

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6077018号
(P6077018)

(45) 発行日 平成29年2月8日 (2017.2.8)

(24) 登録日 平成29年1月20日 (2017.1.20)

(51) Int. Cl.	F I
G 0 6 F 9/45 (2006.01)	G 0 6 F 9/44 3 2 2 E
G 0 6 F 9/38 (2006.01)	G 0 6 F 9/38 3 7 0 C

請求項の数 20 (全 33 頁)

(21) 出願番号	特願2014-558964 (P2014-558964)	(73) 特許権者	595020643
(86) (22) 出願日	平成25年2月27日 (2013.2.27)		クォアルコム・インコーポレイテッド
(65) 公表番号	特表2015-513737 (P2015-513737A)		Q U A L C O M M I N C O R P O R A T E D
(43) 公表日	平成27年5月14日 (2015.5.14)		アメリカ合衆国、カリフォルニア州 9 2
(86) 国際出願番号	PCT/US2013/028029		1 2 1 - 1 7 1 4、サン・ディエゴ、モア
(87) 国際公開番号	W02013/130614		ハウス・ドライブ 5 7 7 5
(87) 国際公開日	平成25年9月6日 (2013.9.6)	(74) 代理人	100108855
審査請求日	平成28年2月4日 (2016.2.4)		弁理士 蔵田 昌俊
(31) 優先権主張番号	61/603,771	(74) 代理人	100109830
(32) 優先日	平成24年2月27日 (2012.2.27)		弁理士 福原 淑弘
(33) 優先権主張国	米国 (US)	(74) 代理人	100103034
(31) 優先権主張番号	13/777,663		弁理士 野河 信久
(32) 優先日	平成25年2月26日 (2013.2.26)	(74) 代理人	100075672
(33) 優先権主張国	米国 (US)		弁理士 峰 隆司
早期審査対象出願		最終頁に続く	

(54) 【発明の名称】 異種CPU-GPU計算のための実行モデル

(57) 【特許請求の範囲】

【請求項 1】

異種計算のための方法であって、

プロセッサを用いて、プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプライントポロジを受信することと、

前記プロセッサを用いて、前記実行モデルの前記パイプライントポロジがグラフィックス処理ユニット (GPU) に実装されるプラットフォームに依存した方法を指示する命令を生成することであって、前記実行モデルの前記パイプライントポロジが前記GPUに実装される前記プラットフォームに依存する方法は、前記GPUのプラットフォームに基づき、前記パイプライントポロジは、消費するカーネルによって消費されるデータを生成する、生成するカーネルを識別し、前記命令は、拡大係数に基づき、前記生成するカーネルによって生成される最大のデータ量は、前記生成するカーネルが受信するデータ量および前記拡大係数に基づく、生成することと、

前記プロセッサを用いて、前記GPUに前記命令を送信することと、
を備える、方法。

【請求項 2】

命令を生成することは、

前記命令を生成するために前記実行モデルの前記パイプライントポロジを指定するコマンドリストをコンパイルすることを備える請求項 1 に記載の方法。

【請求項 3】

前記コマンドリストをコンパイルすることは、

前記実行モデルの前記パイプラインポロジが前記GPUに実装される前記プラットフォームに依存する方法を定義する前記命令を生成するために前記GPUの構成情報に少なくとも基づいて前記コマンドリストをコンパイルすることを備える請求項2に記載の方法。

【請求項4】

前記GPUの前記構成情報は、

前記GPU内のプログラマブル計算ユニット数、及び

前記GPU内のデータレーン数

のうちの1つ以上を備える請求項3に記載の方法。

10

【請求項5】

前記コマンドリストをコンパイルすることは、

前記実行モデルの前記パイプラインポロジにおいて識別されたバッファのサイズ、及び

前記実行モデルの前記パイプラインポロジにおいて識別されたカーネルの重要性のうちの1つ以上に少なくとも基づいて前記コマンドリストをコンパイルすることを備える請求項2に記載の方法。

【請求項6】

前記パイプラインポロジを受信することは、

前記パイプラインポロジを形成するために相互に接続された1つ以上のカーネル及び1つ以上のバッファを示すコマンドリストを受信することを備える請求項1に記載の方法。

20

【請求項7】

装置であって、

グラフィックス処理ユニット(GPU)と、

プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプラインポロジのインディケーションを受信することと、

前記実行モデルの前記パイプラインポロジが前記GPUに実装されるプラットフォームに依存した方法を指示する命令を生成することであって、前記実行モデルの前記パイプラインポロジが前記GPUに実装される前記プラットフォームに依存した方法は、前記GPUのプラットフォームに基づき、前記パイプラインポロジは、消費するカーネルによって消費されるデータを生成する、生成するカーネルを識別し、前記命令は、拡大係数に基づき、前記生成するカーネルによって生成される最大のデータ量は、前記生成するカーネルが受信するデータ量および前記拡大係数に基づく、生成することと、

30

前記GPUに前記命令を送信することと

を行うように構成されたプロセッサと、

を備える、装置。

【請求項8】

前記命令を生成するために、前記プロセッサは、

前記実行モデルの前記パイプラインポロジを指定するコマンドリストをコンパイルように構成される請求項7に記載の装置。

40

【請求項9】

前記コマンドリストをコンパイルするために、前記プロセッサは、

前記実行モデルの前記パイプラインポロジが前記GPUに実装される前記プラットフォームに依存した方法を定義する前記命令を生成するために前記GPUの構成情報に少なくとも基づいて前記コマンドリストをコンパイルするように構成される請求項8に記載の装置。

【請求項10】

前記GPUの前記構成情報は、

前記GPU内のプログラマブル計算ユニット数、及び

50

前記GPU内のデータレーン数

のうちの1つ以上を備える請求項9に記載の装置。

【請求項11】

前記コマンドリストをコンパイルするために、前記プロセッサは、

前記実行モデルの前記パイプラインポロジーにおいて識別されたバッファのサイズ、及び

前記実行モデルの前記パイプラインポロジーにおいて識別されたカーネルの重要性のうちの1つ以上に少なくとも基づいて前記コマンドリストをコンパイルするように構成される請求項8に記載の装置。

【請求項12】

前記パイプラインポロジーを受信するために、前記プロセッサは、

前記パイプラインポロジーを形成するために相互に接続された1つ以上のカーネル及び1つ以上のバッファを示すコマンドリストを受信するように構成される請求項7に記載の装置。

【請求項13】

前記装置は、

メディアプレーヤー、

セットトップボックス、

無線ハンドセット、

デスクトップコンピュータ

ラップトップコンピュータ

ゲームコンソール

ビデオ会議装置、及び

タブレットコンピューティングデバイス、

のうちの1つを備える請求項7に記載の装置。

【請求項14】

非一時的なコンピュータによって読み取り可能な記憶媒体であって、

1つ以上のプロセッサによって実行されたときに、

プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプラインポロジーを受信することと、

前記実行モデルの前記パイプラインポロジーがグラフィックス処理ユニット(GPU)に実装されるプラットフォームに依存した方法を指示する命令を生成することであって、前記実行モデルの前記パイプラインポロジーが前記GPUに実装される前記プラットフォームに依存した方法は、前記GPUのプラットフォームに基づき、前記パイプラインポロジーは、消費するカーネルによって消費されるデータを生成する、生成するカーネルを識別し、前記命令は、拡大係数に基づき、前記生成するカーネルによって生成される最大のデータ量は、前記生成するカーネルが受信するデータ量および前記拡大係数に基づく、生成することと、

前記GPUに前記命令を送信することと

を前記1つ以上のプロセッサに行わせる命令を格納した、非一時的なコンピュータによって読み取り可能な記憶媒体。

【請求項15】

命令を生成することを前記1つ以上のプロセッサに行わせる前記命令は、

前記命令を生成するために前記実行モデルの前記パイプラインポロジーを指定するコマンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる命令を備える請求項14に記載の非一時的なコンピュータによって読み取り可能な記憶媒体。

【請求項16】

前記コマンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる前記命令は、

前記実行モデルの前記パイプラインポロジーが前記GPUに実装される前記プラット

10

20

30

40

50

フォームに依存した方法を定義する前記命令を生成するために前記GPUの構成情報に少なくとも基づいて前記コマンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる命令を備える請求項15に記載の非一時的なコンピュータによって読み取り可能な記憶媒体。

【請求項17】

前記GPUの前記構成情報は、

前記GPU内のプログラマブル計算ユニット数、及び

前記GPU内のデータレーン数

のうちの1つ以上を備える請求項16に記載の非一時的なコンピュータによって読み取り可能な記憶媒体。

10

【請求項18】

前記コマンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる前記命令は、

前記実行モデルの前記パイプラインポロジーにおいて識別されたバッファのサイズ、及び

前記実行モデルの前記パイプラインポロジーにおいて識別されたカーネルの重要性のうちの1つ以上に少なくとも基づいて前記コマンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる命令を備える請求項15に記載の非一時的なコンピュータによって読み取り可能な記憶媒体。

【請求項19】

装置であって、

グラフィックス処理ユニット(GPU)と、

プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプラインポロジーを受信するための手段と、

前記実行モデルの前記パイプラインポロジーが前記GPUに実装されるプラットフォームに依存した方法を指示する命令を生成するための手段であって、前記実行モデルの前記パイプラインポロジーが前記GPUで実装される前記プラットフォームに依存した方法は、前記GPUのプラットフォームに基づき、前記パイプラインポロジーは、消費するカーネルによって消費されるデータを生成する、生成するカーネルを識別し、前記命令は、拡大係数に基づき、前記生成するカーネルによって生成される最大のデータ量は、前記生成するカーネルが受信するデータ量および前記拡大係数に基づく、生成するための手段と、

20

30

前記GPUに前記命令を送信するための手段と、

を備えるプロセッサと、

を備える、装置。

【請求項20】

命令を前記生成するための手段は、

前記命令を生成するために前記実行モデルの前記パイプラインポロジーを指定するコマンドリストをコンパイルするための手段を備える請求項19に記載の装置。

【発明の詳細な説明】

40

【技術分野】

【0001】

本出願は、ここにおける引用によってその内容全体が組み入れられている米国仮特許出願第61/603,771号(出願日:2012年2月27日)の利益を主張するものである。

【0002】

本開示は、アプリケーションの実行に関するものである。本開示は、より具体的には、様々な処理ユニットにおけるアプリケーションの実行に関するものである。

【背景技術】

【0003】

50

グラフィックス処理ユニット（GPU）は、グラフィックス処理に加えての目的のために使用されている。例えば、グラフィックスに関連しないアプリケーションは、GPUの大規模な並列性を利用することによって上昇した速度で実行することができる。この結果、追加のグラフィックスに関連しない処理機能を提供し、汎用GPU（GGPU）と呼ばれるGPUが得られている。例えば、GGPUは、1つ以上のシェーダコア（shader core）を含み、シェーダコアは、グラフィックスに関連しないアプリケーションと同様に、グラフィックス関連のアプリケーションも実行するように構成される。

【発明の概要】

【0004】

概して、本開示は、計算パイプラインを実装するための実行モデルを生成するための技法に関するものである。例えば、多くのデータ処理アルゴリズムは、計算パイプラインとして表すことができ、パイプラインの1つのユニットがデータを受信及び処理し、さらなる処理のためにパイプラインの他のユニットのために処理されたデータを出力する。本開示において説明される技法は、様々なタイプの並列コンピューティングデバイスにおいて効率的に実行することができるような形で計算パイプラインを表すことを考慮することができ、その一例がグラフィックス処理ユニット（GPU）である。例えば、それらの技法は、プラットフォームから独立した方法（platform-independent manner）で（例えば、計算パイプラインを実装することになるコンピューティングデバイスに依存しないで）計算パイプラインを表すことができる。

【0005】

計算パイプラインを表すことで、それらの技法は、利用可能なコンピューティングデバイス、例えば、利用可能なGPU、さらには中央処理装置（CPU）、専用の命令を生成するためにプラットフォームに依存するコンパイルを利用することができる。例えば、計算パイプラインは、プラットフォームから独立した方法で定義することができ、コンパイラは、計算パイプラインが実装されるべきプラットフォーム専用の命令を生成する。

【0006】

一例では、本開示は、異種計算（heterogeneous computing）のための方法について説明する。その方法は、プロセッサを用いて、プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプライントポロジを受信することを含む。その方法は、プロセッサを用いて、実行モデルのパイプライントポロジがグラフィックス処理ユニット（GPU）に実装されるプラットフォームに依存した方法を指示する命令を生成することを含む。この例では、実行モデルのパイプライントポロジがGPUに実装されるプラットフォームに依存した方法は、GPUのプラットフォームに基づく。その方法は、プロセッサを用いて、GPUに命令を送信することを含む。

【0007】

一例では、本開示は、装置について説明する。その装置は、グラフィックス処理ユニット（GPU）と、プロセッサと、を含む。プロセッサは、プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプライントポロジのインディケーション（indication）を受信するように構成される。プロセッサは、実行モデルのパイプライントポロジがGPUに実装されるプラットフォームに依存した方法を指示する命令を生成するようにも構成される。この例では、実行モデルのパイプライントポロジがGPUに実装されるプラットフォームに依存した方法は、GPUのプラットフォームに基づく。プロセッサは、GPUに命令を送信するようにも構成される。

【0008】

一例では、本開示は、1つ以上のプロセッサによって実行されたときに、プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプライントポロジを受信することを1つ以上のプロセッサに行わせる命令が格納されているコンピュータによって読み取り可能な記憶媒体について説明する。命令は、実行モデルのパイプライントポロジがグラフィックス処理ユニット（GPU）に実装されるプラットフォームに依存した方法を指示する命令を生成することも1つ以上のプロセッサに行わせる。こ

の例では、実行モデルのパイプラインポロジがGPUに実装されるプラットフォームに依存した方法は、GPUのプラットフォームに基づく。命令は、GPUに命令を送信することも1つ以上のプロセッサに行わせる。

【0009】

一例では、本開示は、装置について説明する。その装置は、グラフィックス処理ユニット(GPU)と、プロセッサと、を含む。プロセッサは、プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプラインポロジを受信するための手段を含む。プロセッサは、実行モデルのパイプラインポロジがGPUに実装されるプラットフォームに依存した方法を指示する命令を生成するための手段も含む。この例では、実行モデルのパイプラインポロジがGPUに実装されるプラットフォームに依存した方法は、GPUのプラットフォームに基づく。プロセッサは、GPUに命令を送信するための手段も含む。

10

【0010】

1つ以上の例の詳細が添付された図面及び以下の説明において示される。これらの説明及び図面から、及び請求項からその他の特徴、目的、及び利点が明確になるであろう。

【図面の簡単な説明】

【0011】

【図1】実行モデルの例を示した概念図である。

【図2】本開示において説明される1つ以上の例によるデバイスの例を示したブロック図である。

20

【図3】本開示において説明される1つ以上の例による技法例を示したフローチャートである。

【図4】図2のデバイスをより詳細に例示したブロック図である。

【発明を実施するための形態】

【0012】

グラフィックス処理ユニット(GPU)は、データを並列で素早くかつ効率的に処理するように構成することができる。開発者は、GPU上で実行するアプリケーションの形でデータ処理アルゴリズムを開発することができる。例えば、GPUは、1つ以上のアプリケーションを実行するように構成されるシェダプロセッサを含むことができる。これらのアプリケーションの例は、シェダプログラム、例えば、バーテックスシェダ(vertex shader)(頂点シェダ)、ハルシェダ(hull shader)、フラグメントシェダ(fragment shader)、幾何シェダ(geometry shader)、及びグラフィックス処理に関連するその他の該アプリケーション、を含む。

30

【0013】

さらに、幾つかのアプリケーション開発者は、GPUの大規模な並列性を利用し、グラフィックスに関連しないアプリケーションをGPUで実行するのが有益であるとみなすであろう。例えば、GPUによって提供される処理の並列性は、並列行列演算がグラフィックス処理に関連していないときでさえも、それらの行列演算を実行するのに適することができる。グラフィックスに関連しないアプリケーションのその他の例は、並列演算の素早い実行が有益であることができる流体力学又は線形代数に関連する技法を含む。

40

【0014】

該グラフィックスに関連しないアプリケーションを実行することが可能なGPUは、汎用GPU(GPGPU)であるとみなすことができる。例えば、GPUがグラフィックスに関連しないアプリケーションを実行中であるときには、GPUは、GPGPUとして機能している。ほとんどのGPUは、GPGPUとして機能するように構成することができる。

【0015】

例示を目的として、本開示は、GPGPUとして機能しているGPUに関して技法を説明する。しかしながら、それらの技法は、GPUがGPGPUとして機能している(すな

50

わち、グラフィックスに関連しないアプリケーションを実行している）事例には限定されず、それらの技法は、GPUがグラフィックス関連のアプリケーションを実行している事例にも適用することができる。さらに、本開示において説明される技法は、あらゆるタイプの処理ユニット、例えば、中央処理装置（CPU）、アクセラレータ、又はその他のカスタムデバイスによって実装することができる。それらの技法は、GPUに関して説明されが、それらの技法は、その他のタイプの処理ユニットにも拡張可能であることが理解されるべきである。

【0016】

GPU内のシェーダプロセッサは、複数のシェーダコア（これらのコアはグラフィックス関連及びグラフィックスに関連しない、の両方のアプリケーションに関する命令を実行できることを示すためにプログラマブル計算ユニットとも呼ばれる）を含むことができる。プログラマブル計算ユニットの各々は、そのプログラマブル計算ユニットによって実行される命令に関して予約されているローカルメモリ、及びそれらの命令の実行によって生成されたデータ、例えば、命令の実行中に生成される即座の結果、を含むことができる。プログラマブル計算ユニットのローカルメモリは、その他のプログラマブル計算ユニットによってはアクセス不能であることができる。幾つかの例では、GPUで実行されるべき異なるアプリケーションは、異なるプログラマブル計算ユニットによって実行することができる。

【0017】

本開示において説明される技法では、グラフィックス関連のアプリケーションはシェーダと呼ばれ、グラフィックスに関連しないアプリケーションは、カーネルと呼ばれる。例えば、シェーダ（すなわち、グラフィックス関連のアプリケーション）の例は、パーティクルシェーダと、フラグメントシェーダと、幾何シェーダと、含み、ただしこれらに限定されない。カーネル（すなわち、グラフィックスに関連しないアプリケーション）の例は、行列演算、流体力学、画像処理動作、映像処理動作、等を行うためのアプリケーションを含む。

【0018】

さらに、カーネルは、GPUによって実行されるアプリケーションのみに必ずしも限定する必要はなく、GPUの固定機能ユニット（fixed-function unit）（すなわち、プログラミング不能なユニット）も含む。例示のみを目的として、本開示において説明される技法は、GPUで実行されるアプリケーションであるカーネルに関して説明される。例えば、それらの技法は、GPUがGPGPUとして機能するためにGPUのシェーダプロセッサで実行するグラフィックスに関連しないアプリケーションに関して説明される。

【0019】

カーネルは、複数のワークグループ、タスク、又はスレッドを含むことができる（これらはすべて、本開示では同義語として用いられる）。例えば、スレッドは、カーネルのその他のスレッドから独立して実行することができるカーネルの命令の組であることができる。幾つかの例では、カーネルを実行するためには、プログラマブル計算ユニットのうちの1つ以上がカーネルの1つ以上のスレッドを各々実行することができる。例えば、第1のプログラマブル計算ユニットは、カーネルの第1のスレッドを実行することができ、第2のプログラマブル計算ユニットは、同じカーネルの第2のスレッドを実行することができる。幾つかの例では、1つのプログラマブル計算ユニットが1つのカーネルの1つ以上のスレッドを実行することができ、他のプログラマブル計算ユニットは、他のカーネルの1つ以上のスレッドを実行する。幾つかの例では、2つの組み合わせが可能である（すなわち、幾つかのプログラマブル計算ユニットが同じカーネルの異なるスレッドを実行中であり、他方、幾つかのその他のプログラマブル計算ユニットが異なるカーネルのスレッドを実行中である）。

【0020】

GPUは、大規模な並列性を処理のために提供する一方で、開発者、例えば、カーネル

10

20

30

40

50

の開発者は、は、様々なタイプのGPUにおいてパイプライン方式で効率的に実行するカーネルを開発するのは困難であるとみなすであろう。パイプライン方式でカーネルを実行することは、1つのカーネルによって生成されたデータが他のカーネルによって消費されるような形でカーネルを実行することを意味する。他の例として、パイプライン方式でカーネルを実行することは、同じカーネルの他のスレッドによって消費されるべきデータを生成するカーネルのスレッドを実行することを意味する。この開示では、データを生成するスレッドは、生成するスレッド(producer thread)と呼ぶことができ、データを受信するスレッドは、消費するスレッド(consumer thread)と呼ぶことができる。

【0021】

10

幾つかの例では、生成するスレッド及び消費するスレッドは、同じカーネルのスレッドであることができる。幾つかの例では、生成するスレッド及び消費するスレッドは、異なるカーネルのスレッドであることができる。これらの例では、生成するスレッドを含むカーネルは生成するカーネルと呼ぶことができ、消費するスレッドを含むカーネルは消費するカーネルと呼ぶことができる。

【0022】

一例として、データ処理アルゴリズム、例えば、画像処理又は映像処理、を実装するために、開発者は、複数のカーネルを開発することができ、各カーネルは、全体的アルゴリズムの一部を実装する。第1のカーネルは、第2のカーネルによる消費のために処理されるべきデータ(例えば、グラフィックスに関連しないデータ)を受信し、データを処理し、及びデータを出力することができる。第2のカーネルは、第3のカーネルによる消費のために第1のカーネルによって出力されたデータを受信し、データをさらに処理し、データを出力し、以下同様である。

20

【0023】

この例では、第1、第2、及び第3のカーネルがパイプラインを形成すると考えることができ、第1のカーネル(例えば、生成するカーネル)は、第2のカーネル(例えば、第1のカーネルの観点からの消費するカーネル)によって消費されるべきデータを生成する。第2のカーネルは、第3のカーネルによって消費されるべきデータを生成する。この例では、第2のカーネルは、第3のカーネルの観点からの生成するカーネルであり、第3のカーネルは、消費するカーネルである。この方法により、GPUは、パイプライン方式で第1、第2、及び第3のカーネルを実行することができる。

30

【0024】

幾つかの例では、パイプライン方式でカーネルを実行することは、カーネルを逐次に行う(例えば、第1のカーネルを実行し、第2のカーネルを実行することによって後続され、第3のカーネルを実行することによって後続され、以下同様である)ことを意味することができる。しかしながら、本開示において説明される技法は、そのようには限定されない。幾つかの例では、パイプライン方式でカーネルを実行することは、並列して(同時に又は時間が重なる状態で)カーネルを実行することを意味することができる。例えば、GPUは、第2のカーネルが第1のカーネルに関する消費するカーネルであり、第3のカーネルが第2のカーネルに関する消費するカーネルである場合でも第1、第2、及び第3のカーネルのうちの2つ以上を同時に実行することができる。

40

【0025】

開発者は、データ処理アルゴリズムを実装するためにパイプライン方式で実行するカーネルを開発することができるが、開発者は、様々なタイプのGPUにまたがってカーネルの最適な実行を保証することはできないであろう。例えば、開発者は、プロセッサで実行する命令を書くことができる。これらの命令は、いつカーネルを実行すべきかをGPUに命令することをプロセッサに行わせることができる。上述されるように、カーネルは、1つ以上の計算ユニットで実行することができる。しかしながら、開発者は、特定のCPIで利用可能な計算ユニットの数、より一般的には、特定のGPUの並列処理能力、を知らないことがある。

50

【 0 0 2 6 】

この場合は、開発者は、GPUの処理能力を開発者が知らないため、いつカーネルを実行すべきかを予め決定することができないことがある。この結果、開発者は、各々が異なるタイプのGPU専用の異なる命令を書くことになる。例えば、開発者は、第1のGPUタイプ専用の、プロセッサで実行する命令の第1の組を書くことができる。例えば、第1のGPUタイプが3つの計算ユニットを含む場合は、命令の第1の組は、3つの計算ユニットを有するGPUでカーネルを実行することができる方法を定義することができる。開発者は、第2のGPUタイプ専用の、プロセッサで実行する、命令の第2の組も書くことができる。例えば、第2のGPUタイプが4つの計算ユニットを含む場合は、命令の第2の組は、4つの計算ユニットを有するGPUでカーネルを実行することができる方法を定義することができる。

10

【 0 0 2 7 】

幾つかの例では、異なるGPUタイプに関する命令を書くのではなく、開発者は、1つのみのタイプのGPU（例えば、推定される最悪の事態のシナリオでのGPU）に関する命令を書くことができる。これらの例では、1つのタイプのみのGPUがデータ処理アルゴリズムを効率的に実装することができ、その他のGPUタイプは、データ処理アルゴリズムを効率的に実装することができない。

【 0 0 2 8 】

換言すると、開発者がGPUで効率的な方法で実行することをカーネルに行わせる命令を書くことができるプラットフォームから独立した方法は存在することができない。むしろ、開発者は、その他のGPUタイプで非効率的に実行する（例えば、GPUタイプに依存しない）一般的命令を書くことができる。他方、開発者は、結果的には非ポータブルな（non-portable）命令になるプラットフォームに依存する命令を書くことができる。例えば、開発者は、異なるGPUタイプの各々に関して別々の命令を書かなければならないことがあり、それは過度に厄介であろう。

20

【 0 0 2 9 】

本開示において説明される技法は、プラットフォームから独立した方法で（すなわち、異種の計算のために）データ処理アルゴリズムを実装するためにカーネルを効率的に実行することを考慮する。本開示において説明される技法では、異種計算は、プラットフォームから独立した方法での計算を意味する。より詳細に説明されるように、本開示において説明される技法により、開発者は、データ処理アルゴリズムを実装するためのカーネルに関するパイプライン実行モデルを指定する。

30

【 0 0 3 0 】

パイプライン実行モデルを指定するために、開発者は、パイプラインのトポロジーを定義することができる。パイプラインのトポロジーは、相互に接続されたカーネル及びバッファを含む実行グラフであるとみなすことができる。例えば、第1のカーネルが第2のカーネルによって消費されるべきデータを生成する場合。開発者は、第1のカーネルがバッファ、例えば、先入れ先出し（FIFO）バッファ、に結合され、バッファが第2のカーネルに結合されるような形で、トポロジーを定義することができる。この例では、実行グラフは、第1のカーネルがバッファにデータを出力し、第2のカーネルがバッファからデータを受信することを示すことができる。

40

【 0 0 3 1 】

トポロジーを定義することに加えて、開発者は、実行モデルの一部としてトポロジーの特徴を定義することもできる。一例として、開発者は、各カーネルに関する拡大係数（amplification factor）を定義することができる。拡大係数は、カーネルが受信する所定の量のデータに関してカーネルが生成する最大のデータ量を示すことができる。例えば、拡大係数がカーネルに関して5であり、カーネルが2つのデータパケットを受信する場合は、カーネルが生成する最大のデータ量は、10のデータパケットである。

【 0 0 3 2 】

50

他の例として、開発者は、バッファのサイズを定義することができる。例えば、開発者は、バッファの幅（例えば、バッファの記憶場所内に格納することができるデータの量）及びバッファの長さ（例えば、バッファ内の記憶場所の数）を定義することができる。

【0033】

この方法により、開発者は、データ処理アルゴリズムに関するプラットフォームから独立した実行モデルを定義することができる。例えば、開発者は、データ処理アルゴリズムが実装される特定のGPUを説明する必要がある。むしろ、各GPUに関する実行モデルは、同じであることができる。

【0034】

本開示において説明される技法は、開発者がバウンドされた方法で（bounded manner）実行モデルを定義するのを可能にすることができる。例えば、開発者は、いずれのカーネルが必要であるか、いずれのカーネルが生成するカーネルを形成するか及びいずれのカーネルが消費するカーネルを形成するかを完全に定義することができる。バウンドされた方法で実行モデルを定義することは、静的実行モデル（例えば、実行前に定義されるそれ）であるとみなすことができる。

【0035】

バウンドされた方法で実行モデルを定義することは、バウンドされない方法で（unbounded manner）実行モデルを定義することと比較して計算効率の利得を可能にすることができる。実行モデルのバウンドされない定義では、開発者は、実行前には、必要になるカーネル数、いずれのカーネルが生成するカーネルになるか及びいずれのカーネルが消費するカーネルになるかを定義することができない（すなわち、カーネル間の相互接続を定義することができない）。この結果、バウンドされた実行モデルと比較して、バウンドされない実行モデルの性能が最適でなくなる可能性がある。

【0036】

例えば、本開示において説明される技法では、プロセッサは、実行モデルを受信することができ、及び実行モデルをコンパイルしてGPUによって処理することができるオブジェクトコード（すなわち、バイナリコード）にすることができる。コンパイルステップは、プラットフォームに依存するステップであることができる。例えば、プロセッサは、データ処理アルゴリズムが実装されるGPUの処理能力を示す情報で予め構成することができる。一例として、プロセッサは、GPU内の計算ユニット数を示す情報で予め構成することができる。

【0037】

コンパイルステップでは、プロセッサは、メタスケジューラ（meta-scheduler）に関する命令を生成することができる。メタスケジューラは、GPUで実行するソフトウェアであることができ又はGPU内のハードウェアであることができる。メタスケジューラに関する命令は、実行モデルが実行される方法を定義することができる。この例では、実行モデルはバウンド（bound）することができ（例えば、カーネル数及びカーネルの相互接続が知られている）、プロセッサは、GPUの処理能力を示す情報で予め構成することができるため、コンパイラは、GPUが実行モデルを実装する方法を最適化するメタスケジューラに関する命令を定義することができる。バウンドされない実行モデルの場合は、カーネル数及びそれらの各々の相互接続は知ることができず、コンパイラは、GPUでの実行モデルの実行を適切に最適化することができない。

【0038】

図1は、実行モデルの例を示した概念図である。例えば、図1は、実行モデル10を示す。開発者は、データ処理アルゴリズムを実装するための実行モデル10を定義することができる。例えば、開発者は、画像処理、映像処理、線形代数演算、又は流体力学を計算するためのアルゴリズムを実装するために実行モデル10を定義することができる。概して、開発者は、グラフィックス処理ユニット（GPU）によって提供される大規模な並列計算効率を利用するデータ処理アルゴリズムを実装するために実行モデル10を定義することができ、グラフィックスに関連しない目的を含む。GPUがグラフィックスに関

10

20

30

40

50

連しないアルゴリズムを実装中である例では、GPUは、汎用GPU (GPGPU) のように機能しているとみなすことができる。

【 0 0 3 9 】

例示されるように、実行モデル 1 0 は、バッファ 1 2 A 乃至 1 2 D と、カーネル 1 4 A 乃至 1 4 C と、を含む。幾つかの例では、図 1 に示されるバッファ及びカーネルよりも多い又は少ないそれらが存在することができる。バッファ 1 2 A 乃至 1 2 D の例は、先入れ先出し (FIFO) バッファとリングバッファとを含み、ただしこれらに限定されない。

【 0 0 4 0 】

カーネル 1 4 A 乃至 1 4 C の例は、実行モデル 1 0 が実装するように定義される全体的データ処理アルゴリズムの少なくとも一部分を実装する、開発者によって開発されたアプリケーションを含む。開発者は、カーネル 1 4 A 乃至 1 4 C を開発するためにあらゆるプログラミング言語を利用することができる。

【 0 0 4 1 】

開発者が実行モデル 1 0 を定義することができる様々な方法が存在することができる。一例として、開発者は、コンピューティングデバイス、例えば、デスクトップコンピュータ又はラップトップコンピュータ、において実行モデルを定義することができる。開発者は、グラフィカルユーザインタフェース (GUI) を提示するコンピューティングデバイスでアプリケーションを実行することができる。開発者は、図 1 において例示される方法でバッファ 1 2 A 乃至 1 2 D 及びカーネル 1 4 A 乃至 1 4 C を相互接続するために GUI を利用することができる。さらに、開発者は、バッファ 1 2 A 乃至 1 2 D 及びカーネル 1 4 A 乃至 1 4 D の特徴を定義するために GUI を利用することができる。

【 0 0 4 2 】

他の例として、開発者は、特定のアプリケーション処理インタフェース (API) によりコマンドを利用して実行モデルを定義することができる。該 API の例は、Microsoft (登録商標) による DirectX (登録商標)、Khronos グループによる OpenGL (登録商標)、及び Khronos グループによる OpenCL (登録商標) を含む。しかしながら、本開示の態様は、DirectX、OpenGL 又は OpenCL API には限定されず、開発済みの、現在開発中の、又は将来開発予定のその他のタイプの API にまで拡張することができる。さらに、本開示において説明される技法は、API により機能することは要求されない。

【 0 0 4 3 】

例えば、コマンドは、開発者が実行モデルを定義中であることを示すコマンドを含むことができる。コマンドは、バッファ 1 2 A 乃至 1 2 D 及びカーネル 1 4 A 乃至 1 4 C が実行モデル 1 0 に属することを開発者が定義するのを可能にし、及びバッファ 1 2 A 乃至 1 2 D 及びカーネル 1 4 A 乃至 1 4 C が相互接続される方法を定義するコマンドも含むことができる。

【 0 0 4 4 】

いずれの例でも (すなわち、GUI に基づく又はコマンドに基づく)、開発者が実行モデル 1 0 を定義したコンピューティングデバイスは、実行モデル 1 0 を変換することができ、実行モデル 1 0 のトポロジを指定するコマンドリストを含む。例えば、例示されるように、カーネル 1 4 A は、バッファ 1 2 A からデータを受信し、データを処理し、バッファ 1 2 B 及び 1 2 C にデータを格納する。カーネル 1 4 B は、バッファ 1 2 B からデータを受信し、データを処理し、及びバッファ 1 2 D にデータを格納する。カーネル 1 4 C は、バッファ 1 2 D 及び 1 2 C からデータを受信し、データを処理する。

【 0 0 4 5 】

この方法により、バッファ 1 2 A 乃至 1 2 D 及びカーネル 1 4 A 乃至 1 4 C は、計算パイプラインとして構成される。例えば、カーネル 1 4 A は、カーネル 1 4 B 及び 1 4 C に関する生成するカーネルである。カーネル 1 4 B は、カーネル 1 4 A に関する消費するカーネル及びカーネル 1 4 C に関する生成するカーネルである。カーネル 1 4 C は、カーネル 1 4 A 及び 1 4 B の両方に関する消費するカーネルである。

【 0 0 4 6 】

理解を助けるために、図 1 は、実行モデル 1 0 のパイプライントポロジーを例示するとみなすことができる。例えば、開発者は、実行モデル 1 0 のパイプライントポロジーを定義する実行グラフを定義するとみなすことができる。この実行グラフでは、カーネル 1 4 A 乃至 1 4 C は、バッファ 1 2 A 乃至 1 2 D と相互接続されるノードであるとみなすことができる。

【 0 0 4 7 】

幾つかの例では、開発者は、異なる実行モデルを相互接続することもできる。例えば、データ処理アルゴリズムに関して 1 つの実行モデルを定義する代わりに、開発者は、複数の実行モデルを開発することができ、各実行モデルがデータ処理アルゴリズムの一部分を実装する。これらの例では、各々の実行モデル内のカーネルは、全体的データ処理アルゴリズムの一部分の小部分を実装することができる。開発者は、カーネル 1 4 A 乃至 1 4 C 及びバッファ 1 2 A 乃至 1 2 D を相互接続するのと同様の方法で実行モデルを相互接続することができる。例えば、開発者は、バッファ 1 2 A を他の実行モデルに相互接続すること及び / 又はカーネル 1 4 C を他の実行モデルに相互接続することができる。

【 0 0 4 8 】

複数の実行モデルを定義するのが有益であることができる。より詳細に説明されるように、プロセッサは、実行モデル、例えば、実行モデル 1 0、をコンパイルしてオブジェクトコードにし、その結果得られたオブジェクトコードを格納することができる。実行モデル 1 0 が複数の実行モデルのうちの 1 つである例では、プロセッサは、実行モデル 1 0 を再コンパイルする必要がない。換言すると、実行モデルは、全体的なデータ処理アルゴリズムに関するビルディングブロックであるとみなすことができ又はデータ処理アルゴリズムの全体を定義することができる。これで、共通して使用される実行モデルは、実行モデルが使用されるすべての事例に関して再コンパイルする必要がない。

【 0 0 4 9 】

実行モデル 1 0 のトポロジーを定義することに加えて、開発者は、バッファ 1 2 A 乃至 1 2 D 及びカーネル 1 4 A 乃至 1 4 C の特徴を定義することもできる。開発者は、GUI 又は上述されるコマンドに基づくフォーマットを用いて特徴を定義することができる。開発者は、バッファ 1 2 A 乃至 1 2 D 内の記憶場所の数（すなわち、バッファ 1 2 A 乃至 1 2 D の長さ）を定義することができる。開発者は、バッファ 1 2 A 乃至 1 2 D の各記憶場所内に格納することができるデータの量（すなわち、バッファ 1 2 A 乃至 1 2 D の幅）を定義することもできる。

【 0 0 5 0 】

幾つかの例では、開発者は、バッファ 1 2 A 乃至 1 2 D の次元を定義することができる。例えば、幾つかの画像処理技法、例えば、畳み込み、は、ピクセルのブロック（例えば、 7×7 のピクセルのブロック）において生じる。これらの例では、ピクセルをブロック形態で格納するためにバッファ 1 2 A 乃至 1 2 D は二次元バッファであることが有益であることができる。例えば、ピクセルのブロックが 7×7 のピクセルのブロックである場合は、バッファ 1 2 A 乃至 1 2 D のうちの 1 つ以上は、4 9 の記憶場所を有する線形バッファではなく、 7×7 の記憶場所を有する状態で（すなわち、二次元バッファとして）構成することができる。

【 0 0 5 1 】

カーネル 1 4 A 乃至 1 4 C に関しては、開発者は、一例として、拡大係数を定義することができる。拡大係数は、カーネルが消費する所定のデータ量に関してそのカーネルが生成する最大のデータ量を示すことができる。例えば、カーネル 1 4 B に関する拡大係数が 2 であり、カーネル 1 4 B が 5 つのデータパケットをバッファ 1 2 B から受信する場合は、カーネル 1 4 B が生成する最大のデータ量は 1 0 のデータパケットである。他の例として、カーネル 1 4 A 乃至 1 4 C のうちの 1 つ以上に関して、開発者は、カーネルが（例えば、受信されたデータ量から独立して）生成する最大データ量を定義することもできる。

【 0 0 5 2 】

さらに他の例として、開発者は、相対的重要性をカーネル 1 4 A 乃至 1 4 C に割り当てることができる。例えば、重要性は、カーネル 1 4 A 乃至 1 4 C のうちのいずれが割り込まれないで実行すべきかを示すことができ、カーネル 1 4 A 乃至 1 4 C のうちのより重要なそれらは割り込まれないで実行し、カーネル 1 4 A 乃至 1 4 C のうちの重要性がより低いそれらは割り込まれないで又は割り込まれた状態で実行することができる（すなわち、その他の実行のために解放するために実行が断続的に休止される）。

【 0 0 5 3 】

カーネル 1 4 A 乃至 1 4 C 及びバッファ 1 2 A 乃至 1 2 D の特徴は、例示する目的で説明されており、限定するものであるとはみなされるべきでない。開発者が上述されるすべての特徴例を定義する必要はない。例えば、開発者は、バッファ 1 2 A 乃至 1 2 D のサイズ（例えば、長さ及び幅）を定義することができ、カーネル 1 4 A 乃至 1 4 C の特徴は定義することができない。これらの例では、カーネル 1 4 A 乃至 1 4 C が生成又は消費するデータの量は、バッファ 1 2 A 乃至 1 2 D のサイズが既に定義されているため重要ではない。他の例として、開発者は、拡大係数又はカーネル 1 4 A 乃至 1 4 C が生成する最大データ量を定義することができ、バッファ 1 2 A 乃至 1 2 D の特徴は定義しない。これらの例では、実行モデル 1 0 をコンパイルするプロセッサは、拡大係数及び / 又はカーネル 1 4 A 乃至 1 4 C が生成する最大データ量に基づいてバッファ 1 2 A 乃至 1 2 D のサイズを決定することができる。さらに、これらの例では、実行モデル 1 0 をコンパイルするプロセッサは、バッファ 1 2 A 乃至 1 2 D が一次元（すなわち、線形）バッファ又は多次元バッファのいずれであるべきかを決定することができる。

【 0 0 5 4 】

概して、開発者又はプロセッサは、バッファ 1 2 A 乃至 1 2 D が大きくなりすぎないようにする一方で、“行き詰まった”状況を回避するためにカーネル 1 4 A 乃至 1 4 C 及びバッファ 1 2 A 乃至 1 2 D の特徴を決定することができる。行き詰まった状況は、消費するカーネルが、データを受信すべきバッファ内に格納されていないデータを期待したとき又は生成するカーネルが、消費するカーネルがデータを消費するよりも高速でデータを格納しているためバッファがデータでオーバーフローするときに発生する可能性がある。行き詰まった状況では、カーネルは、“停滞し”（hang）、カーネルが行き詰まった状況で停滞しないようにするための追加措置を講じないかぎりデータ処理アルゴリズムの実装を停止する。幾つかの例では、行き詰まりが発生したときに停滞を回避するために追加のタスクを実装するように GPU を構成するよりも、行き詰まりが発生しないようにカーネル 1 4 A 乃至 1 4 C 及びバッファ 1 2 A 乃至 1 2 D の特徴を定義したほうが良いであろう。

【 0 0 5 5 】

開発者又はプロセッサが行き詰まりを緩和するために相対的に大きいサイズのバッファ 1 2 A 乃至 1 2 D を定義することが可能である。しかしながら、バッファ 1 2 A 乃至 1 2 D のサイズが不必要に大きい場合は、プロセッサは、必要なメモリスペースよりもはるかに大きいメモリスペースをバッファ 1 2 A 乃至 1 2 D のために予約する可能性があり、その結果、非効率的なメモリの使用になってしまうおそれがある。

【 0 0 5 6 】

従って、バッファ 1 2 A 乃至 1 2 D 及び / 又はカーネル 1 4 A 1 4 D の特徴を定義することによって、開発者は、メモリが効率的に使用される一方で行き詰まりの機会が低下されるような形で実行モデル 1 0 を定義することができる。この方法により、開発者は、完全にバウンドされた静的な実行モデル 1 0 を定義することができる。例えば、開発者は、実行モデル 1 0 を実装するために必要なカーネル及びバッファの数が動的に（すなわち、実装中に）定義されるのではなく、実装前に、実行モデル 1 0 を実装するために必要なカーネル 1 4 A 乃至 1 4 C 及びバッファ 1 2 A 乃至 1 2 D の数、及びバッファ 1 2 A 乃至 1 2 D 及びカーネル 1 4 A 乃至 1 4 C の特徴を定義することができる。

【 0 0 5 7 】

開発者は、データ処理アルゴリズムを実装する GPU を含むデバイスに実行モデル 1 0

10

20

30

40

50

を格納することができる。例えば、開発者は、実行モデル 10 のパイプラインポロジを指定するコマンドリストを、GPUを含むデバイスに格納することができる。幾つかの例では、開発者が実行モデル 10 をデバイスに格納するのではなく、デバイスのユーザがデバイスでの格納のために実行モデル 10 をダウンロードすることができる。概して、データ処理アルゴリズムを実装する GPU を含むデバイスが実行モデル 10 を格納する方法は、本開示で説明される技法の制約にはならない。換言すると、データ処理アルゴリズムが実装されるべき GPU を含むデバイスに実行モデル 10 (例えば、実行モデル 10 のコマンドのリスト) を格納するためにあらゆる技法を利用することができる。例えば、開発者が実行モデル 10 を定義したコンピューティングデバイスがデータ処理アルゴリズムを実装すべき GPU を含む同じコンピューティングデバイスであることさえも可能である。

10

【0058】

より詳細に説明されるように、GPUを含むデバイスは、プロセッサを含むこともできる。プロセッサは、実行モデル 10 を受信して実行モデル 10 をコンパイルし、GPUが実行モデル 10 によって定義されたデータ処理アルゴリズムを実装するために実行すべきオブジェクトコードにする。本開示において説明される技法により、プロセッサは、GPUの処理能力を説明する実行モデル 10 をコンパイルすることができる。

【0059】

従って、開発者は、プラットフォームから独立した方法で(すなわち、データ処理アルゴリズムを実装する GPU のタイプを考慮せずに)実行モデル 10 のパイプラインポロジを定義することができる。GPUを含むデバイスのプロセッサは、GPUが実行モデル 10 を実装すべき方法を定義する実行モデル 10 に基づいて命令を生成することができる。例えば、プロセッサは、実行モデル 10 をコンパイルし、コンパイルの一部として命令を生成することができる。実行モデル 10 のコンパイル中には、プロセッサは、GPUの処理能力を説明することができる。この方法により、プロセッサは、GPUでの実装のために実行モデル 10 を最適化することができる。これは、開発者が、結果的に非常にポータブルな実行モデル(すなわち、異種計算のために異なるタイプの GPU で効率的に実装することができるモデル)が得られる柔軟で理解しやすい方法で実行モデル 10 を開発するのを可能にすることができる。開発者は、GPUのプラットフォーム専用の又は実装によって定義される挙動に関わる必要がない。

20

【0060】

さらに、バウンドされた方法で実行モデル 10 を定義することは、実行モデル 10 の組み込み式のデバッグを可能にすることができる。一例として、実行モデル 10 を定義することは、上述されるように、行き詰まりの機会を低減させることができる。さらに、バウンドされない実行モデルに関しては、開発者が第 2 のカーネルによって消費されるデータを出力するように第 1 のカーネルを定義し、実装中に第 1 のカーネルを実行させるように不注意に第 2 のカーネルを定義し、従って、第 1 のカーネルが第 2 のカーネルによって生成されたデータを消費するようにする可能性がある。これは、実際には、無限のループを作り出す。バウンドされない実行モデルの場合は、実装まで該状況を決定する方法がない。

30

【0061】

しかしながら、バウンドされた実行モデル 10 を用いることで、開発者は、該無限ループを作り出すのを簡単に回避することができる。例えば、開発者は、該無限ループを作り出したことを GUI で見ることができる。他の例として、アプリケーションが実行モデル 10 のパイプラインポロジを変換して実行モデル 10 のコマンドリスト内に入れたときには、アプリケーションは、該無限ループが存在するかどうかを決定することができる。

40

【0062】

図 2 は、本開示において説明される 1 つ以上の例によるデバイスの例を示したブロック図である。例えば、図 2 は、デバイス 16 を例示する。デバイス 16 の例は、ビデオデバイス、例えば、メディアプレーヤー、セットトップボックス、無線ハンドセット、例えば

50

、携帯電話、パーソナルデジタルアシスタント（PDA）、デスクトップコンピュータ、ラップトップコンピュータ、ゲームプレイコンソール、ビデオ会議装置、タブレットコンピューティングデバイス、等を含み、ただしこれらに限定されない。デバイス16は、図2において例示されるコンポーネントに加えてのそれらを含むことができる。

【0063】

例示されるように、デバイス16は、集積回路（IC）18と、グローバルメモリ20と、を含む。グローバルメモリ20は、デバイス16のためのメモリとみなすことができる。例えば、グローバルメモリ20は、IC18の外部に存在することができ及びIC18及びグローバルメモリ20は、システムバス36を介して通信することができる。グローバルメモリ20は、1つ以上のコンピュータによって読み取り可能な記憶媒体を備えることができる。グローバルメモリ20の例は、ランダムアクセスメモリ（RAM）、又は希望されるプログラムコードを搬送又は格納するために使用することができ及びコンピュータ又はプロセッサによってアクセスすることができるその他のあらゆる媒体を含み、ただしこれらに限定されない。

10

【0064】

幾つかの態様では、グローバルメモリ20は、本開示においてプロセッサ22及びGPU24に起因する機能を実行することをプロセッサ22及び/又はグラフィックス処理ユニット（GPU）24に行わせる命令を含むことができる。従って、グローバルメモリ20は、実行されたときに様々な機能を実行することを1つ以上のプロセッサ（例えば、プロセッサ22及びGPU24）に行わせる命令が格納されているコンピュータによって読み取り可能な記憶媒体であることができる。

20

【0065】

グローバルメモリ20は、幾つかの例では、非一時的な記憶媒体であるとみなすことができる。用語“非一時的な”は、記憶媒体が搬送波又は伝搬される信号において具現化されないことを示すことができる。しかしながら、用語“非一時的な”は、グローバルメモリ20が移動不能である又はその内容が静的であることを意味するとは解釈されるべきでない。一例として、グローバルメモリ20は、データ16から取り外して他のデバイスに移動させることができる。他の例として、グローバルメモリ20に実質的に類似するグローバルメモリは、デバイス16内に挿入することができる。幾つかの例では、非一時的な記憶媒体は、（例えば、RAMにおいて）経時で変化することが可能なデータを格納することができる。

30

【0066】

IC18は、プロセッサ22と、グラフィックス処理ユニット（GPU）24と、を含む。IC18は、追加のコンポーネント、例えば、グローバルメモリ20と通信するためのインタフェースユニット、グローバルメモリ20内のメモリを管理するためのユニット、及びその他の処理ユニット、例えば、ディスプレイプロセッサ、を含むことができる。IC18は、プロセッサ22及びGPU24を収納又は形成するあらゆるタイプの集積回路であることができる。例えば、IC18は、チップパッケージ内の処理チップであるとみなすことができる。

【0067】

プロセッサ22及びGPU24は、単一のIC18の一部として例示されているが、本開示の態様はそうには限定されない。幾つかの例では、プロセッサ22及びGPU24は、異なる集積回路（すなわち、異なるチップパッケージ）に収納することができる。

40

【0068】

プロセッサ22及びGPU24の例は、デジタル信号プロセッサ（DSP）、汎用マイクロプロセッサ、特定用途向け集積回路（ASIC）、フィールドプログラマブルロジックアレイ（FPGA）、又はその他の同等の集積回路又はディスクリート論理回路を含み、ただし、これらに限定されない。幾つかの例では、GPU24は、グラフィックス処理に適する大規模な並列処理能力をGPU24に提供する集積回路及び/又はディスクリート論理回路を含む専用ハードウェアであることができる。幾つかの例では、GPU24は

50

、汎用処理を含むこともでき、汎用の処理タスク（すなわち、グラフィックスに関連しないタスク）を実装するときには汎用GPU（GGPU）と呼ぶことができる。

【0069】

プロセッサ22は、ホストと時々呼ばれ、デバイス16の中央処理装置（CPU）であることができる。プロセッサ22は、様々なタイプのアプリケーションを実行することができる。アプリケーションの例は、ウェブブラウザ、電子リーダー、電子メールアプリケーション、スプレッドシート、ビデオゲーム、映像再生、音声再生、ワードプロセッシング、表示のためのビュー可能なオブジェクトを生成するその他のアプリケーション、又はその他のタイプのアプリケーションを含む。グローバルメモリ20は、1つ以上のアプリケーションの実行のための命令を格納することができる。

10

【0070】

幾つかの例では、プロセッサ22は、処理タスク、例えば、大規模な並列演算を要求するタスク、をGPU24に委託することができる。一例として、グラフィックス処理は、大規模な並列演算を要求し、プロセッサ22は、該グラフィックス処理タスクをGPU24に委託することができる。幾つかの例では、プロセッサ22は、グラフィックス処理に関連しないタスクをGPU24に委託することができる。例えば、データ処理アルゴリズム、例えば、行列演算、画像処理、及び映像処理、は、並列演算を要求し、GPU24は、プロセッサ22と比較して該演算を実装するのにより適している。

【0071】

タスクを実装するために、GPU24は、1つ以上のアプリケーションを実行するように構成することができる。例えば、グラフィックスに関連する処理の場合は、GPU24は、バーテックスシェーダ、フラグメントシェーダ、及び幾何シェーダ、等のアプリケーションを実行することができる。グラフィックスに関連しない処理に関しては、GPU24は、該処理（例えば、グラフィックス関連処理又はグラフィックスに関連しない処理）に関して設計されたアプリケーションを実行することができる。いずれの例（例えば、グラフィックス関連処理又はグラフィックスに関連しない処理）に関しても、プロセッサ22は、以下においてより詳細に説明されるように、1つ以上のアプリケーションを実行するようにGPU24に命令することができる。

20

【0072】

プロセッサ22は、特定のアプリケーション処理インタフェース（API）によりGPU24と通信することができる。例えば、プロセッサ22は、命令、例えば、APIを利用して1つ以上のアプリケーションを実行するようにGPU24に命令する命令、をGPU24に送信することができる。該APIの例は、Microsoft（登録商標）によるDirectX（登録商標）、KhronosグループによるOpenGL（登録商標）、及びKhronosグループによるOpenCL（登録商標）を含む。しかしながら、本開示の態様は、DirectX、OpenGL又はOpenCL APIには限定されず、開発済みの、現在開発中の、又は将来開発予定のその他のタイプのAPIにまで拡張することができる。さらに、本開示において説明される技法は、APIにより機能することは要求されず、プロセッサ22及びGPU24は、あらゆる通信技法を利用することができる。

30

40

【0073】

一例として、グラフィックス関連のアプリケーションに関しては、プロセッサ22は、OpenGL APIを用いてGPUと通信することができる。グラフィックスに関連しないアプリケーションに関しては、プロセッサ22は、OpenCL APIを用いてGPU24と通信することができる。繰り返すと、本開示において説明される技法は、プロセッサ22がOpenGL及び/又はOpenCL APIを用いてGPU24と通信することは必ずしも要求しない。

【0074】

GPU24が実行するグラフィックス関連のアプリケーションは、シェーダと呼ぶことができ、GPU24が実行するグラフィックスに関連しないアプリケーションは、カーネ

50

ルと呼ぶことができる。例えば、グローバルメモリ 20 は、シェーダ及びカーネルの命令を格納することができ、プロセッサ 14 で実行するコンパイラは、シェーダ及びカーネルの命令を GPU 16 での実行のためのオブジェクトコードに変換することができる。他の例として、グローバルメモリ 20 は、GPU 16 が取り出して実行するシェーダ及びカーネルのオブジェクトコードを格納することができる。

【0075】

図 2 において例示されるように、グローバルメモリ 20 は、実行モデル (図 1) を格納する。例えば、グローバルメモリ 20 は、実行モデル 10 のパイプライントポロジを定義するコマンドのリストを格納することができる。上述されるように、開発者は、カーネル 14 A 乃至 14 C を含むために実行モデル 10 のパイプライントポロジを定義しておくことができる。従って、グローバルメモリ 20 は、カーネル 14 A 乃至 14 C のソースコードを格納することができる。代替として又はさらに加えて、グローバルメモリ 20 は、カーネル 14 A 乃至 14 C の予めコンパイルされたソースコード (すなわち、カーネル 14 A 乃至 14 C のオブジェクトコード) を格納することができる。開発者がより多くの又はより少ないカーネル又は異なるカーネルを含むように実行モデルを定義した場合は、グローバルメモリ 20 は、それらのカーネルに関するソースコード及び / 又はオブジェクトコードを格納することができる。カーネル 14 A と 14 C との間の楕円形は、カーネル 14 B もグローバルメモリ 20 に含まれることを示す。

【0076】

この例では、プロセッサ 22 は、実行モデル 10 のパイプライントポロジをグローバルメモリ 20 から取り出すことができる。実行モデル 10 のパイプライントポロジに基づき、プロセッサ 22 は、実行モデル 10 がバッファ 12 A 乃至 12 D を含むことを決定することができる。この例では、プロセッサ 22 は、バッファ 12 A 乃至 12 D に関してグローバルメモリ 20 内で記憶場所を予約することができる。例えば、実行モデル 10 の一部は、バッファ 12 A 乃至 12 D の特徴、例えば、バッファ 12 A 乃至 12 D のサイズ、を含むことができる。この例では、プロセッサ 22 は、バッファ 12 A 乃至 12 D の特徴に基づいてグローバルメモリ 20 内の記憶場所を予約することができる。バッファ 12 A と 12 D との間の楕円形は、バッファ 12 B 及び 12 C もグローバルメモリ 20 に含まれることを示す。

【0077】

他の例として、実行モデル 10 の一部は、カーネル 14 A 乃至 14 C の特徴、例えば、拡大係数及び / 又はカーネル 14 A 乃至 14 C が生成する最大データ量、を含むことができる。この例では、プロセッサ 22 は、カーネル 14 A 乃至 14 C の特徴に基づいてグローバルメモリ 20 内の記憶場所を予約することができる。例えば、拡大係数及び / 又はカーネル 14 A 乃至 14 C が生成する最大データ量を示す値に基づいて、プロセッサ 22 は、バッファ 12 A 乃至 12 D に関する該当するサイズを決定し、決定されたサイズに基づいてグローバルメモリ 20 内の記憶場所を予約することができる。

【0078】

プロセッサ 22 は、バッファ 12 A 乃至 12 D に関するグローバルメモリ 20 内の記憶場所を予約するとして説明されるが、本開示の態様はそのようには限定されないことが理解されるべきである。幾つかの例では、IC 18 又は GPU 24 は、実行モデル 10 によって定義されるデータ処理アルゴリズムを実装するために GPU 24 によって使用されるバッファを管理するように構成される管理ユニット (図 2 には示されない) を含むことができる。これらの例では、プロセッサ 22 は、バッファ 12 A 乃至 12 D のサイズに関して管理ユニットを命令することができ、管理ユニットは、バッファ 12 A 乃至 12 D に関してグローバルメモリ内の記憶場所を予約するように構成することができる。

【0079】

この管理ユニットは、その他の機能、例えば、バッファ 12 A 乃至 12 D 内に格納されたキャッシュバック (cache-backing) データ及び / 又は IC 18 又は GPU 24 のキャッシュ内のカーネル 14 A 乃至 14 C の命令、を実行するように構成す

10

20

30

40

50

ることができる。この管理ユニットは、バッファ 1 2 A 乃至 1 2 D の各々の 1 つに格納されるデータ量を示す情報を格納することもできる。この管理ユニットは、G P U 2 4 での実行時にカーネル 1 4 A 乃至 1 4 C 間でのデータ転送を管理するように構成することができる。例えば、図 1 において例示されるように、実行モデル 1 0 のパイプラインポロジは、カーネル 1 4 A がバッファ 1 2 B に出力し、カーネル 1 4 B がバッファ 1 2 B からデータを受信することを示す。管理ユニットは、カーネル 1 4 A によって生成されたデータのバッファ 1 2 B 内での格納、及びカーネル 1 4 B によるバッファ 1 2 B からのデータの取り出し、そして幾つかの例では、バッファ 1 2 B に格納されるデータ量の格納を管理するように構成することができる。管理ユニットに関する技法は、"GRAPHICS PROCESSING UNIT BUFFER MANAGEMENT" (グラフィックス処理ユニットバッファ管理) という題名を有し、引用によってその内容全体が組み入れられている同時係属米国特許出願第 1 3 / 7 4 7 , 9 4 7 号 (出願日: 2 0 1 3 年 1 月 2 3 日) においても記述される。

10

【 0 0 8 0 】

管理ユニットの利用は、例示することを目的として提供されるものであり、限定するものであるとはみなされるべきでない。例えば、管理ユニットは、プロセッサ 2 2 以外のユニットがグローバルメモリ 2 0 内のバッファ 1 2 A 乃至 1 2 D に関する記憶場所を予約する 1 つの方法例として説明される。しかしながら、本開示の態様は、そのようには限定されず、プロセッサ 2 2、又は G P U 2 4 でさえも、該機能を果たすように構成することができる。例えば、G P U 2 4 がバッファ 1 2 A 乃至 1 2 D 内にデータを格納するときには、G P U 2 4 は、G P U 2 4 がバッファ 1 2 A 乃至 1 2 D に格納したデータ量を格納するようにも構成することができる。説明を容易にするため、技法は、プロセッサ 2 2 又は G P U 2 4 が該機能を実行することに関して説明される。

20

【 0 0 8 1 】

本開示において説明される技法により、プロセッサ 2 2 は、実行モデル 1 0 によって定義されたパイプラインポロジのインディケーション (i n d i c a t i o n) (例えば、一例として、実行モデル 1 0 のコマンドリスト) を受信することができる。プロセッサ 2 2 は、G P U 2 4 がパイプラインポロジを実装する方法を定義する命令を生成することができる。

【 0 0 8 2 】

例えば、例示されるように、プロセッサ 2 2 は、コンパイラ 2 8 を実行することができる。コンパイラ 2 8 は、プロセッサ 2 2 内で形成されていないことを示すために破断線で示される。むしろ、グローバルメモリ 2 0 は、コンパイラ 2 8 のオブジェクトを格納することができ、プロセッサ 2 2 が取り出して実行する。

30

【 0 0 8 3 】

コンパイラ 2 8 は、G P U 2 4 が実行するオブジェクトコード、及び G P U 2 4 が実行モデル 1 0 を実装する方法を定義する命令を生成するために実行モデル 1 0 (例えば、実行モデル 1 0 のコマンドリスト) をコンパイルするように構成することができる。G P U 2 4 が実行モデル 1 0 を実装する方法を定義する命令を生成するためのコンパイルの一部として、コンパイラ 2 8 は、G P U 2 4 の処理能力を説明する (a c c o u n t f o r) ことができる。

40

【 0 0 8 4 】

例えば、例示されるように、グローバルメモリ 2 0 は、G P U 構成 3 2 を格納することができる。G P U 構成 3 2 は、G P U 2 4 の処理能力を定義する又は示す構成情報であることができる。一例として、G P U 構成 3 2 は、G P U 2 4 内のプログラマブル計算ユニット数を示すことができる。上述されるように、カーネルは、G P U 内の 1 つ以上のプログラマブル計算ユニットで実行する。

【 0 0 8 5 】

他の例として、G P U 構成 3 2 は、G P U 2 4 が並列でデータを処理することが可能な方法を示すことができる。例えば、G P U 2 4 は、単一プログラム多データ (S P M D) プログラミングモデル又は単一命令多データ (S I M D) プログラミングモデルを実装す

50

るように構成することができる。一例として、GPU 24がSIMDプログラミングモデルに関して構成される場合は、GPU構成32は、SIMDプログラミングモデル（例えば、8レーンSIMD）を実装するためのGPU 24内のレーン（lane）の数を示す構成情報を含むことができる。

【0086】

GPU構成32は、上述される情報の追加の又は上述される情報と異なるGPU 24の構成情報を含むことができる。概して、GPU構成32は、GPU 24の処理能力を記述する構成情報を含むことができる。

【0087】

さらに、GPU構成32はグローバルメモリ20に格納されるとして例示されるが、本開示の態様は、そのようには限定されない。幾つかの例では、IC 18内のレジスタ又はキャッシュがGPU構成32を格納することができる。これらの例の両方において、プロセッサ22は、グローバルメモリ20からではなく、レジスタからGPU構成32の情報を読み取ることができる。幾つかの例では、プロセッサ22をGPU構成32で予め構成することさえも可能である。

【0088】

コンパイラ28は、実行モデル10をコンパイルするために、実行モデル10の情報に加えて、GPU構成32の情報を利用することができる。コンパイルの結果は、GPU 24が実行するオブジェクトコード、及び、GPU 24が実行モデル10を実装する方法に関する命令であることができる。例えば、オブジェクトコードに加えて、コンパイラ28の出力は、プロセッサ22がグローバルメモリ20に格納するメタ-スケジューラ命令34であることができる。メタ-スケジューラ命令34は、より詳細に説明されるように、GPU 24が実行モデル10を実装する方法を指示する、メタ-スケジューラ30に関する命令であることができる。

【0089】

例えば、メタ-スケジューラ命令34は、GPU 24がカーン処理ネットワーク（Kahn Processing Network）（KPN）に類似する実行モデル10を実装するように指示することができる。例えば、KPNは、データを含むチャネルを決定し、そのチャネルに関するコンシューマ（consumer）を識別し、データのうちの一部の量に関してコンシューマを実行し、全データが処理されるまでこれらのステップを繰り返す。実行モデル10のトポロジーは、（KPNのプロセスに類似する）カーネル及び（KPNのプロセスに類似する）バッファを定義することができる。この方法により、実行モデル10は、バッファの各々に関する消費するカーネルを示す。より詳細に説明されるように、実行モデル10を実装する際には、GPU 24は、データを含むバッファ12A乃至12Dのうちの1つを識別することができ、及び、バッファ12A乃至12Dのうちの識別されたそのデータを消費することになる消費するカーネル（例えば、カーネル14A乃至14Cのうちの1つ）を実行することができる。KPNに関する説明は、単に例示する目的で及び理解を助けるために提供されたものであることが理解されるべきである。本開示において説明される技法は、KPNのそれらに限定される又はKPNのそれらと同一であるとはみなされるべきではない。

【0090】

それらの技法により、メタ-スケジューラ命令34は、デバイスをターゲットにしたバイナリ（device target binary）であることができる。換言すると、実行モデル10はプラットフォームから独立する（すなわち、GPU専用でない）ことができる一方で、メタ-スケジューラ命令34は、プラットフォームに依存する（すなわち、GPU 24専用である）。例えば、コンパイラ28は、GPU 24での実行モデル10の実装を最適化するためにGPU構成32からの情報を利用することができる。

【0091】

一例として、コンパイラ28は、カーネル14A乃至14CがGPU 24のプログラマブル計算ユニットで実行する時間を決定するために情報、例えば、GPU 24のプログラ

10

20

30

40

50

マブル計算ユニット数、を利用することができる。例えば、上述されるように、実行モデル 10 によって定義されるデータ処理アルゴリズムのパイプライン実装は、カーネル 14 A 乃至 14 C のうちの 1 つ以上の並列の実行（例えば、同時）又はカーネル 14 A 乃至 14 C のうちの 1 つ以上の順次の実行（例えば、1 つずつ）を含むことができる。この例では、コンパイラ 28 は、カーネル 14 A 乃至 14 C が、数多くの利用可能なプログラマブル計算ユニットが存在しない場合は順次で実行すべきであること又は数多くの利用可能なプログラマブル計算ユニットが存在する場合は並列で実行すべきであることを指示するメタ - スケジューラ命令 34 を生成することができる。

【0092】

他の例として、コンパイラ 28 は、カーネル 14 A 乃至 14 C の一部のそれらがその他よりも重要であることを示す実行モデル 10 からの情報を利用することができる。例えば、カーネル 14 B がカーネル 14 C よりも重要であると仮定する。この例では、コンパイラ 28 は、カーネル 14 C が幾つかの割り込み有りで行う実行することになった場合でもカーネル 14 B は割り込みなしで行うようにするメタ - スケジューラ命令 34 を生成するために GPU 構成 34 からの情報を利用することができる。

【0093】

例えば、プログラマブル計算ユニットが 1 つのカーネルのスレッドを実行することから他のカーネルのスレッドを実行することに切り換わり、次に戻ることが可能である。この場合は、プログラマブル計算ユニットが切り換わったカーネルは割り込まれたとみなすことができ、プログラマブル計算ユニットがそのカーネルのスレッドを実行するために切り換わって戻るまで休止状態にあるとみなすことができる。幾つかの例では、コンパイラ 28 は、重要なカーネルのスレッドを実行中のプログラマブル計算ユニットが他のカーネルのスレッドの実行に切り換わらないように指示するメタ - スケジューラ命令 34 を生成することができる。

【0094】

さらに他の例として、コンパイラ 28 は、GPU 24 内のデータレーン (data lane) の数を示す GPU 構成 34 からの情報を利用することができる。例えば、GPU 構成 34 は、GPU 24 が 8 レーンの SIMD GPU であることを示すと仮定する。この例では、コンパイラ 28 は、消費するカーネルに結合されたバッファ内に少なくとも 8 つのエントリが存在するまで GPU 24 が消費するカーネルを実行すべきでないことを指示するメタ - スケジューラ命令 34 を生成することができる。例えば、図 1 において例示されるように、カーネル 14 B は、カーネル 14 A にとっての消費するカーネルであり、バッファ 12 B からデータを受信する。GPU 24 は 8 レーン SIMD GPU であると仮定すると、コンパイラ 28 は、少なくとも 8 つのデータ項目がバッファ 12 B 内に格納されるまで GPU 24 がカーネル 14 B を実行すべきでないことを指示するメタ - スケジューラ命令 34 を生成することができる。

【0095】

追加の例として、コンパイラ 28 は、バッファ 12 A 乃至 12 D のサイズを説明することができる。例えば、コンパイラ 28 は、行き詰まりの変化が最小限になるようにカーネル 14 A 乃至 14 C をいつ実行すべきかを決定するためにバッファ 12 A 乃至 12 D のサイズに関する情報を利用することができる。この例では、コンパイラ 28 は、行き詰まりが発生しないようにカーネル 14 A 乃至 14 C が実行する順序を指示するメタ - スケジューラ命令 34 を生成することができる。

【0096】

幾つかの例では、コンパイラ 28 は、GPU 24 に実行モデル 10 を実装する際に誤りが存在するかどうかを決定するように構成することができる。例えば、コンパイラ 28 は、実行モデル 10 の一部であるが、カーネル 14 A 乃至 14 D のいずれの 1 つにも結合されていないバッファ 12 A 乃至 12 D が存在するかどうかを決定するように構成することができる。他の例として、コンパイラ 28 は、カーネル 14 A 乃至 14 D のうちのいずれか 1 つが存在していないバッファにアクセスしようとしているか、又はバッファ内のアウ

10

20

30

40

50

トオブバウンド (o u t - o f - b o u n d s) 記憶場所にアクセスしようとしているかを決定するように構成することができる。コンパイラ 2 8 がコンパイル時に実行モデル 1 0 の機能を検証するため、コンパイラ 2 8 は、実行モデル 1 0 の機能を検証することを G P U 2 4 に行わせる命令をメタ - スケジューラ命令 3 4 に含める必要がない。

【 0 0 9 7 】

プロセッサ 2 2 が、コンパイラ 2 8 を介して、メタ - スケジューラ命令 3 4 を生成した後は、プロセッサ 2 2 は、メタ - スケジューラ命令 3 4 を実行のために取り出すように G P U 2 4 に命令することができる。例示されるように、G P U 2 4 は、メタ - スケジューラ 3 0 を含む。メタ - スケジューラ 3 0 は、G P U 2 4 内のハードウェア、G P U 内のハードウェアで実行するファームウェア、又は G P U 2 4 内のハードウェアで実行するソフトウェアであることができる。メタ - スケジューラ 3 0 は、メタ - スケジューラ命令 3 4 の命令を実行するように構成することができる。

10

【 0 0 9 8 】

本開示において説明される技法では、メタ - スケジューラ 3 0 は、G P U 2 4 のいずれのプログラマブル計算ユニットがカーネル 1 4 A 乃至 1 4 C のいずれのスレッドを何時に実行すべきかを決定するように構成することができる。換言すると、メタ - スケジューラ 3 0 は、実行モデル 1 0 のパイプラインポロジによって定義されたデータ処理アルゴリズムを実装することを G P U 2 4 に行わせるために G P U 2 4 でのカーネル 1 4 A 乃至 1 4 C の実行をスケジューリングするように構成することができる。本開示において説明される技法により、メタ - スケジューラ 3 0 は、メタ - スケジューラ命令 3 4 に基づいて G P U 2 4 でカーネル 1 4 A 乃至 1 4 C を実行するスケジュールを決定することができる。

20

【 0 0 9 9 】

例えば、メタ - スケジューラ命令 3 4 は、カーネル 1 4 A 乃至 1 4 C のうちの 1 つ以上を並列で又は順次実行すべきかを指示することができる。この例では、メタ - スケジューラ 3 0 は、並列又は順次の実行を達成するためにいずれのプログラマブル計算ユニットがカーネル 1 4 A 乃至 1 4 C のスレッドを実行すべきかを決定することができる。他の例として、メタ - スケジューラ命令 3 4 は、カーネル 1 4 A 乃至 1 4 C の重要性を示すことができる。この例では、メタ - スケジューラ 3 0 は、重要なカーネルを実行するプログラマブル計算ユニットがカーネルを実行中に割り込まれないようにいずれのプログラマブル計算ユニットがカーネル 1 4 A 乃至 1 4 C のスレッドを実行すべきかを決定することができる。他の例として、メタ - スケジューラ命令 3 4 は、G P U 2 4 の S I M D 又は S P M D 能力に基づいてカーネル 1 4 A 乃至 1 4 C のうちの 1 つがいつ実行すべきかを指示することができる。この例では、メタ - スケジューラ 3 0 は、メタ - スケジューラ命令 3 4 の命令に基づいて、いずれのプログラマブル計算ユニットが何時にスレッドを実行するかを決定することができる。メタ - スケジューラ 3 0 は、行き詰まりを回避するためにカーネル 1 4 A 乃至 1 4 C がいつ実行すべきかのタイミングを指示するメタ - スケジューラ命令 3 4 内の命令を利用することもできる。

30

【 0 1 0 0 】

繰り返すと、コンパイラ 2 8 は、メタ - スケジューラ命令 3 4 を生成する際に、G P U 構成 3 2 内の情報に基づいて、G P U 2 4 の計算能力を考慮に入れておくことができる。従って、本開示において説明される技法は、メタ - スケジューラ 3 0 がメタ - スケジューラ 3 4 によって示された方法でカーネル 1 4 A 乃至 1 4 C を実行するためにプログラマブル計算ユニットを適切に割り当てることができることをある程度のレベルで保証する。例えば、上述されるように、コンパイラ 2 8 は、メタ - スケジューラ命令 3 4 を生成するために、幾つかの説明例として、幾つかの要因、例えば、G P U 2 4 内でのプログラマブル計算ユニットの数、バッファ 1 2 A 乃至 1 2 D のサイズ、G P U 2 4 の S I M D 又は S P M D 能力、又はカーネル 1 4 A 乃至 1 4 C の重要性、を説明することができる。G P U 2 4 のメタ - スケジューラ 3 0 がカーネル 1 4 A 乃至 1 4 C がどのように実行されるべきかを決定するためにメタ - スケジューラ命令 3 4 を利用するときには、カーネル 1 4 A 乃至

40

50

14CはGPU24で効率的に実行することある程度保証することができる。

【0101】

この方法により、プロセッサ22は、実行モデル10のパイプラインポロジーを示すインディケーションを受信することができ、実行モデル10は、プラットフォームから独立した形で定義される。プロセッサ22は、コンパイラ28を実行することができ、それは、実行モデル10のパイプラインポロジーがどのようにしてGPU24上に実装されるべきかをプラットフォーム専用の方法で定義するメタ-スケジューラ命令34を生成するためのGPU24の処理能力を説明する。GPU24のメタ-スケジューラ30は、メタ-スケジューラ命令34を受信し、メタ-スケジューラ命令34の命令に基づいて、い

10

【0102】

メタ-スケジューラ命令34は、各プラットフォーム専用であるため、メタ-スケジューラ命令34は、GPU24での実行モデル10のパイプラインポロジーの最適な実装になるようにカーネル14A乃至14Cが実行されるべき方法を定義することができる。例えば、メタ-スケジューラ命令34がGPU24と異なるタイプのGPUによって使用される場合は、この異なるタイプのGPUは、メタ-スケジューラ命令34がGPU24のプラットフォーム専用であるため、実行モデル10のパイプラインポロジーを効率的に実装することができない。

【0103】

20

幾つかの例では、メタ-スケジューラ30は、実行ポリシーを実装するように構成することができる。例えば、すべての例においてコンパイラ28が、カーネル14A乃至14Cがいつ実行されるべきかを正確に定義する必要がないことがある。むしろ、コンパイラ28が生成するメタ-スケジューラ命令34は、カーネル14A乃至14CをGPU24で実行すべきことを指示することができ及びいずれのカーネル14A乃至14Cが生成するカーネルであり、いずれが消費するカーネルであることを示すことができる。

【0104】

これらの例では、メタ-スケジューラ30は、カーネル14A乃至14Cがいつ実行すべきかを示す実行ポリシーを実装するように構成することができる。実行ポリシーの一例は、バッファ12A乃至12Dのうちのいずれがデータを格納するかをメタ-スケジューラ30が決定し、データを格納するバッファ12A乃至12Dからデータを受信するカーネル14A乃至14Cのうちの1つ以上を実行する。例えば、メタ-スケジューラ30は、バッファ12A乃至12Dをラウンドロビン方式で検査することができ、及び、データを格納するバッファ12A乃至12Dからデータを消費する全カーネルを実行することができる。

30

【0105】

実行ポリシーの他の例として、メタ-スケジューラ30は、カーネル14A乃至14Cの重要性に基づいていずれのバッファ12A乃至12Dがデータを格納するかを決定することができる。カーネル14A乃至14Cの重要性は、カーネル14A乃至14Cの優先度によって定義することができる。メタ-スケジューラ30は、最初に、カーネル14A乃至14Dの最高の優先度を有するカーネルがデータを受信するバッファ12A乃至12Dのバッファを検査することができる。そのバッファがデータを格納している場合は、メタ-スケジューラ30は、最高の優先度を有するカーネルを実行することができる。次に、メタ-スケジューラ30は、カーネル14A乃至14Dの次に最高の優先度を有するカーネルがデータを受信するバッファ12A乃至12Dのバッファを検査することができる。そのバッファがデータを格納している場合は、メタ-スケジューラ30は、次に最高の優先度を有するカーネルを実行することができ、以下同様である。

40

【0106】

実行ポリシーのさらに他の例として、メタ-スケジューラ30は、バッファ12A乃至12Dのうちのいずれが最も多くの量のデータを格納するかを決定することができる。メタ

50

- スケジューラ 30 は、最も多くの量のデータを格納するバッファ 12 A 乃至 12 D のバッファからデータを受信するカーネル 14 A 乃至 14 C のカーネルを実行することができる。一例として、CPU 24 がバッファ 12 A 乃至 12 D のうちの 1 つにデータを書き込むときには、GPU 24 は、GPU 24 がデータを書き込んだバッファ 12 A 乃至 12 D のうちの 1 つに格納されるデータの量を示す情報を格納することができる。これらの例では、メタ - スケジューラ 30 は、GPU 24 が格納し、バッファ 12 A 乃至 12 D 内のデータの量を示す情報に基づいてバッファ 12 A 乃至 12 D のうちのいずれの 1 つが最も多くの量のデータを格納するかを決定するように構成することができる。

【0107】

しかしながら、メタ - スケジューラ 30 を実行ポリシーを実装するように予め構成する必要がない。むしろ、コンパイラ 28 が、コンパイラ 28 によって生成されるメタ - スケジューラ命令 34 の命令の一部としてメタ - スケジューラ 30 の実行ポリシーを決定することが可能である。

10

【0108】

さらに、開発者がメタ - スケジューラ 30 の実行ポリシーを定義することが可能である。例えば、開発者は、実行モデル 10 の一部としてメタ - スケジューラ 30 の実行ポリシーを定義することができ、コンパイラ 28 は、開発者によって定義された実行ポリシーに関してメタ - スケジューラ 30 に命令するメタ - スケジューラ命令 34 の命令を生成するためにこの開発者によって定義された実行ポリシーを利用することができる。

20

【0109】

しかしながら、開発者が実行ポリシーを定義しないほうが適切である場合がある。例えば、開発者が実行ポリシーを定義した場合は、実行ポリシーは、プラットフォームから独立した方法では適切に機能しないことがある。すべての GPU タイプに関して適切に機能し、決定するアプリケーションにおいて同じ機能上の結果（例えば、同じアプリケーションに関する異なる GPU タイプにわたって同じ結果）を生み出す実行ポリシーを開発者が開発するのは困難であろう。

【0110】

概して、開発者は、カーネル 14 A 乃至 14 C が実行モデル 10 によって定義されたデータ処理アルゴリズムを実装するために適切に実行するかぎりは実行ポリシーには特別の関心を有さないであろう。従って、開発者が実行ポリシーを定義することができないかどうか、及びコンパイラ 28 が実行ポリシーを決定することは関係ないであろう。

30

【0111】

上述される例においては、コンパイラ 28 は、GPU で実行されるオブジェクトコードを生成するために及びメタ - スケジューラ命令 34 を生成するために実行モデル 10 をコンパイルした。幾つかの例では、コンパイラ 28 は、GPU 24 で実行されるオブジェクトコードをグローバルメモリ 20 に格納することができる。実行モデル 10 が複数の実行モデルのうちの 1 つの実行モデルである場合は、コンパイラ 28 は、実行モデル 10 を再コンパイルする必要がない。換言すると、複数の実行モデルを結合することによって、及び幾つかの例では、グローバルメモリ 20 に格納された実行モデルのオブジェクトコードを結合することによって、大きなデータ処理アルゴリズムを生成することが可能である。例えば、データ処理アルゴリズムは、グローバルメモリ 20 に格納された実行モデルのオブジェクトコードから生成することができる。データ処理アルゴリズムの該生成は、GPU 24 が FPG A 又は埋め込まれたデバイスである例にとって有用であることができる。

40

【0112】

本開示において説明される技法では、メタ - スケジューラ 30 に関してメタ - スケジューラ命令 34 を取り出すように GPU 24 に命令することに加えて、プロセッサ 22 は、GPU 24 がカーネル 14 A 乃至 14 C を実行するために必要な追加情報を GPU 24 に提供することができる。例えば、カーネル 14 A 乃至 14 C は、バッファ 12 A 乃至 12 D からのデータに加えて、機能するために追加情報（例えば、引数）を要求することができる。プロセッサ 22 は、該追加情報を GPU 24 に提供することができる。

50

【 0 1 1 3 】

プロセッサ 2 2 は、実行モデル 1 0 を実装するように G P U 2 4 に命令することができる。G P U 2 4 は、メタ - スケジューラ命令 3 4 の命令及びプロセッサ 2 2 がコンパイラ 2 8 によるコンパイルプロセスの一部として生成したオブジェクトコードに基づいて実行モデル 1 0 を実装することができる。幾つかの例では、G P U 2 4 は、プロセッサ 2 2 との同期化なしで実行モデル 1 0 を実装することができる。

【 0 1 1 4 】

さらに、幾つかの例では、コンパイラ 2 8 及びメタ - スケジューラ 3 0 をデバッグモードで構成することが可能である。例えば、開発者が実行モデル 1 0 を開発した後に、開発者は、リリース前に G P U での実行モデルの実装を試験することを希望することができる。試験に関しては、開発者は、デバイス、例えば、デバイス 1 6、に実行モデル 1 0 をロードし、G P U 2 4 で実行モデル 1 0 を試験することができる。試験の一環として、開発者は、デバッグモードを利用することができる。デバッグモードでは、コンパイラ 2 8 は、バッファ 1 2 A 乃至 1 2 D の記憶場所の範囲を単一の記憶場所にまで狭める（例えば、N D r a n g e サイズを最小に小さくする）メタ - スケジューラ命令 3 4 を生成することができる。メタ - スケジューラ命令 3 4 は、カーネル 1 4 A 乃至 1 4 C のうちの 1 つのカーネルのみが一度に実行するように指示することもできる。

【 0 1 1 5 】

デバッグモードでは、開発者は、データがバッファ 1 2 A 乃至 1 2 D 内に格納されている方法、及びカーネル 1 4 A 乃至 1 4 C の各々の 1 つが G P U 2 4 で実行されている方法を追跡することができる。これは、開発者がカーネル 1 4 A 乃至 1 4 C 内の問題又は実行モデル 1 0 内の問題に対処するのを可能にすることができる。

【 0 1 1 6 】

上述されるように、コンパイラ 2 8 は、メタ - スケジューラ命令 3 4 を生成することができる。以下は、実行モデル 1 0 に関するメタ - スケジューラ命令 3 4 の擬似コード例である。幾つかの例では、メタ - スケジューラ命令 3 4 を生成する能力を開発者に与えるのではなく、コンパイラ 2 8 がメタ - スケジューラ命令 3 4 を生成するのが有益であることができる。例えば、開発者がメタ - スケジューラ命令 3 4 を生成した場合は、実行モデル 1 0 はポータブルになることができず、混乱が生じ、さらにデバッグが困難なユーザエラーが生じる可能性がある。

【 0 1 1 7 】

以下の擬似コードでは、F 1 は、バッファ 1 2 A を指し示し、F 2 は、バッファ 1 2 B を指し示し、F 3 は、バッファ 1 2 C を指し示し、F 4 は、バッファ 1 2 D を指し示す。K 1 は、カーネル 1 4 A を指し示し、K 2 は、カーネル 1 4 B を指し示し、K 3 は、カーネル 1 4 C を指し示す。

10

20

30

【数 1】

```

while (true)
{
    //Node 1
    if (K1.F1.size( ) > 0)
    {
        maxSafeNdRangeSizeK1 = min(K1.F1.size( ),
            ((K1.F3.maxSize( ) - K1.F3.size( ))/K1.F3.ampFactor( )),
            ((K1.F2.maxSize( ) - K1.F2.size( ))/K1.F2.ampFactor( )));
        enqueue (K1, maxSafeNdRangeSizeK1);
    }
    //Node 2
    ...
    //Node 3
    ...
    If (F1.size( ) + F2.size( ) + F3.size( ) + F4.size( ) == 0)
        && (K3 is done))
    {
        Exit //Finished execution of graph
    }
}

```

【0118】

図3は、本開示において説明される1つ以上の例による技法例を示したフローチャートである。例えば、図3は、異種計算（例えば、プラットフォームから独立した形での計算）に関する技法を例示する。例示を容易にするために、図2が参照される。

【0119】

図3において例示されるように、プロセッサ22は、プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデル10のパイプラインポロジのインディケーションを受信することができる（38）。例えば、実行モデル10のパイプラインポロジのインディケーションは、実行モデル10の開発者によって作成されたコマンドリストであることができる。データ処理アルゴリズムのプラットフォームから独立したに関する定義は、実行モデル10がGPUの特定のプラットフォームに基づいて設計されていない（例えば、データ処理アルゴリズムを実装するGPUのタイプから独立している）ことを意味する。

【0120】

プロセッサ22のコンパイラ28は、メタ-スケジューラ命令34を生成するためにパイプラインポロジを指定するコマンドリストをコンパイルすることができる（40）。メタ-スケジューラ命令34は、GPU24が実行モデル10のパイプラインポロジ

ーを実装するプラットフォームに依存した方法を指示することができる。GPU 24がパイプラインポロジーを実装するプラットフォームに依存した方法は、メタ-スケジューラ命令が(例えば、GPU 24のGPUタイプに基づいて)GPU構成32によって示されるGPU 24の特定のプラットフォームに基づくことを意味する。プロセッサ22は、実行モデル10のパイプラインポロジーを実装するようにGPU 24に命令するための命令を送信することができる(42)。

【0121】

プロセッサ22が、コンパイラ28を介して、メタ-スケジューラ命令34を生成することができる様々な方法が存在する。一例として、コンパイラ28は、実行モデル10のパイプラインポロジーがGPU 24に実装されるプラットフォームに依存する方法を定義する命令を生成するためにGPU 24の構成情報に少なくとも基づいてコマンドリストをコンパイルすることができる。GPU構成32は、GPU 24の構成情報を提供することができる。例えば、構成情報は、GPU 24内のプログラマブル計算ユニット数を含むことができる。構成情報は、GPU 24内のデータレーン数(すなわち、GPU 24のSIMD又はSPMD能力)を含むことができる。

【0122】

幾つかの例では、コンパイラ28は、実行モデル10で提供された情報に基づいてコマンドリストをコンパイルすることができる。例えば、コンパイラ28は、メタ-スケジューラ命令34を生成するために実行モデル10のパイプラインポロジーで識別されたバッファ(例えば、バッファ12A乃至12D)のサイズに基づいてコマンドリストをコンパイルすることができる。他の例として、コンパイラ28は、メタ-スケジューラ命令34を生成するために実行モデル10のパイプラインポロジーで識別されたカーネル(例えば、カーネル14A乃至14C)の重要性に基づいてコマンドリストをコンパイルすることができる。

【0123】

幾つかの要因、例えば、GPU 24におけるプログラマブル計算ユニット数、GPU 24におけるデータレーン数、バッファ12A乃至12Dのサイズ、及びカーネル14A乃至14Cの重要性、等を利用するコンパイラ28は、例示を目的として提供されるものであり、限定するとはみなされるべきでないことが理解されるべきである。さらに、コンパイラ28は、要因を単独で又はあらゆる組み合わせで利用することができる。例えば、コンパイラ28は、メタ-スケジューラ命令34を生成する際にこれらの要因のうちの1つのみを利用する必要はない。むしろ、コンパイラ28は、メタ-スケジューラ命令34を生成するためにこれらの要因のうちの1つ、これらの要因のうちの1つ以上、及びこれらの要因のあらゆる組み合わせを利用することができる。

【0124】

図4は、図2のデバイスをさらに詳細に例示したブロック図である。例えば、図4は、デバイス16をさらに例示する。デバイス16の例は、無線デバイス、携帯電話、パーソナルデジタルアシスタント(PDA)、ビデオディスプレイを含むビデオゲームコンソール、モバイルビデオ会議装置、ラップトップコンピュータ、デスクトップコンピュータ、テレビセットトップボックス、ダフレットコンピューティングデバイス、電子書籍リーダー、等を含み、ただしこれらには限定されない。デバイス16は、プロセッサ22と、GPU 24と、グローバルメモリ20と、ディスプレイ44と、ユーザインタフェース46と、トランシーバモジュール48と、を含むことができる。プロセッサ22及びGPU 24は、図4において例示されるように、共通のIC 18内に収納することができ、又は別々に収納することができる。さらに、例示されるように、プロセッサ22は、メタ-スケジューラ命令34を生成するためにコンパイラ28を実行することができ、GPU 24は、メタ-スケジューラ命令34の命令を実装するように構成されたメタ-スケジューラ30を含む。

【0125】

デバイス16は、明確化のために図4には示されていない追加のモジュール又はユニッ

10

20

30

40

50

トを含むことができる。例えば、デバイス 16 は、デバイス 16 がモバイル無線電話である例において電話通信を有効にするためのスピーカーとマイクとを含むことができ、これらはいずれも図 4 には示されていない。さらに、デバイス 16 内に示される様々なモジュール及びユニットは、デバイス 16 のすべての例において必要であるわけではない。例えば、ユーザインタフェース 46 及びディスプレイ 44 は、デバイス 16 がデスクトップコンピュータである例ではデバイス 16 の外部に存在することができる。他の例として、ユーザインタフェース 46 は、ディスプレイ 44 がモバイルデバイスのタッチ感応式又はプレゼンス感応式 (presence-sensitive) のディスプレイである例ではディスプレイ 44 の一部であることができる。

【0126】

グローバルメモリ 20、プロセッサ 22、GPU 24、コンパイラ 28、及びメタスケジューラ 30 は、グローバルメモリ 20、プロセッサ 22、GPU 24、コンパイラ 28、及びメタスケジューラ 30 に類似するものであり、図 4 に関してはこれ以上は説明されない。ユーザインタフェース 46 の例は、トラックボールと、マウスと、キーボードと、その他のタイプの入力デバイスとを含み、ただしこれらに限定されない。ユーザインタフェース 46 は、タッチ画面であることもでき、ディスプレイ 44 の一部として組み入れることができる。トランシーバモジュール 48 は、デバイス 16 と他のデバイス又はネットワークとの間の無線又は有線通信を可能にするための回路を含むことができる。トランシーバモジュール 48 は、変調器と、復調器と、増幅器と、有線又は無線通信のためのその他の該回路と、を含むことができる。ディスプレイ 44 は、液晶ディスプレイ (LCD)、陰極線管 (CRT) ディスプレイ、プラズマディスプレイ、タッチ感応式ディスプレイ、プレゼンス感応式ディスプレイ、又は他のタイプのディスプレイデバイスを備えることができる。

【0127】

1 つ以上の例では、説明される機能は、ハードウェア、ソフトウェア、ファームウェア、又はそれらのあらゆる組み合わせにおいて実装することができる。ソフトウェアにおいて実装される場合は、それらの機能は、コンピュータによって読み取り可能な媒体において 1 つ以上の命令又はコードとして格納又は送信すること及びハードウェアに基づく処理ユニットによって実行することができる。コンピュータによって読み取り可能な媒体は、コンピュータによって読み取り可能な記憶媒体を含むことができ、それは、有形な媒体、例えば、データ記憶媒体、又は、例えば、通信プロトコルにより、1 つの場所から他へのコンピュータプログラムの転送を容易にするあらゆる媒体を含む通信媒体、に対応する。このように、コンピュータによって読み取り可能な媒体は、概して、(1) 非一時的である有形なコンピュータによって読み取り可能な記憶媒体又は (2) 通信媒体、例えば、信号又は搬送波、に対応することができる。データ記憶媒体は、本開示において説明される技法の実装のために命令、コード及び/又はデータ構造を取り出すために 1 つ以上のコンピュータ又は 1 つ以上のプロセッサによってアクセスすることができるあらゆる利用可能な媒体であることができる。コンピュータプログラム製品は、コンピュータによって読み取り可能な媒体を含むことができる。

【0128】

一例により、及び制限することなしに、該コンピュータによって読み取り可能な記憶媒体は、希望されるプログラムコードを命令又はデータ構造の形態で格納するために使用することができる及びコンピュータによってアクセス可能である RAM、ROM、EEPROM、CD-ROM 又はその他の光学ディスク記憶装置、磁気ディスク記憶装置、又はその他の磁気記憶デバイス、フラッシュメモリ、又はその他のいずれかの媒体を備えることができる。さらに、どのような接続も、コンピュータによって読み取り可能な媒体であると適切に呼ばれる。例えば、命令が、同軸ケーブル、光ファイバケーブル、より対線、デジタル加入者ライン (DSL)、又は無線技術、例えば、赤外線、無線、及びマイクロ波、を用いてウェブサイト、サーバ、又はその他の遠隔ソースから送信される場合は、該同軸ケーブル、光ファイバケーブル、より対線、DSL、又は無線技術、例えば赤外線、無線

10

20

30

40

50

、及びマイクロ波、は、媒体の定義の中に含まれる。しかしながら、コンピュータによって読み取り可能な記憶媒体およびデータ記憶媒体は、接続、搬送波、信号、又はその他の一時的な媒体は含まず、代わりに、非一時的な、有形の記憶媒体を対象とすることが理解されるべきである。ここにおいて用いられるときのディスク(disk及びdisc)は、コンパクトディスク(CD)(disc)と、レーザディスク(disc)と、光ディスク(disc)と、デジタルバーサタイルディスク(DVD)(disc)と、フロッピー(登録商標)ディスク(disk)と、Blu-ray(登録商標)ディスク(disc)と、を含み、ここで、diskは、通常は磁気的にデータを複製し、discは、レーザを用いて光学的にデータを複製する。上記の組み合わせも、コンピュータによって読み取り可能な媒体の適用範囲内に含まれるべきである。

10

【0129】

命令は、1つ以上のプロセッサ、例えば、1つ以上のデジタル信号プロセッサ(DSP)、汎用マイクロプロセッサ、特定用途向け集積回路(ASIC)、フィールドプログラマブルロジックアレイ(FPGA)又はその他の同等の集積回路又はディスクリット論理回路によって実行することができる。従って、ここにおいて用いられる場合の用語“プロセッサ”は、上記の構造又はここにおいて説明される技法の実装に適するあらゆるその他の構造のうちのいずれかを意味することができる。さらに、幾つかの態様では、ここにおいて説明される機能は、符号化および復号のために構成された専用のハードウェア及び/又はソフトウェアモジュール内において提供されること、又は組み合わされたコーデック内に組み入れることができる。さらに、技法は、1つ以上の回路又は論理素子内に完全に実装することが可能である。

20

【0130】

本開示の技法は、無線ハンドセット、集積回路(IC)又は一組のIC(例えば、チップセット)を含む非常に様々なデバイス又は装置内に実装することができる。本開示では、開示される技法を実施するように構成されたデバイスの機能上の態様を強調するために様々なコンポーネント、モジュール、又はユニットが説明されるが、異なるハードウェアユニットによる実現は必ずしも要求しない。むしろ、上述されるように、様々なユニットは、適切なソフトウェア及び/又はファームウェアと関係させて、ハードウェアユニット内において結合させること又は上述されるように1つ以上のプロセッサを含む相互運用的なハードウェアユニットの集合によって提供することができる。

30

【0131】

様々な例が説明されている。これらの及びその他の例は、以下の請求項の範囲内である。
以下に、本願出願の当初の特許請求の範囲に記載された発明を付記する。

【C1】

異種計算のための方法であって、

プロセッサを用いて、プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプラインポートロギーを受信することと、

前記プロセッサを用いて、前記実行モデルの前記パイプラインポートロギーがグラフィックス処理ユニット(GPU)に実装されるプラットフォームに依存する方法を指示する命令を生成することであって、前記実行モデルの前記パイプラインポートロギーが前記GPUに実装される前記プラットフォームに依存する方法は、前記GPUのプラットフォームに基づくことと、

40

前記プロセッサを用いて、前記GPUに前記命令を送信することと、を備える、方法。

【C2】

命令を生成することは、前記命令を生成するために前記実行モデルの前記パイプラインポートロギーを指定するコマンドリストをコンパイルすることを備えるC1に記載の方法。

【C3】

前記コマンドリストをコンパイルすることは、

前記実行モデルの前記パイプラインポートロギーが前記GPUに実装される前記プラットフ

50

フォームに依存する方法を定義する前記命令を生成するために前記GPUの構成情報に少なくとも基づいて前記コマンドリストをコンパイルすることを備えるC2に記載の方法。

[C4]

前記GPUの前記構成情報は、

前記GPU内のプログラマブル計算ユニット数、及び

前記GPU内のデータレーン数

のうちの1つ以上を備えるC3に記載の方法。

[C5]

前記コマンドリストをコンパイルすることは、

前記実行モデルの前記パイプラインポロジーにおいて識別されたバッファのサイズ、及び

前記実行モデルの前記パイプラインポロジーにおいて識別されたカーネルの重要性

のうちの1つ以上に少なくとも基づいて前記コマンドリストをコンパイルすることを備えるC2に記載の方法。

[C6]

前記パイプラインポロジーを受信することは、

前記パイプラインポロジーを形成するために相互に接続された1つ以上のカーネル及び1つ以上のバッファを示すコマンドリストを受信することを備えるC1に記載の方法。

[C7]

装置であって、

グラフィックス処理ユニット(GPU)と、

プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプラインポロジーのインディケーションを受信し、

前記実行モデルの前記パイプラインポロジーが前記GPUに実装されるプラットフォームに依存した方法を指示する命令を生成し、及び

前記GPUに前記命令を送信するように構成されたプロセッサと、を備え、前記実行モデルの前記パイプラインポロジーが前記GPUに実装される前記プラットフォームに依存した方法は、前記GPUのプラットフォームに基づく、装置。

[C8]

前記命令を生成するために、前記プロセッサは、

前記実行モデルの前記パイプラインポロジーを指定するコマンドリストをコンパイルように構成されるC7に記載の装置。

[C9]

前記コマンドリストをコンパイルするために、前記プロセッサは、

前記実行モデルの前記パイプラインポロジーが前記GPUに実装される前記プラットフォームに依存した方法を定義する前記命令を生成するために前記GPUの構成情報に少なくとも基づいて前記コマンドリストをコンパイルするように構成されるC8に記載の装置

。

[C10]

前記GPUの前記構成情報は、

前記GPU内のプログラマブル計算ユニット数、及び

前記GPU内のデータレーン数

のうちの1つ以上を備えるC9に記載の装置。

[C11]

前記コマンドリストをコンパイルするために、前記プロセッサは、

前記実行モデルの前記パイプラインポロジーにおいて識別されたバッファのサイズ、及び

前記実行モデルの前記パイプラインポロジーにおいて識別されたカーネルの重要性

のうちの1つ以上に少なくとも基づいて前記コマンドリストをコンパイルするように構成されるC8に記載の装置。

10

20

30

40

50

[C 1 2]

前記パイプラインポロジを受信するために、前記プロセッサは、
前記パイプラインポロジを形成するために相互に接続された1つ以上のカーネル及び
1つ以上のバッファを示すコマンドリストを受信するように構成されるC 7に記載の装置
。

[C 1 3]

前記装置は、
メディアプレーヤー、
セットトップボックス、
無線ハンドセット、
デスクトップコンピュータ
ラップトップコンピュータ
ゲームコンソール
ビデオ会議装置、及び
タブレットコンピューティングデバイス、
のうちの1つを備えるC 7に記載の装置。

10

[C 1 4]

コンピュータによって読み取り可能な記憶媒体であって、
1つ以上のプロセッサによって実行されたときに、
プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパ
イプラインポロジのインディケーションを受信し、
前記実行モデルの前記パイプラインポロジが前記G P Uに実装されるプラットフォーム
に依存した方法を指示する命令を生成し、及び
前記G P Uに前記命令を送信することを前記1つ以上のプロセッサに行わせる命令が格納
されており、前記実行モデルの前記パイプラインポロジが前記G P Uに実装される前
記プラットフォームに依存した方法は、前記G P Uのプラットフォームに基づく、コンピ
ュータによって読み取り可能な記憶媒体。

20

[C 1 5]

命令を生成することを前記1つ以上のプロセッサに行わせる前記命令は、
前記命令を生成するために前記実行モデルの前記パイプラインポロジを指定するコマ
ンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる命令を備えるC
1 4に記載のコンピュータによって読み取り可能な記憶媒体。

30

[C 1 6]

前記コマンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる前記命
令は、
前記実行モデルの前記パイプラインポロジが前記G P Uに実装される前記プラットフ
ォームに依存した方法を定義する前記命令を生成するために前記G P Uの構成情報に少な
くとも基づいて前記コマンドリストをコンパイルすることを前記1つ以上のプロセッサに
行わせる命令を備えるC 1 5に記載のコンピュータによって読み取り可能な記憶媒体。

40

[C 1 7]

前記G P Uの前記構成情報は、
前記G P U内のプログラマブル計算ユニット数、及び
前記G P U内のデータレーン数
のうちの1つ以上を備えるC 1 6に記載のコンピュータによって読み取り可能な記憶媒体
。

[C 1 8]

前記コマンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる前記命
令は、
前記実行モデルの前記パイプラインポロジにおいて識別されたバッファのサイズ、及
び

50

前記実行モデルの前記パイプラインポロジにおいて識別されたカーネルの重要性のうち、1つ以上に少なくとも基づいて前記コマンドリストをコンパイルすることを前記1つ以上のプロセッサに行わせる命令を備えるC 1 5に記載のコンピュータによって読み取り可能な記憶媒体。

[C 1 9]

装置であって、

グラフィックス処理ユニット（GPU）と、

プラットフォームから独立した方法でデータ処理アルゴリズムを定義する実行モデルのパイプラインポロジを受信するための手段と、

前記実行モデルの前記パイプラインポロジが前記GPUに実装されるプラットフォームに依存した方法を指示する命令を生成するための手段であって、前記実行モデルの前記パイプラインポロジが前記GPUで実装される前記プラットフォームに依存した方法は、前記GPUのプラットフォームに基づく手段と、

前記GPUに前記命令を送信するための手段と、を備える、装置。

[C 2 0]

命令を生成するための前記手段は、

前記命令を生成するために前記実行モデルの前記パイプラインポロジを指定するコマンドリストをコンパイルするための手段を備えるC 1 9に記載の装置。

10

【図 1】

図 1

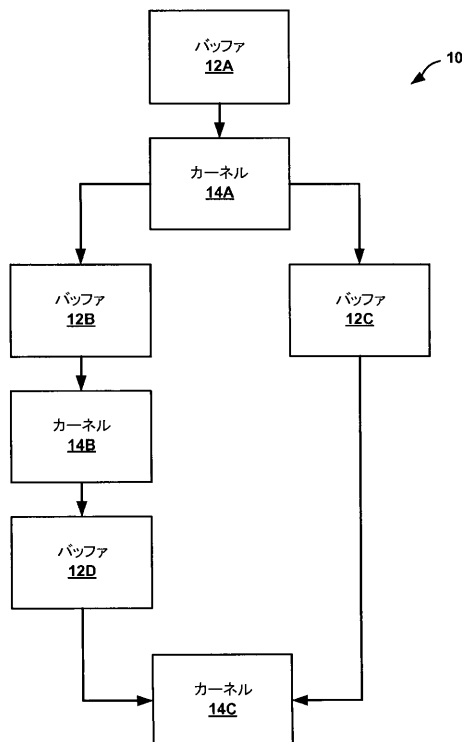


FIG. 1

【図 2】

図 2

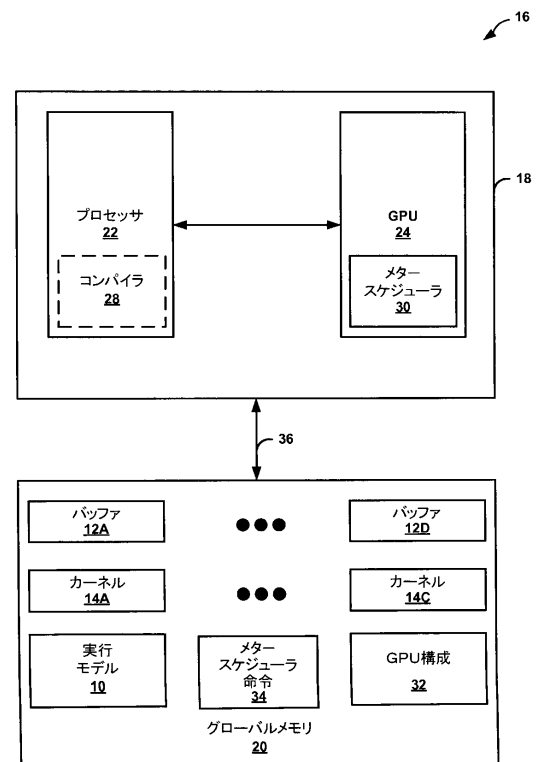


FIG. 2

【図 3】

図 3

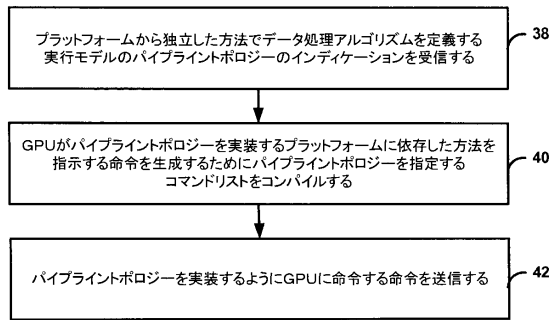


FIG. 3

【図 4】

図 4

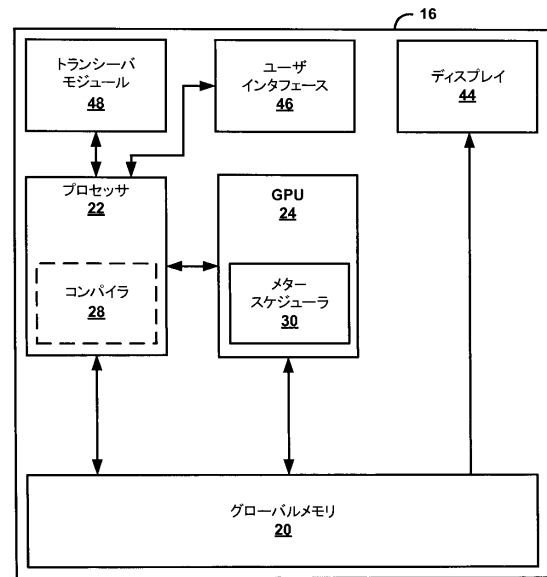


FIG. 4

 フロントページの続き

- (74)代理人 100153051
弁理士 河野 直樹
- (74)代理人 100140176
弁理士 砂川 克
- (74)代理人 100158805
弁理士 井関 守三
- (74)代理人 100179062
弁理士 井上 正
- (74)代理人 100124394
弁理士 佐藤 立志
- (74)代理人 100112807
弁理士 岡田 貴志
- (74)代理人 100111073
弁理士 堀内 美保子
- (72)発明者 ボウルド、アレクセイ・バイ .
アメリカ合衆国、カリフォルニア州 9 2 1 2 1、サン・ディエゴ、モアハウス・ドライブ 5 7
7 5
- (72)発明者 トーゼブスキー、ウィリアム・エフ .
アメリカ合衆国、カリフォルニア州 9 2 1 2 1、サン・ディエゴ、モアハウス・ドライブ 5 7
7 5

審査官 坂庭 剛史

- (56)参考文献 特開 2 0 1 1 - 2 0 4 2 0 9 (J P , A)
特開 2 0 1 1 - 1 7 5 6 2 4 (J P , A)
特開 2 0 0 4 - 1 7 1 2 3 4 (J P , A)
Andrei Hagiescu et al. , Automated architecture-aware mapping of streaming application
onto GPUs , 2011 IEEE International Parallel & Distributed Processing Symposium , 米国 ,
IEEE , 2 0 1 1 年 5 月 1 6 日 , pp.467-478 , ISBN:978-1-61284-372-8
道浦 悌、大野和彦、佐々木敬泰、近藤利夫 , G P G P U におけるデータ転送自動化コンパイラ
の設計 , 情報処理学会研究報告 2 0 1 1 (平成 2 3) 年度 2 [C D - R O M] , 日本 , 一
般社団法人情報処理学会 , 2 0 1 1 年 8 月 1 5 日 , Vol.2011-HPC-130 , No.17 , p p . 1 ~ 9
, ISSN 1884-0930

- (58)調査した分野(Int.Cl. , D B 名)
G 0 6 F 9 / 4 5
G 0 6 F 9 / 3 8