

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6081492号  
(P6081492)

(45) 発行日 平成29年2月15日 (2017.2.15)

(24) 登録日 平成29年1月27日 (2017.1.27)

(51) Int. Cl.

F I

G 0 6 F 15/167 (2006.01)

G 0 6 F 15/167 6 1 0 B

G 0 6 F 9/38 (2006.01)

G 0 6 F 9/38 3 7 0 C

G 0 6 T 1/20 (2006.01)

G 0 6 T 1/20 B

請求項の数 23 (全 35 頁)

(21) 出願番号	特願2014-554817 (P2014-554817)	(73) 特許権者	595020643
(86) (22) 出願日	平成25年1月24日 (2013.1.24)		クゥアルコム・インコーポレイテッド
(65) 公表番号	特表2015-513715 (P2015-513715A)		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/US2013/022900		1 2 1 - 1 7 1 4、サン・ディエゴ、モア
(87) 国際公開番号	W02013/112692		ハウス・ドライブ 5 7 7 5
(87) 国際公開日	平成25年8月1日 (2013.8.1)	(74) 代理人	100108855
審査請求日	平成27年12月28日 (2015.12.28)		弁理士 蔵田 昌俊
(31) 優先権主張番号	61/591,733	(74) 代理人	100109830
(32) 優先日	平成24年1月27日 (2012.1.27)		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	13/747,947		弁理士 野河 信久
(32) 優先日	平成25年1月23日 (2013.1.23)	(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】 グラフィックス並列処理ユニットに関するバッファ管理

(57) 【特許請求の範囲】

【請求項 1】

データ処理動作をパイプライン方式で実行するための方法であって、

グラフィックス処理ユニット (GPU) のシェーダプロセッサの第1のプログラマブル計算ユニットにおいて第1のスレッドを実行することであって、前記シェーダプロセッサは前記第1のプログラマブル計算ユニットを含む複数のプログラマブル計算ユニットを含むことと、

前記 GPU の前記シェーダプロセッサの前記複数のプログラマブル計算ユニットの第2のプログラマブル計算ユニットにおいて第2のスレッドを実行することと、

前記 GPU を含む集積回路 (IC) 内の管理ユニットを用いて、前記第1のスレッドの前記実行によって生成されたデータを前記複数のプログラマブル計算ユニットによって共有される前記 IC の外部の集積グローバルメモリ内のバッファ内に格納することの要求を前記第1のプログラマブル計算ユニットから受信することであって、前記第1のスレッドの前記実行によって生成された前記データは、前記第2のスレッドを実行する前記第2のプログラマブル計算ユニットによって消費されることになり、前記バッファは、先入れ先出し (FIFO) バッファを備えることと、

前記管理ユニットを用いて、前記第1のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の記憶場所を決定することと、

前記 IC を用いて、前記第1のスレッドの前記実行によって生成された前記データを前記バッファ内の前記決定された記憶場所に格納することと、を備える、方法。

10

20

**【請求項 2】**

前記管理ユニットを用いて、前記バッファの状態情報を前記 IC 内に格納することをさらに備え、

前記バッファの前記状態情報は、前記バッファの開始アドレス、前記バッファの終了アドレス、生成されたデータが格納されるべき前記バッファ内のアドレス、及びデータが取り出されるべき前記バッファ内のアドレスのうちの 1 つ以上を含み、

前記バッファ内の前記記憶場所を決定することは、前記バッファの前記格納された状態情報に基づいて前記第 1 のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の前記記憶場所を決定することを備える請求項 1 に記載の方法。

**【請求項 3】**

前記管理ユニットを用いて、前記第 1 のスレッドの前記実行によって生成された前記データの少なくとも一部を取り出すことの要求を前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットから受信することと、

前記管理ユニットを用いて、前記第 1 のスレッドの前記実行によって生成される前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるかどうかを決定することと、をさらに備える請求項 1 に記載の方法。

**【請求項 4】**

前記第 2 のプログラマブル計算ユニットから前記要求を受信することは、前記第 1 のスレッドの前記実行によって生成されたデータを格納することの前記要求を前記第 1 のプログラマブル計算ユニットから受信すると同時に、受信する前に、又は受信した後に前記第 2 のプログラマブル計算ユニットから前記要求を受信することを備える請求項 3 に記載の方法。

**【請求項 5】**

前記第 2 のスレッドによって要求された前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能でないときに、

前記管理ユニットを用いて、第 3 のスレッドを実行するように前記第 2 のプログラマブル計算ユニットに指示することと、

前記管理ユニットを用いて、前記第 2 のスレッドによって要求された前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに前記第 2 のプログラマブル計算ユニットに指示することと、

前記管理ユニットを用いて、前記第 2 のスレッドによって要求された前記データが前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに前記第 2 のスレッドによって要求された前記データを消費するために前記第 2 のスレッドを実行するように前記第 2 のプログラマブル計算ユニットに指示することと、をさらに備える請求項 3 に記載の方法。

**【請求項 6】**

前記管理ユニットを用いて、前記第 2 のスレッドによって要求された前記データに加えて前記集積グローバルメモリからデータを取り出すことと、

前記管理ユニットを用いて、前記第 2 のスレッドによって要求された前記データに加えての前記データを前記 IC 内のキャッシュに格納することと、をさらに備える請求項 3 に記載の方法。

**【請求項 7】**

前記第 1 のスレッドを実行することは、カーネルの生成するスレッドを実行することを備え、前記第 2 のスレッドを実行することは、前記カーネルの消費するスレッドを実行することを備える請求項 1 に記載の方法。

**【請求項 8】**

前記第 1 のスレッドを実行することは、生成するカーネルの前記第 1 のスレッドを実行

10

20

30

40

50

することを備え、前記第2のスレッドを実行することは、消費するカーネルのスレッドを実行することを備える請求項1に記載の方法。

【請求項9】

前記GPUは、前記管理ユニットを含み、前記FIFOバッファは、リングバッファを備える、請求項1に記載の方法。

【請求項10】

前記バッファ内の前記記憶場所を決定することは、前記バッファ内において前記データが格納されるべき前記記憶場所を前記第1のスレッドが示さずに、前記第1のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の前記記憶場所を決定することを備える請求項1に記載の方法。

10

【請求項11】

装置であって、

複数のプログラマブル計算ユニットによって共有される集積グローバルメモリであって、前記集積グローバルメモリはバッファを含み、前記バッファは、先入れ先出し(FIFO)バッファを備える、前記集積グローバルメモリと、

集積回路(IC)であって、

グラフィックス処理ユニット(GPU)であって、

前記複数のプログラマブル計算ユニットを備え、

前記複数のプログラマブル計算ユニットの第1のプログラマブル計算ユニットは、第1のスレッドを実行するように構成され、

20

前記複数のプログラマブル計算ユニットの第2のプログラマブル計算ユニットは、第2のスレッドを実行するように構成された前記GPUと、

前記第1のスレッドの前記実行によって生成されたデータを前記集積グローバルメモリ内の前記バッファ内に格納することの要求を前記第1のプログラマブル計算ユニットから受信し、及び

前記第1のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の記憶場所を決定するように構成された管理ユニットと、を備える、集積回路(IC)とを備え、前記第1のスレッドの前記実行によって生成された前記データは、前記第2のスレッドを実行する前記第2のプログラマブル計算ユニットによって消費され、

前記ICは、前記第1のスレッドの前記実行によって生成された前記データを前記バッファ内の前記決定された記憶場所に格納するように構成される、装置。

30

【請求項12】

前記管理ユニットは、前記バッファの状態情報を前記IC内に格納するように構成され、

前記バッファの前記状態情報は、前記バッファの開始アドレス、前記バッファの終了アドレス、生成されたデータが格納されるべき前記バッファ内のアドレス、及びデータが取り出されるべき前記バッファ内のアドレスのうちの1つ以上を含み、

前記管理ユニットは、前記バッファの前記格納された状態情報に基づいて前記第1のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の前記記憶場所を決定するように構成される請求項11に記載の装置。

40

【請求項13】

前記管理ユニットは、

前記第1のスレッドの前記実行によって生成された前記データの少なくとも一部を取り出すことの要求を前記第2のスレッドを実行する前記第2のプログラマブル計算ユニットから受信し、及び

前記第1のスレッドの前記実行によって生成される前記データが、前記第2のスレッドを実行する前記第2のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるかどうかを決定するように構成される請求項11に記載の装置。

【請求項14】

前記管理ユニットは、前記第1のスレッドの前記実行によって生成されたデータを格納

50

することの前記要求を前記第 1 のプログラマブル計算ユニットから受信すると同時に、受信する前に、又は受信した後に前記第 2 のプログラマブル計算ユニットから前記要求を受信するように構成される請求項 1 3 に記載の装置。

【請求項 1 5】

前記管理ユニットは、

前記第 2 のスレッドによって要求された前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能でないときに、第 3 のスレッドを実行するように前記第 2 のプログラマブル計算ユニットに指示し、

前記第 2 のスレッドによって要求された前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに前記第 2 のプログラマブル計算ユニットに指示し、及び

前記第 2 のスレッドによって要求された前記データが前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに前記第 2 のスレッドによって要求された前記データを消費するために前記第 2 のスレッドを実行するように前記第 2 のプログラマブル計算ユニットに指示するように構成される請求項 1 3 に記載の装置。

【請求項 1 6】

前記管理ユニットは、

前記第 2 のスレッドによって要求された前記データに加えて前記グローバルメモリからデータを取り出し、及び

前記第 2 のスレッドによって要求された前記データに加えて前記データを前記 IC 内のキャッシュに格納するように構成される請求項 1 3 に記載の装置。

【請求項 1 7】

前記第 1 のスレッドは、カーネルの生成するスレッドを備え、前記第 2 のスレッドは、前記カーネルの消費するスレッドを備える請求項 1 1 に記載の装置。

【請求項 1 8】

前記第 1 のスレッドは、生成するカーネルのスレッドを備え、前記第 2 のスレッドは、消費するカーネルのスレッドを備える請求項 1 1 に記載の装置。

【請求項 1 9】

前記 GPU は、前記管理ユニットを含み、前記 FIFO バッファはリングバッファを備える、請求項 1 1 に記載の装置。

【請求項 2 0】

前記管理ユニットは、前記第 1 のスレッドが前記バッファ内で前記データが格納されるべき前記記憶場所を示さずに前記第 1 のスレッドの実行によって生成された前記データが格納されるべき前記バッファ内の前記記憶場所を決定するように構成される請求項 1 1 に記載の装置。

【請求項 2 1】

前記装置は、映像デバイス、セットトップボックス、無線ハンドセット、パーソナルデジタルアシスタント、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、ビデオ会議ユニット、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 1 1 に記載の装置。

【請求項 2 2】

装置であって、

バッファを含む複数のプログラマブル計算ユニットによって共有される集積グローバルメモリであって、前記バッファは、先入れ先出し (FIFO) バッファを備える集積グローバルメモリと、

集積回路 (IC) であって、

グラフィックス処理ユニット (GPU) であって、

第 1 のスレッドを実行するための手段と、

第2のスレッドを実行するための手段と、

前記第1のスレッドの前記実行によって生成されたデータを前記集積グローバルメモリ内の前記バッファ内に格納することの要求を前記第1のスレッドを実行するための前記手段から受信するための手段であって、前記第1のスレッドの前記実行によって生成された前記データは、前記第2のスレッドを実行するための前記手段によって消費される手段と、

前記第1のスレッドを実行するための前記手段によって生成された前記データが格納されるべき前記バッファ内の記憶場所を決定するための手段と、

前記第1のスレッドの前記実行によって生成された前記データを前記バッファ内の前記決定された記憶場所に格納するための手段と、を備えるグラフィックス処理ユニット（GPU）、を備える集積回路（IC）と、を備える、装置。

10

【請求項23】

コンピュータによって読み取り可能な非一時的な記憶媒体であって、  
実行されたときに、

グラフィックス処理ユニット（GPU）のシェーダプロセッサの第1のプログラマブル計算ユニットにおいて第1のスレッドを実行し、ここにおいて、前記シェーダプロセッサは前記第1のプログラマブル計算ユニットを含む複数のプログラマブル計算ユニットを含む、前記第1のスレッドを実行し、

前記GPUの前記シェーダプロセッサの前記複数のプログラマブル計算ユニットの第2のプログラマブル計算ユニットにおいて第2のスレッドを実行し、

20

前記GPUを含む集積回路（IC）内の管理ユニットを用いて、前記第1のスレッドの前記実行によって生成されたデータを前記複数のプログラマブル計算ユニットによって共有される前記ICの外部の集積グローバルメモリ内のバッファ内に格納することの要求を前記第1のプログラマブル計算ユニットから受信し、

前記管理ユニットを用いて、前記第1のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の記憶場所を決定し、及び

前記ICを用いて、前記第1のスレッドの前記実行によって生成された前記データを前記バッファ内の前記決定された記憶場所に格納することを1つ以上のプロセッサに行わせる命令が格納されており、前記第1のスレッドの前記実行によって生成された前記データは、前記第2のスレッドを実行する前記第2のプログラマブル計算ユニットによって消費されることになり、前記バッファは、先入れ先出し（FIFO）バッファを備える、コンピュータによって読み取り可能な非一時的な記憶媒体。

30

【発明の詳細な説明】

【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61/591,733号（出願日：2012年1月27日）の利益を主張するものである。

【0002】

本開示は、メモリアクセス管理に関するものである。本開示は、より具体的には、グラフィックス処理デバイス（GPU）におけるメモリアクセス管理に関するものである。

40

【背景技術】

【0003】

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（GP GPU）と呼ばれるGPUが得られている。例えば、GP GPU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。

50

## 【発明の概要】

## 【0004】

概して、本開示は、グローバルメモリ内に存在しておりグラフィックス処理ユニット（GPU）のためのデータを格納するバッファを管理するための技法に関するものである。例えば、GPUを含む集積回路（IC）チップは、パイプライン管理ユニットを含む。パイプライン管理ユニットは、グローバルメモリ内の1つ以上のバッファの状態情報を維持するように構成することができる。GPUで実行中のアプリケーションがグローバルメモリ内のバッファにアクセスするときには、グローバルメモリ内のバッファの状態情報は、ICチップ内部において入手可能である。このように、GPUは、グローバルメモリ内のバッファの状態情報を決定するためにオフチップメモリアクセスを行う必要がない。

10

## 【0005】

一例では、本開示は、データ処理動作をパイプライン方式で実行するための方法について説明する。その方法は、グラフィックス処理ユニット（GPU）のシェーダプロセッサの第1のプログラマブル計算ユニット（programmable compute unit）において第1のスレッドを実行することと、GPUのシェーダプロセッサの第2のプログラマブル計算ユニットにおいて第2のスレッドを実行することと、を含む。その方法は、GPUを含む集積回路（IC）内の管理ユニットを用いて、第1のスレッドの実行によって生成されたデータをIC外部のグローバルメモリ内のバッファに格納することの要求を第1のプログラマブル計算ユニットから受信することを含む。この例では、第1のスレッドの実行によって生成されたデータは、第2のスレッドを実行する第2のプログラマブル計算ユニットによって消費される。さらに、この例では、バッファは、先入れ先出し方式（FIFO）のバッファ及びリングバッファのうちの1つを備える。その方法は、管理ユニットを用いて、第1のスレッドの実行によって生成されたデータが格納されるべきバッファ内の記憶場所（location）を決定することと、ICを用いて、第1のスレッドの実行によって生成されたデータをバッファ内の決定された記憶場所に格納することと、を含む。

20

## 【0006】

一例では、本開示は、装置について説明する。その装置は、バッファを含むグローバルメモリを含む。この例では、バッファは、先入れ先出し方式（FIFO）のバッファ及びリングバッファのうちの1つを備える。その装置は、グラフィックス処理ユニット（GPU）と管理ユニットとを含む集積回路（IC）も含む。GPUは、第1のスレッドを実行するように構成された第1のプログラマブル計算ユニットと、第2のスレッドを実行するように構成された第2のプログラマブル計算ユニットと、を含む。管理ユニットは、第1のスレッドの実行によって生成されたデータをグローバルメモリ内のバッファに格納することの要求を第1のプログラマブル計算ユニットから受信するように構成される。この例では、第1のスレッドの実行によって生成されたデータは、第2のスレッドを実行する第2のプログラマブル計算ユニットによって消費される。管理ユニットは、第1のスレッドの実行によって生成されたデータが格納されるべきバッファ内の記憶場所を決定するようにも構成される。この例では、ICは、第1のスレッドの実行によって生成されたデータをバッファ内の決定された記憶場所に格納するように構成される。

30

40

## 【0007】

一例では、本開示は、装置について説明する。その装置は、グローバルメモリと、集積回路（IC）と、を含む。グローバルメモリは、バッファを含む。この例では、バッファは、先入れ先出し方式（FIFO）のバッファ及びリングバッファのうちの1つを備える。ICは、第1のスレッドを実行するための手段と、第2のスレッドを実行するための手段と、を備えるグラフィックス処理ユニット（GPU）を含む。ICは、第1のスレッドの実行によって生成されたデータをグローバルメモリ内のバッファ内に格納することの要求を第1のスレッドを実行するための手段から受信するための手段も含む。この例では、第1のスレッドの実行によって生成されたデータは、第2のスレッドを実行するための手段によって消費される。ICは、第1のスレッドを実行するための手段によって生成され

50

たデータが格納されるべきバッファ内の記憶場所を決定するための手段と、第1のスレッドの実行によって生成されたデータをバッファ内の決定された記憶場所に格納するための手段と、も含む。

【0008】

一例では、本開示は、実行されたときに、グラフィックス処理ユニット(GPU)のシェーダプロセッサの第1のプログラマブル計算ユニットにおいて第1のスレッドを実行すること、及びGPUのシェーダプロセッサの第2のプログラマブル計算ユニットにおいて第2のスレッドを実行することを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、GPUを含む集積回路(IC)内の管理ユニットを用いて、第1のスレッドの実行によって生成されたデータをICの外部のグローバルメモリ内のバッファ内に格納することの要求を第1のプログラマブル計算ユニットから受信することも1つ以上のプロセッサに行わせる。この例では、第1のスレッドの実行によって生成されたデータは、第2のスレッドを実行する第2のプログラマブル計算ユニットによって消費される。さらに、この例では、バッファは、先入れ先出し方式(FIFO)のバッファ及びリングバッファのうちの1つを備える。命令は、管理ユニットを用いて、第1のスレッドの実行によって生成されたデータが格納されるべきバッファ内の記憶場所を決定すること、及び、ICを用いて、第1のスレッドの実行によって生成されたデータをバッファ内の決定された記憶場所に格納することも1つ以上のプロセッサに行わせる。

【0009】

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。その説明と図面から、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】

【0010】

【図1】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。

【図2】グラフィックス処理ユニット(GPU)及びグローバルメモリをさらに詳細に例示したブロック図である。

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

【図5】図1のデバイスをさらに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】

【0011】

グラフィックス処理ユニット(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ー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、シェーダプログラム、例えば、バーテックスシェーダ(vertex shader)(頂点シェーダ)、ハルシェーダ(hull shader)、フラグメントシェーダ(fragment shader)、幾何シェーダ(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。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。

【0012】

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(GPU)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、GP GPUとして機能している。ほとんどすべてのGPUは、GP GPUとして機能するように構成することができる。

#### 【0013】

例示を目的として、本開示は、GP GPU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GP GPUとして機能している（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を実行している）事例には限定されない。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。

10

#### 【0014】

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、スレッド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。

20

#### 【0015】

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。

30

#### 【0016】

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GP GPU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。

#### 【0017】

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ことができる。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

40

50



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)。

【0018】

概して、GPUは、単一プログラム多データ(SPMD)プログラミングモデルを実装するように構成することができる。SPMDプログラミングモデルでは、GPUは、(例えば、スレッドとして)複数のプログラマブル計算ユニットでカーネルを実行することができ、各プログラマブル計算ユニットは、それ自体のデータに関する機能を実行する。さらに、SPMDプログラミングモデルでは、プログラマブル計算ユニットは、現在の命令がプログラマブル計算ユニットによって実行されていることを示す各々のプログラムカウンタを含む。

10

【0019】

GPUは、大規模な並列性を処理のために提供する一方で、GPUは、パイプライン方式でカーネルを実行するのにはあまり適していない。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ことになる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。本開示では、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。

20

【0020】

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。

【0021】

例えば、パイプライン方式でカーネルを実行することは、第1のスレッド(例えば、カーネルの生成するスレッド)が第2のスレッド(例えば、同じカーネル又は異なるカーネルの消費するスレッド)によって消費されるデータを生成することと考えることができる。(第1のスレッドに関する消費者であった)この第2のスレッドは、第3のスレッドに関する生成するスレッドであることができる(例えば、第2のスレッドが、第3のスレッドによって消費されるデータを生成する)。第3のスレッドは、第1及び第2のスレッドを含むカーネルと異なるカーネルに関するスレッドであることができ又は第1及び第2のスレッドを含むカーネルのうちの1つに関するスレッドであることができる。この例では、第1、第2、及び第3のスレッドは、処理パイプラインを形成するとみなすことができる。

30

【0022】

パイプライン方式でカーネルを実行することは、カーネル又はスレッドを逐次に(例えば、次々に)実行することを要求するものであるとは解釈されるべきでない。例えば、上例では、GPUが第1、第2、及び第3のスレッドのうちの2つ以上を並列で(例えば、同時に)実行することが可能である。しかしながら、GPUがスレッドを逐次で実行することも可能であり、この場合も依然として、パイプライン方式でカーネルを実行するとみなすことができる。

40

【0023】

カーネルの生成するスレッドを実行するプログラマブル計算ユニットは、生成されたデータをグローバルメモリ(すなわち、GPUを含む集積回路(IC)の外部のオフチップ、システムメモリ)に出力する必要がある。ここで、グローバルメモリは、例えば、システムバスを介してアクセス可能であることができる。同じカーネル又は異なるカーネルの消費するスレッドを実行する他のプログラマブル計算ユニットは、グローバルメモリから

50

生成されたデータを受信する必要がある。さらに詳細に説明されるように、既存のGPUに関しては、グローバルメモリの管理は、計算、時間、及び/又は電力の点で非効率的になり、その結果、パイプライン方式でカーネルを実行するときに性能不良になる可能性がある。

#### 【0024】

本開示は、計算、時間、及び電力の点で効率的なグローバルメモリの管理のための技法について説明する。より詳細に説明されるように、GPUを含む集積回路(IC)は、パイプライン管理ユニット(PMU)を含むことができる。代替として、GPU自体がPMUを含むことができる。PMUは、消費されるべき生成データを格納するグローバルメモリの状態情報を管理するように構成することができる。例えば、プロセッサ又はGPU自体は、グローバルメモリに格納されるべきプログラマブル計算ユニットによって生成されたデータをグローバルメモリに格納するときの記憶場所をグローバルメモリ内において予約することができる。グローバルメモリ内のこれらの予約された記憶場所は、複数のバッファとみなすことができる。幾つかの例では、複数のバッファは、リングバッファ又は先入れ先出し(FIFO)バッファを形成することができる。リングバッファは、FIFOバッファの一例であるとみなすことができる。

10

#### 【0025】

PMUは、オフチップグローバルメモリ内のバッファの状態情報を示す情報を、IC又はGPUの内部に(例えば、オンチップの内部キャッシュメモリ内に)格納することができる。一例として、PMUは、グローバルメモリ内のバッファの開始アドレス及び終了アドレスを示す情報を格納することができる。他の例として、PMUは、生成されたデータが格納されるべき複数のバッファ内にバッファのアドレスを、及び消費されるべきデータが読み取られる複数のバッファ内にバッファのアドレスを格納することができる。さらに他の例として、データを必要とする消費する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データを必要としない消費するカーネルのその他のスレッドの実行を続けることができるように、生成するカーネルがデータの生成を完了したかどうかを示す情報を格納することができる。

20

#### 【0026】

本開示において説明される技法では、PMUは、生成するスレッドによって生成されたデータをバッファ内に格納することの要求を受信することができ、及び、消費するスレッドによる消費のために生成するスレッドによって生成されたデータをバッファから取り出すことの要求を受信することができる。PMUは、バッファの格納された状態情報に基づいて、生成するスレッドの実行によって生成されたデータが格納されるべきバッファ内の記憶場所を決定することができ、及び、バッファの格納された状態情報に基づいて消費するバッファによって消費されるべきデータが取り出されるべき記憶場所を決定することができる。

30

#### 【0027】

GPUを含むICの内部又はGPU自体の内部に格納された情報を用いてグローバルメモリの状態情報を管理することによって、本開示において説明される技法は、GPUがグローバルメモリにアクセスする必要がある回数を最小限にすることができる。例えば、PMUは、GPUを含むICの外部の該情報にアクセスすることによってデータが格納される又は取り出されるアドレスを決定する必要がある。GPUがグローバルメモリにアクセスする必要がある回数を最小限にすることは、電力消費量を低減させ、システムバス帯域幅負荷を低減させ、及びレーテンシーを短縮することができる。

40

#### 【0028】

さらに、以下においてより詳細に説明されるように、既存のGPUでは、カーネルがグローバルメモリを管理する命令を含む必要がある。GPUは、該グローバルメモリ管理命令を実行するクロックサイクルを浪費することがあり、計算上非効率的である可能性がある。PMUがグローバルメモリの状態情報を管理することで、カーネルは、グローバルメモリ管理命令を含める必要がなく、その結果、カーネル命令の複雑さが低下し、さらに、

50

実行する必要があるカーネル命令が少なくなる。このように、本開示において説明される技法は、計算効率を促進させることができる。

【0029】

図1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。デバイス10の例は、映像デバイス、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、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、ゲームコンソール、ビデオ会議ユニット、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10は、図1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。

10

【0030】

例示されるように、デバイス10は、集積回路(IC)12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10のためのメモリであるとみなすことができる。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)、電氣的消去可能プログラマブル読み取り専用メモリ(EEPROM)、フラッシュメモリ、又は、希望されるプログラムコードを命令及び/又はデータ構造の形態で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可能であ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。

20

【0031】

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14及びGPU16に帰する機能を果たすことをプロセッサ14及び/又はGPU16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14及びGPU16)に行わせる命令を格納し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。

【0032】

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た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が取り外し不能であるか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、デバイス10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を、デバイス10に挿入することができる。幾つかの例では、非一時的記憶媒体は、(例えば、RAMにおいて)経時で変化する可能性があるデータを格納することができる。

30

【0033】

IC12は、プロセッサ14と、グラフィックス処理ユニット(GPU)16と、パイプライン管理ユニット(PMU)18と、を含む。IC12は、プロセッサ14、GPU16、及びPMU18を収納する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、IC12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。PMU18は、IC12の一部を形成するハードウェアユニットであることができ又はGPU16内のハードウェアであることができる。PMU18は、IC12内又はGPU内のハードウェアで実行されるソフトウェアであることが可能である。例示及び説明の目的上、技法は、ハードウェアユニットであるPMU18に関して説明される。

40

【0034】

プロセッサ14、GPU16、及びPMU18は、単一のIC12の一部として例示されているが、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、プロセッサ14及びGPU16は、異なる集積回路(すなわち、異なるチップパッケージ)内に収納す

50

ることができる。これらの例では、PMU 18は、GPU 16と同じ集積回路内に収納することができる。幾つかの例では、PMU 18は、GPU 16の一部として形成することができる。一例として、プロセッサ14及びGPU 16は、同じ集積回路（すなわち、同じチップパッケージ）内に収納することができ、PMU 18は、GPU 16内に形成することができる。他の例として、プロセッサ14及びGPU 16は、異なる集積回路（すなわち、異なるチップパッケージ）内に収納することができ、PMU 18は、GPU 16内に形成することができる。

#### 【0035】

プロセッサ14、GPU 16、及びPMU 18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ロジック回路を含み、ただしこれらに限定されない。幾つかの例では、GPU 16及びPMU 18は、以下においてより詳細に説明されるように、グラフィックス処理に適する大規模な並列処理能力をGPU 16に提供し及び管理するグローバルメモリ20をPMU 18に提供する集積回路及び／又はディスクリートロジック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、GPU 16は、汎用処理も含むことができ、汎用処理タスク（例えば、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。

#### 【0036】

プロセッサ14は、ホストと時々呼ばれ、デバイス10の中央処理装置（CPU）であることができる。プロセッサ14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、又は、表示のために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。

#### 【0037】

幾つかの例では、プロセッサ14は、処理タスク、例えば、大規模な並列動作を要求するタスク、をGPU 16に委ねることができる。一例として、グラフィックス処理は、大規模な並列動作を要求し、プロセッサ14は、該グラフィックス処理タスクをGPU 16に委ねることができる。幾つかの例では、プロセッサ14は、グラフィックス処理に関連しないタスクをGPU 16に委ねることができる。例えば、行列演算は、並列動作を要求し、GPU 16のほうがプロセッサ14と比較して該動作を実装するのに適している。

#### 【0038】

タスクを実装するために、GPU 16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、グラフィックス関連の処理に関しては、GPU 16は、パーティックスシェーダ、フラグメントシェーダ、及び幾何学シェーダ、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、GPU 16は、該処理用に設計されたアプリケーション（例えば、行列演算を実装するためのアプリケーション又は流体力学に関するアプリケーション）を実行することができる。いずれの例に関しても（例えば、グラフィックス関連の処理又はグラフィックスに関連しない処理）、プロセッサ14は、1つ以上のアプリケーションを実行するようにGPU 16に命令することができる。

#### 【0039】

プロセッサ14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GPU 16と通信することができる。例えば、プロセッサ14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GPU 16に命令する命令、をGPU 16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

10

20

30

40

50

、本開示の態様は、Direct X、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14及びGPU16は、あらゆる通信技法を利用することができる。

#### 【0040】

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14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14は、OpenCL APIを用いてGPU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14がOpenGL及び/又はOpenCL APIを用いてGPU16と通信することは必ずしも要求しない。

10

#### 【0041】

GPU16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、GPU16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GPU16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、GPU16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。

20

#### 【0042】

シェーダの例は、グラフィックス関連の処理のためのバーテックスシェーダと、フラグメントシェーダと、幾何シェーダと、を含む。カーネルの例は、(例えば、線形代数又は流体力学に関する)グラフィックス処理に関連しないアプリケーションを含む。追加例として、カーネルは、画像処理及び映像処理に関するアプリケーションを含む。

#### 【0043】

GPU16は、シェーダプロセッサを含むことができ、シェーダプロセッサは、シェーダ及びカーネルを実行することができる。例えば、GPU16のシェーダプロセッサは、1つ以上のシェーダコア(プログラマブル計算ユニットと呼ばれる)を含むことができ、1つ以上のプログラマブル計算ユニットの各々がカーネルを実行することができる。

30

#### 【0044】

カーネルは、GPU16で実行するアプリケーションとして説明されるが、カーネルはそのように限定的であるとはみなされるべきでない。カーネルのその他の例は、GPU16の固定機能ユニットを含む。例えば、GPU16は、プログラマブル計算ユニットと、固定機能ユニットと、を含む。プログラマブル計算ユニットは、アプリケーションを実行することによって機能上の柔軟性を提供することができる。固定機能ユニットは、機能上の柔軟性を提供しないハードウェアユニットであることができ、特定の目的のために設計することができる。概して、用語カーネルは、グラフィックスに関連しない目的のためにデータを受信し、そのデータを処理し、及びそのデータを出力とするアプリケーション又はハードウェアユニットを意味する。しかしながら、例示を目的として、本開示において説明される技法は、カーネルが固定機能ユニットである例にまで拡張可能であるという理解の下でカーネルがアプリケーションである例を用いて説明される。

40

#### 【0045】

本開示において説明される技法では、1つのプログラマブル計算ユニットがカーネルの全命令を実行するのではなく、複数のプログラマブル計算ユニットがカーネルの一部を実行するのが可能である。カーネルの一部は、ワークグループ、タスク、又はスレッドと呼ぶことができる(すべて同意語である)。例えば、カーネルのワークグループ、タスク、又はスレッドは、そのカーネルのその他のワークグループ、タスク、又はスレッドから独立して実行することができる命令の組である。

#### 【0046】

50

幾つかの例では、1つ以上のプログラマブル計算ユニットの第1の組がカーネルのスレッドを実行することができ、1つ以上のプログラマブル計算ユニットの第2の組がカーネルのスレッドを実行することができる。幾つかの場合は、プログラマブル計算ユニットの第1の組及びプログラマブル計算ユニットの第2の組が実行するスレッドは、同じカーネルのスレッドであることができる。幾つかの場合は、プログラマブル計算ユニットの第1の組及びプログラマブル計算ユニットの第2の組が実行するスレッドは、異なるカーネルのスレッドであることができる。これらの例のいずれでも、スレッドのうちの1つは、生成されたデータをスレッドの他の1つに出力する必要がある。換言すると、GPU16は、パイプライン方式でカーネルを実行することができる。

【0047】

10

上述されるように、パイプライン方式でカーネルを実行することは、1つのスレッドによって生成されたデータが他のスレッドによって消費され、この他のスレッドによって生成されたデータがさらに他のスレッドによって消費され、以下同様であるような形でカーネルを実行することを意味することができる。これらの例では、スレッドは、異なるカーネル又は同じカーネルのスレッドであることができ、又は、幾つかのスレッドは異なるカーネルに関するものであり、その他のスレッドは同じカーネルに関するものであることができる。これらの例では、カーネルは、データが生成及び消費されるパイプラインを形成するとみることができる。例えば、同じカーネル又は異なるカーネルの第1、第2、及び第3のスレッドが1つのパイプを形成することができ、その中で、第1のスレッドがデータを生成し、処理のために第2のスレッドによって消費されるためにそのデータを送信する。第2のスレッドは、受信されたデータを処理してデータを生成し、生成されたデータを処理のために第3のスレッドに送信し、以下同様である。

20

【0048】

この例では、第1のスレッドは、生成するスレッドと呼ぶことができ、第2のスレッドは、第1のスレッドにとっての消費するスレッド、第3のスレッドにとっての生成するスレッドと呼ぶことができ、第3のスレッドは、消費するスレッドと呼ぶことができる。第1、第2、及び第3のスレッドが異なるカーネル（例えば、それぞれ、第1、第2、及び第3のカーネル）に関するものである例では、第1のカーネルは生成するカーネルと呼ぶことができ、第2のカーネルは、第1のカーネルにとっての消費するカーネル、第3のカーネルにとっての生成するカーネルと呼ぶことができ、第3のカーネルは、消費するカーネルと呼ぶことができる。

30

【0049】

既存のGPUでは、パイプライン方式でカーネルを実行することは、計算上及び電力上非効率的である可能性がある。例えば、プログラマブル計算ユニットの各々は、プログラマブル計算ユニットによって実行される命令を格納するための、処理されるデータを格納するための、及び生成されるデータを格納するためのローカルメモリを含むことができ、生成することができる中間結果を含む。しかしながら、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アクセスすることができない。

【0050】

40

従って、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ために、GPU16は、プログラマブル計算ユニットのローカルメモリに格納された生成データを取り出し、その生成データをグローバルメモリ20に格納することができる。グローバルメモリ20にデータを格納することは、データをオフチップで格納すると呼ぶことができ、その理由は、グローバルメモリ20は、GPU16を収納する集積回路の外部に（すなわち、IC12の外部に）存在するためである。GPU16は、グローバルメモリ20に格納されたデータを取り出し、取り出されたデータを他のプログラマブル計算ユニットのローカルメモリ内にローディングすることができる。

【0051】

説明するための例として、第1のプログラマブル計算ユニットが生成するカーネルのス

50

レッドを実行中であると仮定する。この例では、第1のプログラマブル計算ユニットは、生成するカーネルのスレッドの実行によって生成されたデータを第1のプログラマブル計算ユニットのローカルメモリに格納することができる。GPU16は、第1のプログラマブル計算ユニットのローカルメモリから生成されたデータを取り出し、生成されたデータをグローバルメモリ20内に格納することができる。

【0052】

この例において、第2のプログラマブル計算ユニットが消費するカーネルのスレッドを実行中であると仮定する。この例では、GPU16は、生成するカーネルによって生成されたデータをグローバルメモリ20から取り出し、第2のプログラマブル計算ユニットのローカルメモリ内にデータをローディングすることができる。これで、消費するカーネルは、第2のプログラマブル計算ユニットのローカルメモリ内に格納されたデータを消費することができる。

10

【0053】

上例において、GPU16は、第2のプログラマブル計算ユニットが第1のプログラマブル計算ユニットのローカルメモリにアクセスしないため、生成するカーネルによって生成されたデータをグローバルメモリ20内に格納する必要がある。このように、グローバルメモリ20は、後続して消費されることになる生成されたデータの中間的なストレージとして機能する。

【0054】

概して、生成されたデータがグローバルメモリ20に格納される方法及び/又はグローバルメモリ20からデータが取り出される方法を管理することは、処理上及び計算上非効率的になる可能性がある。一例として、非効率的ではあるが、データがグローバルメモリ20に格納される方法をカーネルが管理することが可能である。例えば、カーネルは、データを格納するための又は格納されたデータを取り出すためのグローバルメモリ20内のアドレス(例えば、ポインタ)を決定することをプログラマブル計算ユニットの算術論理ユニット(ALU)に行わせる命令を含むことが可能である。

20

【0055】

他の例として、グローバルメモリ20は、原子カウンタを格納することができる。原子カウンタの値は、消費されるデータを入手可能であるかどうかを示すことができる。例えば、生成するカーネルは、グローバルメモリ20に格納された原子カウンタの現在値を読み取る命令を含むことができる。生成するカーネルは、生成するカーネルが格納したデータの量に基づいて原子カウンタの値を変更する命令と、原子カウンタの変更された値をグローバルメモリ20内に再び書き込む命令と、も含むことができる。

30

【0056】

消費するカーネルは、グローバルメモリ20に格納された原子カウンタの値を定期的に検査する命令を含むことができる。原子カウンタの値が十分に大きいときには、消費するカーネルは、消費されるべきデータが入手可能であると決定することができる。例えば、原子カウンタの値がXであり、さらに、生成するカーネルがNの量のデータを生成したと仮定する。この例では、消費するカーネルは、原子カウンタの値を定期的に検査することを消費する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に行わせる命令を含むことができる。原子カウンタの値がX+Nであるとプログラマブル計算ユニットが決定したときには、プログラマブル計算ユニットは、格納されたデータを消費のためにグローバルメモリ20から取り出すようにGPU16に要求することができる。

40

【0057】

この方法により、ソフトウェア(すなわち、カーネルの命令)を用いてカーネルをパイプライン方式で実行することが可能である。しかしながら、カーネル内の命令を用いてパイプライン方式でカーネルを実行することが非効率的である様々な理由が存在する。例えば、データを格納すべきアドレス又はグローバルメモリ20内のどの場所にデータが格納されているかを決定する命令をカーネル内に含むことは、不必要に電力を消費すること、及び、グローバルメモリ20内のアドレスを決定するための命令を処理してクロックサイ

50

クルを浪費することをプログラマブル計算ユニットのA L Uに要求することになる。

【 0 0 5 8 】

さらに、原子カウンタの値を定期的に検査することは、オフチップで（すなわち、グローバルメモリ20内の）情報にアクセスするようにG P U 1 6に要求することになる。原子カウンタの値をグローバルメモリ20から読み取ること及び原子カウンタの変更された値をグローバルメモリ20に書き込むことは、望ましくない量の電力を消費する可能性がある。さらに、例示されるように、I C 1 2は、メモリバス24を介してグローバルメモリ20に結合される。メモリバス24が処理できるデータの量には帯域幅上の限界がある。従って、G P U 1 6が原子カウンタの値を読み取る及び書き込むことができるときには遅延が存在することがある。

10

【 0 0 5 9 】

さらに、データが消費するカーネルによって消費されるために入手可能であるときは不明であるため、消費するカーネルを実行するプログラマブル計算ユニットは、データが消費のために入手可能であるかどうかを決定するために原子カウンタの値を検査することをG P U 1 6に定期的に行わせることになる。原子カウンタの値を定期的に検査することは、消費するカーネルのスレッドを“スピニング”（*s p i n n i n g*）した状態にとどまらせることになる可能性がある。例えば、原子カウンタの読み取り値が、データは消費のためにまだ完全には入手可能でないことを示す場合は、プログラマブル計算ユニットは、プログラマブル計算ユニットが原子カウンタの値を再度検査するまで、消費するカーネルのスレッドの実行を休止することができる。データが依然として入手可能でない場合は、プログラマブル計算ユニットは再度待機し、データが入手可能であるかどうかを再度検査することをG P U 1 6に行わせる。この例では、消費するカーネルのスレッドは、消費されるべきデータがグローバルメモリ20において入手可能でない時間の間ビジー - 待機状態にとどまることになる。換言すると、スピニング中には、プログラマブル計算ユニットは、どのような機能も実行していることができず、データの消費を遅延させるおそれがある。

20

【 0 0 6 0 】

（例えば、原子カウンタの値を読み取ることによって）データが入手可能であるかどうかをプログラマブル計算ユニットが決定する頻度が高い場合は、G P U 1 6は、グローバルメモリ20に格納された原子カウンタの値をあまりにも頻繁に読み取りすぎることによって電力を浪費する可能性がある。データが入手可能であるかどうかをプログラマブル計算ユニットが決定する頻度が低い場合は、データが入手可能な時点とG P U 1 7 6がデータを取り出す時点の間に時間の浪費が生じ、それもデータの消費を遅延させる。

30

【 0 0 6 1 】

さらに、グローバルメモリ20が原子カウンタを格納する上記の技法のうちの一部において、1つのカーネルが原子カウンタの値を読みとり中、変更中、及び書き込み中であるときには、その他のカーネルは、原子カウンタの値を読み取ること、変更すること、又は書き込むことを許容されない。該事例では、2つの生成するスレッドがグローバルメモリ20での格納のために同時にデータを出力する必要があるときに、それらのスレッドのうちの1つはデータを出力することができるが、他方のスレッドは、原子カウンタにアクセスできないためデータを出力することができない。該事例では、格納へのアクセスが拒否されたスレッドは、原子カウンタへのアクセスが可能になるまでスピンするおそれがあり、原子カウンタへのアクセスが可能になった時点で、格納へのアクセスを拒否されたスレッドが、グローバルメモリ20にデータを出力することができる。2つの消費するスレッドが同時にデータにアクセスしようと試みたときにも同じことが起きるおそれがある。

40

【 0 0 6 2 】

本開示において説明される技法は、上述される技法と比較して、G P U 1 6がより効率的にパイプライン方式でカーネルを実行するのを可能にすることができる。より詳細に説明されるように、パイプライン管理ユニット（P M U）18は、様々なスレッドによって生成されるデータ及び様々なスレッドによって消費されることになるデータの状態情報を

50



格納するように構成することができる。この方法により、GPU 16は、オフチップで連続的に情報にアクセスし、データがどこに格納されているか及びデータがいつ消費のために入手可能であることを示す必要がない。むしろ、PMU 18は、該情報を内部に（すなわち、IC 12内に）格納することができる。

【0063】

例示されるように、グローバルメモリ20は、バッファ22A乃至22N（総称してバッファ22と呼ばれる）を含むことができる。バッファ22は、グローバルメモリ20内の格納のための記憶場所であることができる。バッファ22の例は、先入れ先出し（FIFO）バッファ又はリングバッファを含む。

【0064】

プロセッサ14は、グローバルメモリ20内に常駐するバッファ数を定義し及びグローバルメモリ20内の記憶場所を予約するように構成することができる。例えば、プロセッサ14は、バッファ22の開始及び終了記憶場所（すなわち、開始及び終了アドレス）を定義することができる。プロセッサ14は、GPU 16のシェーダプロセッサ内に常駐するプログラマブル計算ユニット数に基づいてグローバルメモリ20内に常駐するバッファ数を定義することができる。一例として、プロセッサ14は、各プログラマブル計算ユニットに関して1つ以上の入力バッファ22（すなわち、プログラマブル計算ユニットで実行中のカーネルによって消費されるデータを格納する1つ以上のバッファ）及び各プログラマブル計算ユニットに関してゼロ以上の出力バッファ22（すなわち、GPU 16のプログラマブル計算ユニットで実行中のカーネルによって生成されたデータを格納するゼロ以上のバッファ）が存在するような形でグローバルメモリ20内に常駐するバッファ数を定義することができる。

【0065】

さらに、プロセッサ14は、バッファのサイズを定義するように構成することができる。例えば、プロセッサ14は、バッファ22の各々の中における格納用記憶場所の数（例えば、バッファ22の長さ）を定義するように構成することができる。プロセッサ14は、格納用記憶場所の各々に格納することができるデータの量（例えば、バッファ22の幅）を定義することもできる。幾つかの例では、プロセッサ14は、バッファ22にデータを予めポピュレート（pre-populate）することができる。

【0066】

幾つかの例では、プロセッサ14は、最低数のバッファ22を定義するように構成することができる。一例として、プロセッサ14は、最低数128のバッファ22を定義するように構成することができる。バッファ22の最低数が128であることは、例示を目的とするものであり、限定するものであるとは解釈されるべきでない。バッファ22の最低数は、128よりも多いこと又は少ないことができる。幾つかの例では、バッファ22の最低数に関する要求は必要がない。

プロセッサ14は、バッファ22の状態を決定するために様々な命令を実行するようにも構成することができる。例えば、プロセッサ14は、バッファ22に格納されたデータをIC 12又はGPU 16のバッファ内にコピーする命令及びIC 12又はGPU 16のバッファ内に格納されたデータをバッファ22内にコピーする命令を実行することができる。プロセッサ14は、バッファ22に格納されたデータの量を定義する命令、及び（例えば、バッファが壊れていないことを確認するために）バッファ22の長さ及び幅を確認する命令も実行することができる。プロセッサ14がバッファ22の状態を決定するのを可能にする命令の該実行は、すべての例で要求されるわけではないが、カーネルの開発者がGPU 16ではなくプロセッサ14で命令を実行することによってバッファ22の状態を決定するのに役立つことができる。

【0067】

幾つかの例では、プロセッサ14は、バッファ22に関する拡大係数（amplification factor）を定義するように構成することができる。拡大係数は、バッファ22のうちの1つに格納するためにカーネルのスレッドによって生成することがで

10

20

30

40

50

きる要素の最大数を示すことができる。拡大係数は、データを格納すべきバッファ 22 のうちの 1 つがすべての生成されたデータを格納することができない状況にとって必要になるであろう。この結果、バッファ 22 内での不十分な格納スペースに起因してカーネルの実行を停止させることがあり、立ち往生状態になる可能性がある（例えば、カーネルが実行状態にまったく戻らない）。

【 0 0 6 8 】

該立ち往生が生じる機会を最小限にするために、プロセッサ 14 は、グローバルメモリ 20 の大きな部分を予約することができる（例えば、ほとんどのあらゆるタイプのデータを格納する上で十分な大きさの長くて幅広いバッファ 22 を定義する）。これは、幾つかの場合は適切に機能することができるが、グローバルメモリ 20 の大きな部分を予約する  
10  
のが可能でないその他に関しては適切でない。幾つかの事例では、開発者は、カーネルがあまりにも多すぎるデータを生成せず、それによって立ち往生の機会を最小にするような形でカーネルを開発することができる。

【 0 0 6 9 】

プロセッサ 14 がバッファ 22 を定義するとして説明されているが、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、プロセッサ 14 以外の処理ユニットがバッファ 22 を定義するように構成することができる。幾つかの例では、GPU 16 がバッファ 22 を定義するのが可能である。しかしながら、説明を簡単にするために、これらの技法は、プロセッサ 14 がバッファ 22 を定義するとして説明される。

【 0 0 7 0 】

プロセッサ 14 は、バッファ 22 の情報をパイプライン処理ユニット（PMU）18 に送信することができる。例えば、PMU 18 は、バッファ 22 の数、バッファ 22 の開始及び終了アドレス、バッファ 22 の長さ、幅を示す情報、及び、プロセッサ 14 がバッファ 22 に関して決定したその他の情報を受信することができる。PMU 18 は、バッファ 22 の該状態情報を IC 12 内のレジスタ内に格納することができる。プロセッサ 14 から  
20  
のバッファ 22 の情報が存在することで、PMU 18 は、プログラマブル計算ユニットで実行中のカーネルのスレッドがデータを生成及び消費しているときにバッファ 22 の状態情報を管理するように構成することができる。

【 0 0 7 1 】

例えば、カーネルのスレッドを実行するプログラマブル計算ユニットがデータを生成して生成されたデータを出力後は、PMU 18 は、データを受信し、データが格納されるべきアドレスを決定することができる。例えば、PMU 18 は、バッファ 22 のうちのいずれにデータを格納すべきかを決定することができる。バッファ 22 がリングバッファ又は FIFO バッファである例では、PMU 18 は、バッファ 22 の初めと終わりを識別する  
30  
ポイントに関する情報を格納することができる。リングバッファの場合は、PMU 18 は、有効なデータの初め及び有効なデータの終わりを識別するポイントに関する情報も格納することができる。

【 0 0 7 2 】

従って、生成されたデータが格納されるべきアドレス又は消費のためにデータが取り出されるアドレスを決定することをプログラマブル計算ユニットに行わせる命令を含むカーネルではなく、PMU 18 を、生成されたデータが格納されるべきアドレス又は消費のためにデータが取り出されるアドレスを決定するように構成することができる。この方法により、GPU 16 は、クロックサイクルを浪費せず、さらに、プログラマブル計算ユニットの ALU は、データが格納されるべき又はデータが取り出されるべきアドレスを決定するの  
40  
のに処理電力を浪費することがない。

【 0 0 7 3 】

さらに、PMU 18 は、消費されるべきデータの消費準備が整っている時を決定するように構成することができる。例えば、グローバルメモリ 20 が原子カウンタを格納するのではなく、PMU 18 が IC 12 内にローカルで（例えば、IC 12 内のローカルキャッシュメモリ内のレジスタ内に）原子カウンタを格納することができる。一例として、生成  
50

するスレッドを実行するプログラマブル計算ユニットがデータを出力時に、P M U 1 8 は、内部に格納された原子カウンタの値を読み取り、生成されたデータの量に基づいて原子カウンタの値を変更し、原子カウンタの変更された値をI C 1 2 内に書き込むことができる。この例では、消費スレッドを実行するプログラマブル計算ユニットが原子カウンタの値を読み取るときに、G P U 1 6 は、オフチップのグローバルメモリ20にアクセスすることによって原子カウンタの値を決定する必要がない。代わりに、P M U 1 8 は、原子カウンタの値を提供することができる。

#### 【0074】

幾つかの例では、原子カウンタの値をローカルで格納するP M U 1 8 は、スピニングを低減させることができる。例えば、消費するスレッドを実行するプログラマブル計算ユニットは、消費するスレッドによって消費されるべきデータの要求を出力することができる。この例では、P M U 1 8 は、（例えば、ローカルで格納された原子カウンタの値に基づいて）消費されるべきデータが入手可能であるかどうかを決定することができる。

#### 【0075】

データが消費のためにまだ入手可能でないとP M U 1 8 が決定した場合は、P M U 1 8 は、プログラマブル計算ユニットがまだ入手可能でないデータに依存しない（例えば、同じカーネルの又は可能な場合は異なるカーネルの）異なるスレッドに切り換わるべきであることをプログラマブル計算ユニットに示すことができる。換言すると、P M U 1 8 は、プログラマブル計算ユニットがその他のスレッドの実行を続けることができるようにするためまだ入手可能でないデータを必要とする消費スレッドはスリープ状態にすべきであることを示すことができる。データが入手可能であることが、原子カウンタのローカルで格納された値に基づいてP M U 1 8 によって決定された時点で、P M U 1 8 は、プログラマブル計算ユニットが現在入手可能なデータを用いて消費スレッドを実行することができるようにスリープ中のスレッドに戻る（すなわち、スレッドを起こす）ようにプログラマブル計算ユニットに命令することができる。この方法により、データが消費のためにまだ入手可能でないと、消費スレッドを実行するプログラマブル計算ユニットが、ビジー-待機状態にとどまるのではなく、カーネルのその他のスレッドを実行することができる。

#### 【0076】

他の例として、異なるプログラマブル計算ユニットで実行中の同じカーネルの2つの生成するスレッドが、バッファ22のうちの同じそれに同時にデータを書き込むことを試みたときには、P M U 1 8 は、生成するスレッドのうちの1方にアクセスを許容し、他方の生成するスレッドに対してはアクセスを拒否することができる。この例では、P M U 1 8 は、アクセスを拒否されたスレッドを実行しているプログラマブル計算ユニットに対してカーネルのその他のスレッドを実行するように命令することができる。バッファ22への書き込みアクセスが可能になったことがP M U 1 8 によって決定された時点で、P M U 1 8 は、アクセスが拒否されたスレッドを実行していたプログラマブル計算ユニットに対して、バッファ22への書き込みアクセスが現在は可能であることを示すことができる。この方法により、アクセスを拒否されたスレッドを実行していたプログラマブル計算ユニットが追加のスレッドを実行することができる。

#### 【0077】

同様に、2つの消費するスレッドが、バッファ22のうちの同じそれから同時にデータを読み取ることを試みたときには、P M U 1 8 は、消費するスレッドのうちの1方にアクセスを許容し、他方の消費するスレッドに対してはアクセスを拒否することができる。2つのスレッドが同時に書き込んでいる例と同様に、2つのスレッドが同時に読み取っているこの例では、P M U 1 8 は、アクセスを拒否されたスレッドを実行しているプログラマブル計算ユニットに対してその他のスレッドを実行するように命令することができる。バッファ22への読み取りアクセスが可能になったことがP M U 1 8 によって決定された時点で、P M U 1 8 は、アクセスが拒否されたスレッドを実行していたプログラマブル計算ユニットに対して、バッファ22への読み取りアクセスが現在は可能であることを示すことができる。この方法により、アクセスを拒否されたスレッドを実行していたプログラマ

ブル計算ユニットが追加のスレッドを実行することができる。

【0078】

この方法により、グローバルメモリ20内のバッファ22を定義するプロセッサ14、及びグローバルメモリ20内のバッファ22の状態を管理するPMU18は、GPU16によるパイプライン方式でのカーネルの効率的な実行を可能にすることができる。一例として、PMU18は、パイプライン方式でカーネルを実行するために必要なオフチップアクセスの数を最小限にすることができる。他の例として、PMU18は、データが格納されるべきアドレス又はデータが取り出されるべきアドレスを決定することができるため、GPU16は、該アドレスを決定するためにカーネル内で命令を実行することによって該アドレスを決定するのに電力及びクロックサイクルを浪費しないことができる。換言すると、PMU18は、データを格納する又は取り出す場所を決定する命令をスレッドが含むことなしにデータが格納されるべきアドレス又は取り出されるべきアドレスを決定することができる。さらに、PMU18は、プログラマブル計算ユニットがスピニングせずにカーネルのスレッドを実行するのを可能にすることができる。例えば、生成するカーネルからのデータをまだ入手可能でないときには、PMU18は、消費するカーネルのその他のスレッド（例えば、生成するカーネルからのデータを要求しないスレッド）が実行するのを可能にすることができる。

10

【0079】

図2は、グラフィックス処理ユニット（GPU）及びグローバルメモリをさらに詳細に例示したブロック図である。例えば、図2は、図1のGPU16及びグローバルメモリ20をさらに詳細に例示したものである。例示されるように、GPU16は、シェーダプロセッサ26と、固定機能ユニット30と、パイプライン管理ユニット（PMU）18と、キャッシュ34と、スケジューラ40と、レジスタ44と、を含む。幾つかの例では、レジスタ44は、キャッシュ34の一部であることができる。図2において示される例では、PMU18は、GPU16内で形成されるとして例示される。しかしながら、上述されるように、PMU18は、GPU16の外部で、及びGPU16と同じ集積回路内で形成することができる。

20

【0080】

シェーダプロセッサ26は、プログラマブル計算ユニット28A乃至28N（総称してプログラマブル計算ユニット28と呼ばれる）を含むことができ、それらは、シェーダコアとみなすことができる。固定機能ユニット30は、固定機能計算ユニット32A乃至32N（総称して固定機能計算ユニット32と呼ばれる）を含む。シェーダプロセッサ26及び固定機能ユニット30は、プログラマブル計算ユニット28及び固定機能ユニット32のうちの1つ以上を含むことができる（例えば、例示されるそれらよりも多い又は少ない）。

30

【0081】

プログラマブル計算ユニット28は、上述されるように機能することができる。例えば、プログラマブル計算ユニット28は、グラフィックスに関連するアプリケーション及びグラフィックスに関連しないアプリケーションの両方（例えば、シェーダ及びカーネル）を実行することができる。例えば、プログラマブル計算ユニット28は、デバイス言語（例えば、OpenCL（C言語）で書かれたカーネルを実行することができる。上述されるように、プログラマブル計算ユニット28のうちの1つは、即座の結果を格納するための及びプログラマブル計算ユニット28で実行中のカーネルのスレッド間で共有するためのローカルメモリを含むことができる。プログラマブル計算ユニット28の各々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28によってアクセスすることができない。幾つかの例では、プログラマブル計算ユニット28のうちの1つが、プログラマブル計算ユニット28のうちの他の1つがカーネルのスレッドを実行すべき時間をスケジューリングすることが可能である。

40

【0082】

幾つかの例では、プログラマブル計算ユニット28のうちの1つは、プログラマブル計

50

算ユニット 28 のうちの 1 つ以上のその他のそれらにデータを送信することができる。例えば、パイプライン方式でカーネルを実行するために、生成するスレッドを実行する、プログラマブル計算ユニット 28 のうちの第 1 のそれは、プログラマブル計算ユニット 28 の第 2 のそれにデータ（例えば、グラフィックスに関連しないデータ）を出力することができる。上述されるように、プログラマブル計算ユニット 28 のうちの送信するそれ（例えば、生成するスレッドを実行するプログラマブル計算ユニット）は、バッファ、例えば、グローバルメモリ 20 のバッファ 22 のうちの 1 つ、にデータを格納することができ、プログラマブル計算ユニット 28 のうちの受信するそれ（例えば、消費するスレッドを実行するプログラマブル計算ユニット）は、グローバルメモリ 20 のバッファ 22 のうちの 1 つからデータを取り出すことができる。

10

#### 【0083】

図 2 において示されるように、幾つかの例では、GPU 16 は、内部キャッシュ 34 を含むことができる。しかしながら、キャッシュ 34 は、GPU 16 の内部に限定されるのではなく、IC 12 の内部に存在することができる。幾つかの例では、生成されたデータをオフチップで（例えば、グローバルメモリ 20 内に）格納するのではなく、GPU 16 が GPU 16 又は IC 12 の内部にデータを格納することが可能である。例えば、プログラマブル計算ユニット 28 のうちの送信するそれは、図 2 の例では GPU 16 内に存在するキャッシュ 34 内のバッファ 36 A 乃至 36 N（総称してバッファ 36 と呼ばれる）のうちの 1 つ以上の中にデータを格納することができるが、IC 12 内及び GPU 16 の外部に存在することが可能である。プログラマブル計算ユニット 28 の受信するそれは、キャッシュ 34 内のバッファ 36 からデータを取り出すことができる。キャッシュ 34 内のバッファは、グローバルメモリ 20 内のバッファ 22 のキャッシュバックされた（cache-backed）バッファであることができる。換言すると、グローバルメモリ 20 のバッファ 22 は、消費するスレッドによって消費されるべき生成するスレッドによって生成された完全なデータを格納することができ、バッファ 36 は、グローバルメモリ 20 からデータにアクセスすることと比較して素早いアクセスを行うために生成されたデータの一部を格納するキャッシュとして機能することができる。

20

#### 【0084】

キャッシュ 34 内のバッファは、バッファ 22 に類似することができる。例えば、バッファ 36 は、FIFO バッファ又はリングバッファであることができる。キャッシュ 34 は、オフチップメモリ（例えば、グローバルメモリ 20 のバッファ 22）にアクセスすることに関連するメモリレーテンシー及び電力消費を回避するためにバッファ 36 を含めるのが望ましいであろう。しかしながら、格納のために利用可能なスペースは限られているためバッファ 36 のみを利用するのは実際的でない。この方法により、バッファ 36 内にデータの一部を格納してバッファ 22 内へのスピルオーバー（spill over）を考慮するのが可能である。

30

#### 【0085】

バッファ 36 及びバッファ 22 は、GPU 16 がパイプライン方式でカーネルを実行するのを可能にすることができる。例えば、バッファ 36 及びバッファ 22 は、プログラマブル計算ユニット 28 間での通信を考慮したデータ構造であるとみなすことができる。バッファ 36 及びバッファ 22 は、プログラマブル計算ユニットで実行中のカーネルが出力することができる最低データ量よりも多くのデータ（例えば、2 つ以上のデータユニット）を格納するように構成することができる。この方法により、プログラマブル計算ユニット 28 のうちの 1 つで実行するカーネルのスレッドは、バッファ 36 及びバッファ 22 に格納されており及びプログラマブル計算ユニット 28 のうちの他の 1 つで実行する他のカーネルのスレッドに消費のために渡すことができる可変の量のデータを生成することができる。

40

#### 【0086】

固定機能計算ユニット 32 は、固定された機能を提供することができ及び（限定しない一例として）ハードウェアユニットとして形成することができる。固定機能計算ユニット

50

32は、デバイス言語を用いて書かれた特定の組み込み式カーネルを実行するとみなすことができる。例えば、プログラマブル計算ユニット28は、機能上の柔軟性を提供する一方で、固定機能計算ユニット32は、各々の機能上の柔軟性を制限することができる。例えば、固定機能計算ユニット32は、特定のグラフィックス機能を提供するラスタライゼーションユニット(rasterization unit)、プリミティブ(primitive)アセンブリユニット、ビューポート変換ユニット、及びその他の該ユニットを含むことができる。

#### 【0087】

幾つかの例では、固定機能計算ユニット32は、各々の特定の機能を実行するためにハードワイヤすることができる。さらに、固定機能計算ユニット32は、固定機能計算ユニット32の他の1つがいつ実行すべきかをスケジューリングすることも可能である。さらに、幾つかの事例では、GPU16が固定機能計算ユニット32のうちの特定のそれを含まない場合は、利用不能な固定機能計算ユニットの機能を実行するカーネルを開発することが可能である。換言すると、カーネルは、利用不能な固定機能計算ユニットの固定機能の挙動をエミュレートすることができる。例えば、固定機能テッセレータ(tessellator)を利用可能でない場合は、開発者は、テッセレータの固定機能挙動をエミュレートするテッセレーションカーネルを開発し、プログラマブル計算ユニット28のうちの1つ以上においてカーネルを実行することができる。

#### 【0088】

幾つかの例では、GPU16は、スケジューラ40を含むことができる。スケジューラ40は、様々なプログラマブル計算ユニット28及び固定機能ユニット32にスレッド及び動作を割り当てることができる。例えば、スケジューラ40は、その他のプログラマブル計算ユニット28が利用不足である一方でいずれも利用過剰にならないようにするためにプログラマブル計算ユニット28によって実施されるタスクを負荷均衡化することができる。スケジューラ40は、ハードウェアまたはハードウェアで実行するソフトウェアとして実装することができる。

#### 【0089】

図2において、グローバルメモリ20は、バッファ42A乃至42N(総称してバッファ42と呼ばれる)を含むことができ、キャッシュ34は、バッファ38A乃至38N(総称してバッファ38と呼ばれる)を含むことができる。バッファ38は、必ずしもすべての例で存在する必要はなく、バッファ42に格納されたコマンドに関するキャッシュバックされた格納を提供するための任意選択のオンチップキャッシュとして形成することができる。バッファ42及びバッファ38は、コマンド待ち行列とみなすことができる。すべてのプログラマブル計算ユニット28に関してコマンド待ち行列(例えば、バッファ42及びバッファ38のうちの1つ)及び各タイプの固定機能計算ユニット32に関して1つの待ち行列が存在することができる。バッファ42及びバッファ38は、ゼロ以上のエントリを格納することができる。

#### 【0090】

バッファ42及び任意選択のオンチップバッファ38は、プログラマブル計算ユニット28及び固定機能計算ユニット32に関する仕事量のスケジューリングを援助することができる。例えば、バッファ42は、様々なタスクを実行することをプログラマブル計算ユニット28及び固定機能計算ユニット32に命令するコマンドを格納することができる。例えば、バッファ42内の各エントリは、カーネルのスレッドを実行することを利用可能な1つ以上のプログラマブル計算ユニット28に行わせるための情報を格納すること、及びカーネル引数値に関する情報および依存性情報を格納することができる。幾つかの例では、カーネルのスレッド間での依存性は、1つ以上のプログラマブル計算ユニット28がカーネルを実行する前に満たす必要がある。

#### 【0091】

バッファ22は、プロセッサ14(図1)及びGPU16の両方によってアクセス可能である。一例として、プロセッサ14は、上述される様々なAPIにより呼を用いてバッ

10

20

30

40

50

ファ22にアクセスすることができる。GPU16は、プログラマブル計算ユニット28で実行されたカーネルに基づいてバッファ22にアクセスすることができる。例えば、カーネルは、生成されたデータをグローバルメモリ20内に格納するための機能付きで開発することができる。

【0092】

例示されるように、GPU16は、パイプライン管理ユニット(PMU)18も含むことができる。上述されるように、PMU18は、グローバルメモリ20内でのバッファ22の状態を管理することができる。さらに、PMU18は、キャッシュ34内でのバッファ36の状態を管理することができる。

【0093】

例えば、PMU18は、バッファ22及びバッファ36の長さ及び幅を格納することによってバッファ22及びバッファ36の状態を管理することができ、生成されたデータを格納するために利用可能なバッファ22及びバッファ36の数を含む。一例として、PMU18は、プログラマブル計算ユニット28で実行するカーネルの前にバッファ22を割り当てることができる、及び、カーネルの実行終了時点でバッファ22の割り当てを解除することができる。

【0094】

他の例として、PMU18は、ヘッダポインタ、現在のオフセット、最大深さ、等に関する情報をオンチップレジスタ44に格納することができる。幾つかの例では、PMU18は、テクスチャパラメータがグラフィックス処理において格納される方法と類似の方法でバッファ22及びバッファ36の状態情報を格納することができる。

【0095】

バッファ22は、いずれのバッファ22にデータを格納し又はいずれのバッファ22からデータを取り出すかを決定し、バッファ内でデータを格納すべき又はデータを取り出すべき格納記憶場所を決定し(例えば、アドレスを決定し)、及びプログラマブル計算ユニット28の異なるそれらがデータ崩壊を生じさせるバッファからの情報へのアクセスを試みないようにするという意味で管理を要求することができる。PMU18に、該管理のタスクを負わせることができる。例えば、PMU18を含むGPU16又はPMU18を含むGPU16を含むICの場合、バッファ22の管理は、ICの外部ではなくGPU16を含むIC内にローカル化することができる。この結果、低減された電力消費量、及びプログラマブル計算ユニット28で実行するカーネルの効率的な実行を達成させることができる。

【0096】

一例として、PMU18は、原子カウンタをレジスタ44内に格納することができる。レジスタ44は、キャッシュ34の一部、又はGPU16又はIC12内のその他のメモリの一部であることができる。原子カウンタは、プログラマブル計算ユニット28のうちの1つに関するアクセスが可能であるかどうか(例えば、読み取るためのデータが入手可能であるかどうか又は2つ以上のカーネルが同じバッファ22から同時に書き込むこと又は読み取ることを試みているかどうか)を示すことができる。原子カウンタに基づき、PMU18は、2つのスレッドが同時にデータを書き込むのを試みた場合に発生する可能性があるバッファ22のデータ崩壊を回避するためにプログラマブル計算ユニット28のうちの1つにアクセスを適切に許容し、プログラマブル計算ユニット28のうちのその他のそれらにアクセスを拒否することができる。幾つかの例では、PMU18がプログラマブル計算ユニット28のうちの1つにアクセスを拒否したときには、PMU18は、アクセス(例えば、スレッド)を要求するタスクがスリープ状態になるのを許容し、及び、プログラマブル計算ユニット28のうちの拒否されたそれがその他のタスク(例えば、スレッド)の実行を継続するのを許容することができる。プログラマブル計算ユニット28のうちの拒否されたそれへのアクセスが可能になったときに、PMU18はそのタスクをウェイクアップさせ、さらなる実行のためにそのタスクにデータを提供することができる。この方法により、プログラマブル計算ユニット28は完全にアイドル状態にならず、プログラ

10

20

30

40

50

マブル計算ユニット 28 のその他のタスクを実行することができる。

【0097】

幾つかの例では、グローバルメモリ 20 の 1 つのバッファ 22 からデータを取り出す必要があるときには、PMU 18 は、必要なデータ以上のデータを取り出すことができる。例えば、PMU 18 は、要求されたデータの開始位置及び終了位置を決定することができる。しかしながら、PMU 18 は、要求されたデータの決定された終了位置以降にバッファ 22 内に格納される追加データを取り出すことができる。PMU 18 は、バッファ 36 において格納スペースが利用可能であると PMU 18 が決定したときに該追加データを取り出すことができる。上述されるように、PMU 18 は、グローバルメモリ 20 内のバッファ 22、及びキャッシュ 34 内のバッファ 36 の両方を管理することができる。PMU 18 は、取り出されたデータをキャッシュ 34 に格納することができる。この方法により、追加データが必要なときに既に GPU 16 内に該データが存在している。追加データ（例えば、要求されたデータに加えてのデータ）をバッファ 36 に格納することは、GPU 16 が（例えば、グローバルメモリ 20 から）オフチップでデータにアクセスしなければならない回数をさらに減らすことができる。

10

【0098】

データにアクセスするために、プログラマブル計算ユニット 28 は、バッファにアクセスするためのポインタを利用することができる（例えば、カーネルは、ポインタを用いてデータにアクセスするように開発することができる）。幾つかの例では、PMU 18 は、プログラマブル計算ユニット 28 がデータに適切にアクセスできるようにポインタ情報を維持することができる。例えば、プログラマブル計算ユニット 28 は、バッファ 22 に関する情報を要求する専用命令を PMU 18 に出力することができる。該命令は、バッファ内の要素数、バッファ内に格納されているデータ量（例えば、バッファの幅）、情報の格納場所に関する情報、及びその他の該情報を含むことができる。この方法により、プログラマブル計算ユニット 28 がバッファ 22 に適切にアクセスようにすることは、GPU 16 を収納する IC の内部で行うことができ、それは、GPU 16 を収納する IC 外部でのアクセスを減らすことができる。

20

【0099】

一例として、データが壊れるか失われないようにするために、生成するカーネルは、バッファ 22 の範囲（開始点及び終了点）を問い合わせる命令を含むように開発することができる。この例では、生成するカーネルを実行中であるプログラマブル計算ユニット 28 のうちの 1 つは、バッファ 22 の範囲の問い合わせを PMU 18 に出力することができる。PMU 18 は、（プロセッサ 14 がバッファ 22 を定義したときにプロセッサ 14 から情報を受信することによって）バッファ 22 の範囲の情報をレジスタ 44 に格納しておくことができる。PMU 18 は、バッファ 22 の範囲の結果を生成するカーネルに戻すことができる。

30

【0100】

他の例として、パイプライン方式でカーネルを実行するために、幾つかの例では、パイプライン内でのデータの順序を維持する必要がある。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ことになると仮定する。しかしながら、この場合は、第 3 のカーネルも、第 1 及び第 2 のカーネルが実行中に実行中であることが可能である。この場合は、第 1 のカーネルによって生成されたデータ及び第 3 のカーネルによって生成されたデータの順序を変更することが可能であり、その結果第 2 のカーネルが不正確なデータを消費することになる可能性がある。

40

【0101】

適切な順序設定を保証するために、幾つかの例では、バッファ 22 がアクセスのために利用可能であるかどうかを示す原子カウンタに加えて、PMU 18 が追加の原子カウンタをレジスタ 44 に格納することができる。これらの追加の原子カウンタは、デバイス - 原子カウンタと呼ぶことができる。例えば、バッファ 22 のうちの各々の 1 つと関連付けられたデバイス - 原子カウンタが存在することができる。さらに、PMU 18 又はスケジュー

50



ーラ 40 は、各カーネルの各スレッドによって生成されたデータを格納すべきバッファ 22 内の相対的位置を定義するトークンをそのスレッドに割り当てるように構成することができる。スレッドに関するこのトークンは、デバイス - 原子カウンタの現在値であることができる。

#### 【0102】

例えば、PMU 18 は、0 のトークン値を有するデータを最初に消費する第 1 の消費スレッドを割り当て、1 のトークン値を有するデータを 2 番目に消費する第 2 の消費スレッドを割り当て、以下同様である。これらの消費スレッドの各々は、デバイス - 原子カウンタの値を PMU 18 に要求することができる。デバイス - 原子カウンタの現在値が消費スレッドのトークン値に等しい場合は、消費するスレッドはデータを消費することができる。そうでない場合は、消費するスレッドはデータを消費することができない。

10

#### 【0103】

トークン値がデバイス - 原子カウンタの値と等しい消費スレッドがデータを消費した後は、PMU 18 は、デバイス - 原子カウンタの値を更新することができる。幾つかの例では、消費するスレッドが消費するデータ量は固定させることができ、PMU 18 は、固定された量のデータがバッファ 22 から取り出された後にデバイス - 原子カウンタの値を更新することができる。しかしながら、幾つかの例では、消費するスレッドが消費するデータの量は、固定することができない。これらの例では、消費するスレッドがデータの受信を終了した後に、消費するスレッドは、次の消費するスレッドがデータを消費することができるように PMU 18 がデバイス - 原子カウンタの値を増加すべきであることを PMU 18 に示すことができる。この方法により、デバイス - 原子カウンタは、PMU 18 がその値をレジスタ 44 に格納すること及び更新することができ、データが消費されるべき順序が守られ、順番を狂わせてデータを受信すべきでない消費するスレッドが順番外でデータを受信しないように保証することができる。

20

#### 【0104】

他の例として、PMU 18 は、立ち往生の機会を最小限にするためにレジスタ 44 に情報を格納することができる。例えば、上述されるように、プロセッサ 14 は、バッファ 22 のうちの 1 つ内に格納するためにカーネルのスレッドによって生成することができる要素の最大数を示すバッファ 22 に関する拡大係数を定義するように構成することができる。カーネルが拡大係数によって定義されるよりも多いデータを生成する場合は、カーネルは、立ち往生する（実行を停止する）可能性がある。プロセッサ 14 は、拡大係数の値を PMU 18 に提供することができ、PMU 18 は、拡大係数の値をレジスタ 44 内に格納することができる。

30

#### 【0105】

幾つかの例では、立ち往生する機会を最小限にするために、開発者は、拡大係数の値を要求する命令をカーネル内に含めることができる。カーネルを実行中のプログラマブル計算ユニット 28 のうちの 1 つは、拡大係数の値の要求を PMU 18 に出力することができる。代わりに、PMU 18 は、カーネルを実行中のプログラマブル計算ユニット 28 のうちの 1 つに拡大係数の値を示すことができる。カーネルのスレッドによって生成されるデータの量が拡大係数よりも大きくなるとプログラマブル計算ユニット 28 が決定した場合は、プログラマブル計算ユニット 28 は、生成されたデータの量が拡大係数と等しくなった時点でカーネルの実行を停止することができ、及び、既に生成されていたデータが消費された時点でカーネルの残りのスレッドの実行をスケジューリングすることができる。

40

#### 【0106】

立ち往生を最小限にするための上記の技法に加えて又はその代わりに、PMU 18 は、生成されたデータが拡大係数と等しくなるまで PMU 18 がバッファ 22 にデータを格納することができる事前スケジューリングを実装することができる。PMU 18 は、残りのデータはバッファ 36 に格納することができる。換言すると、PMU 18 は、バッファ 22 にデータを格納する要求が“安全な”範囲内にあるようにすること、バッファ 22 にデータを格納する要求が代わりにバッファ 36 に格納されるようにすることができる。

50

## 【 0 1 0 7 】

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。図 3 において例示されるように、プログラマブル計算ユニット 2 8 のうちの 1 つは、G P U 1 6 のシェーダプロセッサ 2 6 においてカーネルの 1 つ以上のスレッドを実行することができる ( 4 6 )。I C 1 2 内又は G P U 1 6 内にある P M U 1 8 は、カーネルの 1 つ以上のスレッドに関して、I C 1 2 の外部に存在するグローバルメモリ 2 0 にデータを格納する又は I C 1 2 の外部に存在するグローバルメモリ 2 0 からデータを取り出す要求をプログラマブル計算ユニット 2 8 のうちの 1 つから受信することができる ( 4 8 )。

## 【 0 1 0 8 】

10

P M U 1 8 は、データの格納又は取り出しを要求したプログラマブル計算ユニット 2 8 のうちの 1 つに関してアクセスが許容可能であるかどうかを決定することができる ( 5 0 )。アクセスが不可能 ( 5 0 のいいえ ) である場合は、プログラマブル計算ユニット 2 8 のうちの 1 つは、カーネルの追加スレッドを実行することができる ( 5 2 )。この例では、P M U 1 8 は、アクセスが可能になったときにプログラマブル計算ユニットのうちの 1 つに示すことができる。

## 【 0 1 0 9 】

アクセスが可能である ( 5 0 のはい ) 場合は、プログラマブル計算ユニット 2 8 のうちの 1 つは、グローバルメモリ 2 0 内のバッファ ( 例えば、バッファ 2 2 のうちの 1 つ ) 内においてデータが格納されるか又は取り出される記憶場所を決定することができる ( 5 2 )。例えば、P M U 1 8 は、データが格納されるか又は取り出されるグローバルメモリ 2 0 内の記憶場所 ( すなわち、アドレス ) を決定することができる ( 5 4 )。決定された記憶場所に基づき、G P U 1 6 は、グローバルメモリ 2 0 内のバッファ 2 2 のうちの 1 つ内の決定された記憶場所にデータを格納する又は決定された記憶場所からデータを取り出すことができる ( 5 6 )。

20

## 【 0 1 1 0 】

幾つかの例では、バッファ 2 2 のうちの 1 つ内の記憶場所を決定するために、P M U 1 8 は、グローバルメモリ 2 0 内においてデータが格納されるべき又は取り出されるべき記憶場所をカーネルの 1 つ以上のスレッドが示さずに記憶場所を決定することができる。この方法により、カーネルは、データを格納すべき又はデータが取り出されるべきグローバルメモリ 2 0 内の記憶場所を決定するための命令を含む必要がない。

30

## 【 0 1 1 1 】

幾つかの例では、P M U 1 8 は、要求されたデータに加えてデータを取り出すことができる。これらの例では、P M U 1 8 は、追加データをキャッシュ 3 4 に格納することができる。幾つかの例では、P M U 1 8 は、バッファ 2 2 の状態情報をプロセッサ 1 4 から受信することができる。これらの例では、P M U 1 8 は、受信された状態情報に基づいてデータが格納又は取り出されるべきバッファ 2 2 のうちのそれら内の記憶場所を決定することができる。

## 【 0 1 1 2 】

図 4 は、本開示において説明される 1 つ以上の例による他の技法例を示したフローチャートである。例示されるように、G P U 1 6 のシェーダプロセッサ 2 6 の第 1 のプログラマブル計算ユニット ( 例えば、プログラマブル計算ユニット 2 8 のうちの 1 つ ) は、第 1 のスレッドを実行することができる ( 5 8 )。G P U 1 6 のシェーダプロセッサ 2 6 の第 2 のプログラマブル計算ユニット ( 例えば、プログラマブル計算ユニット 2 8 のうちの他の 1 つ ) は、第 2 の異なるスレッドを実行することができる ( 6 0 )。

40

## 【 0 1 1 3 】

G P U 1 6 を含む I C 1 2 内に存在する P M U 1 8 は、第 1 のスレッドの実行によって生成されたデータを、I C 1 2 の外部にあるグローバルメモリ 2 0 内のバッファ ( 例えば、バッファ 2 2 のうちの 1 つ ) 内に格納する要求を第 1 のプログラマブル計算ユニットから受信することができる ( 6 2 )。この例では、第 1 のスレッド ( 例えば、生成するスレ

50

ッド)の実行によって生成されたデータは、第2のスレッド(例えば、消費するスレッド)を実行する第2のプログラマブル計算ユニットによって消費されることになる。さらに、バッファは、先入れ先出し(FIFO)バッファ及びリングバッファのうちの1つであることができ、リングバッファは、FIFOバッファの一例である。

#### 【0114】

PMU18は、第1のスレッドの実行によって生成されたデータが格納されるべきバッファ内の記憶場所を決定することができる(64)。IC12は、第1のスレッドの実行によって生成されたデータをバッファ内の決定された記憶場所に格納することができる(66)。第1のスレッドの実行によって生成されたデータをバッファ内の決定された記憶場所に格納するIC12は、データを格納するIC、データを格納するGPU16、及び/又はデータを格納するPMU18を含むことが理解されるべきである。換言すると、データを格納するIC12は、IC12又はデータを格納するIC12内のいずれかのコンポーネントを意味する。

#### 【0115】

幾つかの例では、PMU18は、バッファ22の状態情報をIC12内(例えば、レジスタ44内)に格納することができる。PMU18は、バッファ22の該状態情報をプロセッサ14から受信することができる。バッファ22の状態情報は、バッファ22の開始アドレス、バッファ22の終了アドレス、生成されたデータが格納されるべきバッファ22内のアドレス、及びデータが取り出されるべきバッファ内のアドレスのうちの1つ以上を含むことができる。これらの例では、PMU18は、バッファ22の格納された状態情報に基づいて第1のスレッドの実行によって生成されたデータが格納されるためのバッファ内の記憶場所を決定することができる。さらに、幾つかの例では、PMU18は、第1のスレッドがバッファ内でデータが格納されるべき記憶場所を示さずに第1のスレッドの実行によって生成されたデータを格納するためのバッファ内の記憶場所を決定することができる。

#### 【0116】

PMU18は、第1のスレッドの実行によって生成されたデータの少なくとも一部を取り出すことの要求を第2のスレッドを実行する第2のプログラマブル計算ユニットから受信することもできる。PMU18は、第1のスレッドの実行によって生成されるデータが、第2のスレッドを実行する第2のプログラマブル計算ユニットによる消費のために取り出すために入手可能であるかどうかを決定することができる。幾つかの例では、PMU18は、第1のスレッドの実行によって生成されたデータを格納することの要求を第1のプログラマブル計算ユニットから受信すると同時に、受信する前に、又は受信した後に第2のプログラマブル計算ユニットからの要求を受信することができる。

#### 【0117】

第2のスレッドによって要求されたデータが、第2のスレッドを実行する第2のプログラマブル計算ユニットによる消費のために取り出し可能でないときには、PMU18は、第3のスレッドを実行するように第2のプログラマブル計算ユニットに指示することができる。PMU18は、第2のスレッドによって要求されたデータが、第2のスレッドを実行する第2のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに第2のプログラマブル計算ユニットに示すこともできる。PMU18は、第2のスレッドによって要求されたデータが第2のスレッドを実行する第2のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに、第2のスレッドによって要求されたデータを消費するために第2のスレッドを実行するように第2のプログラマブル計算ユニットに指示することができる。

#### 【0118】

幾つかの事例では、第1のスレッドは、カーネルの生成するスレッドであることができ、第2のスレッドは、同じカーネルの消費するスレッドであることができる。幾つかの事例では、第1のスレッドは、生成するカーネルのスレッドであることができ、第2のスレッドは、消費するカーネルのスレッドであることができる。

## 【 0 1 1 9 】

図 5 は、図 1 のデバイスをより詳細に例示したブロック図である。例えば、図 5 は、デバイス 10 をさらに示す。デバイス 10 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 ( P D A )、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ユニット、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、タブレットコンピューティングデバイス、電子ブックリーダー、等を含み、ただしこれらに限定されない。デバイス 10 は、プロセッサ 14 と、 G P U 16、グローバルメモリ 20 と、ディスプレイ 68 と、ユーザインタフェース 70 と、トランシーバモジュール 72 と、を含むことができる。示される例では、 P M U 18 は、 G P U 16 内に形成される。幾つかの例では、 P M U 18 は、 G P U 16 を収納する同じ I C (すなわち、 I C 12 ) 内で形成することができる。同じく例示されるように、 C P U 16 は、 I C 12 内に常駐する。しかしながら、プロセッサ 14 も、 I C 12 内に常駐することができる。

10

## 【 0 1 2 0 】

デバイス 10 は、明確化を目的として図 4 には示されていない追加のモジュール又はユニットを含むことができる。例えば、デバイス 10 は、デバイス 10 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、これらのいずれも図 4 には示されていない。さらに、デバイス 10 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 10 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 70 及びディスプレイ 68 は、デバイス 10 がデスクトップコンピュータである例ではデバイス 10 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 70 は、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 ( p r e s e n c e - s e n s i t i v e ) ディスプレイである例ではディスプレイ 68 の一部であることができる。

20

## 【 0 1 2 1 】

図 4 のプロセッサ 14、 G P U 16、 P M U 18、及びグローバルメモリ 20 は、図 1 のプロセッサ 14、 G P U 16、 P M U 18、及びグローバルメモリ 20 と同様であることができる。ユーザインタフェース 70 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 70 は、タッチ画面であることもでき、及びディスプレイ 68 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 72 は、デバイス 10 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回路を含むことができる。トランシーバモジュール 72 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 68 は、液晶ディスプレイ ( L C D )、陰極線管 ( C R T ) 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、又は他のタイプの表示装置を備えることができる。

30

## 【 0 1 2 2 】

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

40

50

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。

【0123】

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、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、通常は磁氣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。

【0124】

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。

【0125】

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、コーデック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。

【0126】

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] データ処理動作をパイプライン方式で実行するための方法であって、

グラフィックス処理ユニット(GPU)のシェーダプロセッサの第1のプログラマブル

10

20

30

40

50

計算ユニットにおいて第 1 のスレッドを実行することと、

前記 G P U の前記シェーダプロセッサの第 2 のプログラマブル計算ユニットにおいて第 2 のスレッドを実行することと、

前記 G P U を含む集積回路 ( I C ) 内の管理ユニットを用いて、前記第 1 のスレッドの前記実行によって生成されたデータを前記 I C の外部のグローバルメモリ内のバッファ内に格納することの要求を前記第 1 のプログラマブル計算ユニットから受信することとであって、前記第 1 のスレッドの前記実行によって生成された前記データは、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによって消費されることになり、前記バッファは、先入れ先出し ( F I F O ) バッファ及びリングバッファのうちの 1 つを備えることと、

10

前記管理ユニットを用いて、前記第 1 のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の記憶場所を決定することと、

前記 I C を用いて、前記第 1 のスレッドの前記実行によって生成された前記データを前記バッファ内の前記決定された記憶場所に格納することと、を備える、方法。

[ C 2 ] 前記管理ユニットを用いて、前記バッファの状態情報を前記 I C 内に格納することをさらに備え、

前記バッファの前記状態情報は、前記バッファの開始アドレス、前記バッファの終了アドレス、生成されたデータが格納されるべき前記バッファ内のアドレス、及びデータが取り出されるべき前記バッファ内のアドレスのうちの 1 つ以上を含み、

前記バッファ内の前記記憶場所を決定することは、前記バッファの前記格納された状態情報に基づいて前記第 1 のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の前記記憶場所を決定することを備える C 1 に記載の方法。

20

[ C 3 ] 前記管理ユニットを用いて、前記第 1 のスレッドの前記実行によって生成された前記データの少なくとも一部を取り出すことの要求を前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットから受信することと、

前記管理ユニットを用いて、前記第 1 のスレッドの前記実行によって生成される前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるかどうかを決定することと、をさらに備える C 1 に記載の方法。

[ C 4 ] 前記第 2 のプログラマブル計算ユニットから前記要求を受信することは、前記第 1 のスレッドの前記実行によって生成されたデータを格納することの前記要求を前記第 1 のプログラマブル計算ユニットから受信すると同時に、受信する前に、又は受信した後に前記第 2 のプログラマブル計算ユニットから前記要求を受信することを備える C 3 に記載の方法。

30

[ C 5 ] 前記第 2 のスレッドによって要求された前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能でないときに、

前記管理ユニットを用いて、第 3 のスレッドを実行するように前記第 2 のプログラマブル計算ユニットに指示することと、

前記管理ユニットを用いて、前記第 2 のスレッドによって要求された前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに前記第 2 のプログラマブル計算ユニットに示すことと、

40

前記管理ユニットを用いて、前記第 2 のスレッドによって要求された前記データが前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに前記第 2 のスレッドによって要求された前記データを消費するために前記第 2 のスレッドを実行するように前記第 2 のプログラマブル計算ユニットに指示することと、をさらに備える C 3 に記載の方法。

[ C 6 ] 前記管理ユニットを用いて、前記第 2 のスレッドによって要求された前記データに加えて前記グローバルメモリからデータを受信することと、

50

前記管理ユニットを用いて、前記第2のスレッドによって要求された前記データに加えての前記データを前記IC内のキャッシュに格納することと、をさらに備えるC3に記載の方法。

[C7] 前記第1のスレッドを実行することは、カーネルの生成するスレッドを実行することを備え、前記第2のスレッドを実行することは、前記カーネルの消費するスレッドを実行することを備えるC1に記載の方法。

[C8] 前記第1のスレッドを実行することは、生成するカーネルの前記第1のスレッドを実行することを備え、前記第2のスレッドを実行することは、消費するカーネルのスレッドを実行することを備えるC1に記載の方法。

[C9] 前記GPUは、前記管理ユニットを含むC1に記載の方法。

10

[C10] 前記バッファ内の前記記憶場所を決定することは、前記第1のスレッドの前記実行によって生成された前記データが前記バッファ内において格納されるべき前記記憶場所を前記第1のスレッドが示さずに前記データが格納されるべき前記バッファ内の前記記憶場所を決定することを備えるC1に記載の方法。

[C11] 装置であって、

バッファを含むグローバルメモリであって、前記バッファは、先入れ先出し(FIFO)バッファ及びリングバッファのうちの1つを備えるグローバルメモリと、

集積回路(IC)であって、

グラフィックス処理ユニット(GPU)であって、

第1のスレッドを実行するように構成された第1のプログラマブル計算ユニットと

20

、  
第2のスレッドを実行するように構成された第2のプログラマブル計算ユニットと  
、を備えるGPUと、

前記第1のスレッドの前記実行によって生成されたデータを前記グローバルメモリ内の前記バッファ内に格納することの要求を前記第1のプログラマブル計算ユニットから受信し、及び

前記第1のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の記憶場所を決定するように構成された管理ユニットと、を備える、集積回路(IC)とを備え、前記第1のスレッドの前記実行によって生成された前記データは、前記第2のスレッドを実行する前記第2のプログラマブル計算ユニットによって消費され、

30

前記ICは、前記第1のスレッドの前記実行によって生成された前記データを前記バッファ内の前記決定された記憶場所に格納するように構成される、装置。

[C12] 前記管理ユニットは、前記バッファの状態情報を前記IC内に格納するように構成され、

前記バッファの前記状態情報は、前記バッファの開始アドレス、前記バッファの終了アドレス、生成されたデータが格納されるべき前記バッファ内のアドレス、及びデータが取り出されるべき前記バッファ内のアドレスのうちの1つ以上を含み、

前記管理ユニットは、前記バッファの前記格納された状態情報に基づいて前記第1のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の前記記憶場所を決定するように構成されるC11に記載の装置。

40

[C13] 前記管理ユニットは、

前記第1のスレッドの前記実行によって生成された前記データの少なくとも一部を取り出すことの要求を前記第2のスレッドを実行する前記第2のプログラマブル計算ユーザから受信し、及び

前記第1のスレッドの前記実行によって生成される前記データが、前記第2のスレッドを実行する前記第2のプログラマブル計算ユニットによる消費のための取り出したために入手可能であるかどうかを決定するように構成されるC11に記載の装置。

[C14] 前記管理ユニットは、前記第1のスレッドの前記実行によって生成されたデータを格納することの前記要求を前記第1のプログラマブル計算ユニットから受信すると同時に、受信する前に、又は受信した後に前記第2のプログラマブル計算ユニットから

50

前記要求を受信するように構成される C 1 3 に記載の装置。

[ C 1 5 ] 前記管理ユニットは、

前記第 2 のスレッドによって要求された前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能でないときに、第 3 のスレッドを実行するように前記第 2 のプログラマブル計算ユニットに指示し、

前記第 2 のスレッドによって要求された前記データが、前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに前記第 2 のプログラマブル計算ユニットに示し、及び

前記第 2 のスレッドによって要求された前記データが前記第 2 のスレッドを実行する前記第 2 のプログラマブル計算ユニットによる消費のための取り出しのために入手可能であるときに前記第 2 のスレッドによって要求された前記データを消費するために前記第 2 のスレッドを実行するように前記第 2 のプログラマブル計算ユニットに指示するように構成される C 1 3 に記載の装置。

[ C 1 6 ] 前記管理ユニットは、

前記第 2 のスレッドによって要求された前記データに加えて前記グローバルメモリからデータを取り出し、及び

前記第 2 のスレッドによって要求された前記データに加えての前記データを前記 I C 内のキャッシュに格納するように構成される C 1 3 に記載の装置。

[ C 1 7 ] 前記第 1 のスレッドは、カーネルの生成するスレッドを備え、前記第 2 のスレッドは、前記カーネルの消費するスレッドを備える C 1 1 に記載の装置。

[ C 1 8 ] 前記第 1 のスレッドは、生成するカーネルのスレッドを備え、前記第 2 のスレッドは、消費するカーネルのスレッドを備える C 1 1 に記載の装置。

[ C 1 9 ] 前記 G P U は、前記管理ユニットを含む C 1 1 に記載の装置。

[ C 2 0 ] 前記管理ユニットは、前記第 1 のスレッドの前記実行によって生成された前記データが前記バッファ内において格納されるべき前記記憶場所を前記第 1 のスレッドが示さずに前記データが格納されるべき前記バッファ内の前記記憶場所を決定することを備える C 1 1 に記載の装置。

[ C 2 1 ] 前記装置は、映像デバイス、セットトップボックス、無線ハンドセット、パーソナルデジタルアシスタント、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、ビデオ会議ユニット、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 1 1 に記載の装置。

[ C 2 2 ] 装置であって、

バッファを含むグローバルメモリであって、前記バッファは、先入れ先出し ( F I F O ) バッファ及びリングバッファのうちの 1 つを備えるグローバルメモリと、

集積回路 ( I C ) であって、

グラフィックス処理ユニット ( G P U ) であって、

第 1 のスレッドを実行するための手段と、

第 2 のスレッドを実行するための手段と、

前記第 1 のスレッドの前記実行によって生成されたデータを前記グローバルメモリ内の前記バッファ内に格納することの要求を前記第 1 のスレッドを実行するための手段から受信するための手段であって、前記第 1 のスレッドの前記実行によって生成された前記データは、前記第 2 のスレッドを実行するための前記手段によって消費される手段と、

前記第 1 のスレッドを実行するための前記手段によって生成された前記データが格納されるべき前記バッファ内の記憶場所を決定するための手段と、

前記第 1 のスレッドの前記実行によって生成された前記データを前記バッファ内の前記決定された記憶場所に格納するための手段と、を備えるグラフィックス処理ユニット ( G P U ) 、を備える集積回路 ( I C ) と、を備える、装置。

[ C 2 3 ] 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、

実行されたときに、

10

20

30

40

50



グラフィックス処理ユニット（GPU）のシェーダプロセッサの第１のプログラマブル計算ユニットにおいて第１のスレッドを実行し、

前記GPUの前記シェーダプロセッサの第２のプログラマブル計算ユニットにおいて第２のスレッドを実行し、

前記GPUを含む集積回路（IC）内の管理ユニットを用いて、前記第１のスレッドの前記実行によって生成されたデータを前記ICの外部のグローバルメモリ内のバッファ内に格納することの要求を前記第１のプログラマブル計算ユニットから受信し、

前記管理ユニットを用いて、前記第１のスレッドの前記実行によって生成された前記データが格納されるべき前記バッファ内の記憶場所を決定し、及び

前記ICを用いて、前記第１のスレッドの前記実行によって生成された前記データを前記バッファ内の前記決定された記憶場所に格納することを１つ以上のプロセッサに行わせる命令が格納されており、前記第１のスレッドの前記実行によって生成された前記データは、前記第２のスレッドを実行する前記第２のプログラマブル計算ユニットによって消費されることになり、前記バッファは、先入れ先出し（FIFO）バッファ及びリングバッファのうちの１つを備える、コンピュータによって読み取り可能な記憶媒体。

10

【図１】

図１

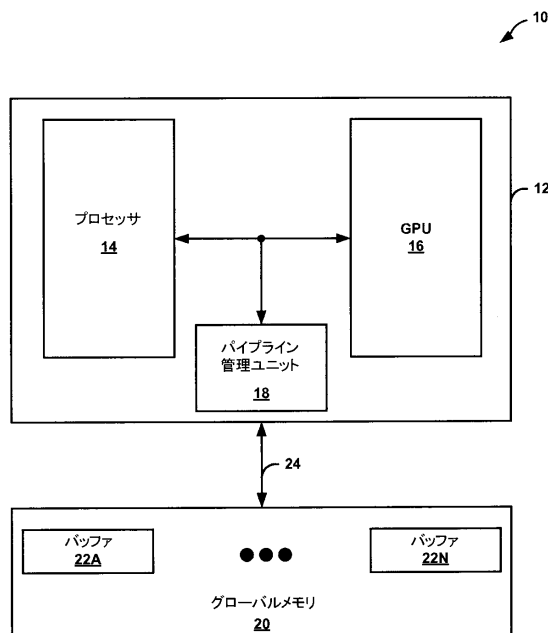


FIG. 1

【図２】

図２

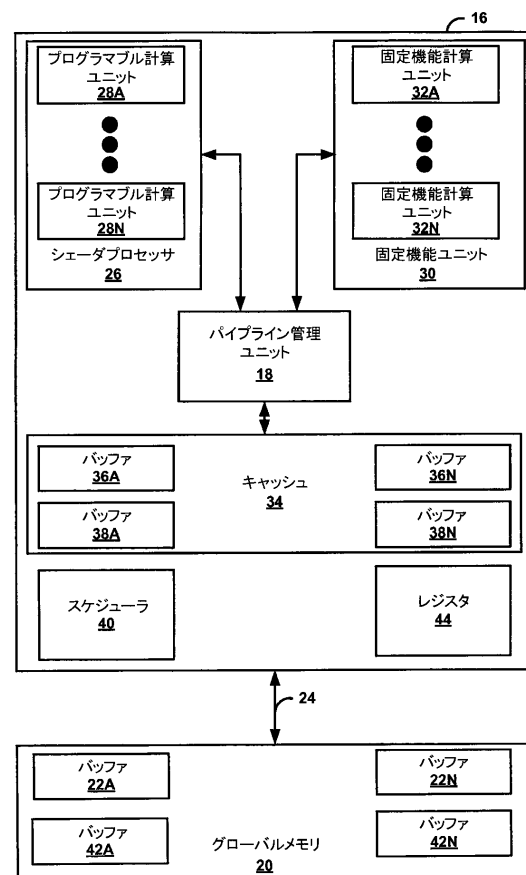


FIG. 2

【図 3】

図 3

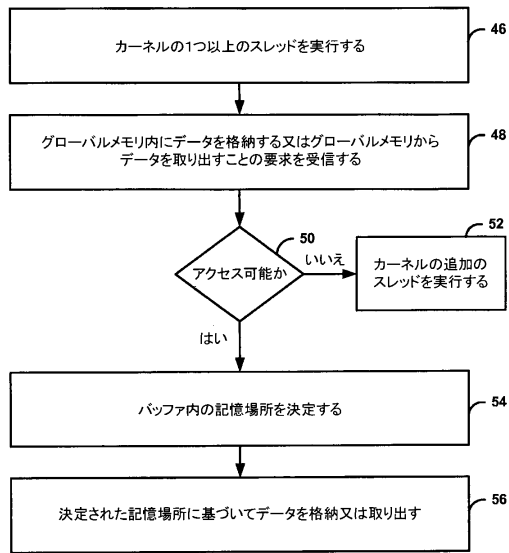


FIG. 3

【図 4】

図 4

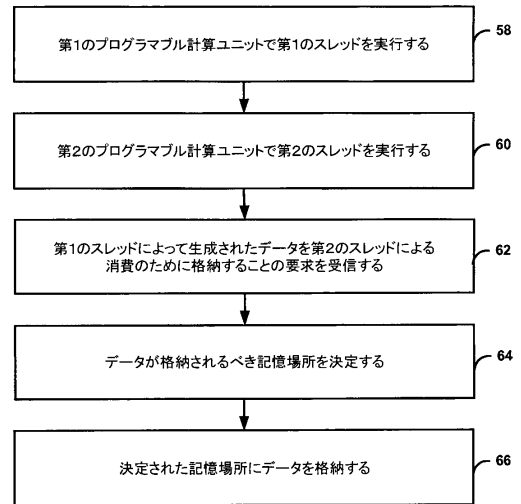


FIG. 4

【図 5】

図 5

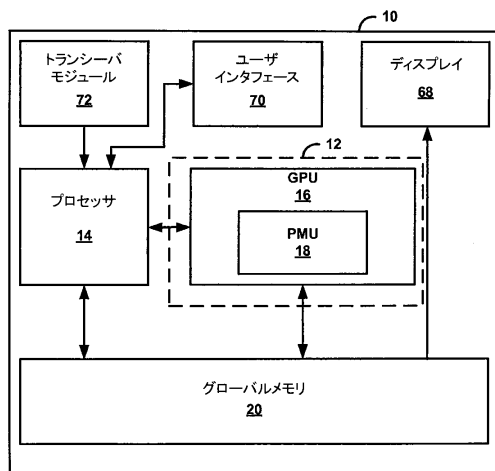


FIG. 5

## 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・ブイ .  
アメリカ合衆国、カリフォルニア州 92121、サン・ディエゴ、モアハウス・ドライブ 57  
75
- (72)発明者 ゴエル、ビネート  
アメリカ合衆国、カリフォルニア州 92121、サン・ディエゴ、モアハウス・ドライブ 57  
75

審査官 清木 泰

(56)参考文献 特開2007-122537(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F15/16 - 15/177  
G06F 9/30 - 9/42  
G06F15/80  
G06T 1/00 - 1/40  
G06T 3/00 - 5/50  
G06T 9/00 - 9/40  
G06T11/00 - 11/40  
G06T15/00 - 17/00  
G06T17/10 - 17/30  
G06F12/08 - 12/12