

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6200824号
(P6200824)

(45) 発行日 平成29年9月20日 (2017.9.20)

(24) 登録日 平成29年9月1日 (2017.9.1)

(51) Int. Cl.		F I	
G06F 15/80	(2006.01)	G06F 15/80	
G06F 9/50	(2006.01)	G06F 9/46	465E
G06F 12/0802	(2016.01)	G06F 12/08	501Z
G06F 12/08	(2016.01)	G06F 12/08	513

請求項の数 4 (全 40 頁)

(21) 出願番号	特願2014-23044 (P2014-23044)	(73) 特許権者	302062931 ルネサスエレクトロニクス株式会社 東京都江東区豊洲三丁目2番24号
(22) 出願日	平成26年2月10日 (2014.2.10)	(74) 代理人	100103894 弁理士 冢入 健
(65) 公開番号	特開2015-149038 (P2015-149038A)	(72) 発明者	京 昭倫 神奈川県川崎市中原区下沼部1753番地 ルネサスエレクトロニクス株式会社内
(43) 公開日	平成27年8月20日 (2015.8.20)	審査官	漆原 孝治
審査請求日	平成28年10月3日 (2016.10.3)		

最終頁に続く

(54) 【発明の名称】 演算制御装置及び演算制御方法並びにプログラム、OpenCLデバイス

(57) 【特許請求の範囲】

【請求項1】

複数の演算素子と、該複数の演算素子に対して設けられた階層の異なる複数のメモリとを有するOpenCL (Open Computing Language) デバイスの前記複数の演算素子による並列演算を制御する演算制御装置であって、

前記複数のメモリのうちの最下位階層のメモリに1つ以上格納されたデータブロックであって、前記並列演算の演算対象としてそのデータが他の階層のメモリに転送されるリードブロックと、前記並列演算後に前記他の階層のメモリから前記最下位階層のメモリに転送される1つ以上のデータブロックであって、前記1つ以上のリードブロックに対する並列演算の演算結果であるライトブロックとに対して夫々設定された属性群を取得して保持する属性群保持部と、

前記属性群保持部により保持された各前記属性群と、前記OpenCLデバイスの構成を示す構成パラメータとに基づいて、夫々の前記リードブロックと前記ライトブロックの転送方式を決定し、決定した転送方式に応じて各前記リードブロックと前記ライトブロックの転送、及び該転送に対応する前記並列演算の制御を行うシナリオ決定部とを有し、

前記属性群は、

前記転送方式を決定するために必要である一方、前記OpenCLデバイスの構成に依存しない属性を複数含むものであり、

当該データブロックを複数のサブブロックに分割して転送するか否か、及び、分割して転送する場合における分割方式を示す割当属性と、

10

20

当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックと共に転送される、該サブブロックの周辺のデータのサイズを示す余白属性と、

当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックがその周辺の他のサブブロックと依存関係が有るか否か、及び、前記依存関係がある場合における全ての依存方向を示す依存属性とを有し、

前記ライトブロックの属性群は、該ライトブロックが既に前記他の階層のメモリに存在し、かつ、前記最下位階層のメモリに転送されると仮定して設定されたものであり、

前記シナリオ決定部は、

第1のカーネルに続いて第2のカーネルを前記OpenCLデバイスに実行させ、かつ、前記第1のカーネルに対応する並列演算の各前記ライトブロックのうちに、前記第2のカーネルに対応する並列演算の前記リードブロックとして使用される継続ライトブロックが含まれており、かつ、前記第1のカーネルに対して設定された前記継続ライトブロックの前記割当属性と、前記第2のカーネルに対して設定された、前記継続ライトブロックに対応する前記リードブロックの前記割当属性とが一致する場合において、該継続ライトブロックが前記最下位階層へ転送されず、前記他の階層のメモリを介して前記第2のカーネルの実行に供されることにより、前記第1のカーネルと前記第2のカーネルの実行をパイプライン化するパイプライン化制御を行い、

前記パイプライン化制御に際して、

前記第1のカーネルに対応する各前記リードブロックの夫々に対して、該リードブロックに対して設定された前記余白属性と前記依存属性に、前記第2のカーネルに対して設定された、前記継続ライトブロックに対応する前記リードブロックの前記余白属性と前記依存属性を夫々論理和加算する、演算制御装置。

【請求項2】

複数の演算素子と、該複数の演算素子に対して設けられた階層の異なる複数のメモリとを有するOpenCL(Open Computing Language)デバイスの前記複数の演算素子による並列演算を制御する演算制御方法であって、

前記複数のメモリのうちの最下位階層のメモリに1つ以上格納されたデータブロックであって、前記並列演算の演算対象としてそのデータが他の階層のメモリに転送されるリードブロックと、前記並列演算後に前記他の階層のメモリから前記最下位階層のメモリに転送される1つ以上のデータブロックであって、前記1つ以上のリードブロックに対する並列演算の演算結果であるライトブロックとに対して夫々設定された属性群を取得して保持する第1のステップと、

保持された各前記属性群と、前記OpenCLデバイスの構成を示す構成パラメータとに基づいて、夫々の前記リードブロックと前記ライトブロックの転送方式を決定し、決定した転送方式に応じて各前記リードブロックと前記ライトブロックの転送、及び該転送に対応する前記並列演算の制御を行う第2のステップとを有し、

前記属性群は、

前記転送方式を決定するために必要である一方、前記OpenCLデバイスの構成に依存しない属性を複数含むものであり、

当該データブロックを複数のサブブロックに分割して転送するか否か、及び、分割して転送する場合における分割方式を示す割当属性と、

当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックと共に転送される、該サブブロックの周辺のデータのサイズを示す余白属性と、

当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックがその周辺の他のサブブロックと依存関係が有るか否か、及び、前記依存関係がある場合における全ての依存方向を示す依存属性とを有し、

前記ライトブロックの属性群は、該ライトブロックが既に前記他の階層のメモリに存在し、かつ、前記最下位階層のメモリに転送されると仮定して設定されたものであり、

前記第2のステップにおいて、

第1のカーネルに続いて第2のカーネルを前記OpenCLデバイスに実行させ、かつ、前記第1のカーネルに対応する並列演算の各前記ライトブロックのうちに、前記第2のカーネルに対応する並列演算の前記リードブロックとして使用される継続ライトブロックが含まれており、かつ、前記第1のカーネルに対して設定された前記継続ライトブロックの前記割当属性と、前記第2のカーネルに対して設定された、前記継続ライトブロックに対応する前記リードブロックの前記割当属性とが一致する場合において、該継続ライトブロックが前記最下位階層へ転送されず、前記他の階層のメモリを介して前記第2のカーネルの実行に供されることにより、前記第1のカーネルと前記第2のカーネルの実行をパイプライン化するパイプライン化制御を行い、

前記パイプライン化制御に際して、

前記第1のカーネルに対応する各前記リードブロックの夫々に対して、該リードブロックに対して設定された前記余白属性と前記依存属性に、前記第2のカーネルに対して設定された、前記継続ライトブロックに対応する前記リードブロックの前記余白属性と前記依存属性を夫々論理和加算する、
演算制御方法。

【請求項3】

複数の演算素子と、該複数の演算素子に対して設けられた階層の異なる複数のメモリとを有するOpenCL(Open Computing Language)デバイスの前記複数の演算素子による並列演算を制御する際に、

前記複数のメモリの中の最下位階層のメモリに1つ以上格納されたデータブロックであって、前記並列演算の演算対象としてそのデータが他の階層のメモリに転送されるリードブロックと、前記並列演算後に前記他の階層のメモリから前記最下位階層のメモリに転送される1つ以上のデータブロックであって、前記1つ以上のリードブロックに対する並列演算の演算結果であるライトブロックとに対して夫々設定された属性群を取得して保持する第1のステップと、

保持された各前記属性群と、前記OpenCLデバイスの構成を示す構成パラメータとに基づいて、夫々の前記リードブロックと前記ライトブロックの転送方式を決定し、決定した転送方式に応じて各前記リードブロックと前記ライトブロックの転送、及び該転送に対応する前記並列演算の制御を行う第2のステップとをコンピュータに実行させ、

前記属性群は、

前記転送方式を決定するために必要である一方、前記OpenCLデバイスの構成に依存しない属性を複数含むものであり、

当該データブロックを複数のサブブロックに分割して転送するか否か、及び、分割して転送する場合における分割方式を示す割当属性と、

当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックと共に転送される、該サブブロックの周辺のデータのサイズを示す余白属性と、

当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックがその周辺の他のサブブロックと依存関係が有るか否か、及び、前記依存関係がある場合における全ての依存方向を示す依存属性とを有し、

前記ライトブロックの属性群は、該ライトブロックが既に前記他の階層のメモリに存在し、かつ、前記最下位階層のメモリに転送されると仮定して設定されたものであり、

前記第2のステップにおいて、

第1のカーネルに続いて第2のカーネルを前記OpenCLデバイスに実行させ、かつ、前記第1のカーネルに対応する並列演算の各前記ライトブロックのうちに、前記第2のカーネルに対応する並列演算の前記リードブロックとして使用される継続ライトブロックが含まれており、かつ、前記第1のカーネルに対して設定された前記継続ライトブロックの前記割当属性と、前記第2のカーネルに対して設定された、前記継続ライトブロックに対応する前記リードブロックの前記割当属性とが一致する場合において、該継続ライトブロックが前記最下位階層へ転送されず、前記他の階層のメモリを介して前記第2のカーネルの実行に供されることにより、前記第1のカーネルと前記第2のカーネルの実行をパイ

10

20

30

40

50

プライン化するパイプライン化制御を行い、

前記パイプライン化制御に際して、

前記第1のカーネルに対応する各前記リードブロックの夫々に対して、該リードブロックに対して設定された前記余白属性と前記依存属性に、前記第2のカーネルに対して設定された、前記継続ライトブロックに対応する前記リードブロックの前記余白属性と前記依存属性を夫々論理和加算する、

プログラム。

【請求項4】

複数の演算素子と、

該複数の演算素子に対して設けられた階層の異なる複数のメモリと、

前記複数の演算素子による並列演算を制御する演算制御部とを備えるOpenCLデバイスであって、

前記演算制御部は、

前記複数のメモリのうちの最下位階層のメモリに1つ以上格納されたデータブロックであって、前記並列演算の演算対象としてそのデータが他の階層のメモリに転送されるリードブロックと、前記並列演算後に前記他の階層のメモリから前記最下位階層のメモリに転送される1つ以上のデータブロックであって、前記1つ以上のリードブロックに対する並列演算の演算結果であるライトブロックとに対して夫々設定された属性群を取得して保持する属性群保持部と、

前記属性群保持部により保持された各前記属性群と、前記OpenCLデバイスの構成を示す構成パラメータとに基づいて、夫々の前記リードブロックと前記ライトブロックの転送方式を決定し、決定した転送方式に応じて各前記リードブロックと前記ライトブロックの転送、及び該転送に対応する前記並列演算の制御を行うシナリオ決定部とを有し、

前記属性群は、

前記転送方式を決定するために必要である一方、前記OpenCLデバイスの構成に依存しない属性を複数含むものであり、

当該データブロックを複数のサブブロックに分割して転送するか否か、及び、分割して転送する場合における分割方式を示す割当属性と、

当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックと共に転送される、該サブブロックの周辺のデータのサイズを示す余白属性と、

当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックがその周辺の他のサブブロックと依存関係が有るか否か、及び、前記依存関係がある場合における全ての依存方向を示す依存属性とを有し、

前記ライトブロックの属性群は、該ライトブロックが既に前記他の階層のメモリに存在し、かつ、前記最下位階層のメモリに転送されると仮定して設定されたものであり、

前記シナリオ決定部は、

第1のカーネルに続いて第2のカーネルを前記OpenCLデバイスに実行させ、かつ、前記第1のカーネルに対応する並列演算の各前記ライトブロックのうちに、前記第2のカーネルに対応する並列演算の前記リードブロックとして使用される継続ライトブロックが含まれており、かつ、前記第1のカーネルに対して設定された前記継続ライトブロックの前記割当属性と、前記第2のカーネルに対して設定された、前記継続ライトブロックに対応する前記リードブロックの前記割当属性とが一致する場合において、該継続ライトブロックが前記最下位階層へ転送されず、前記他の階層のメモリを介して前記第2のカーネルの実行に供されることにより、前記第1のカーネルと前記第2のカーネルの実行をパイプライン化するパイプライン化制御を行い、

前記パイプライン化制御に際して、

前記第1のカーネルに対応する各前記リードブロックの夫々に対して、該リードブロックに対して設定された前記余白属性と前記依存属性に、前記第2のカーネルに対して設定された、前記継続ライトブロックに対応する前記リードブロックの前記余白属性と前記依存属性を夫々論理和加算する、

10

20

30

40

50

OpenCLデバイス。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は並列プロセッサ、具体的にはOpenCLデバイスの演算制御技術に関する。

【背景技術】

【0002】

近年、プロセッサの発熱を抑制する必要性から、プロセッサの動作周波数を上げる代わりに、並列に処理を行うプロセッサ・コア（以下、単に「コア」という）の数を増やすことで性能向上を実現する動きが顕著になっている。複数のコアを有するプロセッサは、マルチコアプロセッサと呼ばれており、マルチコアプロセッサの中の、特にコア数が多いものは、メニコアプロセッサと呼ばれる。本明細書の中で、マルチコアプロセッサとメニコアプロセッサの区別を特にせず、並列に処理を行うコアが複数含まれるプロセッサを概して「並列プロセッサ」という。

10

【0003】

並列プロセッサは、アクセラレータとして様々な分野で用いられている。しかし、メーカ、分野などによって種々のアクセラレータが製造され、アクセラレータ用の言語やフレームワークも様々開発されているため、アクセラレータ間でプログラムコードの移植が困難である。

【0004】

20

この問題を解決するため、並列プロセッサのための標準的なフレームワークとして、OpenCL (Open Computing Language) が定められた (非特許文献1)。ここで、OpenCLの概要を説明する。

【0005】

図17は、非特許文献1におけるFig 3.1に対して符号を追加したものであり、典型的なOpenCLシステムのプラットフォームモデルを示す。

【0006】

図17に示すように、OpenCLシステム10は、Host (ホスト) 12、1つ以上のCompute Device (以下「OpenCLデバイス」という) 14を備える。OpenCLデバイス14は、上述したアクセラレータに該当する。

30

【0007】

夫々のOpenCLデバイス14は、1つ以上のCompute Unit (以下「CU」と略す) 16を有し、夫々のCU16は、1つ以上の演算素子 (Processing Element。以下「PE」と略す) 18を有する。なお、演算素子PE18は、前述したコアに該当する。

【0008】

OpenCLのアプリケーションは、Host 12側で動作するプログラムコードと、OpenCLデバイス14、すなわちアクセラレータ側で動作するプログラムコードとで構成される。Host 12側で動作するプログラムコードは「ホストコード」と呼ばれ、OpenCLデバイス14側で動作するプログラムコードは「カーネル」と呼ばれる。

40

【0009】

Host 12は、API (Application Program Interface) をコールして演算を指示する。OpenCLデバイス14は、指示された演算を実行する。Host 12は、リソースを管理するコンテキストの生成、さらに、OpenCLを通じてデバイス動作を調停するためのコマンドキューの生成を行う。「デバイス動作」は、演算を実行すること、メモリを操作すること、同期をとることなどが含まれる。

【0010】

OpenCLにおいて、カーネルは、work-item (以下、略して「アイテム」とも呼ぶ) として、N次元 (1 ≤ N ≤ 3) のインデックス空間 (Index Space) で実行される。例えば、2次元のインデックス空間として (4, 6) が指定されれば、

50

「4 × 6」の計 24 個のアイテムが実行される。

【0011】

1 アイテムを実行するには、1 P E が利用される。したがって、並列に実行されるアイテムの数と、実在する P E の数とが同一の場合、カーネルは、4 列 6 行計 24 個の P E 上で実行されることになる。

【0012】

なお、並列に実行されるアイテム数に対して実際に存在する P E 数が少ない場合には、実際に存在する P E 上でアイテムの並列実行を繰り返すことが行われる。例えば、2 列 3 行計 6 個の P E しかないときに、上述した (4, 6) のインデックス空間が指定された場合には、該 6 個の P E により、6 個のアイテムの並列実行を 4 回繰り返す必要がある。

10

【0013】

また、OpenCL では、work - group (ワークグループ) という概念も導入されている。ワークグループは、同一の C U 16 上で実行され、かつ、互いに関連するアイテムの集合である。同一のワークグループ内の各アイテムは、同一のカーネルを実行し、該 C U 16 の後述するローカルメモリをシェアする。

【0014】

各ワークグループは、ユニックなグループ ID が割り当てられ、各ワークグループ内のアイテムは、該ワークグループ内においてユニックなローカル ID が割り当てられる。また、アイテムは、ユニックなグローバル ID も割り当てられる。アイテムは、グローバル ID、または、グループ ID とローカル ID の組合せによって識別できる。

20

【0015】

OpenCL デバイス 14 に演算処理を行わせるプロセスは、下記のステップ順で A P I をコールすることにより構成される。

【0016】

<ステップ 1> : 演算処理の参照用のデータ (以下「参照データ」という) やカーネルをホスト 12 から OpenCL デバイス 14 へ転送する。

【0017】

<ステップ 2> : 「カーネル起動コマンド」により OpenCL デバイス 14 上でカーネルの実行を開始させる。

【0018】

<ステップ 3> : OpenCL デバイス 14 におけるカーネルの実行完了後、OpenCL デバイス 14 のメモリ空間から演算処理の結果データをホスト 12 側に転送する。

30

【0019】

図 18 を参照して、メモリ空間を含む OpenCL デバイス 14 の構成を説明する。

図 18 は、非特許文献 1 における Fig 3 . 3 に対して符号を追加したものである。前述したように、OpenCL デバイス 14 は 1 つ以上の C U 16 を備え、夫々の C U 16 は 1 つ以上の P E 18 を有する。

【0020】

上述したステップ 2 におけるカーネルの実行に際して、OpenCL デバイス 14 では、4 つの異なるメモリへのアクセスが生じ得る。この 4 つのメモリは、プライベートメモリ 20、ローカルメモリ 22、グローバルメモリ 32、コンスタントメモリ 34 である。この 4 つのメモリについて、まず、図 19 を参照して、アイテムとワークグループの視点から説明する。なお、図 19 は、非特許文献 1 における Table 3 . 1 である。

40

【0021】

プライベートメモリ 20 は、1 つのアイテムに対応し、該アイテムの実行にのみ用いられる。1 つのアイテムに対応するプライベートメモリ 20 に対して定義された変数は、他のアイテムに対して使用できない。

【0022】

ローカルメモリ 22 は、1 つのグループに対応し、該グループ内の各アイテムによりシェアすることができる。そのため、ローカルメモリ 22 の用途としては、例えば、該グル

50

ープ内の各アイテムによりシェアされる変数をローカルメモリ 22 に割り当てることが挙げられる。

【0023】

グローバルメモリ 32 とコンスタントメモリ 34 は、全てのグループの全てのアイテムからアクセスできる。なお、グローバルメモリ 32 は、アイテムからリードとライトのいずれからもアクセスできるが、コンスタントメモリ 34 は、アイテムからリードアクセスのみができる。以下、グローバルメモリ 32 とコンスタントメモリ 34 をまとめてデバイスメモリ 30 という。

【0024】

アイテムと PE 18 の 1 対 1 の対応関係から、上記 4 つのメモリと、CU 16 及び PE 18 との対応関係は、以下のようになる。

【0025】

プライベートメモリ 20 は、PE 18 と 1 対 1 で対応し、相対応する PE 18 によりのみアクセス可能である。

【0026】

ローカルメモリ 22 は、CU 16 と 1 対 1 で対応し、相対応する CU 16 内の全ての PE 18 によりアクセス可能である。

【0027】

デバイスメモリ 30 は、全ての CU 16 内の全ての PE 18、すなわち OpenCL デバイス 14 内の全ての PE によりアクセス可能である。

【0028】

また、OpenCL デバイス 14 によっては、デバイスメモリ 30 のキャッシュメモリとして機能するキャッシュ 24 がさらに設けられる場合もある。

【0029】

このように、OpenCL デバイス 14 には、階層の異なる複数のメモリが設けられている。これらのメモリは、上位階層にあるほど、PE からのアクセスが高速にできる。デバイスメモリ 30 (最下位)、ローカルメモリ 22 (中位)、プライベートメモリ 20 (最上位) の順で階層が上位になり、同様の順で PE からのアクセスもより高速になる。

【0030】

OpenCL デバイス 14 の性能を十分に引き出すためには、例えば、利用頻度の高いデータを、なるべくより高速なメモリ空間に移動してから参照するようにするなど、デバイスメモリ 30 と、プライベートメモリ 20 / ローカルメモリ 22 との間でのデータ移動について工夫する必要がある。

【0031】

制御方式が OpenCL デバイスと異なる逐次プロセッサの場合においても、グローバルメモリ空間とプライベートメモリ空間との間でのデータ移動が行われる。図 20 に示す逐次プロセッサの例を参照して説明する。

【0032】

図 20 に示す逐次プロセッサ 20 は、演算素子である PE 52、プライベートメモリ 54、グローバルメモリ 56、キャッシュ制御機構 58 を有する。

【0033】

図示のように、逐次プロセッサ 50 の記憶装置は、プライベートメモリ 54 とグローバルメモリ 56 に分けられている。プライベートメモリ 54 は、物理的にオンチップの小容量メモリであり、グローバルメモリ 56 は、物理的にオフチップの大容量のメモリである。

【0034】

逐次プロセッサ 50 では、記憶装置がプライベートメモリ 54 とグローバルメモリ 56 に分けられているが、プライベートメモリ 54 とグローバルメモリ 56 の間に設けられたキャッシュ制御機構 58 により、プライベートメモリ 54 とグローバルメモリ 56 間のデータ移動は自動的に行われ、逐次プロセッサ 50 のユーザは、1 つの大きなメモリ空間し

10

20

30

40

50

が見えない。つまり、逐次プロセッサ50のユーザは、グローバルメモリ56とプライベートメモリ54との間でどのようにデータを移動するかを意図しなくても、PE52に演算処理を行わせるユーザプログラムを容易に開発できる。

【0035】

ところで、並列プロセッサ、特に図18に示すOpenCLデバイス14のように多数のコア(PE)を搭載した場合には、コアの数と同数のプライベートメモリ20が存在し、さらに、CU16の数と同数のローカルメモリ22が存在する。これらのメモリを全て1つのキャッシュ制御機構で統一的に管理するのは、ハードウェアのコストが高く、一般的には実現困難である。

【0036】

他方、キャッシュ制御機構が無いと、OpenCLシステム10のユーザ(以下、単に「ユーザ」と呼ぶ)には、複数のメモリ空間が見えてしまう。前述したように、利用頻度の高いデータを、なるべくより高速なメモリ空間(すなわちより上位階層のメモリ空間)に移動してから参照するようにするなど、より良い性能を追求するためには、演算処理に伴う階層の異なるメモリ間でのデータ移動について、ユーザプログラムで明示的に指示する必要がある。これを正しく実現するためには、ユーザには、上述した各メモリ同士の速度差、容量差、機能差等に関する知識を持つ必要がある。図21を参照して具体例を説明する。

【0037】

図21は、複数のデータブロック(データブロックA~D)から、データブロックA'~B'を得る演算処理を実行する場合を説明するための図である。なお、図21において、ホストからデバイスへのカーネル転送の図示を省略している。また、データブロックA~Dは、上述したステップ1でホスト12からOpenCLデバイス14に転送された参照データであり、グローバルメモリ32に格納される。データブロックA'~B'は、上述したステップ2においてデータブロックA~Dに対して行われた演算の結果であり、グローバルメモリ32に書き込まれ、後に、上述したステップ3でホスト12に転送される。

【0038】

ここで、ステップ2の処理、すなわちカーネルを実行する演算処理を説明する。なお、本明細書において、以下、プライベートメモリについて、複数があり得る場合には、「プライベートメモリ群」という。

【0039】

演算処理の性能を追求しなければ、演算において、プライベートメモリ群/ローカルメモリ22を使わずに、グローバルメモリ32のみを使用する手法が考えられる。この場合、グローバルメモリ32と、プライベートメモリ群/ローカルメモリ22間のデータ転送が無い。

【0040】

この手法は、制御が単純であるが、性能が良くない。演算処理をより良い性能で行うために、上述したように、演算対象のデータをグローバルメモリ32からプライベートメモリ群/ローカルメモリ22に転送してから演算を行い、演算の結果をプライベートメモリ群/ローカルメモリ22に格納してからグローバルメモリ32に転送する手法が用いられる。

【0041】

この手法を用いる場合について、まず、全てのアイテムが同時に並行して実行可能なときの手順(ステップA~C)を説明する。なお、「全てのアイテムが同時に並行して実行可能」とは、PE数が総アイテム数以上であり、かつ、プライベートメモリ群とローカルメモリの容量は、演算対象の全てのデータを格納可能できることなどを意味し、この場合、演算対象のデータをグローバルメモリ32からプライベートメモリ群/ローカルメモリ22への転送、各PE18による演算の並列実行、演算結果をプライベートメモリ群/ローカルメモリ22からグローバルメモリ32への転送が一度しか行われぬ。

【0042】

10

20

30

40

50

<ステップA> : グローバルメモリ32に格納されたデータブロックA~Dをプライベートメモリ群/ローカルメモリ22に転送する。

【0043】

この転送は、例えば、演算対象のデータのうちの、PE18によりのみ使用されるデータをPE18のプライベートメモリに転送し、複数のPE18により共有されるデータをローカルメモリ22に転送することである。

【0044】

なお、以下において、グローバルメモリ32からプライベートメモリ群/ローカルメモリ22へのデータ転送を「リード転送」という。また、データブロックA~Dのような、リード転送されるデータブロックを「リードブロックRB」という。

【0045】

<ステップB> : 各PE18上で演算処理を実行し、演算処理の結果を該PE18がアクセスできるプライベートメモリ/ローカルメモリ22に格納する。

【0046】

<ステップC> : ステップBの演算処理により得られ、プライベートメモリ群/ローカルメモリ22に格納されたデータブロックA'~B'をグローバルメモリ32に転送する。

【0047】

なお、以下において、プライベートメモリ群/ローカルメモリ22からグローバルメモリ32へのデータ転送を「ライト転送」という。また、データブロックA'~B'のような、プライベートメモリ群/ローカルメモリ22に格納され、ライト転送されるデータブロックを「ライトブロックWB」という。

【0048】

該3つのステップの全てについて、ユーザにより作成されたカーネルの中で明示的に指定される必要がある。これらの指定は、演算処理の内容や、OpenCLデバイス14の構成(PE数(=プライベートメモリ数)や、個々のプライベートメモリの容量、ローカルメモリの容量など)に依存する内容が含まれる。

【0049】

例えば、演算対象となるリードブロックRBが複数あり、かつ全てが一度に一つのワークグループ内のプライベートメモリ群/ローカルメモリ22に入りきらないため、それぞれのリードブロックRBをサブブロックに分割しなければならない場合に、該複数のリードブロックRBに対して、ステップAにおいて、サブブロック間の対応付け方法を指定する必要がある。リードブロックRBのサブブロック間の「対応付け方法」は、上記複数のリードブロックRBのサブブロック同士で、どのリードブロックRBのサブブロック同士を、同一のワークグループ内のプライベートメモリ群、または同一ワークグループ内のローカルメモリ22に転送するかを意味する。これは、演算処理の内容、またどのように分割すべきかは、OpenCLデバイス14の構成に依存する。

【0050】

同様に、演算結果として複数のライトブロックWBがある場合には、該複数のライトブロックWBのそれぞれのサブブロックが、どのようなリードブロックRBのサブブロック同士の組合せの下で、演算結果として求められることになるか、という意味での対応付け方法も指定する必要がある。なお、ライトブロックWBの各サブブロックの内容とは、すなわち各ワークグループのプライベートメモリ群またはローカルメモリ22に演算結果として格納されたデータである。そしてライトブロックWBをグローバルメモリ32へ転送するとは、該データをグローバルメモリ32内のライトブロックWBの各サブブロック位置に書き込むことを意味する。リードブロックRBの対応付け方法と同様に、ライトブロックWBの対応付け方法も、演算処理の内容と、OpenCLデバイス14の構成に依存する。

【0051】

上記のように所要のデータブロック全体がワークグループ内のメモリに入りきらない場合以外、PEの総数が、インデックス空間のサイズより小さいなどのときにおいても、全

10

20

30

40

50

てのアイテムが同時に並行して実行することができないため、PEによるアイテムの並列実行を複数回繰り返す必要がある。当然ながら、並列実行の繰り返しに合わせてリード転送とライト転送も繰り返す必要がある。この場合、演算処理の内容とOpenCLデバイス14の構成に応じて、データブロックの分割方法、データブロックを分割して得たサブブロック間の対応付け方法を指定する必要がある。

【0052】

データブロックの「分割方法」は、該データブロックをどのようにしてサブブロックに分割するかを意味する。「サブブロックSB」は、リード転送とライト転送の転送単位である。以下において、リードとライトを区別する必要がある場合については、リードブロックRBを分割して得たサブブロックを「サブリードブロックSRB」といい、ライトブロックWBを分割して得たサブブロック「サブライトブロックSWB」という。

10

【0053】

サブブロックSB間の「対応付け方法」とは、異なるリードブロックあるいはライトブロックに夫々含まれるどのサブブロックSB同士を、同時に同一のプライベートメモリ群、または同一のローカルメモリ22上に存在させるかを意味する。

【0054】

データブロックの分割方法はOpenCLデバイス14の構成に依存し、一方、サブブロックの対応付け方法は演算処理の内容に依存する。分割が必要になると、データブロックを分割しない場合より、さらに指定が複雑である。

【0055】

図22は、OpenCLデバイス14に演算処理を行わせるために、ユーザにより指定する必要のある内容をまとめて示す。

20

【0056】

図示のように、第1の部分は、リード転送用の指定であり、演算処理の内容とOpenCLデバイス14の構成に依存する部分が含まれる。

【0057】

演算処理の内容とOpenCLデバイス14の構成に依存する部分は、例えば、リードブロックRBを分割するか否かの指定(例1)や、分割する場合における分割方法の指定(例2)とサブリードブロックSRB間の対応付け方法の指定(例3)がある。

【0058】

第2の部分は、リードブロックRBまたはサブリードブロックSRBに対する演算処理の指定である。演算処理を指定するための内容であるため、当然演算処理の内容に依存する。さらに、この部分は、リード転送用の指示に合わせる必要があるため、アイテムの並行実行の回数の指定(例4)など、OpenCLデバイス14の構成に依存する内容が含まれる。

30

【0059】

第3の部分は、ライト転送用の指定であり、リード転送用の指示に合わせる必要があるため、必然的に、演算処理の内容とOpenCLデバイス14の構成に依存する部分(例5)が含まれる。

【0060】

このように、より良い性能を追求するために、ユーザは、演算処理の内容と、OpenCLデバイス14の構成に合わせてカーネル(ユーザコード)を開発する必要がある。

40

【0061】

しかし、OpenCLに準拠したデバイスであっても、メーカーが異なれば、各メモリ空間の容量、アクセス速度、アクセス遅延、キャッシュ制御の有無などの点において、千差万別である。そのため、あるOpenCLデバイスにとって、階層の異なるメモリ間のデータ移動が理想的に開発されたユーザコードであっても、他のOpenCLデバイスや、同一シリーズであるものの世代が違うOpenCLデバイスなどにとっては、反って性能劣化を引き起こすものになってしまう可能性がある。つまり、ユーザコードの性能上の可搬性が低い。

50

【 0 0 6 2 】

ある程度の性能可搬性を実現するために、特定の OpenCL デバイスを対象にユーザコードを開発するのではなく、なるべく、既存の多くの種類の OpenCL デバイスの構成を念頭に入れつつ、ユーザコードを作成することが考えられる。しかし、そうした作業は、演算処理の設計者からみれば本質的でない作業のために多大な負担に感じる上に、コードの可読性の低下、そして煩雑度の増大等を引き起こす要因にもなる。

【 0 0 6 3 】

特許文献 1 には、OpenCL デバイスに対して、演算処理に伴う階層の異なる複数のメモリ間でのデータ移動について、ユーザコードの開発者の負担を減らすと共に、ユーザコードの可搬性を高める技術を開示している。

10

【 0 0 6 4 】

該技術では、ユーザコードの開発者は、カーネルの引数として、OpenCL デバイスの演算対象及び演算結果の各データブロック（リードブロックとライトブロック）の夫々について複数の属性を含む属性群を設定し、OpenCL デバイスの演算制御部は、ホストから転送されたカーネルの引数が示す演算対象及び演算結果の各データブロックの属性群と、OpenCL デバイスの構成を示すパラメータとに基づいて、自動的に転送方式を決定すると共に、決定した転送方式でデータの転送と、OpenCL デバイスによる演算を制御する。

【 0 0 6 5 】

なお、転送方式は、主に、デバイスメモリと、ローカルメモリ及びプライベートメモリとの間でどのようにしてデータブロックを転送するかに関するものである。

20

【 0 0 6 6 】

図 2 3 は、特許文献 1 の技術を適用した OpenCL システムにおいて、ユーザがカーネルを開発する際に指定する必要がある内容を示す。図示のように、これらの内容は、属性群の指定とユーザ処理の指定のみであり、いずれも、デバイスの構成に依存しないものである。

【 0 0 6 7 】

図 2 3 に示すように、属性群は、固有属性、演算属性、ポリシ属性のグループに分けられており、各グループの属性は、転送方式を決定するために必要である一方、OpenCL デバイスの構成に依存しないものである。例えば、ポリシ属性には、当該データブロックを複数のサブブロックに分割して転送するか否か、及び、分割して転送する場合における分割方式を示す割当属性や、当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックと共に転送される、該サブブロックの周辺のデータのサイズを示す余白属性や、当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックがその周辺の他のサブブロックと依存関係が有るか否か、及び、依存関係がある場合における全ての依存方向を示す依存属性が含まれる。

30

【 0 0 6 8 】

なお、ライトブロックの属性群は、該ライトブロックが既にローカルメモリまたはプライベートメモリに存在し、かつ、デバイスメモリに転送されると仮定して設定されたものである。

40

【 0 0 6 9 】

特許文献 1 に開示された技術によれば、ユーザコード（カーネル）の開発者は、OpenCL デバイスの構成を知らなくても、カーネルが使用するデータブロック毎に属性群さえ指定すれば、高性能なカーネルを可搬性の高い形で実現できる。それは、OpenCL デバイスの演算制御部は、カーネルが参照する各データブロックの属性群と、OpenCL デバイスの構成とに基づいて、データブロックを該デバイスにとって最適なサイズに適宜に分割して、「リード転送 カーネルによる演算処理 ライト転送」をボディとする反復処理を、分割数に等しい回数だけ繰り返すよう制御可能である。デバイス側の演算制御部の設計については、OpenCL デバイスの構成を熟知した専門家、例えば該 OpenCL デバイスのメーカーの開発者により行えばよい。

50

【先行技術文献】

【特許文献】

【0070】

【特許文献1】特開2013-025547号公報

【非特許文献】

【0071】

【非特許文献1】The OpenCL Specification, Ver: 1.0, Document Revision: 43, Khronos OpenCL Working Group (2009)

【発明の概要】

10

【発明が解決しようとする課題】

【0072】

ここで、連続して実行される2つのカーネルの間に関連性がある場合を考える。例えば、第1のカーネルXは、先に実行され、その演算結果のデータブロックが、後に実行される第2のカーネルYの入力データとして使用されるとする。なお、各データブロックは、サブブロックに分割して転送されるとする。

【0073】

特許文献1のOpenCLシステムの場合、下記の流れになると考える。

<ステップX1>

ホーストは、第1のカーネルXをデバイスに転送する。

20

該第1のカーネルXの引数には、第1のカーネルXによる演算処理に使用される各データブロック(リードブロック)毎に、上述した属性群が含まれている。

【0074】

<ステップX2>

デバイスの演算制御部は、第1のカーネルXの引数に基づいて、グローバルメモリに格納された各リードブロック、及び演算結果となるライトブロックの分割方式を含む転送方式を決定する。そして、各リードブロックをサブリードブロックに分割して、ローカルメモリまたはプライベートメモリへ転送すると共に、PEに演算を行わせる。

【0075】

<ステップX3>

30

デバイスの各PEは、演算制御部の制御に従って演算を行い、サブライトブロックを得る。これらのサブライトブロックは、ローカルメモリまたはプライベートメモリに格納された後、演算制御部の制御に従って、PEからグローバルメモリに転送される。

【0076】

<ステップ4>

全ての処理が終わるまで、演算制御部の制御の元で、ステップX2~ステップX3の処理が繰り返される。最後に、第1のカーネルXの演算結果となるデータブロック(ライトブロック)が得られ、グローバルメモリに格納される。

【0077】

第1のカーネルXの処理が完了すると、下記の流れで第2のカーネルYの処理が行われる。

40

【0078】

<ステップY1>

ホーストは、第2のカーネルYをデバイスに転送する。

該第2のカーネルYの引数には、第2のカーネルYによる演算処理に使用される各データブロック(リードブロック)毎に、上述した属性群が含まれている。

【0079】

第1のカーネルXのライトブロックのうちのあるライトブロックが、第2のカーネルYの演算処理に使用される場合に、第2のカーネルYの引数には、該ライトブロックに対応する、第2のカーネルYのリードブロックについても、属性群が指定されている。以下に

50

において、先に実行されるカーネルのライトブロックであって、後に実行されるカーネルのリードブロックになるデータブロックを「継続ライトブロック」といい、後に実行されるカーネルのリードブロックであって、先に実行されるカーネルのライトブロックになるデータブロックを「継続リードブロック」という。

【0080】

<ステップY2>

デバイスの演算制御部は、第2のカーネルYの引数に基づいて、グローバルメモリに格納された各リードブロック（継続リードブロックを含む）、及び演算結果となるライトブロックの分割方式を含む転送方式を決定する。そして、各リードブロックをサブリードブロックに分割して、ローカルメモリまたはプライベートメモリへ転送すると共に、PEに演算を行わせる。

10

【0081】

<ステップY3>

デバイスの各PEは、演算制御部の制御に従って演算を行い、サブライトブロックを得る。これらのサブライトブロックは、ローカルメモリまたはプライベートメモリに格納された後、演算制御部の制御に従って、PEからグローバルメモリに転送される。

【0082】

<ステップ4>

全ての処理が終わるまで、演算制御部の制御の元で、ステップY2～ステップY3の処理が繰り返される。最後に、第2のカーネルYの演算結果となるデータブロック（ライトブロック）が得られ、グローバルメモリに格納される。

20

【0083】

上述した流れから分かるように、継続ライトブロックがある場合、該継続ライトブロックのサブライトブロックのグローバルメモリへの転送（ライト転送）と、該継続ライトブロックに対応する継続リードブロックのサブリードブロックのローカルメモリ/プライベートメモリへの転送（リード転送）が行われている。

【0084】

継続ライトブロックのサブライトブロックが、既にプライベートメモリ/ローカルメモリにリード転送されたサブリードブロックとして次のカーネルに使用可能であれば、継続ライトブロックをサブリードブロックのローカルメモリ/プライベートメモリからグローバルメモリへ転送せず、そのまま次のカーネルに供するようにすれば、ライト転送もリード転送も削減でき、OpenCLデバイスの処理効率の向上を図ることができる。

30

【0085】

その他の課題と新規な特徴は、本明細書の記述および添付図面から明らかになるであろう。

【課題を解決するための手段】

【0086】

一実施の形態は、OpenCLデバイスの複数の演算素子による並列演算を制御する演算制御方法であり、第1のステップと第2のステップを有する。

【0087】

40

第1のステップは、OpenCLデバイスの最下位階層のメモリ（例えばグローバルメモリ）に1つ以上格納されたリードブロックと、他の階層のメモリ（例えばプライベートメモリやローカルメモリ）から最下位階層のメモリに転送される1つ以上のライトブロックとに対して夫々設定された属性群を取得して保持するステップである。

【0088】

第2のステップは、第1のステップにより保持された各前記属性群と、前記OpenCLデバイスの構成を示す構成パラメータとに基づいて、夫々の前記リードブロックと前記ライトブロックの転送方式を決定し、決定した転送方式に応じて各前記リードブロックと前記ライトブロックの転送、及び該転送に対応する前記並列演算の制御を行うステップである。

50

【 0 0 8 9 】

リードブロックは、並列演算の演算対象として最下位階層のメモリから他の階層のメモリに転送されるデータブロックである。ライトブロックは、並列演算後に他の階層のメモリから最下位階層のメモリに転送されるデータブロックであり、前記1つ以上のリードブロックに対する並列演算の演算結果である。

【 0 0 9 0 】

前記属性群は、前記転送方式を決定するために必要である一方、前記OpenCLデバイスの構成に依存しない属性を複数含むものであり、例えば割当属性、余白属性、依存属性を有する。

【 0 0 9 1 】

割当属性は、当該データブロックを複数のサブブロックに分割して転送するか否か、及び、分割して転送する場合における分割方式を示すものである。

【 0 0 9 2 】

余白属性は、当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックと共に転送される、該サブブロックの周辺のデータのサイズを示すものである。

【 0 0 9 3 】

依存属性は、当該データブロックを複数のサブブロックに分割して転送する場合における、該サブブロックがその周辺の他のサブブロックと依存関係が有るか否か、及び、前記依存関係がある場合における全ての依存方向を示す依存属性を有する。

【 0 0 9 4 】

なお、ライトブロックの属性群は、該ライトブロックが既に前記他の階層のメモリに存在し、かつ、前記最下位階層のメモリに転送されると仮定して設定されたものである。

【 0 0 9 5 】

第2のステップは、第1のカーネルに続いて第2のカーネルをOpenCLデバイスに実行させ、かつ、第1のカーネルに対応する並列演算の各ライトブロックのうちに、第2のカーネルに対応する並列演算のリードブロックとして使用される継続ライトブロックが含まれており、かつ、第1のカーネルに対して設定された継続ライトブロックの割当属性と、第2のカーネルに対して設定された、継続ライトブロックに対応するリードブロック（継続リードブロック）の割当属性とが一致する場合において、該継続ライトブロックが最下位階層へ転送されず、他の階層のメモリを介して第2のカーネルの実行に供されることにより、前記第1のカーネルと前記第2のカーネルの実行をパイプライン化するパイプライン化制御を行うステップをさらに含む。

【 0 0 9 6 】

該パイプライン化制御ステップにおいて、第1のカーネルに対応する各リードブロックの夫々に対して、該リードブロックに対して設定された余白属性と依存属性に、第2のカーネルに対して設定された、継続ライトブロックに対応する継続リードブロックの余白属性と依存属性を夫々論理和加算する。

【 0 0 9 7 】

なお、上記実施の形態の方法を装置やシステムに置き換えて表現したもの、該方法を実行するOpenCLデバイス、該方法をコンピュータに実行せしめるプログラムなども、実施の態様としては有効である。

【 発明の効果 】

【 0 0 9 8 】

前記一実施の形態によれば、OpenCLデバイスの処理効率を向上させることができる。

【 図面の簡単な説明 】

【 0 0 9 9 】

【 図 1 】一実施の形態にかかるOpenCLシステムを示す図である。

【 図 2 】単独カーネルの場合を説明するための具体例を示す図である。

10

20

30

40

50

【図3】リード属性を説明するための図である。

【図4】余白属性を説明するための図である。

【図5】放送属性を説明するための図である。

【図6】割当属性を説明するための図である(その1)

【図7】割当属性を説明するための図である(その2)。

【図8】依存属性を説明するための図である。

【図9】図2に示す各データブロックに対して設定された属性群の例を示す図である。

【図10】図1に示すOpenCLシステムにおける演算ユニットの構成を示すパラメータの例である。

【図11】階層内での対応付けと階層間での対応付けを示す図である。

10

【図12】図2に示す例に対してシナリオ決定部が決定した内容の一部を示す図である。

【図13】関連性のある2つのカーネルを連続して実行する際に行われるパイプライン化制御を説明するための図である(その1)。

【図14】関連性のある2つのカーネルを連続して実行する際に行われるパイプライン化制御を説明するための図である(その2)。

【図15】パイプライン化制御なしとパイプライン化制御ありの場合における各データブロックの転送態様を比較するための図である。

【図16】パイプライン化制御なしとパイプライン化制御ありの場合におけるカーネルの実行態様を比較するための図である。

【図17】OpenCLシステムのプラットフォームモデルを示す図である。

20

【図18】OpenCLデバイスの構成を示す図である。

【図19】OpenCLデバイスにおける各種メモリを説明するための図である。

【図20】逐次プロセッサの例を示す図である。

【図21】特許文献1に開示された技術が解決しようとする課題を説明するための図である(その1)。

【図22】特許文献1に開示された技術が解決しようとする課題を説明するための図である(その2)。

【図23】特許文献1のOpenCLシステムにおいて、ユーザにより指定する必要がある内容を示す図である。

【発明を実施するための形態】

30

【0100】

説明の明確化のため、以下の記載及び図面は、適宜、省略、及び簡略化がなされている。また、様々な処理を行う機能ブロックとして図面に記載される各要素は、ハードウェア的には、CPU、メモリ、その他の回路で構成することができ、ソフトウェア的には、メモリにロードされたプログラムなどによって実現される。したがって、これらの機能ブロックがハードウェアのみ、ソフトウェアのみ、またはそれらの組合せによっていろいろな形で実現できることは当業者には理解されるべきであり、いずれかに限定されるものではない。なお、各図面において、同一の要素には同一の符号が付されており、必要に応じて重複説明は省略されている。

【0101】

40

また、上述したプログラムは、様々なタイプの非一時的なコンピュータ可読媒体(non-transitory computer readable medium)を用いて格納され、コンピュータに供給することができる。非一時的なコンピュータ可読媒体は、様々なタイプの実体のある記録媒体(tangible storage medium)を含む。非一時的なコンピュータ可読媒体の例は、磁気記録媒体(例えばフレキシブルディスク、磁気テープ、ハードディスクドライブ)、光磁気記録媒体(例えば光磁気ディスク)、CD-ROM(Read Only Memory)CD-R、CD-R/W、半導体メモリ(例えば、マスクROM、PROM(Programmable ROM)、EPROM(Erasable PROM)、フラッシュROM、RAM(Random Access Memory))を含む。また、プログラムは、様々なタイプの

50

一時的なコンピュータ可読媒体 (transitory computer readable medium) によってコンピュータに供給されてもよい。一時的なコンピュータ可読媒体の例は、電気信号、光信号、及び電磁波を含む。一時的なコンピュータ可読媒体は、電線及び光ファイバ等の有線通信路、又は無線通信路を介して、プログラムをコンピュータに供給できる。

【0102】

図1は、一実施の形態にかかるOpenCLシステム100を示す。OpenCLシステム100は、OpenCLホスト(以下、単に「ホスト」という)110、OpenCLデバイス(以下、単に「デバイス」という)120を備える。

【0103】

デバイス120は、演算制御部200、演算ユニット140を備える。演算制御部200は、属性群格納部210、カーネル間情報格納部220、シナリオ決定部230を有する。演算ユニット140は、デバイスメモリ150、複数のコンピュータユニット(CU)160、CU160毎に設けられたローカルメモリ170を有する。

【0104】

デバイスメモリ150は、グローバルメモリ152、コンスタントメモリ154を有する。また、夫々のCU160は、複数のPE162と、PE162毎に設けられたプライベートメモリ164を有する。

【0105】

なお、グローバルメモリ152、コンスタントメモリ154、CU160、PE162、プライベートメモリ164、ローカルメモリ170は、通常のOpenCLデバイスにおける同一名称のものと同様であり、ここで詳細な説明を省略する。

【0106】

ホスト110は、デバイス120に演算を行わせる際に、ユーザにより引数の指定がなされた、該演算に対応するカーネルをデバイス120に転送する。

【0107】

デバイス120において、カーネル間情報格納部220は、連続して実行される複数のカーネル間のデータブロックの依存関係を示す情報(「カーネル間データブロック依存情報」という)を格納している。この情報は、ユーザにより直接カーネル間情報格納部220に入力されるようにしてもよいし、上記複数のカーネルのうちの1つ(例えば最も先に実行されるカーネル)の引数に含まれるようにして、該カーネルと共にホスト110からデバイス120に転送され、カーネル間情報格納部220に格納されるようにしてもよい。

【0108】

上記カーネル間データブロック依存情報は、例えば、先に実行されるカーネルのライトブロックのうちに、後に実行されるカーネルのリードブロックとして使用されるものがあるか否か、及びある場合における、継続ライトブロックと継続リードブロックとの対応関係を示すものである。

【0109】

なお、単独カーネルの場合には、カーネル間データブロック依存情報は、「依存なし」を示すことになる。単独カーネルは、他のカーネルと関連性が無く、その演算対象となるリードブロックが、他のカーネルの演算結果となるライトブロックでもなければ、そのライトブロックが他のカーネルのリードブロックとしても使用されないカーネルを意味する。

【0110】

演算制御部200は、ホスト110からのカーネルが示す演算内容、属性群格納部210に格納された各属性群と、カーネル間情報格納部220に格納されたカーネル間データブロック依存情報とに基づいて、演算ユニット140を制御する。演算ユニット140は、演算制御部200の制御に従って演算を行い、演算結果を得る。演算ユニット140が行う演算は、演算対象のデータをデバイスメモリ150からプライベートメモリ164

10

20

30

40

50

／ローカルメモリ 170 への転送（リード転送）と、演算結果のデータをデバイスメモリ 150 への転送（ライト転送）が伴う。

【0111】

通常、演算対象のデータと演算結果のデータのいずれもグローバルメモリ 152 に格納されるため、以下において、「リード転送」と「ライト転送」がグローバルメモリ 152 と、プライベートメモリ 164 / ローカルメモリ 170 との間の転送であるとして説明をする。演算対象のデータがコンスタントメモリ 154 にも格納される場合には、「リード転送」がコンスタントメモリ 154 からプライベートメモリ 164 / ローカルメモリ 170 への転送も含むことを理解されたい。

【0112】

< 単独カーネルの場合 >

分かりやすいように、まず、図 2 に示す具体例を参照しながら、単独カーネルを実行する場合について説明する。

【0113】

この場合、カーネル間情報格納部 220 には、「依存なし」を示すカーネル間データブロック依存情報が格納されることになる、OpenCL システム 100 における処理は、特許文献 1 に開示された OpenCL システムの処理と同様である。

【0114】

なお、以下の説明における用語「次元」は、前述したインデックス空間の次元に対応する。OpenCL では、インデックス空間は、3 次元まで指定することができるが、説明上の便宜のため、2 次元を例にする。

【0115】

図 2 に示すように、ユーザが実現したい処理は、式 (1) に示す処理である。

【数 1】

$$R_{(x, y)} = \sum [P_i * Q_{(x, (y+i))}] \quad (i=0 \sim 8) \quad (1)$$

【0116】

演算対象のデータは、リードブロック P とリードブロック Q に含まれる各データである。演算結果のデータは、ライトブロック R に含まれる各データである。図 2 に示す例では、ライトブロックは、ライトブロック R の 1 つのみである。

【0117】

式 (1) に示す処理を実現するためのカーネルは、当然ながら、ユーザにより作成される。カーネルの作成は、デバイス 120 の各 PE 162 により実行されるプログラムコード（以下、「ユーザ指定処理」という）の作成処理と、引数の指定処理を含む。

【0118】

本実施の形態の OpenCL システム 100 では、OpenCL で定められた各引数以外に、さらに、8 つの属性を含む「属性群」がサポートされている。これらの 8 つの属性は、「サイズ属性」、「リード属性」、「ライト属性」、「余白属性」、「放送属性」、「割当属性」、「階層属性」、「依存属性」であり、ユーザにより、カーネルの引数として、全てのリードブロックとライトブロックに対して指定される。なお、属性の指定に際して、必ずしも 8 つの属性の全てを指定する必要があるとは限らない。

【0119】

上記属性群に含まれる 8 つの属性は、グローバルメモリ 152 と、プライベートメモリ 164 及び / またはローカルメモリ 170 との間で転送される方式（転送方式）を決定するために参照する必要があるパラメータである。以下、特別な説明が無い限り、「転送」は、グローバルメモリ 152 と、プライベートメモリ 164 及び / またはローカルメモリ 170 との間の転送を意味する。

【0120】

10

20

30

40

50

「転送方式」は、下記の内容が含まれる。

(1) 転送方式 1 : 分割の有無

該転送方式 1 は、「分割無し」方式と「分割有り」方式のいずれかである。

【 0 1 2 1 】

リードブロックについて、「分割無し」とは、該リードブロックの全てのデータを一度のリード転送により転送することを意味し、「分割有り」とは、該リードブロックを複数のサブリードブロック S R B に分割し、一度のリード転送により一つのワークグループへ 1 つのサブリードブロック S R B のみを転送することを意味する。

【 0 1 2 2 】

ライトブロックについて、「分割無し」とは、該ライトブロックの全てのデータ（演算結果）を一度のライト転送によりグローバルメモリ 1 5 2 に転送することを意味し、「分割有り」とは、該ライトブロックを複数のサブライトブロック S W B に分割し、一度のライト転送により一つのワークグループから 1 つのサブライトブロック S W B のみを転送することを意味する。

【 0 1 2 3 】

(2) 転送方式 2 : 分割方式

該転送方式 2 は、該データブロックを複数のサブブロックに分割する場合に、どのように分割するかを意味する。

【 0 1 2 4 】

(3) 転送方式 3 : 分配方式（サブリードブロック S R B の対応付け方法）

該転送方式 3 は、リードブロックを対象として、上記転送方式 1 が分割有りの場合に指定される。

【 0 1 2 5 】

すなわち分配方式は、リードブロックが複数あるときには、同じ回のリード転送により転送される、各リードブロックからの 1 つずつのサブリードブロック S R B をどのようにしてプライベートメモリ 1 6 4 / ローカルメモリ 1 7 0 に分配するかを意味する。リードブロックが 1 つしか無い場合には、分割方式は、該サブリードブロック S R B をプライベートメモリ 1 6 4 / ローカルメモリ 1 7 0 に均等に分配する方式になる。

【 0 1 2 6 】

(4) 転送方式 4 : 統合方式（サブライトブロック S W B の対応付け方法）

該転送方式 4 は、ライトブロックを対象として、上記転送方式 1 が分割有りの場合に指定される。

【 0 1 2 7 】

なお、「統合する」とは、グローバルメモリ 1 5 2 における、該ライトブロックの個々のサブライトブロック S W B に割り当てられた領域に書き込むことを意味する。

【 0 1 2 8 】

統合方式は、ライトブロックが複数あるときには、同じ回のライト転送により転送される、プライベートメモリ群 / ローカルメモリ 1 7 0 に格納されている複数のサブライトブロック S W B を、どのようにして夫々のライトブロックに統合するかを意味する。ライトブロックが 1 つしか無い場合には、統合方式は、プライベートメモリ群 / ローカルメモリ 1 7 0 に格納された該ライトブロックの各サブライトブロック S W B のデータを、上記 1 つのライトブロックに統合する方式になる。

【 0 1 2 9 】

上記した「転送方式」は、従来では、ユーザコードにより指定される必要がある。演算処理の内容、及びデバイス 1 2 0（具体的には演算ユニット 1 4 0）の構成に依存するため、ユーザコードが煩雑であり、可搬性を維持するのは極めて困難と言ってよい。

【 0 1 3 0 】

対して、本実施の形態の OpenCL システム 1 0 0 は、演算処理の内容と OpenCL デバイスの構成の双方を考慮しつつその都度、個々のユーザが時間を費やして場当たり的に転送方式を決定してきた処理を、OpenCL デバイスの構成に依存する処理と依存

10

20

30

40

50

しない処理とに分け、さらに、OpenCLデバイスに依存する処理については、OpenCLデバイスに依存しない処理により指定したパラメータと、OpenCLデバイスの構成を示すパラメータとに基づいて自動的に決定するようにしたものである。

【0131】

具体的には、本実施の形態のOpenCLシステム100では、OpenCLデバイスに依存しない処理は、データブロックに対する上記属性群の指定に該当する。これらの属性群は、転送方式を決定する上で必要なパラメータであるものの、OpenCLデバイスの構成に依存しないものである。以下、詳細に説明する。なお、以下において、各属性が「転送方式を決定する上で必要であるものの、OpenCLデバイスの構成に依存しないもの」であることを繰り返せず、このこと以外の要素についてのみを説明する。

10

【0132】

上記8つの属性は、さらに、「固有属性」、「演算属性」、「ポリシ属性」の3種類に分けられる。図2に示す各データブロックの例を参照しながら説明する。

【0133】

「固有属性」は、演算処理の内容及びユーザの意志に関係なく、当該データブロックが持つ固有の属性である。本実施の形態において、「固有属性」は、下記の「サイズ属性」である。

【0134】

<サイズ属性>

この属性は、データブロックのサイズを示すパラメータであり、例えば、次元毎のワード数、及びワード毎のバイト数やビット数である。この属性は、全てのデータブロックに対して必ず指定される。

20

【0135】

そのため、図2に示すリードブロックP、リードブロックQ、及びライトブロックRのサイズ属性は、下記ようになる。

【0136】

リードブロックPは、2次元のデータブロックであり、X方向サイズ L_x とY方向サイズ L_y が、共に3ワードである。なお、リードブロックPは、1ワードのビット数が、8ビットである。そのため、サイズ属性として、リードブロックPに対して、「X方向サイズ L_x : 3、Y方向サイズ L_y : 3、ビット数/ワード : 8」が指定される。

30

【0137】

同様に、リードブロックQに対して、「X方向サイズ L_x : 640、Y方向サイズ L_y : 480、ビット数/ワード : 16」が指定される。

【0138】

また、ライトブロックRに対して、「X方向サイズ L_x : 640、Y方向サイズ L_y : 480、ビット数/ワード : 32」が指定される。

【0139】

「演算属性」は、ユーザの意志には関係しないものの、演算処理の内容に関係する属性である。本実施の形態において、「演算属性」は、下記の「リード属性」、「ライト属性」を含む。演算属性は、各データブロックが1つ以上のサブブロックに分割される仮定の元で指定される。ライトブロックについては、さらに、リードブロックが既にプライベートメモリ群/ローカルメモリ上に存在する仮定が加えられる。なお、データブロックが1つのサブブロックに分割されることは、分割されないことである。

40

【0140】

<リード属性>

この属性は、まず、該データブロックが演算対象のデータ(つまり、リード転送されるデータ)であるか否か、及び演算対象のデータである場合の転送順位を示す。転送順位は、該データブロックの各サブブロックを、どのような順位で転送するかを指定するパラメータである。

【0141】

50

ライトブロックは、リード転送されないで、リード転送されないことを示す「NONE」が指定される。リードブロックは、リード転送されるデータブロックであるので、「リード転送される」として、転送順位が指定される。

【0142】

本実施の形態のOpenCLシステム100において、リード属性として設定される転送順位は、「TOP LEFT」、「BOTTOM RIGHT」、「ランダム」が指定可能である。なお、「ランダム」が指定された場合に限り、別途、転送順位を示す情報を格納した領域へのポインタが指定される。

【0143】

「TOP LEFT」は、左上端のサブブロックから転送することを示し、「BOTTOM RIGHT」は、右下端のサブブロックから転送することを示す。図3は、リード属性として「TOP LEFT」が指定された場合のサブブロックの転送順位を示す。

10

【0144】

図3に示すように、この場合、グローバルメモリ152に格納されたデータブロック(リードブロック)は、転送順位が、左上端のサブブロック1、サブブロック1の右隣のサブブロック2、サブブロック3の右隣のサブブロック3、・・・となっている。

【0145】

図2に示す各データブロックのリード属性を説明する。

式(1)に示す演算処理の内容に基づいて、リードブロックPとリードブロックQは、リード転送されるデータであるため、リード属性として、転送順位が「TOP LEFT」に指定される。

20

【0146】

一方、ライトブロックRは、リード転送されるデータではないため、リード属性が「NONE」に指定される。

【0147】

<ライト属性>

この属性は、まず、該データブロックが演算結果のデータ(つまり、ライト転送されるデータ)であるか否か、及び演算結果のデータである場合の転送順位を示す。転送順位は、各サブライトブロックSWBを、どのような順位で転送するかを指定するパラメータである。

30

【0148】

リードブロックは、ライト転送されないで、ライト転送されないことを示す「NONE」が指定される。ライトブロックは、ライト転送されるデータブロックであるので、「ライト転送される」として、転送順位が指定される。

【0149】

本実施の形態のOpenCLシステム100において、ライト属性として設定される転送順位は、「TOP LEFT」、「BOTTOM RIGHT」、「ランダム」が指定可能である。なお、「ランダム」が指定された場合に限り、別途、転送順位を示す情報を格納した領域へのポインタが指定される。

【0150】

40

リードブロックは、ライト転送されないで、ライト転送されないことを示す「NONE」が指定される。そのため、図2に示すリードブロックPとリードブロックQは、ライト属性として「NONE」が指定される。

【0151】

ライトブロックは、ライト転送されるデータブロックであるので、「ライト転送される」として、上記転送順位が指定される。この転送順位は、該データブロックの各サブブロックを、どのような順位で転送するかを指定するパラメータである。

【0152】

本実施の形態のOpenCLシステム100において、ライト属性として設定される転送順位は、リード属性として設定される転送順位と同様に、「TOP LEFT」と、「

50

「BOTTOM RIGHT」と、「ランダム」が指定可能である。「ランダム」が指定された場合に限り、別途、転送順位を示す情報を格納した領域へのポインタが指定される。

【0153】

ライト属性として設定される転送順位の夫々のパラメータの意義は、リード属性として設定される転送順位の相対応するパラメータと同一であるので、ここで詳細な説明を省略する。

【0154】

図2に示す各データブロックのライト属性を説明する。

式(1)に示す演算処理の内容に基づいて、ライトブロックRは、ライト転送されるデータであるため、ライト属性として、転送順位が「TOP LEFT」に指定される。

10

【0155】

一方、リードブロックPとリードブロックQは、ライト転送されるデータではないため、ライト属性が「NONE」に指定される。

【0156】

「ポリシ属性」は、演算処理の内容と共に、ユーザがどのように転送及び演算処理を実行したいかの意志に関係する属性である。本実施の形態のOpenCLシステム100において、「ポリシ属性」は、「余白属性」、「放送属性」、「割当属性」、「階層属性」、「依存属性」を含む。ポリシ属性も、各データブロックが1つ以上のサブブロックに分割される仮定の元で指定される。ライトブロックについては、さらに、ライトブロックが既にプライベートメモリ20/ローカルメモリ22上に存在する仮定が加えられる。なお、データブロックが1つのサブブロックに分割されることは、分割されないことである。

20

【0157】

<余白属性>

この属性は、リードブロックを対象とするパラメータであり、サブリードブロックSRB内のデータと共に転送される、該サブリードブロックSRBの境界と隣接する該サブリードブロックSRB外のデータの量を示すパラメータである。また、余白属性は、次元毎に指定される。余白属性のパラメータの単位は、ワードである。

【0158】

なお、ライトブロックについては、余白属性を指定できない、または指定したとしても無視されるようになっている。

30

【0159】

図4は、余白属性として、X方向に「1」、Y方向に「2」が指定された場合に該データブロック(リードブロック)の各サブリードブロックSRBの転送時の転送範囲を示す。この場合、該サブリードブロックSRBの転送時の転送範囲は、該サブリードブロックSRB内のデータに加え、該サブリードブロックSRBの左右両側の境界に隣接する1列ずつのデータと、上下両端の境界に隣接する2行ずつのデータが含まれる。なお、図4では、データブロックの上端に位置するサブリードブロックSRBを例としており、該サブリードブロックSRBの上端の境界に隣接するデータが無いため、該サブリードブロックSRBの転送時の転送範囲には、上端の境界に隣接するデータが含まれない。

【0160】

後に詳細に説明するが、本実施の形態のOpenCLシステム100では、各リードブロックについて、1度のリード転送時に1つのCU160に対して1つのサブリードブロックSRBが転送され、1度のリード転送により該CU160のプライベートメモリ/ローカルメモリに転送された各サブブロックサブリードブロックSRBを対象として該CU160により演算が行われる。そして、演算の結果は、夫々のライトブロックの1つのサブブロックとして該CU160からグローバルメモリに転送される。

40

【0161】

例えば、2次元画像に対して3×3のコンボリューション演算を行う場合、注目画素と上下左右の4つの方向において、該注目画素に隣接する1画素が必要である。サブリードブロックSRB内のデータのみが転送されるのでは、該サブリードブロックSRBの最も

50

外側に位置する各画素に対する演算ができない。そのため、この場合、余白属性として、X方向とY方向に共に「1」を指定する必要がある。

【0162】

図2に示す各データブロックの余白属性を説明する。なお、余白属性の指定は、他の属性の指定と関係するため、余白属性のみでは説明しにくい。そのため、ここでは、図2に示す各データブロックに対して指定された余白属性の値のみを示し、それらの意義の詳細については後述する。

【0163】

リードブロックPとリードブロックQは、余白属性の指定対象であるが、リードブロックPに対しては、X方向とY方向のいずれについても余白が「0」に設定される。また、

10

【0164】

リードブロックQに対しては、X方向の余白が「0」、Y方向の余白が「9」に指定される。従って、リードブロックPの転送時には、サブリードブロックSRB内のデータのみが転送される。一方、リードブロックQの転送時には、サブリードブロックSRB内のデータに加え、該サブリードブロックSRBと下端で隣接する9行のデータも転送される。

【0165】

ライトブロックRは、ライト転送されるデータであるため、余白属性の指定対象ではない。すなわち、ライトブロックRの転送時には、当該サブライトブロックSWB内のデータのみが転送される。

20

【0166】

<放送属性>

この属性は、各リードブロックについてはサブリードブロックSRBの転送先、各ライトブロックについてはサブライトブロックSWBの転送元が、プライベートメモリとローカルメモリのいずれになるかを指定するパラメータであり、「ON」と「OFF」のいずれである。例えば、図5に示すように、放送属性の「ON」は、上記転送先または転送元としてローカルメモリを指定し、放送属性の「OFF」は、上記転送先または転送元としてプライベートメモリを指定する。

【0167】

図2に示す各データブロックの放送属性を説明する。なお、放送属性の指定も、他の属性と関係するため、放送属性のみでは説明しにくい。そのため、ここでは、図2に示す各データブロックに対して指定された放送属性の値のみを示し、それらの意義の詳細については後述する。

30

【0168】

リードブロックPは、放送属性が「ON」に指定される。そのため、リードブロックPの転送時に、各サブリードブロックSRBは、ローカルメモリ170に転送される。

【0169】

リードブロックQとライトブロックRは、放送属性が「OFF」に指定される。そのため、リードブロックQの転送時には、各サブリードブロックSRBは、プライベートメモリ164に転送される。また、ライトブロックRのサブライトブロックSWBは、プライベートメモリ群からグローバルメモリ152に転送される。

40

【0170】

<割当属性>

この属性は、サブリードブロックSRBとサブライトブロックSWBをどのようにしてCU160のプライベートメモリ群/ローカルメモリに割り当てる割当方式を示すパラメータである。

【0171】

本実施の形態のOpenCLシステム100では、「縦優先」と「横優先」の2つの割当方式がサポートされる。

【0172】

50

割当属性を詳細に説明する前に、まず、ワークグループサイズ WGs と、1つのアイテムに割り当てるデータ量を説明する。

【0173】

放送属性が「OFF」のデータブロックに対応する「ワークグループサイズ WGs 」は、1つのワークグループ内のアイテム数により表されるデータ量であり、例えばアイテム数が N であれば、ワークグループサイズ WGs は N ワードになる。一方、放送属性が「ON」のデータブロックに対応するワークグループサイズ WGs は、常に1と見なされる。以降、放送属性が「OFF」のデータブロックに対応する「ワークグループサイズ WGs 」を単にワークグループサイズ WGs と呼ぶ。このワークグループサイズ WGs は、後にデバイス120における演算制御部200のシナリオ決定部230により決定されるが、ユーザが、カーネルの作成時に、その最大値と最小値のいずれをとるかを指定することができる。

10

【0174】

また、1つのアイテムに割り当てるデータ量も、後にシナリオ決定部230により決定される。ワークグループサイズ WGs が最小値をとる場合とは、1アイテムが1PEに対応する場合である。それに対し、PEのプライベートメモリを M ($M: 2$ 以上の整数)分割して M アイテムを1PEに対応させることで、ワークグループサイズ WGs を最大 M 倍大きく指定できるOpenCLデバイスの場合、最大値は最小値の M 倍となる。一方、最大値のワークグループサイズ WGs を利用する場合、1PE当たりのプライベートメモリは最小値のワークグループサイズ WGs を利用する場合の M 分の1となるだけである。以下、説明を簡単にするため、シナリオ決定部230は、常に、ワークグループサイズ WGs としてその最小値を採用するものとし、その結果、1アイテムが1PEに対応するので、以降、アイテムという表現の代わりにPEという表現を用いる。

20

【0175】

「縦優先」は、サブブロックのX方向のサイズ($SBsx$)がワークグループサイズ WGs であり、Y方向のサイズ $SBsy$ が1つのアイテムに割り当てるデータ量、すなわち1PEに割り当てるデータ量であるようにデータブロックを割り当てる方式である。

【0176】

割当属性として「縦優先」が指定されたデータブロック内の各サブブロックは、1サブブロックが1ワークグループに対応し、かつ、該サブブロック内の1列のデータが該ワークグループ内の1つのPE162のプライベートメモリ164に対応するように、グローバルメモリ152とプライベートメモリ164間で転送される。

30

【0177】

図6は、割当属性として「縦優先」が指定されたリードブロックの場合の例を示す。この場合、夫々のサブリードブロックSRBは、該サブリードブロックSRB内の全てのデータが同一のワークグループ内のPE162のプライベートメモリ164に格納され、かつ、該サブブロック内の同列のデータが同一のPE162のプライベートメモリ164に格納されるように、グローバルメモリ152からプライベートメモリ164に転送される。

【0178】

「横優先」は、サブブロックのX方向のサイズ $SBsx$ がワークグループサイズ WGs の整数倍になるようにデータブロックを分割する方式である。なお、サブブロックのY方向のサイズ $SBsy$ は、後にシナリオ決定部230により決定される。

40

【0179】

割当属性として「横優先」が指定されたデータブロック内の各サブブロックは、1サブブロックが1つのワークグループに対応し、該サブブロック内の全てのデータが、同一のワークグループに含まれるPE162のプライベートメモリ164に転送される。さらに、サブブロックの各行は、該行をワークグループサイズ WGs の量毎に区切って得た整数個の区切りブロックのデータが、区切りブロック毎に、該ワークグループ内の WGs 個のPE162のプライベートメモリ164に均等に分配される。

50

【 0 1 8 0 】

図7は、割当属性として「横優先」が指定されたリードブロックの場合の例を示す。なお、該例では、サブライトブロックS W BのY方向のサイズS B s yが1ワードとされている。つまり、該リードブロックは、サブリードブロックS R BのX方向のサイズS B s xがワークグループサイズW G sの整数倍であり、行数が1であるように分割される。

【 0 1 8 1 】

図示のように、この場合、該リードブロックの夫々のサブリードブロックS R Bは、該サブリードブロックS R B内の全てのデータが同一のワークグループ内のP E 1 6 2のプライベートメモリ1 6 4に格納されるように転送される。さらに、該サブリードブロックS R Bを行方向にワークグループサイズW G sの量毎に区切って得た整数個の区切りブロックのデータが、区切りブロック毎に、該ワークグループ内のW G s個のP E 1 6 2のプライベートメモリ1 6 4に均等に分配されるように転送される。例えば、データ1から、データaの前のデータまでの複数のデータは、1つの区切りブロックを構成しており、同一のワークグループの複数のP E 1 6 2のプライベートメモリ1 6 4に夫々格納されるように転送される。また、データaから、データjの前のデータまでの複数のデータも、1つの区切りブロックを構成しており、ワークグループの複数のP E 1 6 2のプライベートメモリ1 6 4に夫々格納されるように転送される。

【 0 1 8 2 】

図2に示す各データブロックの割当属性を説明する。なお、割当属性の指定も、他の属性と関係するため、割当属性のみでは説明しにくい。そのため、ここでは、図2に示す各データブロックに対して指定された割当属性の値のみを示し、それらの意義の詳細については後述する。

【 0 1 8 3 】

リードブロックPは、割当属性として「縦優先」が指定されたとする。

リードブロックQとライトブロックRも、割当属性が「縦優先」に指定されたとする。

【 0 1 8 4 】

< 階層属性 >

この属性は、1以上の自然数で指定される階層数である。同一の階層数が指定された複数のリードブロックは、1度のリードあるいはライト転送により、夫々1つのサブブロックS Bが転送される。

【 0 1 8 5 】

図2に示す各データブロックの階層属性を説明する。なお、階層属性の指定も、他の属性と関係するため、階層属性のみでは説明しにくい。そのため、ここでは、図2に示す各データブロックに対して指定された階層属性の値のみを示し、それらの意義の詳細については後述する。

【 0 1 8 6 】

リードブロックPとリードブロックQは、階層属性が「2」と「1」に夫々指定される。また、ライトブロックRの階層属性が「1」に指定される。

【 0 1 8 7 】

< 依存属性 >

この属性は、当該データブロックに対して、サブブロックと、該サブブロックと隣接する8つの他のサブブロックとのデータ依存関係を示すパラメータであり、「依存有り」、「依存無し」を指定可能である。また、「依存有り」については、さらに、3種類の依存関係を指定可能である。

【 0 1 8 8 】

「依存有り」は、R 1型、R 2型、R 3型の3種類がある。図8を参照して説明する。なお、図8において、点線により囲まれる枠は、サブブロックS Bを示し、数字は、転送順位を示す。

【 0 1 8 9 】

図8は、割当属性とリード属性が「縦優先」と「T O P L E F T」に夫々指定されて

いる場合の R 1 型 ~ R 3 型依存関係を示す。

【 0 1 9 0 】

図示のように、この場合、「R 1 型依存」は、注目サブブロック S B 視点から見て上のサブブロック S B と依存関係を有する場合に指定される。「R 2 型依存」は、注目サブブロック S B から見て上、斜め左上、左のサブブロック S B と依存関係を有する場合に指定される。また、「R 3 型依存」は、注目サブブロック S B から見て、隣接する 8 つのサブブロック S B のうちの、下、斜め右下の 2 つを除いた 6 つとデータ依存関係を有する場合に指定される。

【 0 1 9 1 】

図 2 に例示したデータブロックの依存属性を説明する。なお、依存属性の指定も、他の属性と関係するため、依存属性のみでは説明しにくい。そのため、ここでは、図 2 に示す各データブロックに対して指定された依存属性の値のみを示し、それらの意義の詳細については後述する。

10

【 0 1 9 2 】

該例において、リードブロック P、リードブロック Q、ライトブロック R のいずれも、依存属性が「NONE」に指定される。

【 0 1 9 3 】

以上、本実施の形態の OpenCL システム 1 0 0 でサポートされる 8 つの属性を説明した。これらの 8 つの属性は、いずれもデバイス 1 2 0 における演算ユニット 1 4 0 の構成に依存しない。後の説明時に分かりやすいように、図 2 に示す 3 つのデータブロックに対して夫々設定された属性群を図 9 にまとめて示す。

20

【 0 1 9 4 】

本実施の形態の OpenCL システム 1 0 0 において、ホスト 1 1 0 からデバイス 1 2 0 に転送したカーネルの引数には、演算対象と演算結果の夫々のデータブロックに対して指定した上記の属性群が含まれる。デバイス 1 2 0 の演算制御部 2 0 0 における属性群格納部 2 1 0 は、これらの各属性群を格納して、シナリオ決定部 2 3 0 に供する。

【 0 1 9 5 】

シナリオ決定部 2 3 0 は、属性群格納部 2 1 0 に格納された各属性群に基づいて、後述する転送シナリオ（以下、単に「シナリオ」ともいう）を決定すると共に、決定されたシナリオに基づいて、演算ユニット 1 4 0 による演算、及び演算に伴うデータの転送を制御する。なお、シナリオ決定部 2 3 0 は、上記制御に際して、指示セットを演算ユニット 1 4 0 に送信することを繰り返す。指示セットは、リード転送指示 R、演算実行指示 S、ライト転送指示 W を含む。また、指示セットの 1 回の送信は、リード転送指示 R、演算実行指示 S、ライト転送指示 W の順にこれらの指示を送信することを意味する。

30

【 0 1 9 6 】

演算ユニット 1 4 0 において、演算制御部 2 0 0 からのリード転送指示 R に応じて、リード転送が行われる。

【 0 1 9 7 】

次いで、演算制御部 2 0 0 からの演算実行指示 S に応じて、演算ユニット 1 4 0 の P E 1 6 2 により、リード転送指示 R に応じたリード転送でプライベートメモリ 1 6 4 / ローカルメモリ 1 7 0 に格納されたデータに対して演算処理が行われる。演算結果となる各々のデータは、夫々の P E 1 6 2 の対応するプライベートメモリ 1 6 4 / ローカルメモリ 1 7 0 に格納される。

40

【 0 1 9 8 】

そして、演算制御部 2 0 0 からのライト転送指示 W に応じて、プライベートメモリ 1 6 4 / ローカルメモリ 1 7 0 に格納されているデータ（演算結果）は、ライト転送によりグローバルメモリ 1 5 2 へ転送される。

【 0 1 9 9 】

ここで、シナリオ決定部 2 3 0 によるシナリオの決定処理を説明する。

シナリオ決定部 2 3 0 は、属性群格納部 2 1 0 に格納された各属性群と、演算ユニット

50

140の構成を示すパラメータとに基づいて、シナリオを決定する。このシナリオは、転送方式に該当する。

【0200】

シナリオ決定部230の動作を説明する。説明に当たって、演算処理及びデータブロックについては、図2に示す例を使用する。なお、図2に示す各データブロックは、図9に示すように属性群が設定され、属性群格納部210に格納されたとする。また、演算ユニット140の構成を示すパラメータの例を、図10に示す。

【0201】

また、上記において各属性を説明する際に、分かりやすいように、1アイテムが1PEに対応するとした。以下の説明においては、1PEが1以上のアイテムに対応可能であるとする。

10

【0202】

シナリオ決定部230は、下記の規則に従って転送方式の決定と、決定した転送方式に応じた演算処理の制御を行う。

【0203】

<規則1>

シナリオ決定部230は、まず、全てのデータブロックに対して、共通のワークグループサイズWGsを設定すると共に、階層属性が同一である複数のデータブロックが共通の分割数でサブブロックに分割されるように、データブロックの分割サイズ、反復回数を決

20

【0204】

「分割サイズ」とは、サブブロックのサイズを意味し、「分割数」は、1つのデータブロックを分割して得たサブブロックの数を意味する。1となる分割数は、分割しないことを意味する。また、「反復回数」は、該データブロックの全てのサブブロックの転送に必要な転送回数を意味する。

【0205】

データブロックのサイズ、分割サイズ、分割数、ワークグループサイズWGs、反復回数などは、下記の式(2)~(5)に示す関係を有する。

【0206】

データブロックのサイズ = X方向サイズLx × Y方向サイズLy (2)

30

分割サイズ

= サブブロックのX方向のサイズSBsx × Y方向のサイズSBsy (3)

分割数 = データブロックのサイズ ÷ (分割サイズ × ワークグループサイズWGs) (4)

反復回数

= { 分割数 × (X方向サイズLx ÷ ワークグループサイズWGs) } / WG数 (5)

【0207】

シナリオ決定部230は、放送属性が「OFF」であるデータブロックのうちの、リード属性が「NONE」でないデータブロック(リードブロック)またはライト属性が「NONE」でないデータブロック(ライトブロック)のそれぞれに対して、分割サイズと、余白属性が指定されている場合の余白分の総量(余白属性が指定されていない場合には「0」との和が、個々のアイテムのプライベートメモリの利用可能容量の合計値を上回らない、かつ、同一の階層属性を持つデータブロックが同一の分割数に分割される制約を満たすように、ワークグループサイズWGsと分割サイズ(正確にはサブブロックのY方向サイズSBsy)を決定する。なお、上記制約を満たす前提の下、ターゲットOpenCLデバイスの推奨ワークグループサイズWGsがあれば、該ワークグループサイズWGsを採用し、ワークグループサイズWGsの上限と下限が定められていれば、上限と下限により決められる範囲内にワークグループサイズWGsを決定する。

40

【0208】

<規則2>

階層属性が同一であるリードブロックの相対応するサブリードブロックSRBを含むサ

50

ブリードブロックSRB群毎を、同時にプライベートメモリまたはローカルメモリ空間へ転送すると共に、ユーザ指定処理を起動する。なお、転送されるサブリードブロックSRBに対して余白属性により余白が指定される場合に、該余白分のデータも転送する。

【0209】

<規則3>

階層属性が異なるデータブロックの分割数の掛け算となる回数だけ、相対応するサブリードブロックSRBのサブリードブロックSRB群をプライベートメモリまたはローカルメモリ区間に転送した後、ユーザ指定処理を起動する。

【0210】

例えば、階層属性が「1」であるリードブロックの分割数をN、階層属性が「2」であるリードブロックの分割数をMとすると、シナリオ決定部230は、ユーザ指定処理をN×M回だけ呼び出すように動作する。また、各回の呼出しに先立ち、サブリードブロックSRBの1種類の組合せをプライベートメモリまたはローカルメモリ区間に転送する。

【0211】

シナリオ決定部230は、サブブロックの対応付け方法(上述した転送方式における分配方式及び統合方式)を決定してから転送を行う。図12に示すように、同一の階層のデータブロック同士のサブブロック間と、異なる階層のデータブロック同士のサブブロック間とは対応付け方法が異なる。

<規則4>

ユーザ指定処理の各回の起動後に、演算処理の結果となるサブライトブロックSWBをグローバルメモリ区間に転送する。サブライトブロックSWBの転送は、ユーザ指定処理の起動後に行われ、転送方向が「プライベートメモリ及び/またはローカルメモリ空間からグローバルメモリ空間へ」である点を除き、サブリードブロックSRBのときと同様である。

【0212】

図11は、図9に示すリードブロックP、リードブロックQ、ライトブロックRに対して指定された属性群と、図10に示す演算ユニット140の構成を示すパラメータとに基づいて、シナリオ決定部230により決定されたワークグループサイズWGs、縦分割サイズ(サブブロックのY方向のサイズSBy)、反復回数の例を示す。

【0213】

シナリオ決定部230は、まず、演算ユニット140のワークグループサイズWGs(32)を仮決定する。階層属性が「1」である各データブロックの分割数を4とすると、リードブロックQによるプライベートメモリ占有量は、アイテム毎に、「129×2B」の0.258KBとなる。なお、129は、「リードブロックQのY方向サイズLy(480)/分割数(4)」(縦分割サイズ)に、余白属性により指定されたY方向の9を加算して求められた値である。同様に、ライトブロックRによるプライベートメモリ占有量は、ワークアイテムWI毎に、「(480/4)×4B」の0.738KBに求められる。

【0214】

リードブロックQとライトブロックRによるアイテム毎のプライベートメモリ占有量の和が、演算ユニット140の構成を示すパラメータ(図10)におけるアイテム毎のプライベートメモリ容量(1KB)より小さいので、規則1が満たされているとして、リードブロックQとライトブロックRのワークグループサイズWGs、分割サイズが決まる。また、反復回数は、上記式(6)に従って「5」と算出される。

【0215】

また、階層属性が「2」のリードブロックPについては、放送属性が「ON」であるため、ローカルメモリへ割り当てられる。また、そのサイズ(3×3=9)は、演算ユニット140のワークグループWG毎のローカルメモリ容量4KBより小さいので、分割無し(反復回数:1)でローカルメモリに転送される。

【0216】

10

20

30

40

50

最後に、依存属性が「依存有り」のデータブロックが存在する場合のシナリオ決定部 230 の動作をまとめる。以下の説明が分かりやすいように、ここで、各データブロックを夫々の分割サイズに分割した上で対応付けされたサブブロック群のことを、「サブブロック集合」と呼ぶ。

【0217】

シナリオ決定部 230 は、サブブロック集合間でサブブロック同士が依存関係にある場合は、依存元のサブブロック集合に対する処理が行われてから、依存先のサブブロック集合を処理対象とするように、反復の順序を制御する。ここで、サブブロック集合 M0 内のある一つのサブブロック MX が、サブブロック集合 M1 内のある一つのサブブロック MY と依存関係にあり、かつ MX が定義元であるとする、サブブロック集合 M0 と M1 は依存関係にあると呼び、M0 が M1 の依存元、M1 が M0 の依存先と呼ぶ。

10

【0218】

このようにして、OpenCL システム 100 は、単独カーネルについての処理を行う。なお、単独カーネルが複数続いた場合には、OpenCL システム 100 は、上述した処理をカーネル毎に繰り返す。

【0219】

< 関連性のある 2 つのカーネルの場合 >

この 2 つのカーネルを夫々第 1 のカーネル X と第 2 のカーネル Y とする。該 2 つのカーネルに関連性があるとは、先に実行される第 1 のカーネル X の演算結果となるライトブロックのうちに継続ライトブロックがあり、該継続ライトブロックが、後続の第 2 のカーネル Y の演算対象となるリードブロックとして使用されることを意味する。すなわち、第 2 のカーネル Y のリードブロックのうちに、継続リードブロックが含まれている。

20

【0220】

図 13 は、この 2 つのカーネルのリードブロックとライトブロックを示す。これらのデータブロックのうちに、データブロック R は継続ライトブロックであり、データブロック S は継続ライトブロック R に対応する継続リードブロックである。継続ライトブロック R と継続リードブロック S は、同様のデータブロックであるが、第 1 のカーネル X のライトブロックと第 2 のカーネル Y のリードブロックとして異なる属性群が設定されているため、分かりやすいように、別々のデータブロックとして図示すると共に、異なる符号を付与している。

30

【0221】

演算制御部 200 の属性群格納部 210 には、図 13 に示す各データブロックの夫々の属性群が格納されている。これらの属性群は、単独カーネルの場合に説明した属性群と同様である。

【0222】

また、演算制御部 200 のカーネル間情報格納部 220 には、カーネル間データブロック依存情報として、第 1 のカーネル X の後に第 2 のカーネル Y が実行されること、第 1 のカーネル X のライトブロック R が、第 2 のカーネル Y のリードブロック S として使用されることを示す情報が格納される。

【0223】

例として、第 1 のカーネル X が、図 2 に示す例の演算を行うためのカーネルであるとする。そのため、第 1 のカーネル X のリードブロックは、データブロック P 及び Q であり、第 1 のカーネル X のライトブロックは、データブロック R であり、該データブロック R は、後続の第 2 のカーネル Y に使用される継続ライトブロックである。

40

【0224】

また、独立カーネルのときに説明した通りに、リードブロック P の属性群のうちに、割当属性が「縦優先」であり、余白属性が「X : 0 , Y : 0」(余白なし)であり、依存属性が「NONE」(サブリードブロック間に依存性なし)である。

【0225】

また、第 2 のカーネル Y が、第 1 のカーネル X のライトブロック R と同様のデータブロ

50

ックであるデータブロックSを演算して、データブロックUとVを得るものとする。そのため、第2のカーネルYのリードブロックは、データブロックSであり、該データブロックSは、継続リードブロックである。また、第2のカーネルYのライトブロックは、データブロックU及びVである。

【0226】

図示の例では、継続リードブロックSの属性群のうちに、割当属性が「縦優先」であり、余白属性が「X:1, Y:2」であり、依存属性が「R2型」（左のサブブロックと、左上のサブブロックと、上のサブブロックに依存する）である。

【0227】

すなわち、継続リードブロックSをサブブロックに分割して転送する際に、該サブブロックと共に、指定されたその左隣の1列のデータと右隣の1列のデータ、及び上隣の2行のデータと下隣の2行のデータも転送するよう、余白属性により規定されている。また、継続リードブロックSの各サブブロックの転送及び演算の順序については、ターゲットのサブブロックの処理の前に、その依存元の各サブブロック（ここでは、ターゲットのサブブロックの左と、左上と、上との3つのサブブロック）の処理が行われるよう、依存属性により規定されている。

10

【0228】

第2のカーネルYのライトブロックUとVについても属性群が設定されており、それらは、第2のカーネルYの演算内容などに応じてユーザにより設定されたものであり、ここで詳細な説明と図示を省略する。

20

【0229】

本実施の形態において、演算制御部200のシナリオ決定部230は、継続ライトブロックRの割当属性と、継続リードブロックSの割当属性とに基づいて、パイプライン化制御を行うか否かを決定する。具体的には、継続ライトブロックRと継続リードブロックSの割当属性が同一であれば、パイプライン化制御を行うと決定し、継続ライトブロックRと継続リードブロックSの割当属性が異なれば、パイプライン化制御を行わないと決定する。

【0230】

「パイプライン化制御」とは、継続ライトブロックRがグローバルメモリ152に転送されず、プライベートメモリ164またはローカルメモリ170を介して第2のカーネルYの実行に供されることにより、第1のカーネルXと第2のカーネルYの実行をパイプライン化することである。

30

【0231】

パイプライン化制御を行わないと決定した場合に、シナリオ決定部230は、第1のカーネルXと第2のカーネルYを夫々単独カーネルとして、まず第1のカーネルXを実行させ、ライトブロックRの全てのデータが得られてグローバルメモリ152に格納された後に、該ライトブロックRを第2のカーネルYのリードブロックSとして、第2のカーネルYを実行させる。

【0232】

パイプライン化制御を行うと決定した場合に、シナリオ決定部230は、パイプライン化制御を実現するために、先に実行される第1のカーネルXの各データブロックの余白属性と依存属性を変更する。

40

【0233】

具体的には、シナリオ決定部230は、第1のカーネルXのリードブロックPに対して、その余白属性と依存属性に、継続リードブロックSの余白属性と依存属性を夫々論理和加算して新たな余白属性と依存属性とする。

【0234】

また、リードブロックQに対しても同様に、シナリオ決定部230は、その余白属性と依存属性に、継続リードブロックSの余白属性と依存属性を夫々論理和加算して新たな余白属性と依存属性とする。

50

【 0 2 3 5 】

図 1 4 は、シナリオ決定部 2 3 0 により変更後のリードブロック P とリードブロック Q の属性群を示す。図示のように、変更後、リードブロック P は、余白属性が元の「 X : 0 , Y 0 」から「 X : 1 , Y : 2 」になり、依存属性が元の「 NONE 」から「 R 2 型 」になっている。

【 0 2 3 6 】

また、リードブロック Q も、上記変更後、余白属性が元の「 X : 0 , Y 9 」から「 X : 1 , Y : 1 1 」になり、依存属性が元の「 NONE 」から「 R 2 型 」になっている。

【 0 2 3 7 】

そして、シナリオ決定部 2 3 0 は、第 1 のカーネル X の各ライトブロックについては余白属性と依存属性が変更された属性群を使用し、他の各データブロックについては元の属性群を使用して、各データブロックの転送方式を決定すると共に、決定した転送方式に応じて各リードブロックとライトブロックの転送、及び該転送に対応する並列演算の制御を行う。この際、継続ライトブロック R の各サブブロックをグローバルメモリ 1 5 2 に転送せず、プライベートメモリ 1 6 4 またはローカルメモリ 1 7 0 を介して第 2 のカーネル Y の実行に供する点と、第 1 のカーネル X と第 2 のカーネル Y の実行を並列させる点とを除き、単独カーネルの場合と同様である。

【 0 2 3 8 】

図 1 5 の上部は、第 1 のカーネル X と第 2 のカーネル Y の実行に際して、パイプライン化制御が行われない場合における各データブロックの転送態様を示し、図 1 5 の下部は、第 1 のカーネル X と第 2 のカーネル Y の実行に際して、パイプライン化制御が行われた場合における各データブロックの転送態様を示す。図中実線矢印は、グローバルメモリ 1 5 2 を介したデータの受渡しを意味し、点線矢印は、プライベートメモリ 1 6 4 またはローカルメモリ 1 7 0 を介したデータの受渡しを意味する。

【 0 2 3 9 】

図 1 5 に示すように、パイプライン化制御が行われない場合において、第 1 のカーネル X の演算結果となるライトブロック R (継続ライトブロック) は、グローバルメモリ 1 5 2 に格納された後に、第 2 のカーネル Y のリードブロック S (継続リードブロック) としてグローバルメモリ 1 5 2 からプライベートメモリ 1 6 4 またはローカルメモリ 1 7 0 に転送され、第 2 のカーネル Y の実行に供する。

【 0 2 4 0 】

一方、パイプライン化制御が行われた場合には、継続ライトブロック R は、グローバルメモリ 1 5 2 に転送されず、プライベートメモリ 1 6 4 またはローカルメモリ 1 7 0 を介して第 2 のカーネル Y の実行に供されるようになっている。

【 0 2 4 1 】

なお、プライベートメモリ 1 6 4 またはローカルメモリ 1 7 0 を介したデータの受渡しについては、例えば、該データを生成する処理と、該データを用いて演算を行う処理とを同一の P E に割り当て、該 P E は、前の演算により生成してプライベートメモリ 1 6 4 またはローカルメモリ 1 7 0 に格納したデータをそのまま読み出して後の演算に使えばよい。

【 0 2 4 2 】

また、OpenCL デバイスの各 P E がリング状に接続され、各 P E が該リング上の隣接する P E とデータの受渡しが可能である場合には、前の演算を行った P E は、次の演算を行う P E との間の各 P E を介して、前の演算で生成したデータを渡せばよい。

【 0 2 4 3 】

図 1 6 の上部は、パイプライン化制御が行われない場合における第 1 のカーネル X と第 2 のカーネル Y の実行態様を示し、図 1 6 の下部は、パイプライン化制御が行われた場合における第 1 のカーネル X と第 2 のカーネル Y の実行態様を示す。

【 0 2 4 4 】

図 1 6 において、斜線で塗り潰された枠は、第 1 のカーネル X のためのリード転送の期

10

20

30

40

50

間を示し、点で塗り潰された枠は、第1のカーネルXのためのライト転送の期間を示す。また、縦線で塗り潰された枠は、第2のカーネルYのためのリード転送の期間を示し、網状に塗り潰された枠は、第2のカーネルYのためのライト転送の期間を示す。

【0245】

また、白い枠は、第1のカーネルXの実行期間を示し、波線で塗り潰された枠は、第2のカーネルYの実行期間を示す。

【0246】

図16に示すように、パイプライン化制御が行われない場合において、第1のカーネルXと第2のカーネルYは、夫々独立カーネルとして実行され、実行される度に、サブリードブロックのリード転送（グローバルメモリからプライベートメモリ/ローカルメモリへの転送）と、サブライトブロックのライト転送（プライベートメモリ/ローカルメモリからグローバルメモリへの転送）が生じている。

10

【0247】

一方、パイプライン化制御が行われた場合には、第1のカーネルXの各回の実行により生成されたサブライトブロックのライト転送が行われず、第2のカーネルYの各回の実行の演算対象となるサブリードブロック（第1のカーネルXにより生成されたサブライトブロック）のリード転送が行われない。

【0248】

さらに、第2のカーネルYの実行の演算対象となるサブリードブロックのリード転送が行われないため、第2のカーネルYの実行に並行して、第1のカーネルXの次の実行のためのリード転送を行うことができる。

20

【0249】

このように、本実施の形態のOpenCLシステム100におけるデバイス120の演算制御部200は、第1のカーネルXと第2のカーネルYを連続して実行させる際に、第1のカーネルXの各ライトブロックのうちに、第2のカーネルYリードブロックとして使用される継続ライトブロックが含まれており、かつ、上記継続ライトブロックの割当属性と、第2のカーネルYの、該継続ライトブロックに対応するリードブロック（継続リードブロック）割当属性とが一致する場合において、継続ライトブロックのサブブロックがそのまま継続リードブロックのサブブロックとして用いることができるため、パイプライン化制御を行うことにより、第1のカーネルXと第2のカーネルYの実行のパイプライン化を実現している。

30

【0250】

具体的には、シナリオ決定部230は、シナリオ決定部230は、パイプライン化制御に際して、まず、第1のカーネルXの各リードブロックの余白属性に、第2のカーネルYの継続リードブロックの余白属性を論理和加算して、第1のカーネルXの各リードブロックの余白属性を変更すると共に、第1のカーネルXの各リードブロックの依存属性に、第2のカーネルYの継続リードブロックの依存属性を論理和加算することにより、第1のカーネルXの各リードブロックの属性群を変更する。

【0251】

そして、第1のカーネルXと第2のカーネルYの各データブロックの属性群（第1のカーネルXのリードブロックについては、変更後の属性群）に基づいて、各データブロックの転送方式を決定し、決定した転送方式に応じて各データブロックの転送、及び該転送に対応する並列演算の制御を行う。この際、継続ライトブロックの各サブブロックをグローバルメモリに転送せず、プライベートメモリまたはローカルメモリを介して第2のカーネルYの実行に供すると共に、第1のカーネルXと第2のカーネルYの実行を並列させることにより、第1のカーネルXと第2のカーネルYのパイプライン化を実現する。

40

【0252】

図16に示すように、シナリオ決定部230のパイプライン化制御の結果、第1のカーネルXと第2のカーネルYの実行完了までの時間は、パイプライン化制御が行われない場合より「T」の分短縮されており、OpenCLデバイスの処理が高速になっている。

50

【0253】

なお、第1のカーネルXと第2のカーネルYの実行をパイプライン化するために、第1のカーネルXの各リードブロックの余白属性に、第2のカーネルYの継続リードブロックの余白属性を論理和加算しているため、第1のカーネルXのサブリードブロックのサイズと、該サブリードブロックと共に転送されるデータのサイズとの総和が、プライベートメモリまたはローカルメモリの容量を超える可能性がある。

【0254】

これを回避するために、シナリオ決定部230は、第1のカーネルXの継続ライトブロックの割当属性と、第2のカーネルYの継続リードブロックの割当属性とが一致し、かつ、変更後の余白属性に基づいて決定した第1のカーネルXのサブリードブロックの転送時のサイズが予め設定された閾値以下であることを条件に、パイプライン化制御を行うと決定するようにすればよい。なお、この閾値は、プライベートメモリ及び/またはローカルメモリのサイズに応じて予め定めればよい。

10

【0255】

また、分かりやすいように、上記において、関連性のある2つのカーネルのパイプライン化について説明したが、同様の原理で、3以上の関連性のあるカーネルのパイプライン化もできる。

【0256】

すなわち、連続して実行されるN個(N:2以上の正数)のカーネルがあり、実行順が2番目以降の各カーネルは、そのリードブロックに、直前に実行されるカーネルのライトブロックが含まれているときに、シナリオ決定部230は、各K番目のカーネル(1$\leq K \leq N-1$)の継続ライトブロックの割当属性と、該継続ライトブロックに対応する、(K+1)番目のカーネルの継続リードブロックの割当属性とが一致する場合にのみ、パイプライン化制御を行うと決定する。

20

【0257】

パイプライン化制御を行うと決定した場合に、シナリオ決定部230は、1番目のカーネルの各リードブロックの余白属性に、2番目以降の各カーネルの継続リードブロックの余白属性を論理和加算すると共に、1番目のカーネルの各リードブロックの依存属性に、2番目以降の各カーネルの継続リードブロックの依存属性を論理和加算する。

【0258】

その後、シナリオ決定部230は、各データブロックの転送方式の決定、及び決定した転送方式に応じて各データブロックの転送、及び該転送に対応する並列演算の制御に際して、継続ライトブロックの各サブブロックをグローバルメモリに転送せず、プライベートメモリまたはローカルメモリを介して次のカーネルの実行に供すると共に、各カーネルの実行を並列させることにより、これらのカーネルのパイプライン化を実現する。

30

【0259】

パイプライン化制御の対象となるカーネルの数が多いほど、1番目のカーネルのサブリードブロックの転送時のサイズが大きくなる。1番目のカーネルのサブリードブロックの転送時のサイズがプライベートメモリまたはローカルメモリの容量を超えることを避けるために、シナリオ決定部230は、上記N個のカーネルを、連続して実行される1つ以上のカーネルにより構成された複数のグループに分け、グループ毎にパイプライン化制御を行うようにすればよい。グループ分けに際しては、変更後の余白属性に基づいて決定した、該グループ内に最も先に実行されるカーネルのサブリードブロックの転送時のサイズが予め設定された閾値以下になるようにすればよい。

40

【0260】

以上、本発明者によってなされた発明を実施の形態に基づき具体的に説明したが、本発明は既に述べた実施の形態に限定されるものではなく、その要旨を逸脱しない範囲において種々の変更が可能であることはいうまでもない。

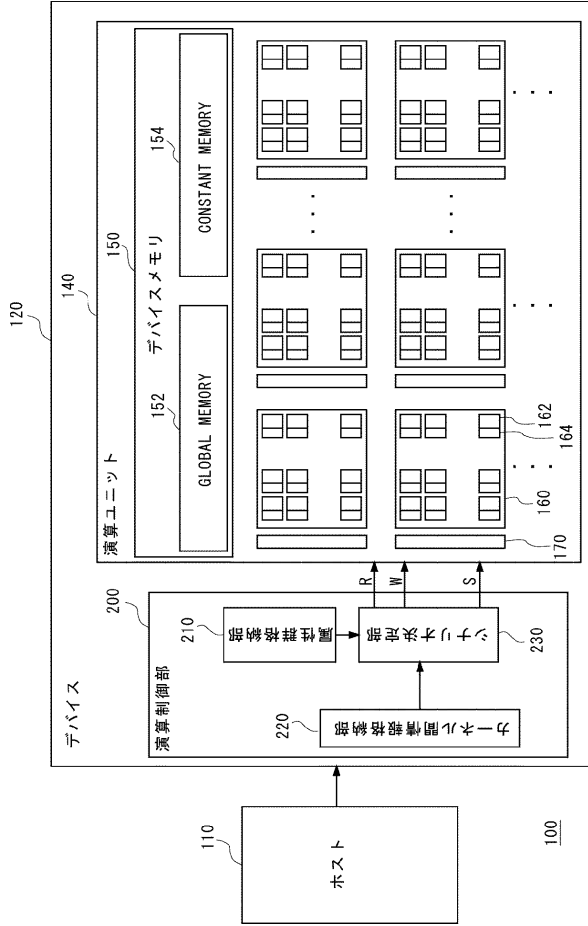
【符号の説明】

【0261】

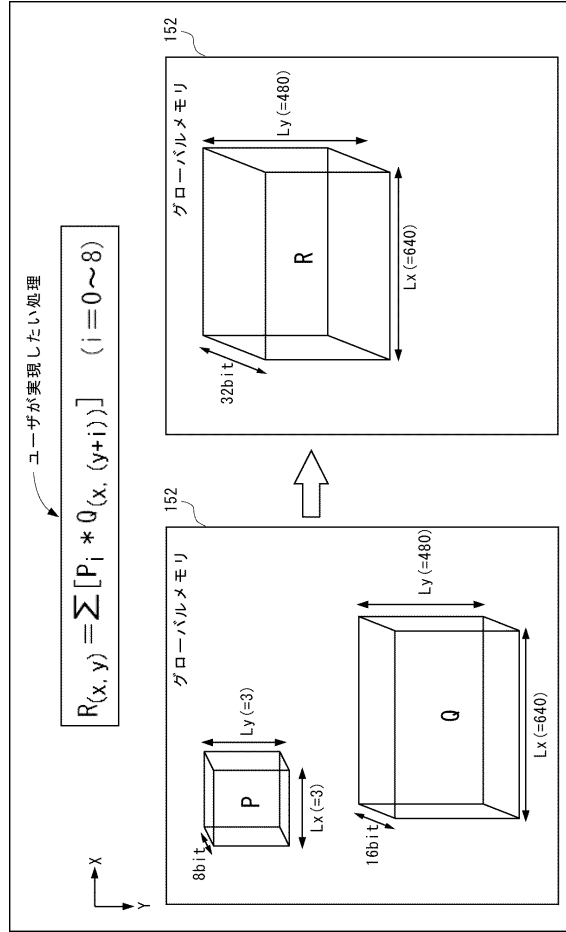
50

1 0	O p e n C L システム	
1 2	ホスト	
1 4	O p e n C L デバイス	
1 6	C U	
1 8	P E	
2 0	プライベートメモリ	
2 2	ローカルメモリ	
2 4	キャッシュ	
3 0	デバイスメモリ	
3 2	グローバルメモリ	10
3 4	コンスタントメモリ	
5 0	逐次プロセッサ	
5 2	P E	
5 4	プライベートメモリ	
5 6	グローバルメモリ	
5 8	キャッシュ制御機構	
1 0 0	O p e n C L システム	
1 1 0	ホスト	
1 2 0	デバイス	
1 4 0	演算ユニット	20
1 5 0	デバイスメモリ	
1 5 2	グローバルメモリ	
1 5 4	コンスタントメモリ	
1 6 0	C U	
1 6 2	P E	
1 6 4	プライベートメモリ	
1 7 0	ローカルメモリ	
2 0 0	演算制御部	
2 1 0	属性群格納部	
2 2 0	カーネル間情報格納部	30
2 3 0	シナリオ決定部	
L x	データブロックのX方向サイズ	
L y	データブロックのY方向サイズ	
S B s x	サブブロックのX方向サイズ	
S B s y	サブブロックのY方向サイズ	
S R B	サブリードブロック	
S W B	サブライトブロック	
W G	ワークグループ	
W G s	ワークグループサイズ	
W I	ワークアイテム	40
R	リード転送指示	
W	ライト転送指示	
S	演算実行指示	

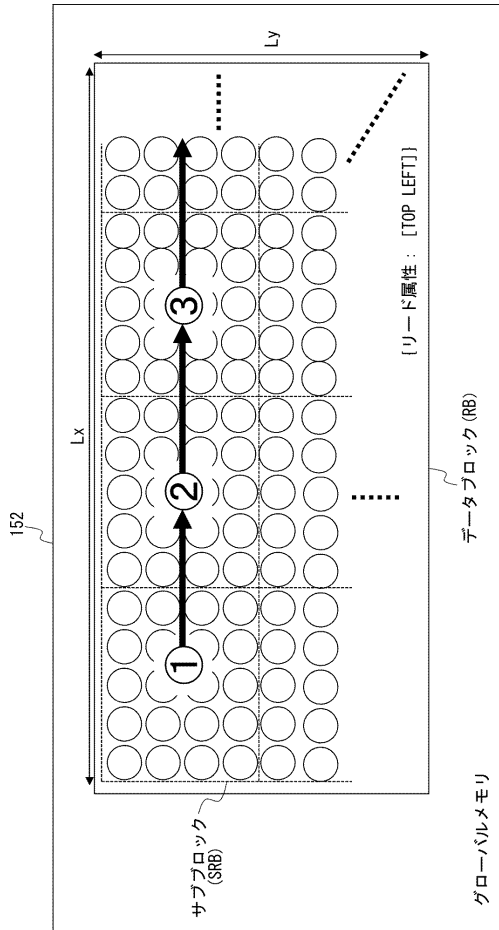
【図1】



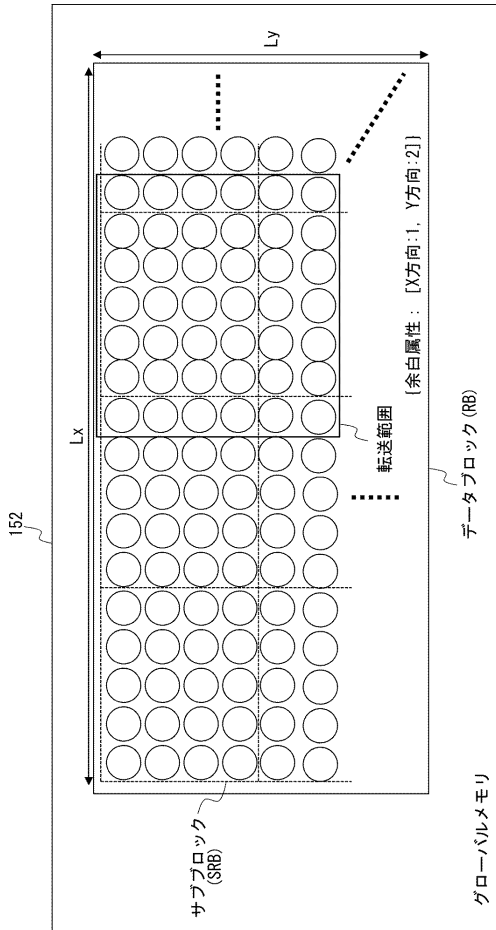
【図2】



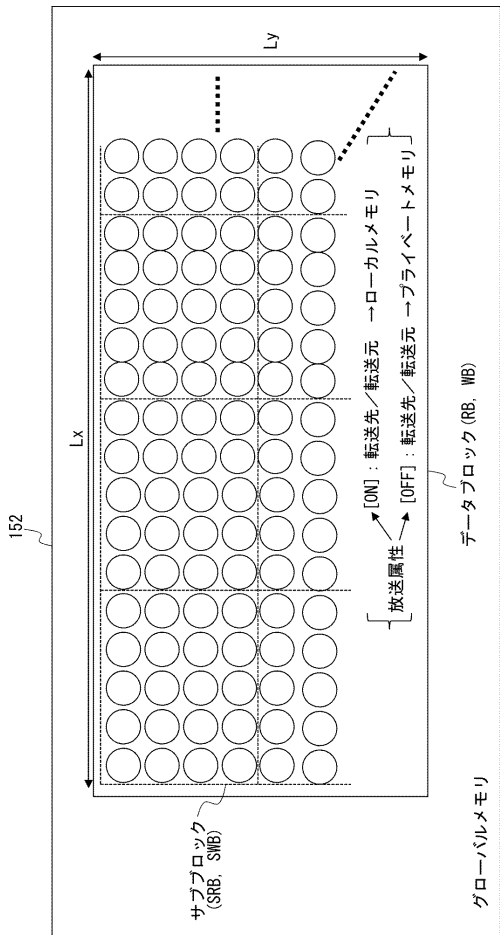
【図3】



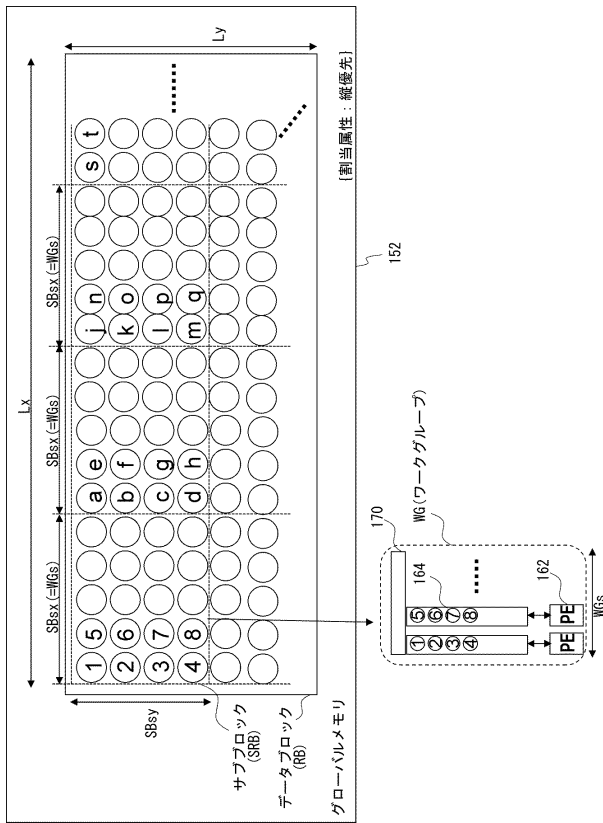
【図4】



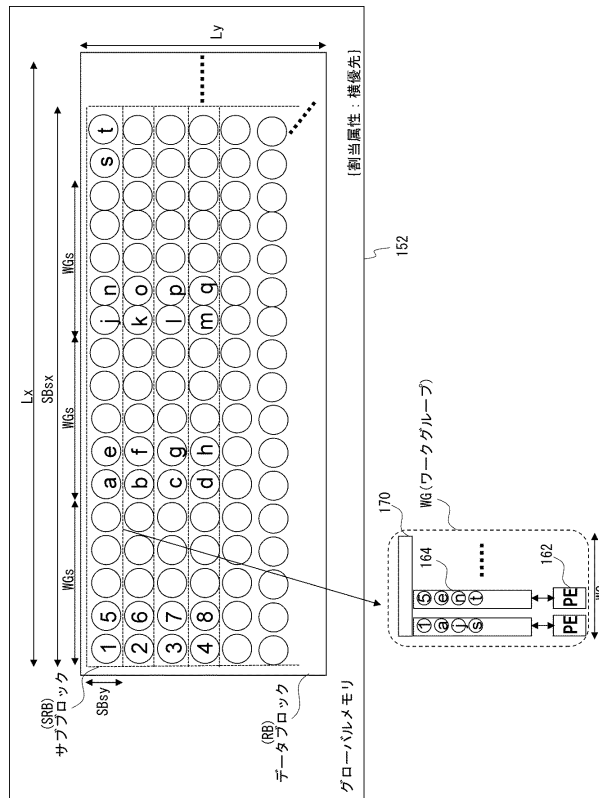
【図 5】



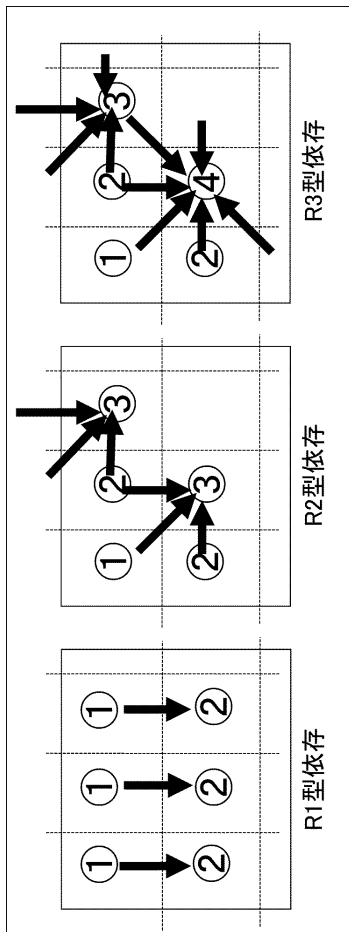
【図 6】



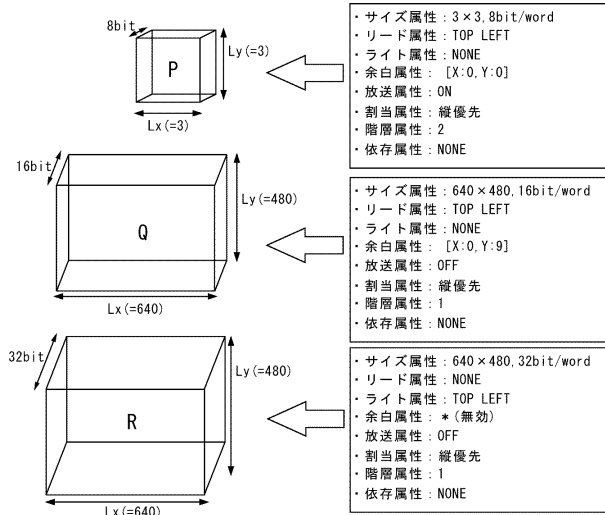
【図 7】



【図 8】



【図9】



【図10】

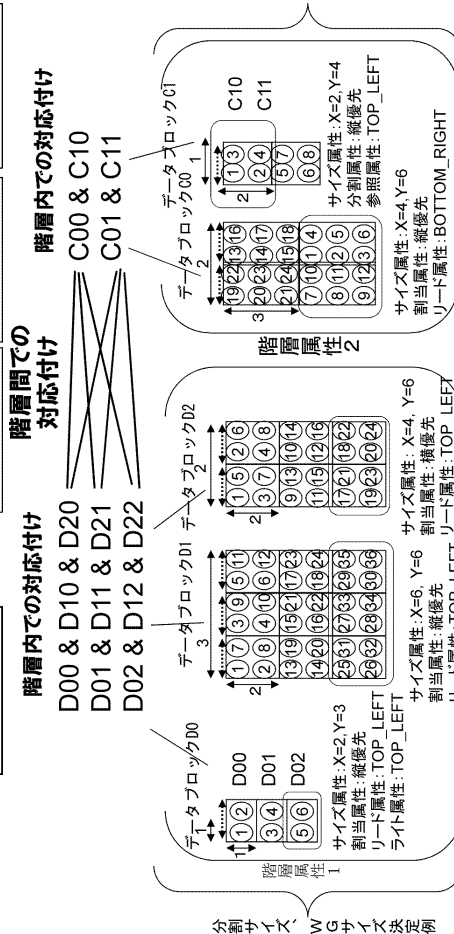
演算ユニットの構成を示すパラメータの例

- ワークアイテム(WI)毎のプライベートメモリ容量: 1KB
- ワークグループ(WG)毎のローカルメモリ容量: 4KB
- WGサイズ(WGs)の最大値: 32
- WG数の最大値: 16

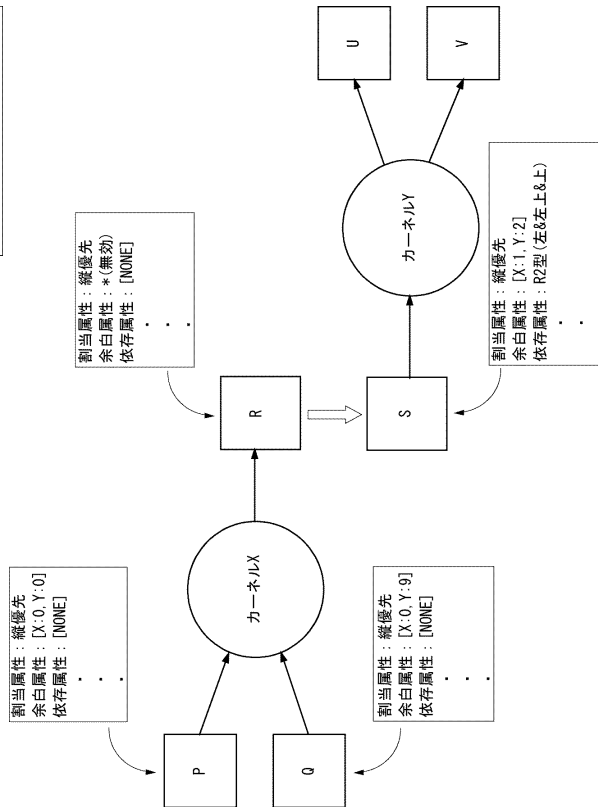
【図12】

- ワークグループサイズWG_s=32
- <階層1>
 - 分割数が4である場合の各データブロックのWI毎のプライベートメモリ占有量(サブブロックのY方向サイズSbsy)
 - リードブロックQ: WI毎129×2B (=0.258KB/WI)
 - ライトブロックR: WI毎120×4B (=0.480KB/WI)
 - 反復回数: 5(=分割数4×(640÷WGサイズ32)÷WG数16)
- <階層2>
 - 分割数: 1
 - リードブロックP
 - 「分割無し」でローカルメモリへ
 - 反復回数: 1

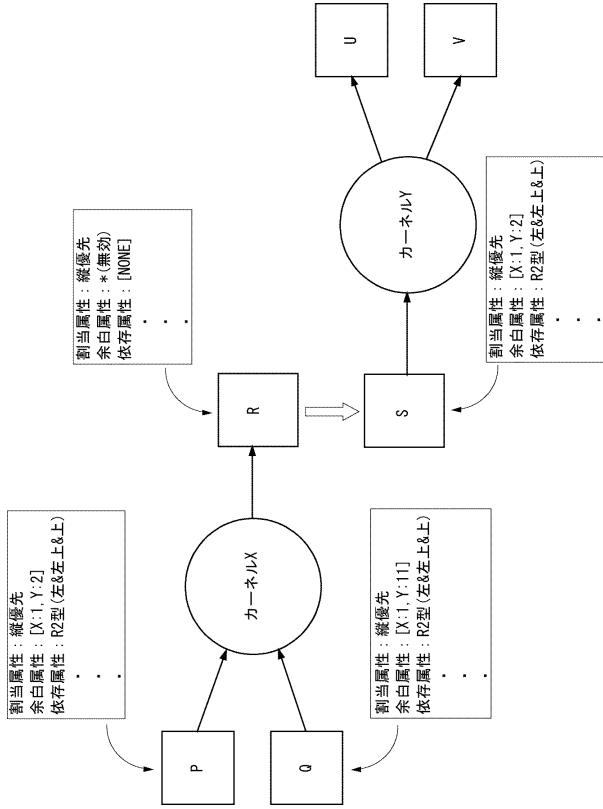
【図11】



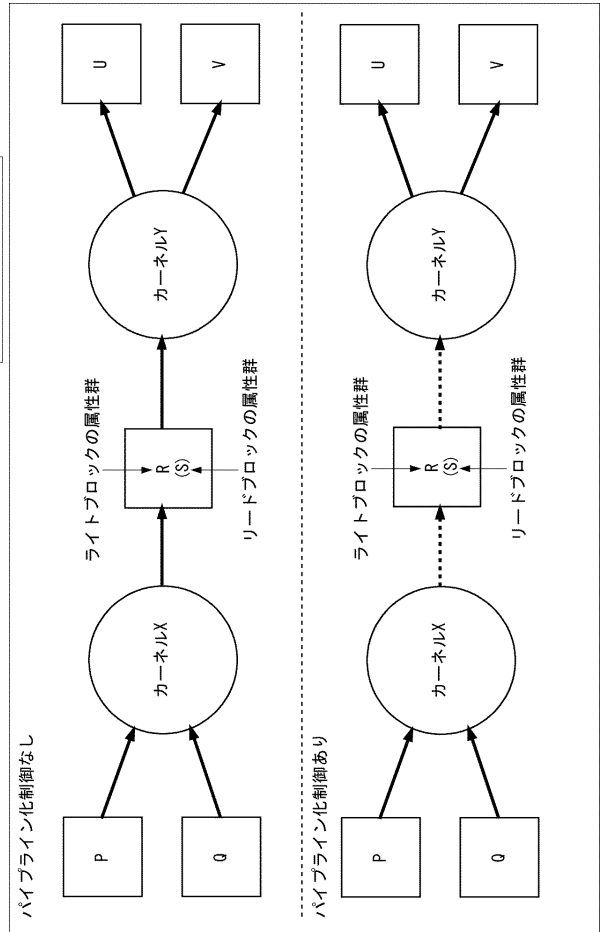
【図13】



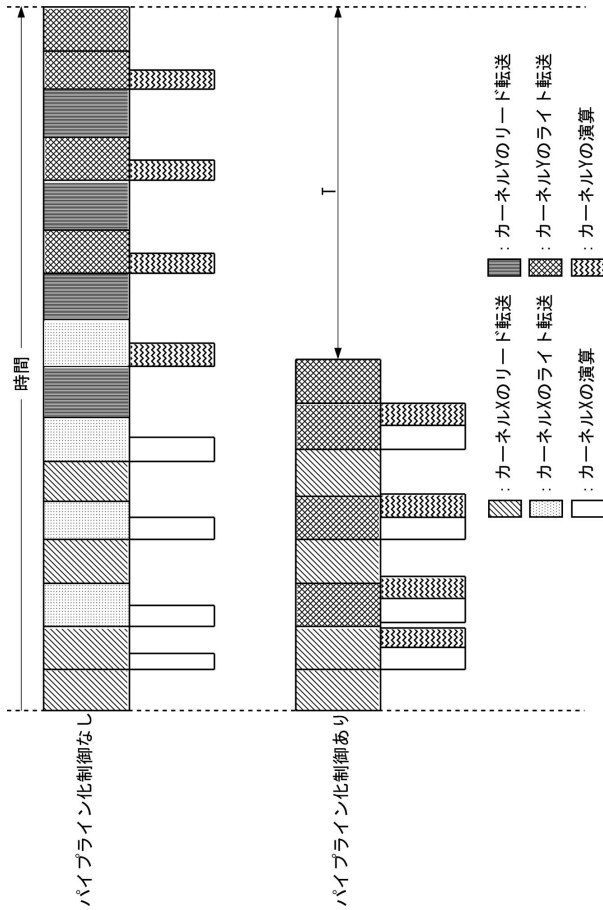
【図14】



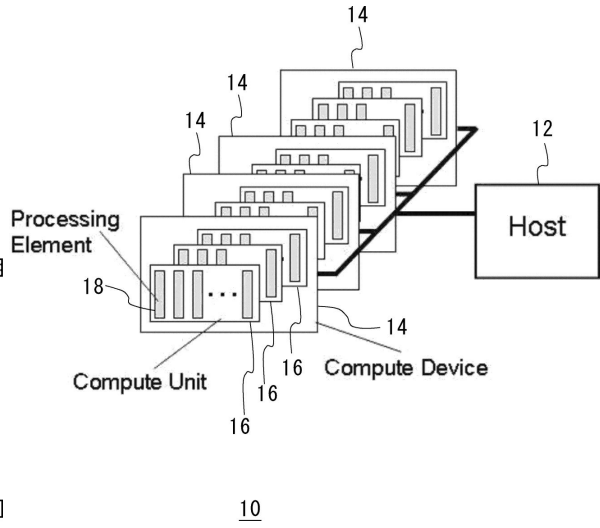
【図15】



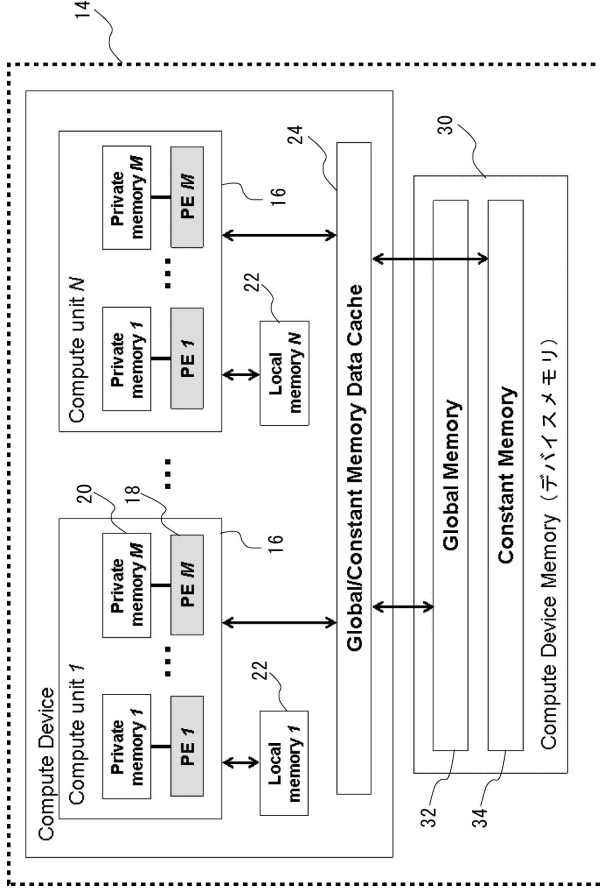
【図16】



【図17】



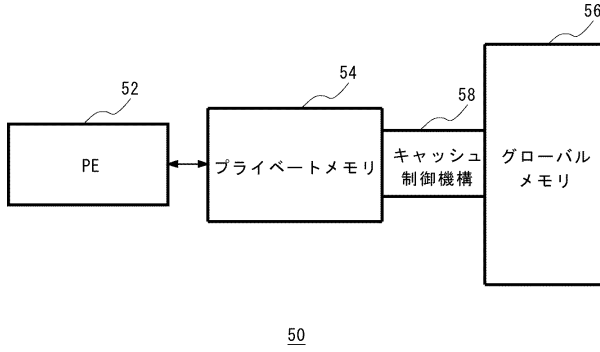
【図18】



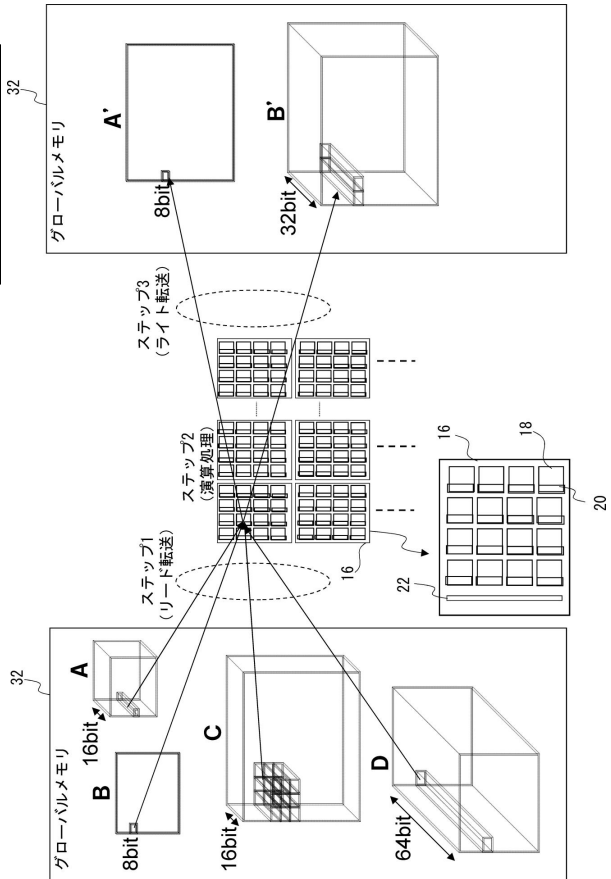
【図19】

	Global	Constant	Local	Private
Host	Dynamic allocation Read / Write access	Dynamic allocation Read / Write access	Dynamic allocation No access	No allocation No access
Kernel	No allocation Read / Write access	Static allocation Read-only access	Static allocation Read / Write access	Static allocation Read / Write access

【図20】



【図21】



【図 2 2】

1. リード転送用の指定

演算処理の内容とデバイスの構成に依存する部分（下記例を参照）を含む。

- ・例1：リードブロックを分割するか否かの指定
- ・例2：リードブロックを分割する場合に、その分割方法の指定
- ・例3：リードブロックを分割する場合に、異なるリードブロックに夫々含まれるサブリードブロック間の対応付け方法の指定

2. リードブロックまたはサブブロックに対する演算処理の指定

デバイスの構成に依存する部分（下記例を参照）を含む。

- ・例4：リード転送用の指定に合わせて、アイテムの並行実行の回数の指定

3. ライト転送用の指定

演算処理の内容とデバイスの構成に依存する部分（下記例を参照）を含む。

- ・例5：リード転送用の指定に合わせて指定

【図 2 3】

1. 属性群の指定

- (1) 固有属性：演算内容及びユーザの意志に関係しない属性
 - (2) 演算属性：演算内容に関係するものの、ユーザの意志に関係しない属性
 - (3) ポリシ属性：演算内容と、ユーザの意志に関係する
- いずれも、デバイスの構成と関係しない。

2. ユーザ処理の指定

演算内容の指定であり、デバイスの構成と関係しない。

フロントページの続き

(56)参考文献 特開2013-25547(JP,A)

特開2008-217134(JP,A)

佐藤 裕幸, グラフィックス処理用プロセッサGPUによるSAR画像再生処理の高速化, 電子情報通信学会技術研究報告, 日本, 社団法人電子情報通信学会, 2010年 4月15日, 第110巻, 第8号, pp.7-12

(58)調査した分野(Int.Cl., DB名)

G06F 15/80

G06F 9/50

G06F 12/08

G06F 12/0802