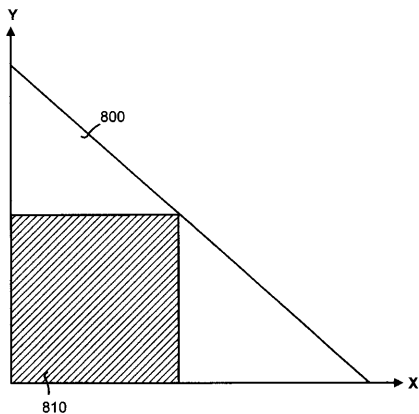
【図6】



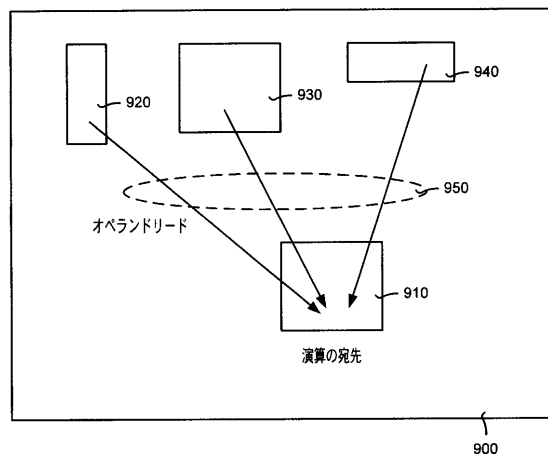
【図7】



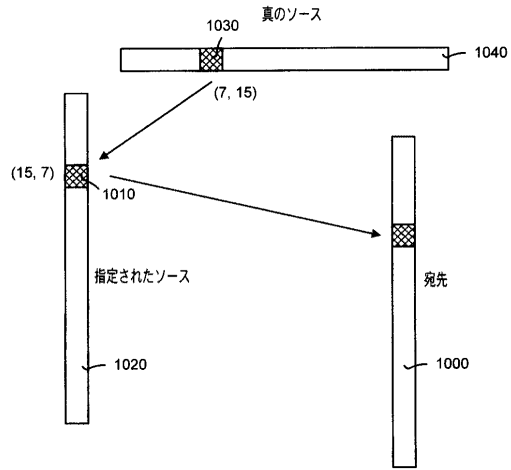
【図8】



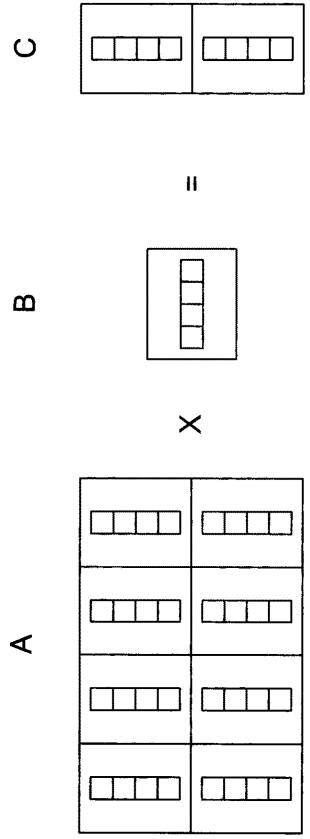
【図9】



【図10】



【図11】



## フロントページの続き

- (72)発明者 イアン エー . バック  
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ  
イクロソフト コーポレーション内
- (72)発明者 パトリス ワイ . シマード  
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ  
イクロソフト コーポレーション内

審査官 長谷川 篤男

- (56)参考文献 Optical Character Recognition on Graphics Hardware , integrative paper , UNCCH, Departme  
nt of Computer Science. , 2002年11月, 第1 - 9頁, URL , [http://www.cs.unc.edu/~  
adyilie/IP/Final.pdf](http://www.cs.unc.edu/~adyilie/IP/Final.pdf)  
Christian-A. Bohn , Kohonen Feature Mapping through Graphics Hardware , Proceedings of I  
nt. Conf. on Comp. Intelligence and Neurosciences , 1998年, 第1-4頁, URL , [http://  
www.fh-wedel.de/~bo/db/pubs/bohn\\_KFMGPU1998.pdf](http://www.fh-wedel.de/~bo/db/pubs/bohn_KFMGPU1998.pdf)

## (58)調査した分野(Int.Cl. , DB名)

G06N 3/00 - 3/063  
JSTPlus (JDreamII)  
IEEE Xplore