

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5996554号
(P5996554)

(45) 発行日 平成28年9月21日(2016.9.21)

(24) 登録日 平成28年9月2日(2016.9.2)

(51) Int. Cl. F I
G O 6 F 9 / 5 4 (2 0 0 6 . 0 1) G O 6 F 9 / 4 6 4 8 O B

請求項の数 38 (全 28 頁)

(21) 出願番号	特願2013-549593 (P2013-549593)	(73) 特許権者	595020643
(86) (22) 出願日	平成24年1月13日 (2012.1.13)		クアアルコム・インコーポレイテッド
(65) 公表番号	特表2014-505946 (P2014-505946A)		QUALCOMM INCORPORATED
(43) 公表日	平成26年3月6日 (2014.3.6)		アメリカ合衆国、カリフォルニア州 92
(86) 国際出願番号	PCT/US2012/021344		121-1714、サン・ディエゴ、モア
(87) 国際公開番号	W02012/097316		ハウス・ドライブ 5775
(87) 国際公開日	平成24年7月19日 (2012.7.19)	(74) 代理人	100108855
審査請求日	平成25年9月17日 (2013.9.17)		弁理士 蔵田 昌俊
審査番号	不服2015-13892 (P2015-13892/J1)	(74) 代理人	100109830
審査請求日	平成27年7月23日 (2015.7.23)		弁理士 福原 淑弘
(31) 優先権主張番号	13/007,333	(74) 代理人	100158805
(32) 優先日	平成23年1月14日 (2011.1.14)		弁理士 井関 守三
(33) 優先権主張国	米国 (US)	(74) 代理人	100194814
			弁理士 奥村 元宏

最終頁に続く

(54) 【発明の名称】 汎用グラフィクス処理装置における計算リソースパイプライン化

(57) 【特許請求の範囲】

【請求項1】

汎用グラフィクス処理装置 (G P G P U) であって、
処理パイプラインのステージとして選択的に動作するように構成された前記 G P G P U の 2 以上のプログラム可能な並行処理装置と、

前記並行処理装置間の移送のためにデータを保持するように構成された前記 G P G P U の 1 以上のローカルメモリバッファであって、前記ローカルメモリバッファの各々が前記処理パイプラインにおける前記並行処理装置のうち少なくとも 2 つの間で直接接続される、 1 以上のローカルメモリバッファと、

ここにおいて、前記データは、前記ローカルメモリバッファを介して、前記並行処理装置の一方から前記並行処理装置の他方へと直接受け渡される、

前記処理パイプライン内のデータシーケンスを保持するように構成された制御装置であって、データセットのデータスレッドのシーケンスを記録するために前記並行処理装置のうち少なくとも 1 つに前記データセットが入るときにシーケンス決定カウンタを実行し、前記シーケンス決定カウンタによって記録されるのと同じシーケンスで前記並行処理装置から前記データセットの前記データスレッドをリリースするために前記並行処理装置のうちの前記少なくとも 1 つから前記データセットが出るときにシーケンスエンフォーシングバリアを実行するように構成された制御装置と

を備える G P G P U 。

【請求項2】

10

20

前記 1 以上のローカルメモリバッファは、前記並行処理装置間の前記データの移送を可能にするハードウェアベースのデータフロー制御メカニズムを含む、請求項 1 に記載の G P G P U。

【請求項 3】

前記 1 以上のローカルメモリバッファは、ハードウェアベースの先入れ先出しバッファ (F I F O)、後入れ先出しバッファ (L I F O S) またはインデックス付けされたバッファのうち少なくとも 1 つを備える、請求項 1 に記載の G P G P U。

【請求項 4】

前記 1 以上のローカルメモリバッファにデータを送信するように前記並行処理装置のうち 1 以上を構成し、前記 1 以上のローカルメモリバッファからデータを受信するように前記並行処理装置のうち 1 以上を構成するように構成された制御装置、

10

をさらに備える請求項 1 に記載の G P G P U。

【請求項 5】

前記制御装置は、前記ローカルメモリバッファにデータを送信し前記ローカルメモリバッファからデータを受信するように前記並行処理装置を構成するために 1 以上のアプリケーションプログラミングインタフェース (A P I) を実行するように構成される、請求項 4 に記載の G P G P U。

【請求項 6】

前記制御装置は、前記処理パイプラインにおいて前の処理装置からのデータ出力を保持するために前記ローカルメモリバッファの各々についての必要な幅を決定し、前記決定された幅を有するように前記ローカルメモリバッファの各々を構成するようにさらに構成された請求項 1 に記載の G P G P U。

20

【請求項 7】

前記制御装置は、1 以上のアプリケーションプログラミングインタフェース (A P I) を実行して、前記ローカルメモリバッファの各々についての前記幅を決定し、前記決定された幅によって前記ローカルメモリバッファの各々を構成し、前記ローカルメモリバッファの各々の深さを決定するように構成される、請求項 6 に記載の G P G P U。

【請求項 8】

前記制御装置は、前記ローカルメモリバッファの各々の深さを決定するようにさらに構成され、前記ローカルメモリバッファの各々は前記深さを前記幅とトレードすることが構成可能である、請求項 6 に記載の G P G P U。

30

【請求項 9】

前記並行処理装置のうちの一つは、前記処理パイプラインの第 1 ステージとして動作し、デバイスメモリからオリジナルデータセットを取り出すように構成される、請求項 1 に記載の G P G P U。

【請求項 10】

前記並行処理装置のうちの一つは、前記処理パイプラインの最終ステージとして動作し、デバイスメモリにパイプライン処理されたデータセットを格納するように構成される、請求項 1 に記載の G P G P U。

【請求項 11】

40

前記並行処理装置のうちの一つは、前記処理パイプラインの中間ステージとして動作し、前記ローカルメモリバッファのうち一方を介して前記処理パイプラインにおいて前記並行処理装置のうちの前ものからデータセットを受信し、前記ローカルメモリバッファのうち他方を介して前記処理パイプラインにおいて前記並行処理装置のうち後続のものに前記データセットを送信するように構成される、請求項 1 に記載の G P G P U。

【請求項 12】

前記並行処理装置のうちの前記少なくとも一つは、前記データセットを処理するために、デバイスメモリから補助データを取り出すように構成される、請求項 11 に記載の G P G P U。

【請求項 13】

50

汎用グラフィクス処理装置（GPGPU）によってデータを処理する方法であって、前記方法は、

処理パイプラインのステージとして選択的に動作するように前記GPGPUの2以上のプログラム可能な並行処理装置を構成することと、

前記並行処理装置間の移送のためのデータを保持するように前記GPGPUの1以上のローカルメモリバッファを構成することであって、なお、前記ローカルメモリバッファの各々は、前記処理パイプラインにおける前記並行処理装置のうち少なくとも2つの間で直接接続され、ここにおいて、前記データは、前記ローカルメモリバッファを介して、前記並行処理装置の一方から前記並行処理装置の他方へと直接受け渡される、構成することと

10

前記処理パイプライン内のデータシーケンスを保持することであって、データセットのデータスレッドのシーケンスを記録するために前記並行処理装置のうちの少なくとも1つに前記データセットが入るときにシーケンス決定カウンタを実行することと、前記シーケンス決定カウンタによって記録されるのと同じシーケンスで前記並行処理装置から前記データセットの前記データスレッドをリリースするために前記並行処理装置のうちの前記少なくとも1つから前記データセットが出るときにシーケンスエンフォーシングバリアを実行することとを備える、前記データシーケンスを保持することと、

を備える、方法。

【請求項14】

前記1以上のローカルメモリバッファは、前記並行処理装置間の前記データの移送を可能にするハードウェアベースのデータフロー制御メカニズムを含む、請求項13に記載の方法。

20

【請求項15】

前記1以上のローカルメモリバッファは、ハードウェアベースの先入れ先出しバッファ（FIFO）、後入れ先出しバッファ（LIFO）またはインデクス付けされたバッファのうちの少なくとも1つを備える、請求項13に記載の方法。

【請求項16】

前記1以上のローカルメモリバッファにデータを送信するように前記並行処理装置の1以上を構成することと、

前記1以上のローカルメモリバッファからデータを受信するように前記並行処理装置の1以上を構成することと、

30

をさらに備える請求項13に記載の方法。

【請求項17】

前記1以上の並行処理装置を構成することは、前記ローカルメモリバッファにデータを送信し前記ローカルメモリバッファからデータを受信するように前記並行処理装置を構成するために1以上のアプリケーションプログラミングインタフェース（API）を実行することを備える、請求項16に記載の方法。

【請求項18】

前記1以上のローカルメモリバッファを構成することは、

前記処理パイプラインにおいて前の処理装置からのデータ出力を保持するために前記ローカルメモリバッファの各々についての必要な幅を決定することと、

40

前記決定された幅を有するように前記ローカルメモリバッファの各々を構成することと、

を備える、請求項13に記載の方法。

【請求項19】

前記1以上のローカルメモリバッファを構成することは、1以上のアプリケーションプログラミングインタフェース（API）を実行して、前記ローカルメモリバッファの各々についての前記幅を決定し、前記決定された幅によって前記ローカルメモリバッファの各々を構成し、前記ローカルメモリバッファの各々の深さを決定することを備える、請求項18に記載の方法。

50

【請求項 20】

前記 1 以上のローカルメモリバッファを構成することは、前記ローカルメモリバッファの各々の深さを決定することをさらに備え、前記ローカルメモリバッファの各々は前記深さを前記幅とトレードすることが構成可能である、請求項 18 に記載の方法。

【請求項 21】

前記 2 以上の並行処理装置を構成することは、前記並行処理装置のうちの 1 つを、前記処理パイプラインの第 1 ステージとして動作し、デバイスメモリからオリジナルデータセットを取り出すように構成することを備える、請求項 13 に記載の方法。

【請求項 22】

前記 2 以上の並行処理装置を構成することは、前記並行処理装置のうちの 1 つを、前記処理パイプラインの最終ステージとして動作し、デバイスメモリにパイプライン処理されたデータセットを格納するように構成することを備える、請求項 13 に記載の方法。

10

【請求項 23】

前記 2 以上の並行処理装置を構成することは、前記並行処理装置のうちの少なくとも 1 つを、前記処理パイプラインの中間ステージとして動作し、前記ローカルメモリバッファのうち一方を介して前記処理パイプラインにおいて前記並行処理装置のうちの前のものからデータセットを受信し、前記ローカルメモリバッファのうち他方を介して前記処理パイプラインにおいて前記並行処理装置のうち後続のものに前記データセットを送信するように構成することを備える、請求項 13 に記載の方法。

【請求項 24】

前記並行処理装置のうちの少なくとも 1 つを構成することは、前記データセットを処理するために、デバイスメモリから補助データを取り出すように前記並行処理装置のうちの前記少なくとも 1 つを構成することを備える、請求項 23 に記載の方法。

20

【請求項 25】

汎用グラフィクス処理装置 (GPGPU) であって、
処理パイプラインのステージとして選択的に動作するように前記 GPGPU の 2 以上のプログラム可能な並行処理装置を構成するための手段と、

前記並行処理装置間の移送のためのデータを保持するように前記 GPGPU の 1 以上のローカルメモリバッファを構成するための手段と、なお、前記ローカルメモリバッファの各々は、前記処理パイプラインにおける前記並行処理装置のうち少なくとも 2 つの間で直接接続され、ここにおいて、前記データは、前記ローカルメモリバッファを介して、前記並行処理装置の一方から前記並行処理装置の他方へと直接受け渡される、構成するための手段と、

30

前記処理パイプライン内のデータシーケンスを保持するための手段であって、データセットのデータスレッドのシーケンスを記録するために前記並行処理装置のうちの少なくとも 1 つに前記データセットが入るときにシーケンス決定カウンタを実行するための手段と、前記シーケンス決定カウンタによって記録されるのと同じシーケンスで前記並行処理装置から前記データセットの前記データスレッドをリリースするために前記並行処理装置のうちの前記少なくとも 1 つから前記データセットが出るときにシーケンスエンフォーシングバリアを実行するための手段とを備える、前記データシーケンスを保持するための手段

40

と
を備える GPGPU。

【請求項 26】

前記 1 以上のローカルメモリバッファは、前記並行処理装置間の前記データの移送を可能にするハードウェアベースのデータフロー制御メカニズムを含む、請求項 25 に記載の GPGPU。

【請求項 27】

前記 1 以上のローカルメモリバッファにデータを送信するように前記並行処理装置の 1 以上を構成するための手段と、

前記 1 以上のローカルメモリバッファからデータを受信するように前記並行処理装置の

50

1 以上を構成するための手段と、
をさらに備える請求項 25 に記載の G P G P U。

【請求項 28】

前記ローカルメモリバッファにデータを送信し前記ローカルメモリバッファからデータを受信するように前記並行処理装置を構成するために 1 以上のアプリケーションプログラミングインタフェース (A P I) を実行するための手段、をさらに備える請求項 27 に記載の G P G P U。

【請求項 29】

前記処理パイプラインにおいて前の処理装置からのデータ出力を保持するために前記ローカルメモリバッファの各々についての必要な幅を決定するための手段と、

前記決定された幅を有するように前記ローカルメモリバッファの各々を構成するための手段と、

をさらに備える請求項 25 に記載の G P G P U。

【請求項 30】

1 以上のアプリケーションプログラミングインタフェース (A P I) を実行して、前記ローカルメモリバッファの各々についての前記幅を決定し、前記決定された幅によって前記ローカルメモリバッファの各々を構成し、前記ローカルメモリバッファの各々の深さを決定するための手段をさらに備える請求項 29 に記載の G P G P U。

【請求項 31】

前記ローカルメモリバッファの各々の深さを決定するための手段をさらに備え、前記ローカルメモリバッファの各々は前記深さを前記幅とトレードすることが構成可能である、請求項 29 に記載の G P G P U。

【請求項 32】

汎用グラフィクス処理装置 (G P G P U) によってデータを処理するための命令を備えるコンピュータ可読媒体であって、前記命令は、実行時に、プログラム可能なプロセッサに、

処理パイプラインのステージとして選択的に動作するように前記 G P G P U の 2 以上のプログラム可能な並行処理装置を構成することと、

前記並行処理装置間の移送のためにデータを保持するように前記 G P G P U の 1 以上のローカルメモリバッファを構成することであって、なお、前記ローカルメモリバッファの各々は、前記処理パイプラインにおける前記並行処理装置のうち少なくとも 2 つの間で直接接続され、ここにおいて、前記データは、前記ローカルメモリバッファを介して、前記並行処理装置の一方から前記並行処理装置の他方へと直接受け渡される、と、

前記処理パイプライン内のデータシーケンスを保持することと、

を行わせ、前記命令は、前記プログラム可能なプロセッサに、データセットのデータスレッドのシーケンスを記録するために前記並行処理装置のうちの少なくとも 1 つに前記データセットが入るときにシーケンス決定カウンタを実行することと、前記シーケンス決定カウンタによって記録されるのと同じシーケンスで前記並行処理装置から前記データセットの前記データスレッドをリリースするために前記並行処理装置のうちの前記少なくとも 1 つから前記データセットが出るときにシーケンスエンフォーシングバリアを実行することとを行わせる、コンピュータ可読媒体。

【請求項 33】

前記 1 以上のローカルメモリバッファは、前記並行処理装置間の前記データの移送を可能にするハードウェアベースのデータフロー制御メカニズムを含む、請求項 32 に記載のコンピュータ可読媒体。

【請求項 34】

前記プログラム可能なプロセッサに、

前記 1 以上のローカルメモリバッファにデータを送信するように前記並行処理装置の 1 以上を構成させる命令と、

前記 1 以上のローカルメモリバッファからデータを受信するように前記並行処理装置の

10

20

30

40

50

1 以上を構成させる命令と、

をさらに備える請求項 3 2 に記載のコンピュータ可読媒体。

【請求項 3 5】

前記プログラム可能なプロセッサに、前記ローカルメモリバッファにデータを送信し前記ローカルメモリバッファからデータを受信するように前記並行処理装置を構成するために 1 以上のアプリケーションプログラミングインタフェース (A P I) を実行させる命令、をさらに備える請求項 3 4 に記載のコンピュータ可読媒体。

【請求項 3 6】

前記プログラム可能なプロセッサに、

前記処理パイプラインにおいて前の処理装置からのデータ出力を保持するために前記ローカルメモリバッファの各々についての必要な幅を決定させる命令と、

前記決定された幅を有するように前記ローカルメモリバッファの各々を構成させる命令と、

をさらに備える請求項 3 2 に記載のコンピュータ可読媒体。

【請求項 3 7】

前記プログラム可能なプロセッサに、1 以上のアプリケーションプログラミングインタフェース (A P I) を実行して、前記ローカルメモリバッファの各々についての前記幅を決定し、前記決定された幅によって前記ローカルメモリバッファの各々を構成し、前記ローカルメモリバッファの各々の深さを決定させる命令、をさらに備える請求項 3 6 に記載のコンピュータ可読媒体。

【請求項 3 8】

前記プログラム可能なプロセッサに、前記ローカルメモリバッファの各々の深さを決定させる命令をさらに備え、前記ローカルメモリバッファの各々は前記深さを前記幅とトレードすることが構成可能である、請求項 3 6 に記載のコンピュータ可読媒体。

【発明の詳細な説明】

【技術分野】

【0001】

[0001] 本開示は、データを処理することに関し、より具体的には、汎用グラフィクス処理装置を使用してデータを処理することに関する。

【背景技術】

【0002】

[0002] 汎用グラフィック処理装置 (G P G P U) は、2 D および 3 D グラフィクスを処理するように元々設計されたグラフィック処理装置の一般化されたバージョンである。G P G P U は、G P U のハイパワー並行処理を、グラフィクス処理を超えて汎用データ処理アプリケーションに拡張する。一例として、G P U は、非グラフィカル計算のためにグラフィクス処理装置に一定のアプリケーションアクセスを与える O p e n C L 仕様にしたがってデータを処理するように構成されうる。「OpenCL Specification, Version 1.1」は 2 0 1 0 年 6 月にリリースされ、公的に入手可能である。

【0003】

[0003] G P G P U は、処理装置間の同期またはデータ共有を許容しない高並列構造 (highly parallel structure) で配置されたプログラム可能な処理装置を含む。代わりに、個々の処理装置は、外部メモリとデータセットのみ交換する。この構造により、G P G P U のためのアプリケーションは、本質的に並行であるものに限定される。G P G P U アーキテクチャは高並列処理されうるので、それらは、パイプラインベース計算の効率的な実装を阻む。この限定は各処理ステージにおける並行処理を使用する 2 D および 3 D グラフィクス処理に及ぶが、ステージ間の計算リソースのパイプライン化を必要とする。

【発明の概要】

【0004】

[0004] 本開示は、パイプラインベースのアプリケーションの効率的な処理を許容するために並行処理装置によって汎用グラフィクス処理装置 (G P G P U) のアーキテクチャを

10

20

30

40

50

拡張するための技法を説明する。例えば、本技法は、並行処理装置間の移送についてのデータを保持するように処理パイプラインのステージとして動作する並行処理装置に接続されたローカルメモリバッファを構成することを含めうる。ローカルメモリバッファは、並行処理装置間の、オンチップ、低電力、直接的なデータ移送を可能にする。ローカルメモリバッファは、並行処理装置間のデータ移送を可能にするためにハードウェアベースのデータフロー制御メカニズムを含めうる。このように、データは、ローカルメモリバッファを介して処理パイプラインにおいて1つの並行処理装置から次の並行処理装置へと直接受け渡され、実際には並行処理装置を一連のパイプラインステージに変換する。ローカルメモリバッファは、処理パイプラインにおける並行処理装置の各々がシステムメモリへの呼び出しを行ない、データを取り出すおよび/または格納する必要性を減らすまたは取り除くことによって、メモリ帯域幅使用量を著しく減らすことができる。

10

【0005】

[0005] 本技法は、いくつかの例では、前の並行処理装置からのデータ出力を保持するためにバッファに必要とされる幅を有するようにローカルメモリバッファの各々を構成することを含めうる。例えば、ローカルメモリバッファは、幅と深さを交換することが構成可能であるハードウェアベースのバッファでありうる。さらに、いくつかの例では、本技法は、処理パイプライン内でデータシーケンスをプリザーブ(preserve)するためにシーケンシングバリア(sequencing barriers)を実行することを含めうる。例えば、データセットのデータスレッドのシーケンスは、並行処理装置にデータセットが入るときに記録され、データセットが処理された後で、データセットのデータスレッドは、記録されるのと同じシーケンスで並行処理装置からリリースされうる。

20

【0006】

[0006] 一例では、本開示は、処理パイプラインのステージとして選択的に動作するように構成された2以上の並行処理装置と、並行処理装置間の移送のためにデータを保持するように構成された1以上のローカルメモリバッファであって、バッファの各々が並行処理装置のうち少なくとも2つの間で接続される、1以上のローカルメモリバッファと、を備えるGPGPUを対象とする。

【0007】

[0007] 別の例では、本開示は、処理パイプラインのステージとして選択的に動作するように2以上の並行処理装置を構成することと；行処理装置間の移送のためにデータを保持するように1以上のローカルメモリバッファを構成することと、なお、バッファの各々は、並行処理装置のうち少なくとも2つの間で接続される；を備えるGPGPUによってデータを処理する方法を対象とする。

30

【0008】

[0008] さらなる例では、本開示は、処理パイプラインのステージとして選択的に動作するように2以上の並行処理装置を構成するための手段と；並行処理装置間の移送のためにデータを保持するように1以上のローカルメモリバッファを構成するための手段と、なお、バッファの各々は、並行処理装置のうち少なくとも2つの間で接続される；を備えるGPGPUを対象とする。

【0009】

[0009] 別の例では、本開示は、GPGPUによってデータを処理するための命令を備えるコンピュータ可読媒体であって、実行時に、プログラム可能なプロセッサに、処理パイプラインのステージとして選択的に動作するように2以上の並行処理装置を構成させ、並行処理装置間の移送のためにデータを保持するように1以上のローカルメモリバッファを構成させる、なお、バッファの各々が並行処理装置のうち少なくとも2つの間で接続される、コンピュータ可読媒体を対象とする。

40

【0010】

[0010] 1以上の例の詳細は、添付図面および下記の詳細な説明で記載されている。他の特徴、目的および利点は、明細書および図面、ならびに、特許請求の範囲から明らかとなるであろう。

50

【図面の簡単な説明】

【0011】

【図1】図1は、処理パイプラインを実装することが構成可能である汎用グラフィックス処理装置（GPGPU）を含むデバイスを図示するブロック図である。

【図2】図2は、並行処理を実行するように構成された並行処理装置を含む従来のGPGPUを図示するブロック図である。

【図3】図3は、処理パイプラインを実装するように構成されたローカルメモリバッファと並行処理装置を含む図1のGPGPUの一例を図示するブロック図である。

【図4】図4は、処理パイプラインのステージとして並行処理装置間でデータを移送するために並行処理装置に接続されたローカルメモリバッファを含むGPGPUの例示的な動作を図示するフローチャートである。

【図5】図5は、GPGPUのローカルメモリバッファと並行処理装置によって実装される処理パイプライン内でデータシーケンスをプリザーブする例示的な動作を図示するフローチャートである。

【詳細な説明】

【0012】

[0016] 本開示は、パイプラインベースのアプリケーションの効率的な処理を許容するために並行処理装置によって汎用グラフィックス処理装置（GPGPU）のアーキテクチャを拡張するための技法を説明する。具体的には、本技法は、並行処理装置間の移送についてのデータを保持するように処理パイプラインのステージとして動作する並行処理装置に接続されたローカルメモリバッファを構成することを含む。ローカルメモリバッファは、並行処理装置間の、オンチップ、低電力、直接的なデータ移送を可能にする。ローカルメモリバッファは、並行処理装置間のデータ移送を可能にするためにハードウェアベースのデータフロー制御メカニズムを含めうる。このように、データは、ローカルメモリバッファを介して処理パイプラインにおいて1つの並行処理装置から次の並行処理装置へと直接受け渡され、実際には並行処理装置を一連のパイプラインステージに変換する。ローカルメモリバッファは、処理パイプラインにおける並行処理装置の各々がシステムメモリへの呼び出しを行ない、データを取り出しおよび/または格納する必要性を減らすまたは取り除くことによって、メモリ帯域幅使用量を著しく減らすことができる。

【0013】

[0017] 図1は、処理パイプライン10を実装することが構成可能である汎用グラフィックス処理装置（GPGPU）6を含むデバイス2を図示するブロック図である。以下でより詳細に説明されるように、GPGPU6の処理パイプライン10は、処理パイプライン10のステージとして動作するように構成された2以上の並行処理装置と、処理パイプライン10を実装するために並行処理装置間の移送のためにデータを保持するように構成された1以上のローカルメモリバッファとを含む。

【0014】

[0018] デバイス2は、データを送信および受信すること、様々なデータ処理アプリケーションをサポートすること、および、ユーザへの提示のために処理されたデータを出力することが可能である。デバイス2の例は、限定されないが、モバイル無線電話、携帯情報端末（PDA）、ビデオゲーミングデバイス、ビデオゲーミングコンソール、ビデオコンファレンシング装置（video conferencing units）、ラップトップコンピュータ、デスクトップコンピュータ、タブレットコンピュータ、テレビセットトップボックス、デジタル記録デバイス、デジタルメディアプレイヤー、および同様なものを含む。

【0015】

[0019] 図1で図示された例では、デバイス2は、ホストプロセッサ4、処理パイプライン10を伴うGPGPU6、ディスプレイ8、スピーカ10、デバイスメモリ12、トランシーバモジュール14、およびユーザ入力デバイス16を含む。他のケースでは、例えば、デバイス2がデスクトップコンピュータである場合、ディスプレイ8、スピーカ10および/またはユーザインタフェース16はデバイス2に外付けであってもよい。ホスト

10

20

30

40

50

プロセッサ 4 と G P G P U 6 は、デジタル信号プロセッサ (D S P)、汎用マイクロプロセッサ、特定用途向け集積回路 (A S I C)、フィールドプログラム可能なゲートアレイ (F P G A)、または他の同等な集積またはディスクリート論理回路を備えうる。

【 0 0 1 6 】

[0020] ホストプロセッサ 4 は 1 以上のアプリケーションを実行しうる。アプリケーションの例は、ウェブブラウザ、e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、オーディオおよびビデオ編集アプリケーション、または、ディスプレイ 8 および / またはスピーカ 1 0 を介したユーザへの提示のためのビジュアルおよび / またはオーディオ出力を生成する他のアプリケーション、を含む。G P G P U 6 はまた 1 以上のアプリケーションを実行しうる。G P G P U 6 は、ホストプロセッサ 4 によって実行されたアプリケーションをサポートしてアプリケーションを実行しうる。具体的には、G P G P U 6 は、ディスプレイ 8 および / またはスピーカ 1 0 を介してユーザへの提示のためのデータを準備するためにアプリケーションを実行しうる。

10

【 0 0 1 7 】

[0021] G P G P U 6 は、グラフィクス処理を超えて汎用データ処理アプリケーションに G P U の高電力並行処理を拡張するグラフィクス処理装置 (G P U) の一般化されたバージョンである。例として、G P G P U 6 は、非グラフィカル計算のために G P U に一定のアプリケーションアクセスを与える O p e n C L 仕様にしたがってデータを処理するように構成されうる。従来の G P G P U は、図 2 を参照して以下でさらに詳細に説明され、パイプラインベースのアプリケーションの効率的な実装を防ぐ高並列構造で並べられるプログラム可能な処理装置を含む。この限定は各処理ステージで並行処理を使用する 2 D および 3 D グラフィクス処理アプリケーションに及ぶが、ステージ間の計算リソースのパイプライン化を必要とする。

20

【 0 0 1 8 】

[0022] パイプラインベースのアプリケーションは、第 1 ステージがオリジナルデータセットを処理するように構成され、第 2 ステージが第 1 ステージの出力を処理するように構成され、第 3 ステージが第 3 ステージの出力を処理するように構成され、アプリケーションに必要とされるステージの数について同様に続くようにステージで処理されるべきデータセットを必要とする。パイプラインベースのアプリケーションの最も効率的な実装は、処理パイプラインにおいてあるステージから次のステージへと直接データセットを受け渡すことである。パイプラインベースのアプリケーションのあまり効率的でない実装は、処理パイプラインにおける各ステージについて、オフチップメモリから前回のステージによって処理されたデータを取り出し、そのあとで、次のステージのためのオフチップメモリに戻って処理されたデータを格納することである。このあまり効率的でない実装はいまだに、データセットが処理パイプラインにおいて各ステージによって正しいシーケンスで処理されることを確実にするシーケンシングメカニズムを必要とする。従来の G P G P U は、処理パイプライン、または、パイプラインベースアプリケーションを実行するのに必要なシーケンシングメカニズムさえも実装するように構成されることができない。

30

【 0 0 1 9 】

[0023] 本開示における技法によれば、また、従来の G P G P U とは異なり、いくつかの例において、G P G P U 6 は、2 D および 3 D グラフィクス処理アプリケーションを含むパイプラインベースのアプリケーションを実行するために処理パイプライン 1 0 を実装することが構成可能である。図 3 を参照して以下でより詳細に説明されるように、G P G P U 6 の処理パイプライン 1 0 は、処理パイプライン 1 0 のステージとして動作するように構成された 2 以上の並行処理装置と、処理パイプライン 1 0 を実装するために並行処理装置間の移送のためにデータを保持するように構成された 1 以上のローカルメモリバッファとを含む。処理パイプライン 1 0 に含まれるローカルメモリバッファは、並行処理装置間の、オンチップ、低電力、直接的なデータ移送を可能にする。このように、データは、ローカルメモリバッファを介して処理パイプライン 1 0 において 1 つの並行処理装置から次の並行処理装置へと直接受け渡され、実際には並行処理装置を一連のパイプラインステージ

40

50

に変換する。処理パイプライン10の実装は、処理パイプライン10における並行処理装置の各々がGPGPU6からオフチップで位置されるデバイスメモリ12への呼び出しを行ない、データを取り出すおよび/または格納する必要性を減らすまたは取り除くことによって、メモリ帯域幅使用量を著しく減らすことができる。

【0020】

[0024] 本開示の技法は、前の並行処理装置からのデータ出力を保持するためにバッファに必要とされる幅を有するように処理パイプライン10内でローカルメモリバッファの各々を構成することを含めうる。例えば、ローカルメモリバッファは、深さを幅と交換することが構成可能であるハードウェアベースのバッファでありうる。さらに、本技法は、処理パイプライン10内でデータシーケンスをプリザブするためにシーケンシングバリアを実行することを含む。例えば、データセットのデータスレッドのシーケンスは、データセットが処理パイプライン10内の並行処理装置に入るときに記録され、データセットが処理された後で、データセットのデータスレッドは、記録されるのと同じシーケンスで並行処理装置からリリースされうる。

10

【0021】

[0025] 例えば、GPGPU6が処理パイプライン10を実装するように構成されるとき、GPGPU6は、ウェブブラウザ、eメール、ビデオゲーム、およびホストプロセッサ4によって実行されるビデオ編集アプリケーションをサポートして、パイプラインベースの2Dおよび3Dグラフィックス処理アプリケーションを実行しうる。別の例として、GPGPU6が処理パイプライン10を実装するように構成されないとき、GPGPU6は、画像ベースの探索アプリケーション、画像記述子生成/抽出、ラジオメトリック画像調整(radiometric image adjustments)、オーディオ処理、およびホストプロセッサ4によって一般的に実行される他の動作のような高並行構造で効率的に動作するアプリケーションを実行しうる。

20

【0022】

[0026] ある場合には、GPGPU6が、パイプラインベースのグラフィック処理アプリケーションをサポートしてアプリケーションを実行しうる。パイプラインベースのグラフィックス処理アプリケーションは、処理パイプライン10を使用するGPGPU6自体によって、または、デバイス2に含まれる別個のGPUによって実行されうる。例えば、GPGPU6は、画像特殊効果アプリケーション、GPUパイプラインのための頂点(vertices)生成、およびGPUパイプラインからのカラーバッファを使用するグラフィックスポスト処理アプリケーションを実行しうる。

30

【0023】

[0027] ディスプレイ8およびスピーカ10は双方とも、デバイス2のための出力デバイスを備える。あるケースでは、ディスプレイ8とスピーカ10は、ユーザにビジュアルおよびオーディオ出力の両方を提示するために一緒に使用されうる。他のケースでは、ディスプレイ8とスピーカ10は、ユーザに出力を提示するために、別々に使用されうる。例として、ディスプレイ8は、液晶ディスプレイ(LCD)、ブラウン管(CRT)ディスプレイ、プラズマディスプレイまたは別のタイプのディスプレイデバイスを備えうる。

【0024】

[0028] ユーザ入力デバイス16は、デバイス2のための1以上のユーザ入力デバイスを備える。例えば、ユーザ入力デバイス16は、トラックボール、マウス、キーボード、マイクロフォン、および/または他のタイプの入力デバイスを含めうる。他の例では、ユーザ入力デバイス16は、タッチスクリーンを備え、ディスプレイ8の一部として組み込まれうる。ユーザは、ユーザ入力デバイス16を介してホストプロセッサ4および/またはGPGPU6によって実行されるべき1以上のアプリケーションを選択しうる。

40

【0025】

[0029] ホストプロセッサ4は、トランシーバモジュール14を介してホストプロセッサ4および/またはGPGPU6によって処理されるべきデータをダウンロードしうる。ホストプロセッサ4はまた、トランシーバモジュール14を介してホストプロセッサ4およ

50

び/またはG P G P U 6によって実行される1以上のアプリケーションをダウンロードしうる。トランシーバモジュール14は、デバイス2と他のデバイスとの間の無線通信または有線通信、またはネットワークを可能にする回路を含めうる。トランシーバモジュール14は、変調器、復調器、増幅器、および有線通信または無線通信のための他の当該回路を含めうる。

【0026】

[0030] デバイスメモリ12は、ホストプロセッサ4および/またはG P G P U 6によって処理されるべきデータを格納し、また、ホストプロセッサ4および/またはG P G P U 6から受信される処理されたデータを格納しうる。さらに、デバイスメモリ12は、ホストプロセッサ4および/またはG P G P U 6によって実行された1以上のアプリケーションを格納しうる。デバイスメモリ12は、1以上のコンピュータ可読記憶媒体を備えうる。デバイスメモリ12の例は、限定されないが、ランダムアクセスメモリ(RAM)、読み出し専用メモリ(ROM)、電子的に消去可能なプログラム可能な読み出し専用メモリ(EEPROM(登録商標))、CD-ROMまたは他の光学ディスクストレージ、磁気ディスクストレージまたは他の磁気ストレージデバイス、フラッシュメモリ、または命令またはデータ構造の形式で所望プログラムコードを搬送または格納するために使用されることができ、また、コンピュータまたはプロセッサによってアクセスされることができ、任意の他の媒体を含む。

10

【0027】

[0031] 図2は、並行処理を実行するように構成された並行処理装置22A-22Dを含む従来のG P G P U 18を図示するブロック図である。いくつかの例では、G P G P U 18は、図1を参照して上述されているデバイス2と実質的に同様なデバイス内に含まれうる。G P G P U 18は、データ配信装置20、並行処理装置22A-22D(「並行処理装置22」)、およびG P G P U 18に外付けのデバイスメモリ26に並行処理装置22を接続するバス24を含む。

20

【0028】

[0032] 従来のG P G P U 18は、2Dおよび3Dグラフィックスを処理するように元々設計されたGPUの一般化バージョンである。G P G P U 18は、GPUの高電力並行処理を、グラフィックス処理を超えて汎用処理アプリケーションに拡張することができる。例として、G P G P U 18は、OpenCL仕様に従ってデータを処理するように構成されうる。OpenCL仕様は、非グラフィカルコンピューティングのためにGPUに一定のアプリケーションアクセスを与える。OpenCL用語では、データスレッドは作業項目(work item)と呼ばれ、データセットは作業グループ(work group)と呼ばれ、処理装置は計算装置(compute units)と呼ばれ、処理装置の集まりは、計算グループ(compute group)と呼ばれる。

30

【0029】

[0033] 一般的なGPUタスクは高度に並行であり、所与の処理装置内で処理されているデータセットのデータスレッド間での情報交換を必要としない。例えば、頂点について計算された値は、異なる頂点について計算された値から独立しており、ピクセルについて計算された値は、異なるピクセルについて計算された値から独立している。GPUの並行性質を模倣するために、G P G P U 18は、高並行構造で配列された並行処理装置22を含むように設計される。

40

【0030】

[0034] G P G P U 18のアーキテクチャは、並行処理装置22間のデータ共有または同期を許容しないほど、高並行である。動作において、データ配信装置20は、並行処理装置22の各々に、デバイスメモリ26に格納されたデータセットを割り当てる。処理中、割り当てられたデータセットのデータスレッドは、並行処理装置22の各々の内で共有され同期されうる。しかしながら、異なるデータセットのデータスレッドは、並行処理装置22間で共有または同期されることができない。代わりに、並行処理装置22の各々は、バス24を介してデバイスメモリ26と割り当てられたデータセットのみ交換する。より

50

具体的には、並行処理装置 22 の各々は、バス 24 を介してデバイスメモリ 26 から処理についての割り当てられたデータセットを取り出し、データセットを処理した後で、バス 24 を介してデバイスメモリ 26 に戻って処理されたデータセットを格納する。

【 0031 】

[0035] G P G P U 18 の並行アーキテクチャは、並行処理装置 22 間のパイプラインベースアプリケーションの効率的な実装を阻む。パイプラインベースのアプリケーションでは、処理装置は、異なる処理タスクについて 1 つのステージから別のステージへとデータが移動することを可能にするためにパイプラインにおいてステージとして接続される。G P G P U 18 におけるパイプラインベースのアプリケーションに対する限定は、2 D および 3 D グラフィクス処理アプリケーションに拡張する、そしてそれは各処理ステージで並行処理を使用するが、ステージ間でのパイプライン化を必要とする。

10

【 0032 】

[0036] したがって、G P G P U 18 のアプリケーションは、本質的に並行であるものに限定される。並行処理装置 22 の各々は算術論理装置 (A L U) のクラスタまたは他の構成可能な論理素子を備えうる。したがって、並行処理装置 22 は、G P G P U 18 によって実行されるアプリケーションに依存して異なる動作を実行することがプログラム可能または構成可能である。G P G P U 18 の高並行構造で効率的に動作するアプリケーションは、画像ベースの探索アプリケーション、画像記述子生成/抽出、ラジオメトリック画像調整 (radiometric image adjustments)、オーディオ処理、およびデジタル信号プロセッサ (D S P) によって一般的に実行される他の動作および同様なものを含めうる。さらに、G P G P U 18 によって実行されるアプリケーションは、画像特殊効果生成、G P U パイプラインのための頂点生成、G P U パイプラインからのカラーバッファを使用してグラフィクスポスト処理動作のようなパイプラインベースのグラフィクス処理アプリケーションとのインタラクションを必要としうる。

20

【 0033 】

[0037] 図 3 は、図 1 の例示的な G P G P U 6 を図示するブロック図であり、処理パイプライン 10 を実装するように構成されたローカルメモリバッファ 44 A - 44 C と並行処理装置 42 A - 42 D とを含む。他の例では、G P G P U 6 は、より多数またはより少数の並行処理装置およびローカルメモリバッファを含めうる。

【 0034 】

[0038] 図 3 の例では、G P G P U 6 は、データ配信装置 40、並行処理装置 42 A - 42 D (「 並行処理装置 42 」) および G P G P U 6 に外付けのデバイスメモリ 12 (図 1) に並行処理装置 42 を接続するバス 46 を含む。従来の G P G P U とは異なり (例えば、図 3 の G P G P U 18)、G P G P U 6 はまた、並行処理装置 42 間で接続されたローカルメモリバッファ 44 A - 44 C (「 ローカルメモリバッファ 44 」) を含む。並行処理装置 42 と並行処理装置 42 間で接続されたローカルメモリバッファ 44 の組み合わせは、処理パイプライン 10 と呼ばれうる。G P G P U 6 はまた、制御装置 30 およびローカルメモリ 38 を含む。ローカルメモリ 38 は、ローカルメモリバッファ 44 に類似したバッファ、レジスタ、または G P G P U 6 のデータを一時的に格納するキャッシュを備えうる。制御装置 30 は、アプリケーションプログラミングインタフェース (A P I) 32、バッファマネージャ 34、およびシーケンスマネージャ 36 を含む。

30

40

【 0035 】

[0039] ローカルメモリバッファ 44 は、並行処理装置 42 間のデータ移送を可能にするハードウェアベースのデータフロー制御メカニズムを含めうる。例えば、ローカルメモリバッファ 44 は、ハードウェアベースの先入れ先出し (F I F O) バッファ、後入れ先出し (L I F O) バッファまたはインデクス付けされたバッファのような他のタイプのハードウェアベースのバッファを備えうる。ローカルメモリバッファ 44 A がハードウェアベースの F I F O を備える場合には、例えば、ローカルメモリバッファ 44 A は、バッファにデータを書き込むスペースがあるときローカルメモリバッファ 44 A へデータを並行処理装置 42 A が送信し、そうでないときには書き込み要求をストールすることを可能にす

50

るデータフロー制御メカニズムを含む。その場合、ローカルメモリバッファ44Aはまた、バッファから読み出すのに利用可能なデータがあるとき、ローカルメモリバッファ44Aからデータを並行処理装置42Bが受信し、そうでないときには読み出し要求をストールすることを可能にするデータフロー制御メカニズムを含む。ローカルメモリバッファ44がハードウェアベースのデータフロー制御メカニズムを含むとき、あまり効率的でないソフトウェアベースのデータフロー制御は、並行処理装置42間のデータの移送を可能にするのに必要ではない。

【0036】

[0040] ローカルメモリバッファ44は、並行処理装置42間の、オンチップ、低電力、直接的なデータ移送を可能にする。ローカルメモリバッファ44は「ローカル」である、なぜならば、それらは、GPGPU6内で、処理装置42と同じチップ上で位置されるからである。このように、データは、ローカルメモリバッファ44を介して処理パイプライン10において並行処理装置42の一方から並行処理装置42の他方へと直接受け渡される。並行処理装置42は、GPGPU6に外付けであるまたはGPGPU6からオフチップに配置されているデバイスメモリ12でデータを繰り返し取り出し格納することを必要としていない。したがって、ローカルメモリバッファ44は、並行処理装置42を一連のパイプラインステージに変換し、GPGPU6内で処理パイプライン10を実装する。

10

【0037】

[0041] 図示された例では、ローカルメモリバッファ44の各々は、処理パイプライン10が純粋に直列なパイプラインであるように連続順で並行処理装置42の2つの間で直接接続される。ローカルメモリバッファ44は、それが2つの並行処理装置42によってのみアクセス可能であるように、それらが接続され並行処理装置42のいずれによってもアドレス可能なバスに、「直接」接続される。例えば、ローカルメモリバッファ44Aは、並行処理装置42Aおよび42Bとの間で直接接続され、ローカルメモリバッファ44Bは、並行処理装置42Bと42Cとの間で直接接続され、ローカルメモリバッファ44Cは、並行処理装置42Cと42Dとの間で直接接続される。

20

【0038】

[0042] 他の例では、メモリバッファ44の各々はまた、連続順でない並行処理装置42のうち1以上に直接接続されうる。この場合、ローカルメモリバッファ44の各々は、クロスバー接続を介して並行処理装置42のいずれかに直接接続されうる。例えば、ローカルメモリバッファ44Aは、並行処理装置42Aがローカルメモリバッファ44Aを介して並行処理装置42B-42Dのいずれかにデータを移送しうるように、クロスバー接続を介して並行処理装置42の各々に直接接続されうる。クロスバー接続の使用は、ローカルメモリバッファ44を並行処理装置42に対してより幅広くアクセス可能にし、純粋に直列ではない処理パイプラインの実装を可能にする。

30

【0039】

[0043] 処理パイプライン10が純粋に直列なパイプラインを備える図示された例では、並行処理装置42は、ローカルメモリバッファ44の次のもの(successive one)にデータを書き込む許可のみを有し、ローカルメモリバッファ44の前のものからデータを読み出す許可のみを有しうる。例えば、並行処理装置42Bは、ローカルメモリバッファ44Aからデータを読み出すことのみ可能であり、ローカルメモリバッファ44Bにデータを書き込むことのみ可能でありうる。処理パイプラインがクロスバー接続を含める場合、並行処理装置42は、ローカルメモリバッファ44のいずれかに読み出し且つ書き込む許可を有しうる。例えば、並行処理装置42Bは、ローカルメモリバッファ44Aで、また、ローカルメモリバッファ44Bで、データを読み出し書き込むことが可能であることがある。

40

【0040】

[0044] 上述されるように、ローカルメモリバッファ44は、FIFOバッファ、LIFOバッファ、またはインデックス付けされたバッファのうちの少なくとも1つを備えうる。ローカルメモリバッファ44に使用されるバッファのタイプは、処理パイプライン10で

50

必要とされるハードウェアベースのデータフロー制御メカニズムのタイプに依存しうる。ローカルメモリバッファ44に使用されるバッファのタイプはまた、ローカルメモリバッファ44が1対1接続またはクロスバー接続を介して並行処理装置42に接続されるかに依存しうる。さらに、クロスバー接続が使用されるとき、制御装置30のバッファマネージャ34は、どの並行処理装置42が所与時間にどのローカルメモリバッファ44にアクセスするかを管理するために、いくらかのメモリ制御を実行する必要があることがある。

【0041】

[0045] 上述されるように、ローカルメモリバッファ44は、1対1またはクロスバー接続のいずれかを介して並行処理装置42の少なくとも2つの間で直接接続されうる。しかしながら、ローカルメモリバッファ44は、並行処理装置42によってアドレス可能なバスでないことがある。このように、ローカルメモリバッファ44の指定されたメモリコントローラは必要でないことがある。具体的には、メモリコントローラは、バスにわたってローカルメモリバッファ44に対して読み出しおよび書き込みコマンドを処理する必要はない。

10

【0042】

[0046] ローカルメモリバッファ44は、並行処理装置42の各々がバス46を介してデバイスメモリ12への呼び出しを行ない、データを取り出すおよび/または格納する必要性を減らすまたは取り除くことによって、メモリ帯域幅使用量を著しく減らすことができる。動作において、並行処理装置42Aは、処理パイプライン10の第1の処理装置として、バス46を介してデバイスメモリ12からオリジナルデータセットを取り出す。データセットは、データ配信装置40によって並行処理装置42Aに割り当てられうる。さらに、並行処理装置42Dは、処理パイプライン10の最終処理装置として、バス46を介してデバイスメモリ12にポストパイプラインデータセットを格納する。並行処理装置42Bおよび42Cは、処理パイプライン10の中間処理装置として、ローカルメモリバッファ44のうち一方を介して並行処理装置42のうち前のものからデータセットを受信し、ローカルメモリバッファ44のうち他方を介して並行処理装置42のうち後続のものにデータセットを送信する。したがって、中間処理装置は、データを取り出しおよび/または格納するためにデバイスメモリ12と相互作用することを必要とされない。いくつかの場合では、中間処理装置は、処理パイプライン10の特定のステージを実行するためにデバイスメモリから補助データを取り出しうる。しかしながら、処理用の主要なデータセットは、ローカルメモリバッファ44を介して処理パイプライン10に沿って直接受け渡される。

20

30

【0043】

[0047] 上述されるように、GPGPU6は、グラフィクス処理を超えて汎用データ処理アプリケーションにGPUの高電力並行処理を拡張するGPUの一般化されたバージョンである。例として、GPGPU6は、非グラフィカル計算のためにグラフィクス処理装置に一定のアプリケーションアクセスを与えるOpenCL仕様にしたがってデータを処理するように構成されうる。OpenCL用語では、データスレッドは作業項目(work item)と呼ばれ、データセットは作業グループ(work group)と呼ばれ、処理装置は計算装置(compute units)と呼ばれ、処理装置の集まりは、計算グループ(compute group)と呼ばれる。

40

【0044】

[0048] 本開示の技法によれば、GPGPU6は、2Dおよび3Dグラフィクス処理アプリケーションを含むパイプラインベースのアプリケーションを実行するために処理パイプライン10を実装することが構成可能である。より具体的には、GPGPU6の制御装置30は、処理パイプラインのステージとして動作するように並行処理装置42を構成する。制御装置30はまた、並行処理装置42間の移送のためのデータを保持するように、並行処理装置42間で接続されたローカルメモリバッファ44を構成する。

【0045】

[0049] 並行処理装置42は、GPGPU6によって実行されるアプリケーションに依存

50

して異なる動作を実行することがプログラム可能または構成可能でありうる。制御装置 30 は、アプリケーションにしたがって動作するように並行処理装置 42 の各々を構成しうる。例えば、並行処理装置 22 の各々は算術論理装置 (ALU) のクラスタまたは他の構成可能な論理素子を備えうる。

【0046】

[0050] ローカルメモリバッファ 44 はまた、GPGPU6 によって実行されるアプリケーションに依存して並行処理装置 42 からの異なるタイプのデータ出力を保持することがプログラム可能または構成可能でありうる。例えば、ローカルメモリバッファ 44 は、ハードウェアベースのバッファを備えうるが、構成可能な態様のセット (a set of configurable aspects) を含めうる。構成可能な態様の 1 つは、並行処理装置 42 からの異なるタイプのデータ出力を適応させるためのローカルメモリバッファ 44 の幅でありうる。例えば、ローカルメモリバッファ 44 は、深さを幅とトレードすることが構成可能でありうる。制御装置 30 のバッファマネージャ 34 は、並行処理装置 42 のうち前のもののデータ出力を保持するためにローカルメモリバッファ 44 の各々に必要とされる幅を決定しうる。バッファマネージャ 34 は、並行処理装置 42 の各々からデータ出力のタイプを認識するので、データを保持するためにローカルメモリバッファ 44 の各々によって必要とされる幅を認識する。バッファマネージャ 34 は、そのあとで、決定された幅を有するようにローカルメモリバッファ 44 の各々を構成しうる。

10

【0047】

[0051] いったん並行処理装置 42 とローカルメモリバッファ 44 が GPGPU6 内で処理パイプライン 10 を実装するように構成されると、並行処理装置 42 は、ローカルメモリバッファ 44 を介してデータを移送しうる。制御装置 30 は、ローカルメモリバッファ 44 にデータを送信するように並行処理装置 42 のうち 1 以上を構成し、ローカルメモリバッファ 44 からデータを受信するように並行処理装置 44 のうち 1 以上を構成しうる。例えば、制御装置 30 は、それぞれ、ローカルメモリバッファ 44 A、44 B、および 44 C にデータを送信するように並行処理装置 42 A、42 B および 42 C を構成しうる。制御装置 30 はまた、それぞれ、ローカルメモリバッファ 44 A、44 B、および 44 C からデータ受信するように並行処理装置 42 B、42 C、および 42 D を構成しうる。

20

【0048】

[0052] ハードウェアベースのフロー制御メカニズムを有するローカルメモリバッファ 44 は、新規 API 32 を導入することによって、OpenCL 規格のような GPGPU6 規格を使用して露出されうる (exposed)。例えば、制御装置 30 は、API 32 の 1 以上を実行して、ローカルメモリバッファ 44 の各々に必要とされる幅を決定し、決定された幅でローカルメモリバッファ 44 の各々を構成し、ローカルメモリバッファ 44 の各々の深さを決定しうる。さらに、制御装置 30 は、API 32 の 1 以上を実行してローカルメモリバッファ 44 にデータを送信しローカルメモリバッファ 44 からデータを受信するように並行処理装置 42 を構成しうる。ローカルメモリバッファ 44 に含まれるハードウェアベースのデータフロー制御メカニズムは、並行処理装置 42 が、さらなるソフトウェアベースのデータフロー制御なしに、ローカルメモリバッファ 44 にデータを送信し、ローカルメモリバッファ 44 からデータを受信することを可能にする。

30

40

【0049】

[0053] さらに GPGPU6 の制御装置 30 は、並行処理装置 42 のうち 1 以上内でデータシーケンスをプリザーブすることによって処理パイプライン 10 内でデータシーケンスをプリザーブしうる。GPGPU6 によって実行されるパイプラインベースのアプリケーション、特に 3D グラフィクスアプリケーションは、処理パイプライン 10 内で一定のシーケンスで処理されるべきデータを必要としうる。データが処理パイプラインの各ステージで処理されるとき、データは、条件、キャッシュヒットまたはミス、および同様なもののような実行課題 (execution issues) に起因してシーケンスを変更しうる。制御装置 30 のシーケンスマネージャ 36 は、並行処理装置 42 の少なくともいくつかの内でデータシーケンスをプリザーブするためにシーケンシングバリアを実行しうる。シーケンシング

50

バリアは、処理パイプライン 10 内の処理速度を減速させることができるので、シーケンスマネージャ 36 は、精確な処理のためにデータシーケンスプリザベーションを必要とするこれらの並行処理装置 42 においてシーケンシングバリアのみを実行しうる。

【0050】

[0054] シーケンスマネージャ 36 によって実行されるシーケンシングバリアは、シーケンス決定カウンタ (SDC) およびシーケンスエンフォーシングバリア (SEB) を含めうる。例えば、シーケンシングバリアは、SDC および SEB について OpenCL 言語に新しい関数呼び出しを追加することによって、OpenCL 規格のような GPGPU 規格を使用して露出されうる (exposed)。

【0051】

[0055] シーケンスマネージャ 36 は、データセットが並行処理装置 42 のいずれか 1 つに入るとき SDC を実行しうる。シーケンスマネージャ 36 は、そのあとで、ローカルメモリ 38 内で受信されたデータセットのデータスレッドのシーケンスを記録することによって SDC 動作を実行する。例えば、シーケンスマネージャ 36 は、データスレッドがデバイスメモリ 12 から受信される順でデータセットの各データスレッドのインデックスを記録しうる。

【0052】

[0056] シーケンスマネージャ 36 は、データセットが並行処理装置 42 のうちの 1 つから出るときに、SEB を実行しうる。シーケンスマネージャ 36 は、そのあとで、SDC によって記録されるのと同じシーケンスで並行処理装置 42 のうちの 1 つからデータセットのデータスレッドをリリースすることによって SEB 動作を実行する。例えば、シーケンスマネージャ 36 は、ローカルメモリ 38 に記録されたデータスレッドインデックスにアクセスし、インデックスが記録された順にしたがって各データスレッドをリリースする。このように、データセットのデータスレッドは、データセットのデータスレッドが並行処理装置 42 のうち現在のものに入るのと同じ順で並行処理装置 42 の後続のものに入るであろう。

【0053】

[0057] 一例では、制御装置 30 は、パイプラインベースの 3D グラフィクス処理アプリケーションを実行するように GPGPU 6 を構成しうる。その場合、制御装置 30 は、3D グラフィクス処理パイプラインのステージとして動作するように並行処理装置 42 を構成しうる。例えば、制御装置 30 は、頂点シェーダとして動作するように並行処理装置 42 A を構成し、トライアングルラステライザとして動作するように並行処理装置 42 B を構成し、フラグメントシェーダとして動作するように並行処理装置 42 C を構成し、ピクセルブレンダとして動作するように並行処理装置 42 D を構成しうる。

【0054】

[0058] 制御装置 30 はまた、3D グラフィクス処理パイプライン 10 を実装するために並行処理装置 42 間の移送のためのデータを保持するようにハードウェアベースのデータフロー制御メカニズムによってローカルメモリバッファ 44 を構成しうる。例えば、制御装置 30 は、頂点シェーダとして動作する並行処理装置 42 A と、トライアングルラステライザとして動作する並行処理装置 42 B との間の移送のためのポスト頂点シェーダの頂点データを保持するようにローカルメモリバッファ 44 A を構成しうる。制御装置 30 は、トライアングルラステライザとして動作する並行処理装置 42 B と、フラグメントシェーダとして動作する並行処理装置 42 C との間の移送のためにプレフラグメントシェーダピクセルデータを保持するようにローカルメモリバッファ 44 B を構成しうる。最後に、制御装置 30 は、フラグメントシェーダとして動作している並行処理装置 42 C とピクセルブレンダとして動作している並行処理装置 42 D との間の移送のためにポストフラグメントシェーダピクセル値を保持するようにローカルメモリバッファ 44 C を構成しうる。

【0055】

[0059] 3D グラフィクス処理アプリケーションを実行するとき、データ配信装置 40 は、頂点シェーダとして動作している並行処理装置 42 A にオリジナル頂点データセットを

10

20

30

40

50

割り当てうる。並行処理装置 4 2 A は、バス 4 6 を介して、デバイスメモリ 1 2 から、割り当てられたオリジナル頂点データセットを取り出す。データセットが入るとき、シーケンスマネージャ 3 6 は、頂点データのシーケンスを記録するために S D C を実行する。並行処理装置 4 2 A はそのあとで頂点シェーディング動作を実行し、ローカルメモリバッファ 4 4 A にポスト頂点シェーダの頂点データを送信する。データセットが並行処理装置 4 2 A から出るとき、シーケンスマネージャ 3 6 は、S D C によって記録されるのと同じシーケンスで頂点データをリリースするために S E B を実行する。このように、頂点データは、頂点シェーダとして動作する並行処理装置 4 2 A に頂点データが入ったのと同じ順で、トライアングルラスライザとして動作する並行処理装置 4 2 B に到達するであろう。

【 0 0 5 6 】

[0060] トライアングルラスライザとして動作する並行処理装置 4 2 B は、ローカルメモリバッファ 4 4 A からポスト頂点シェーダの頂点データを受信する。いくつかの場合においては、並行処理装置 4 2 B はまた、トライアングルラスライズ化動作を実行するためにバス 4 6 を介してデバイスメモリ 1 2 から補助データを取り出しうる。並行処理装置 4 2 B はそのあとでトライアングルラスライズ化動作を実行し、ローカルメモリバッファ 4 4 B にプレフラグメントシェーダピクセルデータを送信する。いくつかの例では、シーケンスマネージャ 3 6 は、頂点データが並行処理装置 4 2 B に入るとき S D C を実行し、ピクセルデータが並行処理装置 4 2 B から出るとき S E B を実行してデータシーケンスをプリザーブしうる。他の例では、シーケンシングバリアは必須ではないので、並行処理装置 4 2 B に対して実行されない。

【 0 0 5 7 】

[0061] 並行処理装置 4 2 C は、フラグメントシェーダを動作し、ローカルメモリバッファ 4 4 B からプレフラグメントシェーダピクセルデータを受信する。データセットが入るとき、シーケンスマネージャ 3 6 は、ピクセルデータのシーケンスを記録するために S D C を実行する。いくつかの場合においては、並行処理装置 4 2 C はまた、フラグメントシェーダ動作を実行するためにバス 4 6 を介してデバイスメモリ 1 2 から補助データを取り出しうる。並行処理装置 4 2 C はそのあとでフラグメントシェーディング動作を実行し、ポストフラグメントシェーダピクセル値をローカルメモリバッファ 4 4 C に送信する。データセットが並行処理装置 4 2 C から出るとき、シーケンスマネージャ 3 6 は、S D C によって記録されるのと同じシーケンスでピクセルデータをリリースするために S E B を実行する。このように、ピクセルデータは、フラグメントシェーダとして動作している並行処理装置 4 2 C にピクセルデータが入ったのと同じ順で、ピクセルブレンダとして動作する並行処理装置 4 2 D に到達するであろう。

【 0 0 5 8 】

[0062] 並行処理装置 4 2 D は、ピクセルブレンダとして動作し、ローカルメモリバッファ 4 4 C からポストフラグメントシェーダピクセル値を受信する。並行処理装置 4 4 D は、ピクセルブレンディング動作を実行し、バス 4 6 を介してデバイスメモリ 1 2 にポストパイプラインデータセットを格納する。いくつかの例では、シーケンスマネージャ 3 6 は、ピクセルデータが並行処理装置 4 2 D に入るとき S D C を実行し、画像データが並行処理装置 4 2 D から出るとき S E B を実行してデータシーケンスをプリザーブしうる。他の例では、シーケンシングバリアは必須ではないので、並行処理装置 4 2 D に対して実行されない。3 D グラフィクス処理アプリケーションの上述された例は、単なる例示であり、開示された技法は、G P G P U 6 において様々なパイプラインベースのアプリケーションを実行するために使用されうる。

【 0 0 5 9 】

[0063] 図 4 は、処理パイプライン 1 0 のステージとしての並行処理装置間でデータを移送するために並行処理装置 4 2 に接続されたローカルメモリバッファ 4 4 を含む G P G P U 6 の例示的な動作を図示するフローチャートである。図示される動作は、図 3 の G P G P U 6 を参照して説明される。

【 0 0 6 0 】

10

20

30

40

50

[0064] G P G P U 6 の制御装置 3 0 は、処理パイプライン 1 0 のステージとして動作するように並行処理装置 4 2 を構成する (5 0)。例えば、制御装置 3 0 は、3 D グラフィックス処理パイプラインのステージとして動作するように並行処理装置 4 2 を構成する。その例では、制御装置 3 0 は、頂点シェーダとして動作するように並行処理装置 4 2 A を構成し、トライアングルラスライザとして動作するように並行処理装置 4 2 B を構成し、フラグメントシェーダとして動作するように並行処理装置 4 2 C を構成し、ピクセルブレンダとして動作するように並行処理装置 4 2 D を構成しうる。

【 0 0 6 1 】

[0065] 制御装置 3 0 はまた、並行処理装置 4 2 間の移送のためのデータを保持するようにローカルメモリバッファ 4 4 を構成し、結果、並行処理装置 4 2 を処理パイプライン 1 0 に変換する (5 2)。ローカルメモリバッファ 4 4 は、並行処理装置 4 2 間のデータ移送を可能にするためにハードウェアベースのデータフロー制御メカニズムを含めうる。例えば、ローカルメモリバッファ 4 4 は、ハードウェアベースの F I F O、L I F O、またはインデクス付けされたバッファを備えうる。ローカルメモリバッファ 4 4 は、並行処理装置 4 2 の少なくとも 2 つの間で直接接続されうる。例えば、3 D グラフィックス処理パイプラインの場合、ローカルメモリバッファ 4 4 A は、頂点シェーダとして動作する並行処理装置 4 2 A と、トライアングルラスライザとして動作する並行処理装置 4 2 B との間で直接接続され、ポスト頂点シェーダの頂点データ (post-vertex shader vertex data) を保持するように構成されうる。ローカルメモリバッファ 4 4 B は、トライアングルラスライザとして動作する並行処理装置 3 2 B と、フラグメントシェーダとして動作する並行処理装置 4 2 C との間で直接接続され、プレフラグメントシェーダのピクセルデータを保持するように構成されうる。最後に、ローカルメモリバッファ 4 4 C は、フラグメントシェーダとして動作する並行処理装置 4 2 C と、ピクセルブレンダとして動作する並行処理装置 4 2 D との間で直接接続され、ポストフラグメントシェーダピクセル値を保持するように構成されうる。

【 0 0 6 2 】

[0066] さらに、制御装置 3 0 のバッファマネージャ 3 4 は、並行処理装置 4 2 のうち前のものからのデータ出力を保持するためにローカルメモリバッファ 4 4 の各々に必要とされる幅を決定しうる (5 4)。バッファマネージャ 3 4 は、並行処理装置 4 2 の各々からデータ出力のタイプを認識するので、データを保持するためにローカルメモリバッファ 4 4 の各々によって必要とされる幅を認識する。バッファマネージャ 3 4 は、そのあとで、決定された幅を有するようにローカルメモリバッファ 4 4 の各々を構成しうる (5 6)。ある場合においては、ローカルメモリバッファ 4 4 は、ハードウェアベースでありうるが、構成可能な態様のセット (a set of configurable aspects) を含む。例えば、ローカルメモリバッファ 4 4 は、深さを幅とトレードすることが構成可能でありうる。

【 0 0 6 3 】

[0067] 例えば、バッファマネージャ 3 4 は、頂点シェーダとして動作している並行処理装置 4 2 A がポスト頂点シェーダの頂点データを出力するということを認識し、ポスト頂点シェーダの頂点データを保持するのに必要とされた幅を有するようにローカルメモリバッファ 4 4 A を構成しうる。バッファマネージャ 3 4 はまた、トライアングルラスライザとして動作している並行処理装置 4 2 B がプレフラグメントシェーダピクセルデータを出力するということを認識し、プレフラグメントシェーダ画素データを保持するのに必要とされた幅を有するようにローカルメモリバッファ 4 4 B を構成しうる。さらに、バッファマネージャ 3 4 は、フラグメントシェーダとして動作している並行処理装置 4 2 C がポストフラグメントシェーダピクセル値を出力するということを認識し、ポストフラグメントシェーダピクセル値を保持するのに必要とされる幅を有するようにローカルメモリバッファ 4 4 C を構成しうる。

【 0 0 6 4 】

[0068] いったん並行処理装置 4 2 とローカルメモリバッファ 4 4 が G P G P U 6 内で処理パイプライン 1 0 を実装するように構成されると、並行処理装置 4 2 は、ローカルメモ

10

20

30

40

50

リバッファ 44 を介して互いの間でデータを移送しうる (58)。より具体的には、制御装置 30 は、ローカルメモリバッファ 44 にデータを送信するように並行処理装置 42 のうち 1 以上を構成し、ローカルメモリバッファ 44 からデータを受信するように並行処理装置 44 のうち 1 以上を構成しうる。例えば、制御装置 30 は、それぞれ、ローカルメモリバッファ 44 A、44 B、および 44 C にデータを送信するように並行処理装置 42 A、42 B、および 42 C を構成しうる。制御装置 30 はまた、それぞれ、ローカルメモリバッファ 44 A、44 B、および 44 C からデータ受信するように並行処理装置 42 B、42 C、および 42 D を構成しうる。

【0065】

[0069] 図 5 は、GPGPU 6 の並行処理装置 42 とローカルメモリバッファ 44 とによって実装される処理パイプライン内でデータシーケンスをプリザブする例示的な動作を図示するフローチャートである。GPGPU 6 の制御装置 30 は、並行処理装置 42 のうち 1 以上内でデータシーケンスをプリザブすることによって処理パイプライン内でデータシーケンスをプリザブしうる。図示される動作は、図 3 の GPGPU 6 の並行処理装置 42 A を参照して説明される。同様な動作が他の並行処理装置 42 のうちのいずれについても実行されうる。

10

【0066】

[0070] 例として、並行処理装置 42 およびローカルメモリバッファ 44 は、3D グラフィクス処理パイプラインを実装するように構成されうる。その例では、並行処理装置 42 A は、頂点シェーダとして動作するように構成され、並行処理装置 42 B は、トライアングルラスタライザとして動作するように構成され、並行処理装置 42 C は、フラグメントシェーダとして動作するように構成され、並行処理装置 42 D は、ピクセルブレンダとして動作するように構成されうる。

20

【0067】

[0071] 処理パイプライン 10 のステージ、例えば頂点シェーダとして動作するように構成された並行処理装置 42 A は、処理のためにデータセットを受信する (62)。例えば、データ配信装置 40 は、頂点データのデータセットを並行処理装置 42 A に割り当て、並行処理装置 42 A は、バス 46 を介してデバイスメモリ 12 から割り当てられたデータセットを受信しうる。データセットが並行処理装置 42 A に入る時に、制御装置 30 のシーケンスマネージャ 36 は、シーケンス決定カウンタ (SDC) を実行する (64)。SDC にしたがって、シーケンスマネージャ 36 は、ローカルメモリ 38 内で受信されたデータセットのデータスレッドのシーケンスを記録する (66)。例えば、シーケンスマネージャ 36 は、データスレッドがデバイスメモリ 12 から受信される順でデータセットの各データスレッドのインデックスを記録しうる。

30

【0068】

[0072] 頂点シェーダとして動作するように構成された並行処理装置 42 A は、そのあとで、ポスト頂点シェーダの頂点データを生成するためにデータセットを処理する (68)。上述されているように、並行処理装置 42 A は、トライアングルラスタライザとして動作するように構成された並行処理装置 42 B にデータセットを移送するために、ローカルメモリバッファ 44 A にポスト頂点シェーダの頂点データを送信するように構成されうる。データセットが並行処理装置 42 A から出る時に、シーケンスマネージャ 36 は、シーケンスエンフォーシングバリア (SEB) を実行する (70)。SEB にしたがって、シーケンスマネージャ 36 は、SDC によって記録されるのと同じシーケンスで並行処理装置 42 A からデータセットのデータスレッドをリリースする (72)。例えば、シーケンスマネージャ 36 は、ローカルメモリ 38 に記録されたデータスレッドインデックスにアクセスし、インデックスが記録された順にしたがって各データスレッドをリリースする。このように、頂点シェーダとして動作するように構成された並行処理装置 42 A に複数の頂点が入ったのと同じ順でトライアングルラスタライザとして動作するように構成された並行処理装置 42 B に複数の頂点が入るであろう。

40

【0069】

50

[0073] 1以上の例では、説明された機能は、ハードウェア、ソフトウェア、ファームウェア、またはそれらのいずれの組み合わせにおいて実装されうる。ソフトウェアで実装される場合には、機能または動作は、非一時的なコンピュータ可読媒体で1以上の命令またはコードとして格納され、ハードウェアベースの処理装置によって実行されうる。コンピュータ可読媒体は、データ記憶媒体のようなタンジブル媒体に対応するコンピュータ可読媒体、または、例えば通信プロトコルにしたがって、1つの場所から別の場所へとコンピュータプログラムの移送を容易にする任意の媒体を含む通信媒体、を含めうる。このように、コンピュータ可読媒体は一般的に、(1)非一時的であるタンジブルコンピュータ可読記憶媒体または(2)信号または搬送波のような通信媒体に対応しうる。データ記憶媒体は、本開示で説明される技法の実装についての命令、コードおよび/またはデータ構造を取り出すために1以上のコンピュータまたは1以上のプロセッサによってアクセスされることができ任意の利用可能な媒体でありうる。コンピュータプログラムプロダクトは、コンピュータ可読媒体を含めうる。

【0070】

[0074] 例として、また限定されないが、そのようなコンピュータ可読媒体は、RAM、ROM、EEPROM、CD-ROMあるいは他の光学ディスクストレージ、磁気ディスクストレージあるいは他の磁気ストレージデバイス、フラッシュメモリ、のようなノンランジトリ媒体、あるいは、命令あるいはデータ構造の形態で所望プログラムコードを格納または搬送するために使用されることができ、また、コンピュータによってアクセスされることができ、任意の他の媒体を備えることができる。また、いずれの接続もコンピュータ可読媒体と適切に名付けられる。例えば、命令がウェブサイト、サーバ、あるいは、同軸ケーブル、光ファイバーケーブル、ツイストペア、デジタル加入者ライン(DSL)、あるいは赤外線、無線、およびマイクロ波のような無線技術を使用している他の遠隔ソース、から送信される場合には、そのときには、同軸ケーブル、光ファイバーケーブル、ツイストペア、DSL、あるいは赤外線、無線、およびマイクロ波のような無線技術は、媒体の定義に含まれる。しかしながら、コンピュータ可読記憶媒体およびデータ記憶媒体は接続、搬送波、信号、または他のランジエント媒体を含まないが、代わりに、ノンランジエント、タンジブル記憶媒体を対象としているということは理解されるべきである。ここに使用されているように、ディスク(disk)とディスク(disc)は、コンパクトディスク(compact disc)(CD)、レーザーディスク(登録商標)(laser disc)、光学ディスク(optical disc)、デジタル汎用ディスク(digital versatile disc)(DVD)、フロッピー(登録商標)ディスク(disk)およびブルーレイ(登録商標)ディスクを含んでおり、「ディスク(disks)」は、大抵、データを磁気で再生し、「ディスク(disks)」は、レーザーで光学的に再生する。上記のものの組み合わせも、コンピュータ可読媒体の範囲内に含まれるべきである。

【0071】

[0075] 命令は、1以上のDSP、汎用マイクロプロセッサ、ASIC、FPGA、または他の同等な集積またはディスクリートな論理回路のような1以上のプロセッサによって実行されうる。したがって、ここで使用される用語「プロセッサ」は、前述の構造のうちの一つかまたはここで説明される技法の実施に適切な任意の他の構造を指す。さらに、いくつかの態様では、ここで説明される機能は、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および/またはソフトウェアのモジュール内で提供されうる、または、組み合わせられたコーデックに組み込まれうる。また、本技法は、1つまたは複数の回路または論理構成要素において十分に実装されることができ。

【0072】

[0076] 本開示の技法は、無線ハンドセット、集積回路(IC)または1セットのIC(例えばチップセット)を含む種々さまざまなデバイスまたは装置で実装されうる。様々なコンポーネント、モジュールまたは装置は、開示された技法を実行するように構成されたデバイスの機能的態様を強調するために本開示で説明されており、異なるハードウェア装置による実現を必ずしも必要としていない。むしろ、上述されているように、様々な装置

10

20

30

40

50

は、コーデックハードウェア装置で組み合わせられ、または、適切なソフトウェアおよび/またはファームウェアと併せて上述されるような1以上のプロセッサを含むインタオペラティブハードウェア装置の集まりによって与えられる。

【0073】

[0077] 様々な例が説明されている。これらおよび他の例は、特許請求の範囲内にある。
以下に本件出願当初の特許請求の範囲を付記する。

[C 1]

汎用グラフィック処理装置 (G P G P U) であって、
処理パイプラインのステージとして選択的に動作するように構成された2以上の並行処理装置と、

10

前記並行処理装置間の移送のためにデータを保持するように構成された1以上のローカルメモリバッファであって、前記バッファの各々が前記並行処理装置のうち少なくとも2つの間で接続される、1以上のローカルメモリバッファと、
を備える G P G P U。

[C 2]

前記1以上のローカルメモリバッファの各々は、前記処理パイプラインにおいて前記並行処理装置のうち前記少なくとも2つの間で直接接続される、[C 1] に記載の G P G P U。

[C 3]

前記1以上のローカルメモリバッファは、前記並行処理装置間の前記データの移送を可能にするハードウェアベースのデータフロー制御メカニズムを含む、[C 1] に記載の G P G P U。

20

[C 4]

前記1以上のローカルメモリバッファは、ハードウェアベースの先入れ先出しバッファ (F I F O)、後入れ先出しバッファ (L I F O S) またはインデックス付けされたバッファのうちの少なくとも1つを備える、[C 1] に記載の G P G P U。

[C 5]

前記1以上のローカルメモリバッファにデータを送信するように前記並行処理装置のうち1以上を構成し、前記1以上のローカルメモリバッファからデータを受信するように前記並行処理装置のうち1以上を構成するように構成された制御装置、
をさらに備える [C 1] に記載の G P G P U。

30

[C 6]

前記制御ユニットは、前記ローカルメモリバッファにデータを送信し前記ローカルメモリバッファからデータを受信するように前記並行処理装置を構成するために1以上のアプリケーションプログラミングインタフェース (A P I) を実行するように構成される、[C 5] に記載の G P G P U。

[C 7]

前記処理パイプラインにおいて前の処理装置からのデータ出力を保持するために前記ローカルメモリバッファの各々についての必要な幅を決定するように、前記決定された幅を有するように前記ローカルメモリバッファの各々を構成するように、構成された制御装置

40

をさらに備える [C 1] に記載の G P G P U。

[C 8]

前記制御装置は、1以上のアプリケーションプログラミングインタフェース (A P I) を実行して、前記ローカルメモリバッファの各々についての前記幅を決定し、前記決定された幅によって前記ローカルメモリバッファの各々を構成し、前記ローカルメモリバッファの各々の深さを決定するように構成される、[C 7] に記載の G P G P U。

[C 9]

前記ローカルメモリバッファの各々は深さを幅とトレードすることが構成可能である、[C 7] に記載の G P G P U。

50

[C 1 0]

前記処理パイプライン内でデータシーケンスをプリザーブする制御装置、をさらに備える [C 1] に記載の G P G P U。

[C 1 1]

前記制御装置は、

前記データセットのデータスレッドのシーケンスを記録するために前記並行処理装置のうちの少なくとも1つにデータセットが入るときにシーケンス決定カウンタを実行するように、

前記シーケンス決定カウンタによって記録されるのと同じシーケンスで前記並行処理装置から前記データセットの前記データスレッドをリリースするために前記並行処理装置のうちの前記少なくとも1つから前記データセットが出るときにシーケンスエンフォーシングバリアを実行するように、

構成される、 [C 1 0] に記載の G P G P U。

[C 1 2]

前記並行処理装置のうちの1つは、前記処理パイプラインの第1ステージとして動作し、デバイスメモリからオリジナルデータセットを取り出すように構成される、 [C 1] に記載の G P G P U。

[C 1 3]

前記並行処理装置のうちの1つは、前記処理パイプラインの最終ステージとして動作し、デバイスメモリにパイプライン処理されたデータセットを格納するように構成される、 [C 1] に記載の G P G P U。

[C 1 4]

前記並行処理装置のうちの少なくとも1つは、前記処理パイプラインの中間ステージとして動作し、前記ローカルメモリバッファのうち一方を介して前記処理パイプラインにおいて前記並行処理装置のうちの前のものからデータセットを受信し、前記ローカルメモリバッファのうち他方を介して前記処理パイプラインにおいて前記並行処理装置のうち後続のものに前記データセットを送信するように構成される、 [C 1] に記載の G P G P U。

[C 1 5]

前記並行処理装置のうちの前記少なくとも1つは、前記データセットを処理するために、デバイスメモリから補助データを取り出すように構成される、 [C 1 4] に記載の G P G P U。

[C 1 6]

汎用グラフィック処理装置 (G P G P U) によってデータを処理する方法であって、前記方法は、

処理パイプラインのステージとして選択的に動作するように2以上の並行処理装置を構成することと；

前記並行処理装置間の移送のためのデータを保持するように1以上のローカルメモリバッファを構成することと、なお、前記バッファの各々は、前記並行処理装置のうち少なくとも2つの間で接続される；

を備える、方法。

[C 1 7]

前記ローカルメモリバッファの各々は、前記処理パイプラインにおいて前記並行処理装置のうち前記少なくとも2つの間で直接接続される、 [C 1 6] に記載の方法。

[C 1 8]

前記1以上のローカルメモリバッファは、前記並行処理装置間の前記データの移送を可能にするハードウェアベースのデータフロー制御メカニズムを含む、 [C 1 6] に記載の方法。

[C 1 9]

前記1以上のローカルメモリバッファは、ハードウェアベースの先入れ先出しバッファ (F I F O) 、 後入れ先出しバッファ (L I F O S) またはインデクス付けされたバッフ

10

20

30

40

50

アのうちの少なくとも1つを備える、[C 1 6]に記載の方法。

[C 2 0]

前記1以上のローカルメモリバッファにデータを送信するように前記並行処理装置の1以上を構成することと、

前記1以上のローカルメモリバッファからデータを受信するように前記並行処理装置の1以上を構成することと、

をさらに備える [C 1 6]に記載の方法。

[C 2 1]

前記1以上の並行処理装置を構成することは、前記ローカルメモリバッファにデータを送信し前記ローカルメモリバッファからデータを受信するように前記並行処理装置を構成するために1以上のアプリケーションプログラミングインタフェース (A P I) を実行することを備える、[C 2 0]に記載の方法。

10

[C 2 2]

前記1以上のローカルメモリバッファを構成することは、

前記前記処理パイプラインにおいて前の処理装置からのデータ出力を保持するために前記ローカルメモリバッファの各々についての必要な幅を決定することと、

前記決定された幅を有するように前記ローカルメモリバッファの各々を構成することと

、

を備える、[C 1 6]に記載の方法。

[C 2 3]

前記1以上のローカルメモリバッファを構成することは、1以上のアプリケーションプログラミングインタフェース (A P I) を実行して、前記ローカルメモリバッファの各々についての前記幅を決定し、前記決定された幅によって前記ローカルメモリバッファの各々を構成し、前記ローカルメモリバッファの各々の深さを決定することを備える、[C 2 2]に記載の方法。

20

[C 2 4]

前記ローカルメモリバッファの各々は深さを幅とトレードすることが構成可能である、[C 2 2]に記載の方法。

[C 2 5]

前記処理パイプライン内でデータシーケンスをプリザーブすること、

をさらに備える [C 1 6]に記載の方法。

30

[C 2 6]

前記データシーケンスをプリザーブすることは、

前記データセットのデータスレッドのシーケンスを記録するために前記並行処理装置のうちの少なくとも1つにデータセットが入るときにシーケンス決定カウンタを実行することと、

前記シーケンス決定カウンタによって記録されるのと同じシーケンスで前記並行処理装置から前記データセットの前記データスレッドをリリースするために前記並行処理装置のうちの前記少なくとも1つから前記データセットが出るときにシーケンスエンフォーシングバリアを実行することと、

40

をさらに備える、[C 2 5]に記載の方法。

[C 2 7]

前記2以上の並行処理装置を構成することは、前記並行処理装置のうちの1つを、前記処理パイプラインの第1ステージとして動作し、デバイスメモリからオリジナルデータセットを取り出すように構成することを備える、[C 1 6]に記載の方法。

[C 2 8]

前記2以上の並行処理装置を構成することは、前記並行処理装置のうちの1つを、前記処理パイプラインの最終ステージとして動作し、デバイスメモリにパイプライン処理されたデータセットを格納するように構成することを備える、[C 1 6]に記載の方法。

[C 2 9]

50

前記 2 以上の並行処理装置を構成することは、前記並行処理装置のうちの少なくとも 1 つを、前記処理パイプラインの中間ステージとして動作し、前記ローカルメモリバッファのうち一方を介して前記処理パイプラインにおいて前記並行処理装置のうちの前のものからデータセットを受信し、前記ローカルメモリバッファのうち他方を介して前記処理パイプラインにおいて前記並行処理装置のうち後続のものに前記データセットを送信するように構成することを備える、[C 1 6] に記載の方法。

[C 3 0]

前記並行処理装置のうちの少なくとも 1 つを構成することは、前記データセットを処理するために、デバイスメモリから補助データを取り出すように前記並行処理装置のうちの前記少なくとも 1 つを構成することを備える、[C 2 9] に記載の方法。

10

[C 3 1]

汎用グラフィック処理装置 (G P G P U) であって、
処理パイプラインのステージとして選択的に動作するように 2 以上の並行処理装置を構成するための手段と；

前記並行処理装置間の移送のためのデータを保持するように 1 以上のローカルメモリバッファを構成するための手段と、なお、前記バッファの各々は、前記並行処理装置のうち少なくとも 2 つの間で接続される；

を備える G P G P U 。

[C 3 2]

前記ローカルメモリバッファの各々は、前記処理パイプラインにおいて前記並行処理装置のうち前記少なくとも 2 つの間で直接接続される、[C 3 1] に記載の G P G P U 。

20

[C 3 3]

前記 1 以上のローカルメモリバッファは、前記並行処理装置間の前記データの移送を可能にするハードウェアベースのデータフロー制御メカニズムを含む、[C 3 1] に記載の G P G P U 。

[C 3 4]

前記 1 以上のローカルメモリバッファにデータを送信するように前記並行処理装置の 1 以上を構成するための手段と、

前記 1 以上のローカルメモリバッファからデータを受信するように前記並行処理装置の 1 以上を構成するための手段と、

をさらに備える [C 3 1] に記載の G P G P U 。

30

[C 3 5]

前記ローカルメモリバッファにデータを送信し前記ローカルメモリバッファからデータを受信するように前記並行処理装置を構成するために 1 以上のアプリケーションプログラミングインタフェース (A P I) を実行するための手段、をさらに備える [C 3 4] に記載の G P G P U 。

[C 3 6]

前記処理パイプラインにおいて前の処理装置からのデータ出力を保持するために前記ローカルメモリバッファの各々についての必要な幅を決定するための手段と、

前記決定された幅を有するように前記ローカルメモリバッファの各々を構成するための手段と、

をさらに備える [C 3 1] に記載の G P G P U 。

40

[C 3 7]

1 以上のアプリケーションプログラミングインタフェース (A P I) を実行して、前記ローカルメモリバッファの各々についての前記幅を決定し、前記決定された幅によって前記ローカルメモリバッファの各々を構成し、前記ローカルメモリバッファの各々の深さを決定すること、をさらに備える [C 3 6] に記載の G P G P U 。

[C 3 8]

前記ローカルメモリバッファの各々は深さを幅とトレードすることが構成可能である、[C 3 6] に記載の G P G P U 。

50

[C 3 9]

前記処理パイプライン内でデータシーケンスをプリザーブするための手段、をさらに備える [C 3 1] に記載の G P G P U。

[C 4 0]

前記データセットのデータスレッドのシーケンスを記録するために前記並行処理装置のうちの少なくとも1つにデータセットが入るときにシーケンス決定カウンタを実行するための手段と、

前記シーケンス決定カウンタによって記録されるのと同じシーケンスで前記並行処理装置から前記データセットの前記データスレッドをリリースするために前記並行処理装置のうちの前記少なくとも1つから前記データセットが出るときにシーケンスエンフォースングバリアを実行するための手段と、

をさらに備える [C 3 9] に記載の G P G P U。

[C 4 1]

汎用グラフィクス処理装置 (G P G P U) によってデータを処理するための命令を備えるコンピュータ可読媒体であって、実行時に、プログラム可能なプロセッサに、
処理パイプラインのステージとして選択的に動作するように2以上の並行処理装置を構成させる；

前記並行処理装置間の移送のためにデータを保持するように1以上のローカルメモリバッファを構成させる、なお、前記バッファの各々は、前記並行処理装置のうち少なくとも2つの間で接続される；

コンピュータ可読媒体。

[C 4 2]

前記ローカルメモリバッファの各々は、前記処理パイプラインにおいて前記並行処理装置のうち前記少なくとも2つの間で直接接続される、 [C 4 1] に記載のコンピュータ可読媒体。

[C 4 3]

前記1以上のローカルメモリバッファは、前記並行処理装置間の前記データの移送を可能にするハードウェアベースのデータフロー制御メカニズムを含む、 [C 4 1] に記載のコンピュータ可読媒体。

[C 4 4]

前記プログラム可能なプロセッサに、
前記1以上のローカルメモリバッファにデータを送信するように前記並行処理装置の1以上を構成させる命令と、

前記1以上のローカルメモリバッファからデータを受信するように前記並行処理装置の1以上を構成させる命令と、

をさらに備える [C 4 1] に記載のコンピュータ可読媒体。

[C 4 5]

前記プログラム可能なプロセッサに、前記ローカルメモリバッファにデータを送信し前記ローカルメモリバッファからデータを受信するように前記並行処理装置を構成するために1以上のアプリケーションプログラミングインタフェース (A P I) を実行させる命令、をさらに備える [C 4 4] に記載のコンピュータ可読媒体。

[C 4 6]

前記プログラム可能なプロセッサに、
前記処理パイプラインにおいて前の処理装置からのデータ出力を保持するために前記ローカルメモリバッファの各々についての必要な幅を決定させる命令と、

前記決定された幅を有するように前記ローカルメモリバッファの各々を構成させる命令と、

をさらに備える [C 4 1] に記載のコンピュータ可読媒体。

[C 4 7]

前記プログラム可能なプロセッサに、1以上のアプリケーションプログラミングインタ

10

20

30

40

50

フェース (API) を実行して、前記ローカルメモリバッファの各々についての前記幅を決定し、前記決定された幅によって前記ローカルメモリバッファの各々を構成し、前記ローカルメモリバッファの各々の深さを決定させる命令、をさらに備える [C 4 6] に記載のコンピュータ可読媒体。

[C 4 8]

前記ローカルメモリバッファの各々は深さを幅とトレードすることが構成可能である、 [C 4 6] に記載のコンピュータ可読媒体。

[C 4 9]

前記プログラム可能なプロセッサに前記処理パイプライン内でデータシーケンスをプリザーブさせる命令、をさらに備える [C 4 1] に記載のコンピュータ可読媒体。

[C 5 0]

前記プログラム可能なプロセッサに、前記データセットのデータスレッドのシーケンスを記録するために前記並行処理装置のうちの少なくとも1つにデータセットが入るときにシーケンス決定カウンタを実行させる命令と、

前記シーケンス決定カウンタによって記録されるのと同じシーケンスで前記並行処理装置から前記データセットの前記データスレッドをリリースするために前記並行処理装置のうちの前記少なくとも1つから前記データセットが出るときにシーケンスエンフォーシングバリアを実行させる命令と、

をさらに備える [C 4 9] に記載のコンピュータ可読媒体。

10

20

【 図 1 】

図 1

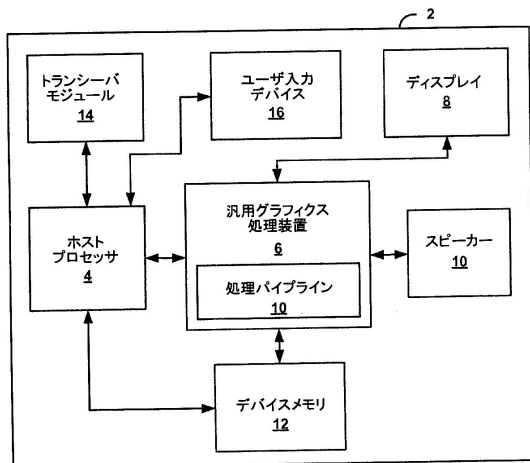


FIG. 1

【 図 2 】

図 2

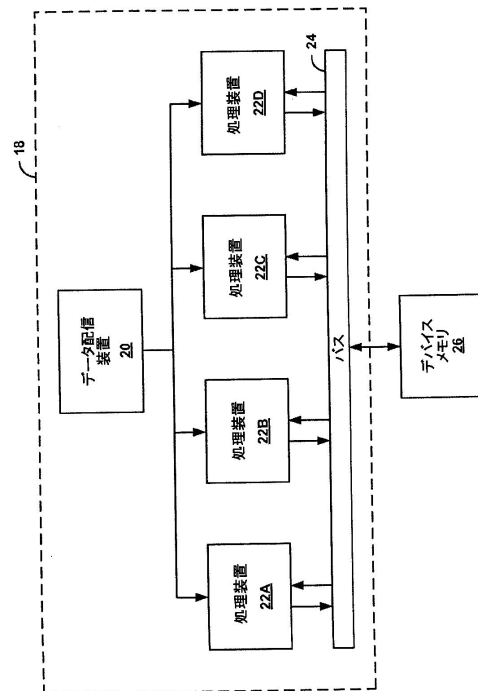


FIG. 2

【 図 3 】

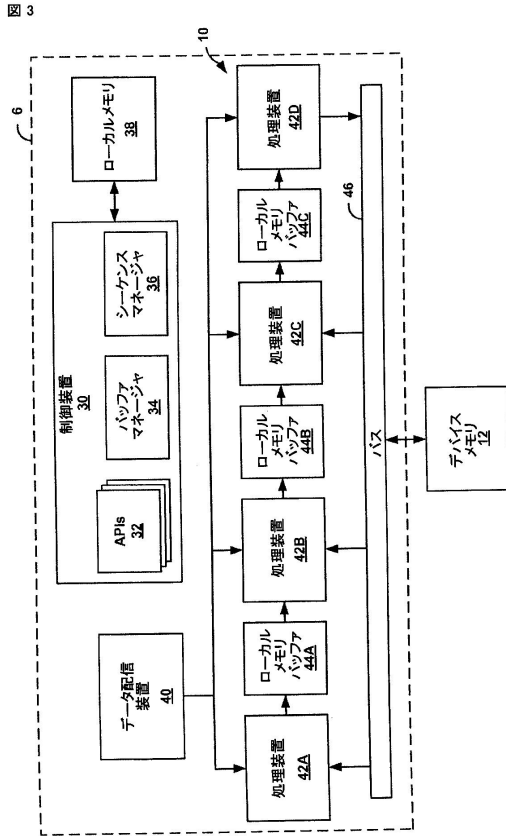


FIG. 3

【 図 4 】

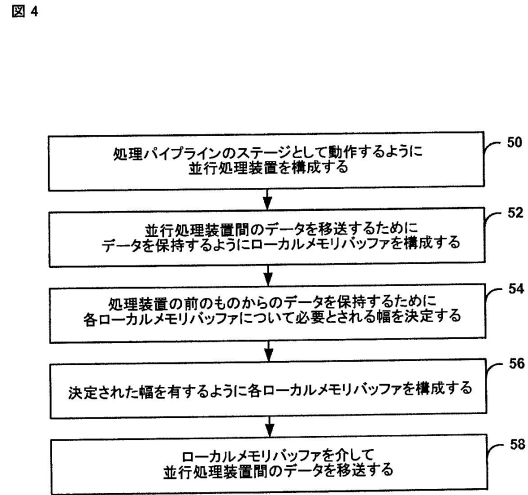


FIG. 4

【 図 5 】

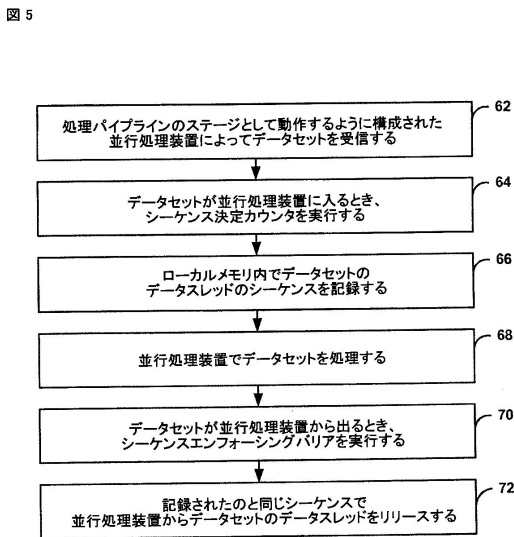


FIG. 5

フロントページの続き

- (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 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 グルバー、アンドリュー
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 クルスティク、アレクサンドラ・エル。
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 シンプソン、ロバート・ジェイ。
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 シャープ、コリン
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 ユ、チュン
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

合議体

審判長 高木 進
審判官 辻本 泰隆
審判官 須田 勝巳

- (56)参考文献 特表2004-515856(JP,A)
特開2008-9697(JP,A)
特開平6-282447(JP,A)
特開2010-287110(JP,A)
特開2010-244096(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G06F9/46-9/54