

(19) 日本国特許庁 (JP)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2019-75101

(P2019-75101A)

(43) 公開日 令和1年5月16日 (2019.5.16)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 12/00 (2006.01)	G06F 12/00 560F	5B033
G06F 9/305 (2006.01)	G06F 12/00 580	5B160
G06F 9/302 (2006.01)	G06F 9/305 A	
	G06F 9/302 A	

審査請求 未請求 請求項の数 25 O L (全 17 頁)

(21) 出願番号	特願2018-173507 (P2018-173507)	(71) 出願人	390019839
(22) 出願日	平成30年9月18日 (2018.9.18)		三星電子株式会社
(31) 優先権主張番号	62/573,390		Samsung Electronics Co., Ltd.
(32) 優先日	平成29年10月17日 (2017.10.17)		大韓民国京畿道水原市靈通区三星路129
(33) 優先権主張国	米国 (US)		129, Samsung-ro, Yeongtong-gu, Suwon-si, Gyeonggi-do, Republic of Korea
(31) 優先権主張番号	15/854,557	(74) 代理人	110000051
(32) 優先日	平成29年12月26日 (2017.12.26)		特許業務法人共生国際特許事務所
(33) 優先権主張国	米国 (US)	(72) 発明者	張 牧 天
			アメリカ合衆国, 95051 カリフォルニア州, サンタクララ, ピアトリノプレイス 2920

最終頁に続く

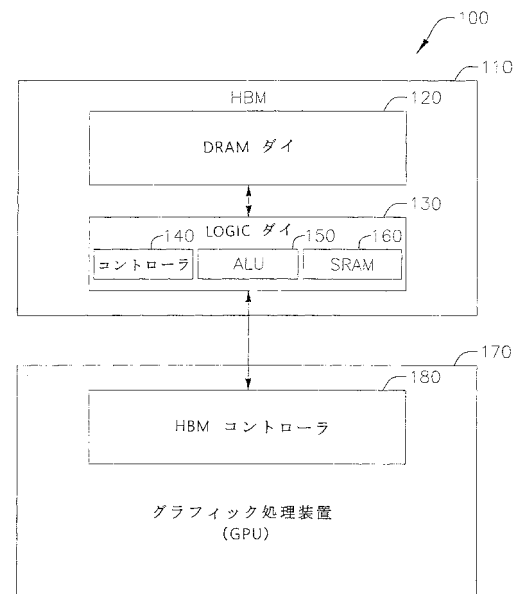
(54) 【発明の名称】 インメモリのコマンド処理方法、これを適用する高帯域幅メモリ (HBM)、及びHBMシステム

(57) 【要約】 (修正有)

【課題】 高帯域幅メモリシステムでインメモリのコマンドを調整するためのシステム、及び方法を提供する。

【解決手段】 高帯域幅メモリ (HBM) システム 100 において、グラフィック処理装置の HBM コントローラ 180 によってファンクション・イン・メモリ命令を HBM 110 に送信する。HBM のロジック部分 130 は、ファンクション・イン・メモリ命令を受信し、ロジック部分が有するコントローラ 140、ALU 150、及び SRAM 160 を用いて命令の実行を調整する。

【選択図】 図 1



【特許請求の範囲】**【請求項 1】**

H B M (H i g h - B a n d w i d t h M e m o r y) システムにおいてインメモリのコマンドを処理する方法であって、

グラフィック処理装置の H B M コントローラによって、ファンクション・イン・メモリ命令を H B M に送信する段階と、

コントローラ、A L U (A r i t h m e t i c L o g i c U n i t)、及び S R A M を含む前記 H B M のロジック部分によって、前記ファンクション・イン・メモリ命令を受信する段階と、

前記コントローラ、前記 A L U、又は前記 S R A M のうちの少なくとも 1 つを用いる前記ロジック部分によって、前記ファンクション・イン・メモリ命令に基づいて、前記ファンクション・イン・メモリ命令の実行を調整する段階と、を有することを特徴とする方法。

10

【請求項 2】

前記ファンクション・イン・メモリ命令の実行を調整する段階は、

前記コントローラによって、前記ファンクション・イン・メモリ命令を、演算及び少なくとも 1 つのデータ位置を含む計算ファンクション・イン・メモリ命令として識別する段階と、

前記コントローラによって、前記少なくとも 1 つのデータ位置に応じて、前記 H B M の D R A M から少なくとも 1 つのデータを検索する段階と、

20

前記コントローラによって、前記少なくとも 1 つのデータ及び前記演算を前記 A L U に提供する段階と、

前記 A L U によって、前記少なくとも 1 つのデータに対して前記演算を実行する段階と、

前記演算の実行結果を前記 D R A M に格納する段階と、を有することを特徴とする請求項 1 に記載の方法。

【請求項 3】

前記演算は、アトミック演算 (a t o m i c o p e r a t i o n) 及びデータ型を含み、

前記少なくとも 1 つのデータ位置は、宛先レジスタ、メモリアドレス、ソースレジスタ、定数、又は参照レジスタのうちの少なくとも 1 つを含むことを特徴とする請求項 2 に記載の方法。

30

【請求項 4】

前記アトミック演算は、A D D、S U B T R A C T、E X C H A N G E、M A X、M I N、I N C R E M E N T、D E C R E M E N T、C O M P A R E - A N D - S W A P、A N D、O R、X O R、又は N O T の関数のうちの少なくとも 1 つを含むことを特徴とする請求項 3 に記載の方法。

【請求項 5】

前記演算は、A L U 演算及び D R A M アクセス演算を含み、

前記少なくとも 1 つのデータ位置は、宛先レジスタ及び少なくとも 1 つのソースレジスタを含むことを特徴とする請求項 2 に記載の方法。

40

【請求項 6】

前記 D R A M アクセス演算は、前記 A L U 演算と対をなすロード命令又は格納命令を含むことを特徴とする請求項 5 に記載の方法。

【請求項 7】

前記ファンクション・イン・メモリ命令の実行を調整する段階は、

前記コントローラによって、前記ファンクション・イン・メモリ命令をソースレジスタ及び宛先レジスタを含むファンクション・イン・メモリ移動命令として識別する段階と、

前記コントローラによって、前記ソースレジスタに応じて、前記 H B M の D R A M から少なくとも 1 つのデータを検索する段階と、

50

前記コントローラによって、前記少なくとも1つのデータを前記宛先レジスタの前記 D R A M に格納する段階と、を有することを特徴とする請求項 1 に記載の方法。

【請求項 8】

前記ファンクション・イン・メモリ命令の実行を調整する段階は、

前記コントローラによって、前記ファンクション・イン・メモリ命令をソースレジスタ及び宛先レジスタのうちの少なくとも1つを含むファンクション・イン・メモリスクラッチパッド命令として識別する段階と、

前記コントローラのタイミングパラメータを、D R A M のタイミングパラメータから S R A M のタイミングパラメータに調整する段階と、

前記コントローラによって、前記 S R A M で前記 S R A M のタイミングパラメータに応じて、前記ファンクション・イン・メモリスクラッチパッド命令を実行する段階と、を有することを特徴とする請求項 1 に記載の方法。

10

【請求項 9】

前記少なくとも1つのデータ位置がグラフィック処理装置のキャッシュを含む場合、前記 H B M コントローラによって、前記グラフィック処理装置における前記ファンクション・イン・メモリ命令の実行を調整する段階を有することを特徴とする請求項 1 に記載の方法。

【請求項 10】

H B M (H i g h - B a n d w i d t h M e m o r y) であって、
D R A M と、

20

コントローラ、A L U (A r i t h m e t i c L o g i c U n i t)、及び S R A M を含んで命令を実行するロジック部分と、を備え、

前記命令は、前記ロジック部分によって実行されるとき、前記ロジック部分がファンクション・イン・メモリ命令に基づいて、前記 D R A M、前記コントローラ、前記 A L U、又は前記 S R A M のうちの少なくとも1つを用いることによって、前記ファンクション・イン・メモリ命令の実行を調整することを特徴とする H B M。

【請求項 11】

前記ファンクション・イン・メモリ命令の実行調整は、

前記コントローラによって、前記ファンクション・イン・メモリ命令を演算及び少なくとも1つのデータ位置を含む計算ファンクション・イン・メモリ命令として識別し、

30

前記コントローラによって、前記少なくとも1つのデータ位置に応じて前記 D R A M から少なくとも1つのデータを検索し、

前記コントローラによって、前記少なくとも1つのデータ及び前記演算を前記 A L U に提供し、

前記 A L U によって、前記少なくとも1つのデータに対して前記演算を実行し、

前記演算の実行結果を前記 D R A M に格納することを特徴とする請求項 10 に記載の H B M。

【請求項 12】

前記演算は、アトミック演算 (a t o m i c o p e r a t i o n) 及びデータ型を含み、

40

前記少なくとも1つのデータ位置は、宛先レジスタ、メモリアドレス、ソースレジスタ、定数、又は参照レジスタのうちの少なくとも1つを含むことを特徴とする請求項 11 に記載の H B M。

【請求項 13】

前記アトミック演算は、A D D、S U B T R A C T、E X C H A N G E、M A X、M I N、I N C R E M E N T、D E C R E M E N T、C O M P A R E - A N D - S W A P、A N D、O R、X O R、又は N O T の関数のうちの少なくとも1つを含むことを特徴とする請求項 12 に記載の H B M。

【請求項 14】

前記演算は、A L U 演算及び D R A M アクセス演算を含み、

50

前記少なくとも1つのデータ位置は、宛先レジスタ及び少なくとも1つのソースレジスタを含むことを特徴とする請求項11に記載のHBM。

【請求項15】

前記DRAMアクセス演算は、前記ALU演算と対をなすロード命令又は格納命令を含むことを特徴とする請求項14に記載のHBM。

【請求項16】

前記ファンクション・イン・メモリ命令の実行調整は、

前記コントローラによって、前記ファンクション・イン・メモリ命令をソースレジスタ及び宛先レジスタを含むファンクション・イン・メモリ移動命令として識別し、

前記コントローラによって、前記ソースレジスタに応じて、前記DRAMから少なくとも1つのデータを検索し、

前記コントローラによって、前記少なくとも1つのデータを前記宛先レジスタの前記DRAMに格納することを特徴とする請求項10に記載のHBM。

【請求項17】

前記ファンクション・イン・メモリ命令の実行調整は、

前記コントローラによって、前記ファンクション・イン・メモリ命令をソースレジスタ及び宛先レジスタのうちの少なくとも1つを含むファンクション・イン・メモリスクラッチパッドの命令として識別し、

前記コントローラのタイミングパラメータをDRAMのタイミングパラメータからSRAMのタイミングパラメータに調整し、

前記コントローラによって、前記SRAMで前記SRAMのタイミングパラメータに応じて前記ファンクション・イン・メモリスクラッチパッド命令を実行することを特徴とする請求項10に記載のHBM。

【請求項18】

HBM(High-Bandwidth Memory)システムであって、

DRAMを含むDRAMダイと、

コントローラ、ALU(Arithmetic Logic Unit)、及びSRAMを含むロジックダイと、

を有するHBMと、

前記HBMから分離して前記コントローラに命令を送信するグラフィック処理装置メモリコントローラと、を備え、

前記コントローラは、

前記命令を受信して、前記命令が一般命令である場合には、前記DRAMダイの前記DRAMに前記命令を伝達し、前記命令がファンクション・イン・メモリ命令である場合には、前記命令の実行を調整し、

前記命令の実行は、

前記コントローラ、前記ALU、及び前記SRAMのうちの少なくとも1つを用いて実行することを特徴とするHBMシステム。

【請求項19】

前記ファンクション・イン・メモリ命令の実行調整は、

前記コントローラによって、前記ファンクション・イン・メモリ命令を演算及び少なくとも1つのデータ位置を含む計算ファンクション・イン・メモリ命令として識別し、

前記コントローラによって、前記少なくとも1つのデータ位置に応じて前記DRAMから少なくとも1つのデータを検索し、

前記コントローラによって、前記少なくとも1つのデータ及び前記演算を前記ALUに提供し、

前記ALUによって、前記少なくとも1つのデータに対して前記演算を実行し、

前記演算の実行結果を前記DRAMに格納することを特徴とする請求項18に記載のHBMシステム。

【請求項20】

10

20

30

40

50

前記演算は、アトミック演算 (a t o m i c o p e r a t i o n) 及びデータ型を含み、

前記少なくとも1つのデータ位置は、宛先レジスタ、メモリアドレス、ソースレジスタ、定数、又は参照レジスタのうちの少なくとも1つを含むことを特徴とする請求項19に記載のHBMシステム。

【請求項21】

前記アトミック演算は、A D D、S U B T R A C T、E X C H A N G E、M A X、M I N、I N C R E M E N T、D E C R E M E N T、C O M P A R E - A N D - S W A P、A N D、O R、X O R、又はN O Tの関数のうちの少なくとも1つを含むことを特徴とする請求項20に記載のHBMシステム。

10

【請求項22】

前記演算は、A L U演算及びD R A Mアクセス演算を含み、

前記少なくとも1つのデータ位置は、宛先レジスタ及び少なくとも1つのソースレジスタを含むことを特徴とする請求項19に記載のHBMシステム。

【請求項23】

前記D R A Mアクセス演算は、前記A L U演算と対をなすロード命令又は格納命令を含むことを特徴とする請求項22に記載のHBMシステム。

【請求項24】

前記ファンクション・イン・メモリ命令の実行調整は、

前記コントローラによって、前記ファンクション・イン・メモリ命令をソースレジスタ及び宛先レジスタを含むファンクション・イン・メモリ移動命令として識別し、

20

前記コントローラによって、前記ソースレジスタに応じて、前記D R A Mから少なくとも1つのデータを検索し、

前記コントローラによって、前記少なくとも1つのデータを前記宛先レジスタの前記D R A Mに格納することを特徴とする請求項18に記載のHBMシステム。

【請求項25】

前記ファンクション・イン・メモリ命令の実行調整は、

前記コントローラによって、前記ファンクション・イン・メモリ命令をソースレジスタ及び宛先レジスタのうちの少なくとも1つを含むファンクション・イン・メモリスクラッチパッド命令として識別し、

30

前記コントローラのタイミングパラメータをD R A MのタイミングパラメータからS R A Mのタイミングパラメータに調整し、

前記コントローラによって、前記S R A Mで前記S R A Mのタイミングパラメータに応じて前記ファンクション・イン・メモリスクラッチパッド命令を実行することを特徴とする請求項18に記載のHBMシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、メモリ帯域幅の管理に関し、より詳しくは、プロセッサ及び高帯域幅メモリ (H B M : H i g h - B a n d w i d t h M e m o r y) のロジックダイ (L o g i c D i e) 上のメモリコントローラを有するHBMシステム、及びHBMの処理方法に関する。

40

【背景技術】

【0002】

HBMは、主にグラフィック処理装置 (G P U) 用高性能メモリとして用られる。HBMは、通常のD R A Mに比べて非常に広いバスを有するという利点がある。現在のHBMシステムの構造 (a r c h i t e c t u r e) は、HBMでバッファとして機能するロジックダイと、多重スタックD R A Mダイ (例：ダイス (d i c e)) と、グラフィック処理装置のHBMコントローラと、を備える。メモリシステムにメモリ内の処理 (例：インメモリプロセッシング) 機能を追加することによって、さらに性能が向上するが、HBMの

50

変更による既存のグラフィック処理装置の環境変化は最小限にすべきである。

【 0 0 0 3 】

なお、上述の内容は、本発明の実施形態における背景技術の理解を助けるためのものであって、先行技術を構成しない技術内容を含む。

【 先行技術文献 】

【 特許文献 】

【 0 0 0 4 】

【 特許文献 1 】 特開 2 0 0 3 - 0 1 5 8 2 4 号 公 報

【 発明の概要 】

【 発明が解決しようとする課題 】

10

【 0 0 0 5 】

本発明は、上記従来技術に鑑みてなされたものであって、本発明の目的は、メモリシステムにファンクション・イン・メモリを導入して性能を向上しつつも、既存のグラフィック処理装置の環境変化を最小限にすることができるインメモリのコマンド処理方法と、このための H B M、及び H B M システムを提供することにある。

【 課題を解決するための手段 】

【 0 0 0 6 】

本発明は、高帯域幅メモリシステムでインメモリのコマンドを調整するためのシステム、及び方法を提供する。

【 0 0 0 7 】

20

上記目的を達成するためになされた本発明の一態様による高帯域幅メモリ (H B M) システムにおいてインメモリ (In - memory) のコマンドを処理するシステム、及び方法は、グラフィック処理装置の H B M コントローラによって、ファンクション・イン・H B M (Function - in - H B M、以下「ファンクション・イン・メモリ」という) 命令を H B M に送信し、前記 H B M のロジック部分 (Logic component) で前記ファンクション・イン・メモリ命令を受信することを含む。前記ロジック部分は、コントローラ、A L U (Arithmetic Logic Unit)、及び S R A M を含む。前記ロジック部分は、前記ファンクション・イン・メモリ命令に基づいて、前記コントローラ、前記 A L U、又は前記 S R A M のうちの少なくとも 1 つを用いることによって、前記ファンクション・イン・メモリ命令の実行を調整する。

30

【 0 0 0 8 】

前記ファンクション・イン・メモリ命令の実行調整は、前記ファンクション・イン・メモリ命令を、計算 (computational) ファンクション・イン・メモリ命令として識別することを含む。前記計算ファンクション・イン・メモリ命令は、演算 (operation) 及び少なくとも 1 つのデータ位置を含む。前記計算ファンクション・イン・メモリ命令の実行調整は、前記少なくとも 1 つのデータ位置に応じて、前記 H B M の D R A M から少なくとも 1 つのデータを検索し、前記少なくとも 1 つのデータ及び前記演算を前記 A L U に提供し、前記 A L U によって前記少なくとも 1 つのデータに対して前記演算を実行し、前記演算の実行結果を D R A M に格納することを有する。

【 0 0 0 9 】

40

前記演算は、アトミック演算 (atomic operation) 及びデータ型であり、前記少なくとも 1 つのデータ位置は、宛先レジスタ、メモリアドレス、ソースレジスタ、定数、又は参照レジスタのうちの少なくとも 1 つを含む。

【 0 0 1 0 】

前記アトミック演算は、A D D、S U B T R A C T、E X C H A N G E、M A X、M I N、I N C R E M E N T、D E C R E M E N T、C O M P A R E - A N D - S W A P、A N D、O R、X O R、又は N O T の関数のうちの少なくとも 1 つを含む。

【 0 0 1 1 】

前記演算は、A L U 演算及び D R A M アクセス演算であり、前記少なくとも 1 つのデータ位置は、宛先レジスタ及び少なくとも 1 つのソースレジスタを含む。

50

【 0 0 1 2 】

前記 D R A M アクセス演算は、前記 A L U 演算と対をなすロード命令又は格納命令を含む。

【 0 0 1 3 】

前記ファンクション・イン・メモリ命令の実行調整は、前記コントローラによって前記ファンクション・イン・メモリ命令をファンクション・イン・メモリ移動命令として識別することを含む。前記ファンクション・イン・メモリ移動命令は、ソースレジスタと、宛先レジスタと、を含む。前記コントローラは、前記ソースレジスタに応じて、前記 H B M の D R A M から前記少なくとも 1 つのデータを検索し、前記少なくとも 1 つのデータを前記宛先レジスタの D R A M に格納する。

10

【 0 0 1 4 】

前記ファンクション・イン・メモリ命令の実行調整は、前記ファンクション・イン・メモリ命令をファンクション・イン・メモリスクラッチパッド命令として識別することを含む。前記ファンクション・イン・メモリスクラッチパッド命令は、ソースレジスタ又は宛先レジスタのうちの少なくとも 1 つを含む。前記ファンクション・イン・メモリスクラッチパッド命令の実行調整は、前記コントローラのタイミングパラメータを D R A M のタイミングパラメータから S R A M のタイミングパラメータに調整し、前記 S R A M で前記 S R A M のタイミングパラメータに応じて、前記ファンクション・イン・メモリスクラッチパッド命令を実行することを含む。

20

【 0 0 1 5 】

前記 H B M システムにおいて、前記インメモリのコマンドを処理するシステム、及び方法は、前記少なくとも 1 つのデータ位置がグラフィック処理装置のキャッシュを含む場合、前記 H B M コントローラによって前記グラフィック処理装置における前記ファンクション・イン・メモリ命令の実行を調整することを含む。

【 発明の効果 】

【 0 0 1 6 】

本発明によれば、既存のグラフィック処理装置の環境変化を最小限にすると共に、メモリシステムの性能を向上させるためにインメモリのコマンドを処理する。

【 図面の簡単な説明 】

【 0 0 1 7 】

30

【 図 1 】 本発明の一実施形態による高帯域幅メモリシステム構造のブロック図である。

【 図 2 】 本発明の一実施形態によるファンクション・イン・メモリの命令セットアーキテクチャからの命令を処理する演算を示すフローチャートである。

【 発明を実施するための形態 】

【 0 0 1 8 】

以下、本発明を実施するための形態の具体例を、図面を参照しながら詳細に説明する。

【 0 0 1 9 】

本発明の技術的思想の特徴、及びこれを達成するための方法は、実施形態に関する以下の具体的な説明及び図面を参照して容易に理解される。以下、図面を参照しながら、本発明の実施形態について詳細に説明する。図面全般に亘って同一の参照符号は、同一の構成要素を示す。なお、本発明は、多様な実施形態を有し、本明細書で説明する実施形態のみに限定されない。本明細書で説明する実施形態は、本発明を全て完全に開示する例として提供され、これを通じて本発明が属する技術分野の通常の知識を有する者に、本発明の態様及び特徴を十分に伝えられる。従って、本発明が属する技術分野の通常の技術者が本発明の態様及び特徴を完全に理解するために、不要なプロセス、構成、及び技法については説明を省略する。特に言及しない限り、図面及び詳細な説明全般に亘って同一の参照符号は、同一の構成要素を示し、同一の構成要素に関する重複説明は省略する。図面に示す構成要素、階層、及び領域は、明確性のために誇張することがある。

40

【 0 0 2 0 】

以下の説明では、様々な実施形態が完全に理解されるように、数多くの特定の細部説明

50

を記載する。しかし、このような特定の細部説明がなくても、様々な実施形態を実施し得ることは明らかである。別の例では、公知の構造、及び装置は、様々な実施形態が不要且つあいまいにならないようにブロック図の形態で示す。

【0021】

何れかの要素、階層、領域、又は構成が、他の要素、階層、領域、又は構成に対して、「上に」、「連結された」、又は「結合された」と記載される場合、これは、直接、他の要素や階層上に、他の要素や階層に連結された、又は他の要素や階層に結合されたことを意味する、或いは1つ以上の媒介要素、媒介階層、媒介領域、又は媒介構成が存在することを意味する。但し、「直接連結された／直接結合された」との記載は、中間の構成なしで他の構成と直接連結されるか、又は結合される1つの構成を示す。一方、「間に」、「間にすぐ」、又は「隣接した」及び「すぐ隣」のように、構成間の関係を説明する他の表現も同様に解釈する。また、何れのかの要素又は階層が、二つの要素又は階層の間にあると記載する場合、これは、単に二つの要素又は階層の間にあることを意味するか、又は1つ以上の媒介要素又は階層が存在することを意味する。

10

【0022】

本明細書で用いる用語は、単に特定の実施形態を説明するためのものであって、本発明を限定するものではない。本明細書で用いる単数形態の用語は、文脈上特に明らかな指示がない限り、複数形態の用語も含む。また、「含む」、「有する」、「備える」等の用語を本明細書で用いる場合、これは、記載する特徴、数字、段階、動作、要素、及び／又は構成の存在を明示するが、1つ以上の他の特徴、数字、段階、動作、要素、構成、及び／又はこれらの集合の存在や付加を排除するものではない。本明細書で用いる「及び／又は」という用語は、挙げられる1つ以上の関連項目の任意の組み合わせ及び全ての組み合わせを含む。

20

【0023】

本明細書で用いる「実質的に」、「約」、「略」、及びこれに類似する用語は、程度 (degree) を示す用語ではなく、近似 (approximation) を示す用語として用いるものであり、本発明が属する技術分野の通常の技術者に認識される測定値又は計算値に内在する偏差を説明するためのものである。本明細書で用いる「約」又は「略」は、当該測定及び特定の量の測定に関する誤差 (即ち、測定システムの限界) を考慮し、本発明の技術分野に属する通常の技術者が決定した特定の値に対する許容可能な偏差範囲内の平均及び記載した値を含む。例えば、「約」は、1つ以上の標準偏差内、又は記載した値の $\pm 30\%$ 、 20% 、 10% 、 5% 内を意味する。また、本発明の実施形態を説明する際に「し得る及び／又はできる」という表現を用いる場合、これは、「本発明の1つ以上の実施形態」を示す。本明細書で用いられる「用いる」及び「用いられる」という用語は、それぞれ「利用する」及び「利用される」という用語と同じ意味である。更に、「例示的な」という用語は、例示又は一例を指称する。

30

【0024】

特定の実施形態が異なって実施される場合、詳細なプロセス順序は、説明している順序とは異なって実行され得る。例えば、連続して説明した二つのプロセスは、実質的に同時に、又は説明した順序とは逆に実行される。

40

【0025】

実施形態及び／又は中間構造の概略図である図面を参照して、様々な実施形態を本明細書で説明する。図示した形状は、例えば、製造技術及び／又は許容誤差の結果により変わる。また、本明細書に開示している特定の構造的又は機能的説明は、単に本発明の技術的思想による実施形態を説明するために例示する。従って、本明細書に開示する実施形態は、説明している領域の特定形態に制限されるものではなく、例えば、製造過程で生じる偏差形状を含む。例えば、矩形で示した注入領域 (implanted region) は、通常円形や湾曲した形態を有し、及び／又は注入領域から非注入領域 (non-implanted region) への二元的変化というよりもその境界部分で注入濃度の傾斜を有する。同様に、注入によって形成された埋込領域は、注入が起きる面と埋込領域と

50

の間の領域にいくらかの注入をもたらす。従って、図に示す領域は、事実上概略的なものであり、その形態は装置の領域の実際の形態を示したものでなく、示した形状に制限しようとするものでもない。

【0026】

本明細書で説明する本発明の実施形態によると、電子/電気装置、及び/又は任意の他の関連装置や構成は、任意の適切なハードウェア、ファームウェア（例えば、特定用途向け集積回路）、ソフトウェア、又はソフトウェア、ファームウェア、及びハードウェアの適切な組み合わせを利用して具現される。例えば、このような装置の様々な構成要素は、1つの集積回路（IC）チップ又は個別のICチップ上に形成される。また、このような装置の様々な構成要素は、フレキシブル印刷回路フィルム（flexible printed circuit film）、テープキャリアパッケージ（TCP：Tape Carrier Package）、印刷回路基板（PCB：Printed Circuit Board）上に具現されるか又は1つの基板（substrate）上に形成される。更に、このような装置の多様な構成要素は、1つ以上のコンピューティング装置の1つ以上のプロセッサで実行され、コンピュータプログラムの命令を実行して、本明細書に記載する多様な機能を行うための他のシステム構成要素と相互作用するプロセス又はスレッド（thread）である。コンピュータプログラムの命令はメモリに格納され、メモリは、例えばランダムアクセスメモリ（RAM：Random Access Memory）又はフラッシュメモリ（例：NANDフラッシュメモリ）装置などの標準メモリ装置を用いるコンピューティング装置で具現される。コンピュータプログラムの命令は、例えばCD-ROM、フラッシュドライブ等のような他の非一時的コンピュータ読み取り可能な媒体に格納される。また、本発明が属する技術分野の通常の技術者には、本発明の技術範囲を逸脱することなく、多様なコンピューティング装置の機能が、単一のコンピューティング装置に結合若しくは統合されるか、又は特定のコンピューティング装置の機能が1つ以上の他のコンピューティング装置に分散されることが認識される。

10

20

30

【0027】

本明細書で用いる技術用語及び科学用語を含む全ての用語は、特に定義しない限り、本発明が属する技術分野の通常の知識を有する者が一般的に理解するものと同様の意味を有する。また、通常用いられる辞典に定義されているような用語は、関連技術及び/又は本明細書の文脈上の意味と一致すると解釈され、本明細書で明らかに定義しない限り、理想的又は過度に形式的な意味として解釈されない。

【0028】

図1は、本発明の一実施形態による高帯域幅メモリシステム構造のブロック図である。

【0029】

図1に示す本発明の実施形態は、ファンクション・イン・メモリ（Function-in-HBM）のHBMシステム100、及びHBM用命令セットアーキテクチャ（ISA：Instruction Set Architecture）の拡張のためのシステムを提供する。

【0030】

HBMシステム100は、HBM110に統合される追加の計算リソースを支援する。例えば、本実施形態において、HBMシステム100は、一部のデータの演算及び移動をインメモリ（in-memory）で実行させ、大容量のスクラッチパッドを提供する。

40

【0031】

HBMシステム100は、グラフィック処理装置（GPU）170に連結された少なくとも1つのHBM110を含む。本実施形態において、HBM110は、DRAM120（例：1つ以上のDRAMダイ（die））と、ロジック部分（logic component）130（例：ロジックダイ（die））と、を含む。ロジック部分130は、コントローラ140、ALU150、及びSRAM160を有し、グラフィック処理装置170は、HBM110とインターフェースするためのHBMコントローラ180を含む。

50

【 0 0 3 2 】

本発明の一実施形態によるコントローラ 1 4 0 は、グラフィック処理装置 1 7 0 からの命令の実行を調整する。命令は、一般命令とファンクション・イン・メモリ命令との両方を含む。例えば、一般命令（ファンクション・イン・メモリ命令ではなく従来のロードファンクション及び格納ファンクション）は、HBMコントローラ 1 8 0 によって送信され、コントローラ 1 4 0 で受信されて、既存の方法で実行される。

【 0 0 3 3 】

また、コントローラ 1 4 0 は、インメモリ・ファンクション（例：ファンクション・イン・メモリ命令）の実行を調整する。例えば、コントローラ 1 4 0 は、データ移動演算（例：ロード／格納の対命令）を実行する。一例として、コントローラ 1 4 0 は、本来複数の一般命令であったファンクション・イン・メモリ命令を実行する。例えば、コントローラ 1 4 0 は、ALU 1 5 0 を利用する計算ファンクション・イン・メモリ命令（例：アトミック命令（Atomic Instructions）及びALU命令）の実行を調整する。この場合、コントローラ 1 4 0 は、DRAM 1 2 0 からデータを検索して、処理のために該当データ（及びALU演算）をALU 1 5 0 に提供することによって、命令の実行を調整する。その結果は、DRAM 1 2 0 に格納されるか、又はグラフィック処理装置 1 7 0 に戻される。一例として、ファンクション・イン・メモリ命令は、ロード命令又は格納命令と対をなす 1 つ以上のALU命令を含む。

【 0 0 3 4 】

他の実施形態として、コントローラ 1 4 0 はスクラッチパッドの読み取り命令及び書き込み命令の実行を調整する。以下では、このようなファンクション・イン・メモリの各類型について詳細に説明する。

【 0 0 3 5 】

本発明の一実施形態によるALU 1 5 0 は、様々な計算の動作（例：単純な計算のコマンド）を実行する。一例として、ALU 1 5 0 は、算術演算、ビット演算、シフト演算等を実行する 3 2 ビットALUである。例えば、ALU 1 5 0 は、ADD、SUBTRACT、EXCHANGE、MAX、MIN、INCREMENT、DECREMENT、COMPARE - AND - SWAP、AND、OR、及びXORの演算を実行する。ALU 1 5 0 は、アトミック演算及び非アトミック演算に利用される。

【 0 0 3 6 】

一実施形態として、コントローラ 1 4 0 は演算を提供し、ALU 1 5 0 へのデータ入力及びALU 1 5 0 からDRAM 1 2 0 へのデータ出力を管理する。また、他の実施形態として、ALU 1 5 0 が直接DRAM 1 2 0 からデータを検索して、DRAM 1 2 0 にデータを格納する。さらに他の実施形態では、コントローラ 1 4 0 がDRAM 1 2 0 からデータを検索して、DRAM 1 2 0 にデータを格納する役割を担う。

【 0 0 3 7 】

本発明の一実施形態によるSRAM 1 6 0 は、低レイテンシのスクラッチパッドとして構成される。一実施形態として、SRAM 1 6 0 は、同一のコマンド／アドレス（CA：Command / Address）及びデータ（DQ）インターフェースをDRAM 1 2 0 と共有し、他の実施形態として、SRAM 1 6 0 は、固有のCA及びDQインターフェースを有する。SRAM 1 6 0 は、DRAM 1 2 0 のアドレス範囲から区別される固有のアドレス範囲を含む。

【 0 0 3 8 】

コントローラ 1 4 0 は、入力される読み取り／書き込み命令のアドレスを用いて、要請がスクラッチパッド演算であるか否かを判定する。他の実施形態として、グラフィック処理装置 1 7 0 は、具体的に指定されたスクラッチパッド命令をコントローラ 1 4 0 に送信する。

【 0 0 3 9 】

SRAM 1 6 0 を利用する場合、コントローラ 1 4 0 は、SRAM 1 6 0 のタイミングパラメータ（例：DRAM 1 2 0 のタイミングパラメータより速いか又は低い

10

20

30

40

50

レイテンシ)に応じて動作するように自身のタイミングパラメータを変える。S R A M 1 6 0の利用は、ユーザ(例:プログラマ)によって指定され、S R A M 1 6 0の空間は、実行時間中に割り当てられる。動作中のスクラッチパッドは、グラフィック処理装置のL 1スクラッチパッドと同様に動作する(例:低レイテンシメモリを提供)。グラフィック処理装置のL 1スクラッチパッドは通常小さい(コア当たり15kB)ため、拡張されたH B Mスクラッチパッド(例:S R A M 1 6 0)はD R A M 1 2 0を利用する場合よりも性能が向上する。

【0040】

本発明の一実施形態によるファンクション・イン・メモリI S Aは、H B M 1 1 0で利用可能な追加のリソースを利用するために提供される。例えば、ファンクション・イン・メモリI S Aは、計算ファンクション・イン・メモリ命令(例:ファンクション・イン・メモリアトミック命令及びファンクション・イン・メモリA L U命令)と、データ移動ファンクション・イン・メモリ命令と、ファンクション・イン・メモリスクラッチパッド命令と、を含む演算を許容するために、既存の命令セットを拡張する。各々のファンクション・イン・メモリ命令は、命令をファンクション・イン・メモリ命令として識別するファンクション・イン・メモリの指定子と、H B Mにより実行される1つ以上の演算と、データ位置(例:レジスタ、メモリ、提供された定数等)と、を含む。一実施形態として、ファンクション・イン・メモリ命令は、`<designator> . <operation> . <data location 1> <data location 2>`の形式に配列される。

【0041】

本実施形態において、命令は、ユーザ(例:プログラマ)、コンパイラ、又はグラフィック処理装置によって、ファンクション・イン・メモリ命令として指定される。例えば、一部のプログラミング言語において、ユーザは、演算を実行する位置(例:インメモリ、グラフィック処理装置、又は中央処理装置)、又は利用するメモリを指定する。他の例として、コンパイラは、H B M 1 1 0で実行するコマンドを識別し、グラフィック処理装置170よりもH B M 1 1 0のファンクション・イン・メモリのコマンドを優先して実行する。さらに他の例として、グラフィック処理装置170は、ソース及び宛先メモリアドレスを分析して、ファンクション・イン・メモリ命令が有効か否かを判定する。例えば、グラフィック処理装置170は、メモリアドレスを分析して、メモリアドレスの少なくとも1つがグラフィック処理装置のキャッシュに位置するか、又はH B Mではない他のメモリに位置するかを判定し、これらの場合、一般命令を実行する(例:非インメモリ命令)。

【0042】

図2は、ファンクション・イン・メモリの命令セットアーキテクチャ(F I M I S A)からの命令を処理する方法を示すフローチャートである。

【0043】

本発明の一実施形態によるグラフィック処理装置170は、ファンクション・イン・メモリ(F I M)命令をH B M 1 1 0に送信する(S 200段階)。グラフィック処理装置170は、要請を処理してH B M 1 1 0に送信し、任意の返還情報を処理するH B Mコントローラ180を含む。本実施形態において、H B Mコントローラ180は、ファンクション・イン・メモリ命令に含まれるメモリアドレスの位置を推定することによって、ファンクション・イン・メモリ命令が適切か否かを検証する。例えば、アドレスが他のH B Mを示すか、又はアドレスの何れかがグラフィック処理装置170のキャッシュ、即ちローカルキャッシュ用である場合、ファンクション・イン・メモリ命令は不適切なものとなる。本実施形態において、ファンクション・イン・メモリ命令の検証は、H B M 1 1 0に命令が送信される前に、初期のグラフィック処理装置パイプライン(G P U P i p e l i n e)段階で行われる。例えば、本実施形態において、グラフィック処理装置のローカルキャッシュコントローラは、ファンクション・イン・メモリ命令が適切か否かを検証する。

【0044】

HBM110は、ロジック部分130でファンクション・イン・メモリ命令を受信する(S210段階)。コントローラ140は、命令を処理して実行を調整する。コントローラ140は、命令がファンクション・イン・メモリ命令なのか否かを検証し、その命令による演算を判定する(S220段階)。例えば、コントローラ140は、命令がALU150を利用する演算命令なのか、移動命令なのか、それともスクラッチパッド命令なのかを判定する。次に、コントローラ140は、命令自体(例:移動命令)を完了することによって、又は必須ロジックハードウェア(例:ALU150又はSRAM160)を用いることによって、命令の実行を調整する(S230段階)。

【0045】

本発明の一実施形態によるISAは計算命令を含む。本実施形態において、コントローラ140は計算命令を受信して、ALU150での計算命令の実行を調整する。計算命令は、ファンクション・イン・メモリアトミック命令と、ファンクション・イン・メモリALU命令と、を含む。

【0046】

本実施形態において、HBM110は、ファンクション・イン・メモリISAを用いてファンクション・イン・メモリアトミック命令を処理する。アトミック命令は、一般に3つの段階、即ち、メモリ位置からデータを読み取る段階と、データに対して関数(例:ADD、SUBTRACT等)を実行する段階と、結果データを再度所定のメモリ位置に書き込む段階と、に分類される。

【0047】

HBM110は、ファンクション・イン・メモリアトミック命令を受信すると、内部で上記3つの段階を全て実行する。従来のアトミック命令と比較すると、ファンクション・イン・メモリアトミック命令は、命令の実行を完了するために、グラフィック処理装置170によってHBM110に送信された追加情報を含む。例えば、グラフィック処理装置170は、アトミック演算、データ型、宛先レジスタ、ソースレジスタ、参照レジスタ(例:COMPARE-AND-SWAP関数のための)、メモリアドレス、及び関数型を含むファンクション・イン・メモリアトミック命令を送信する。

【0048】

本実施形態において、コントローラ140は、グラフィック処理装置のHBMコントローラ180からファンクション・イン・メモリアトミック命令を受信する。コントローラ140は、命令を読み取り、ファンクション・イン・メモリの指定子を用いて読み取られた命令が、ファンクション・イン・メモリ命令か否かを判定する。コントローラ140は、該当命令がファンクション・イン・メモリアトミック命令か否かを判定するために演算を用いる。この演算は、該当関数がファンクション・イン・メモリアトミック命令という信号を送信する以外にも、ALU150により実行される関数型(例:アトミックADD、アトミックCOMPARE-AND-SWAP、アトミックOR等)と、演算が実行されるデータ型(例:符号のある32ビット定数、符号のない32ビット定数等)と、を示す。

【0049】

次に、コントローラ140は、提供されたデータ位置からデータを読み取り、関数と共にALU150にデータを提供する。ALU150は、データに対して関数を実行し、その結果は元のデータ位置に格納される。

【0050】

一例として、ファンクション・イン・メモリアトミック命令の一般フォーマットは、`im.atom.<function>.<data type> <destination register> <memory address> <source register or constant> <reference register>`である。表1は、例示値を有する関数の例(又はファンクション・イン・メモリアトミック命令の例)の一部を示す。

【0051】

10

20

30

40

50

【表 1】

関数	ファンクション・イン・メモリ命令の例
General Format	f i m . a t o m . < f u n c t i o n > < d a t a l o c a t i o n 1 > < d a t a l o c a t i o n 2 >
ADD	f i m . a t o m . a d d . u 3 2 % r 6 , [% r 1 3] , - 1 0
Exchange	f i m . a t o m . e x c h . b 3 2 % r 7 , [% r 1 4] , % r 4
Find Max	f i m . a t o m . m a x . s 3 2 % r 8 , [% r 1 5] , % r 4
Find Min	f i m . a t o m . m i n . s 3 2 % r 9 , [% r 1 6] , % r 4
Increment	f i m . a t o m . i n c . u 3 2 % r 1 0 , [% r 1 7] , 1 7
Decrement	f i m . a t o m . d e c . u 3 2 % r 1 1 , [% r 1 8] , 1 3 7
Compare-and-Swap	f i m . a t o m . c a s . b 3 2 % r 1 3 , [% r 1 9] , % r 1 2 , % r 4
AND	f i m . a t o m . a n d . b 3 2 % r 1 7 , [% r 1 1 0] , % r 1 6
OR	f i m . a t o m . o r . b 3 2 % r 1 9 , [% r 1 1 1] , % r 1 8
XOR	f i m . a t o m . x o r . b 3 2 % r 2 0 , [% r 1 1 2] , % r 4

10

20

【0052】

本実施形態において、ファンクション・イン・メモリ命令ISAは、ファンクション・イン・メモリALU命令のための命令を含む。グラフィック処理装置によって実行される通常の演算は、メモリから必要なデータを引き出し、結果データを格納することに伴うロード命令及び格納命令を必要とする。本実施形態において、ファンクション・イン・メモリALU命令は、既存の関数とこれに伴うロード命令及び格納命令を単一のファンクション・イン・メモリALU命令に圧縮する。一例として、ファンクション・イン・メモリALU命令の一般フォーマットは、ファンクション・イン・メモリの指定子と、ALU演算及びこれと対をなすロード演算/格納演算を含む演算と、少なくとも1つのデータ位置と、を含む。例えば、ファンクション・イン・メモリALU命令は、f i m . < f u n c t i o n > . < l o a d / s t o r e > < d e s t i n a t i o n r e g i s t e r > < s o u r c e r e g i s t e r > < l o a d r e g i s t e r > の形式に配列される。例えば、ALU命令は、表2（ファンクション・イン・メモリALU ロード命令/格納命令の例）に示すように、ロード命令及び/又は格納命令と対をなす。

30

40

【0053】

【表 2】

非ファンクション・イン・メモリ命令	ファンクション・イン・メモリ命令の例
l d \$ r 1 , [\$ r 2] a d d \$ r 3 , \$ r 0 , \$ r 1	f i m . a d d . l d \$ r 3 , \$ r 0 , [\$ r 2]
a d d \$ r 3 , \$ r 0 , \$ r 1 s t [\$ r 2] , \$ r 3	f i m . a d d . s t [\$ r 2] , \$ r 0 , \$ r 1

【0054】

一例として、ファンクション・イン・メモリALU命令の演算は、ファンクション・イ

50

ン・メモリアトミック命令の演算と同様である。例えば、本実施形態において、コントローラ 140 は、グラフィック処理装置の HBM コントローラ 180 からファンクション・イン・メモリ ALU 命令を受信する。コントローラ 140 は、命令を読み取り、ファンクション・イン・メモリの指定子を用いて読み取った命令がファンクション・イン・メモリ命令か否かを判定する。コントローラ 140 は、命令がファンクション・イン・メモリ ALU 命令か否かを判定するため、及び ALU 150 により実行される関数型（例：ADD、EXCHANGE、MIN、MAX、OR 等）を判定するために演算を用いる。

【0055】

次に、コントローラ 140 は、提供されたデータ位置からデータを読み取り、関数と共に ALU 150 にデータを提供する。ALU 150 は、データに対して関数を実行し、その結果は元のデータ位置又は指示されたデータ位置に格納される。

10

【0056】

本実施形態において、HBM 110 及びファンクション・イン・メモリ ISA は、ファンクション・イン・メモリ移動命令のために構成される。移動命令は、対応する格納命令と対をなすロード命令として識別される。ロード及び格納のアドレスが同一の HBM に位置すると、関数はインメモリで実行される。本実施形態において、対をなすロード及び格納の関数は、単一のファンクション・イン・メモリ移動命令に併合される。表 3 は、ロード命令 / 格納命令の例（又はファンクション・イン・メモリの移動命令の例）を示す。

【0057】

【表 3】

20

非ファンクション・イン・メモリ命令	ファンクション・イン・メモリ命令の例
ld \$r1, [\$r2] st [\$r3], \$r1	fim.mov [\$r3], [\$r2]

【0058】

コントローラ 140 は、ファンクション・イン・メモリ移動命令を受信すると、ファンクション・イン・メモリの指定子により、命令がファンクション・イン・メモリ命令であることを認識して、演算が移動演算であることを認識する。次に、コントローラ 140 は、ソースアドレスからデータを読み取り、宛先アドレスにデータを格納する。本実施形態において、コンパイラは、表 3 に示すように、単一のファンクション・イン・メモリ移動命令に併合されるロード命令及び格納命令の対を識別する。本実施形態において、グラフィック処理装置 170 は、移動命令を HBM 110 に送信する前に移動命令を分析し、ファンクション・イン・メモリ移動命令が適切か否かを判定する。例えば、グラフィック処理装置 170（例：HBM コントローラ 180）は、ソースメモリアドレス及び宛先メモリアドレスを分析し、アドレスの何れかがグラフィック処理装置のキャッシュに存在するか否かを判定する。この例において、グラフィック処理装置 170 は、命令を一般のロード命令と格納命令とに分離する。

30

【0059】

本実施形態において、HBM 110 及びファンクション・イン・メモリ ISA は、低レイテンシスクラッチパッドとして SRAM 160 を利用する。コントローラ 140 は、ファンクション・イン・メモリスクラッチパッド命令を識別する。上述のように、SRAM 160 は、DRAM 120 のアドレス範囲から区別された特定のアドレス範囲を含む。コントローラ 140 は、要請が DRAM 120 メモリアドレス又は SRAM 160 メモリアドレスに対応するか否かを識別する。本実施形態において、ISA は、特定のスクラッチパッドのコマンドを含む。例えば、ISA は、ファンクション・イン・メモリ読み取りスクラッチパッドのコマンド（例：FIM.RD__SP）と、ファンクション・イン・メモリ書き込みスクラッチパッドのコマンド（例：FIM.WR__SP）と、を含む。この例では、コントローラ 140 は入力コマンドのメモリアドレスを推定しない。

40

【0060】

50

S R A M 1 6 0 は、D R A M 1 2 0 より低いレイテンシ（即ち、より速く）で動作する。従って、スクラッチパッドは、データのロード時間及び格納時間に対して非決定的に H B M 1 1 0 をレンダリングする。コントローラ 1 4 0 は、D R A M メモリが関連関数を実行する場合、D R A M 1 2 0 のタイミングパラメータに応じて動作し、S R A M メモリの関連関数（例：ファンクション・イン・メモリスクラッチパッド命令）を実行する間は、S R A M 1 6 0 のタイミングパラメータに応じて動作することによって、命令の実行を調整する。従って、ファンクション・イン・メモリスクラッチパッド命令を受信すると、コントローラ 1 4 0 は S R A M 1 6 0 のタイミングパラメータに対応するように自身のタイミングパラメータを調整し、ファンクション・イン・メモリ読み取り／書き込みスクラッチパッドのコマンドを実行する。

10

【 0 0 6 1 】

本実施形態において、ユーザは低レイテンシスクラッチパッドとして S R A M 1 6 0 を利用するデータ構造を設定する。ユーザがスクラッチパッドを用いることによってデータ構造を設定すると、コンパイラは指定子（例：アセンブリの l l s p ）を含むように要請を変換し、その結果、グラフィック処理装置 1 7 0 は H B M 1 1 0 に空間を割り当てる。

【 0 0 6 2 】

以上、上述した本発明の実施形態は、高帯域幅メモリシステム及び命令セットアーキテクチャを提供する。

【 0 0 6 3 】

上述の説明内容は、本発明の一実施形態を説明するためのものであって、本発明を限定するものではない。本発明の技術分野に属する通常の技術者は、本発明の技術範囲から逸脱することなく、多様に変形実施することが可能である。

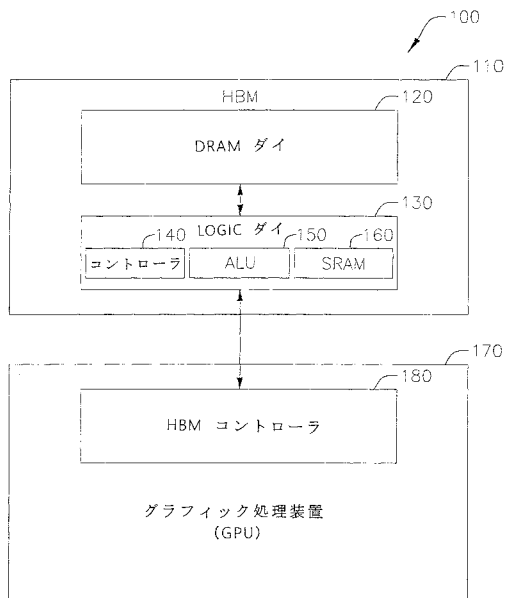
20

【 符号の説明 】**【 0 0 6 4 】**

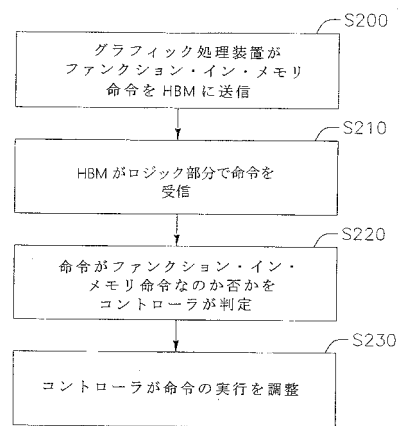
1 0 0	H B M システム
1 1 0	H B M
1 2 0	D R A M
1 3 0	ロジック部分
1 4 0	コントローラ
1 5 0	A L U
1 6 0	S R A M
1 7 0	グラフィック処理装置
1 8 0	H B M コントローラ

30

【図 1】



【図 2】



フロントページの続き

(72)発明者 マラディ, クリシャン テジャ

アメリカ合衆国, 9 5 1 3 5 カリフォルニア州, サンノゼ, ラウトレク ドライブ 4 1 9 6

(72)発明者 牛 迪 民

アメリカ合衆国, 9 4 0 8 7 カリフォルニア州, サニーベール, ホルトハウス テラス 5 2 7

(72)発明者 チェン 宏 忠

アメリカ合衆国, 9 5 0 3 2 カリフォルニア州, ロスガトス, ユニット6 カールトン アベニ
ュー 1 2 0

F ターム(参考) 5B033 BD00

5B160 CB03 GA00