

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7550614号
(P7550614)

(45)発行日 令和6年9月13日(2024.9.13)

(24)登録日 令和6年9月5日(2024.9.5)

(51)国際特許分類		F I	
G 0 6 F	12/00 (2006.01)	G 0 6 F	12/00 5 6 0 F
G 0 6 F	9/34 (2018.01)	G 0 6 F	9/34 3 5 0 A
G 0 6 F	12/06 (2006.01)	G 0 6 F	9/34 3 5 0 B
G 1 1 C	5/04 (2006.01)	G 0 6 F	12/06 5 4 0 E
H 1 0 B	12/00 (2023.01)	G 1 1 C	5/04 2 2 0
請求項の数 22 (全32頁) 最終頁に続く			
(21)出願番号	特願2020-191783(P2020-191783)	(73)特許権者	390019839
(22)出願日	令和2年11月18日(2020.11.18)		三星電子株式会社
(65)公開番号	特開2021-128752(P2021-128752 A)		S a m s u n g E l e c t r o n i c s C o . , L t d .
(43)公開日	令和3年9月2日(2021.9.2)		大韓民国京畿道水原市靈通区三星路 1 2 9
審査請求日	令和5年9月22日(2023.9.22)		1 2 9 , S a m s u n g - r o , Y e o n g t o n g - g u , S u w o n - s i , G y e o n g g i - d o , R e p u b l i c o f K o r e a
(31)優先権主張番号	62/975,577	(74)代理人	110000051
(32)優先日	令和2年2月12日(2020.2.12)		弁理士法人共生国際特許事務所
(33)優先権主張国・地域又は機関	米国(US)	(72)発明者	マラディ, クリシュナ テジャ
(31)優先権主張番号	16/859,829		アメリカ合衆国, 9 5 1 3 5 カリフォルニア州, サンノゼ, ロートレック ド
(32)優先日	令和2年4月27日(2020.4.27)		最終頁に続く
(33)優先権主張国・地域又は機関	米国(US)		
早期審査対象出願			

(54)【発明の名称】 インメモリコンピューティングに対するデータ配置のための方法及びその方法が適用されたメモリモジュール

(57)【特許請求の範囲】

【請求項 1】

メモリモジュールであって、
ダイナミックランダムアクセスメモリ（D R A M）バンクを含むメモリダイと、
ホストプロセッサからオペランド及び命令を受信するメモリコントローラと、を備え、
前記 D R A M バンクは、
複数のページに配列された D R A M セルのアレイと、
算術論理ユニット（A L U）を含むインメモリコンピューティング（I M C）モジュールと、を含み、
前記メモリコントローラは、
前記命令に基づいて、複数のデータレイアウトの中から、前記 D R A M セルのアレイの前記複数のページから前記 I M C モジュールへの前記オペランドの配置を指定するためのデータレイアウトを決定し、
前記決定されたデータレイアウトに従って、前記 D R A M バンクに前記オペランドを供給し、
前記命令に従って、前記オペランドに対して前記 A L U による演算を実行するように、前記 D R A M バンクの前記 I M C モジュールを制御するように構成され、
前記 I M C モジュールは、オペランドレジスタを更に含み、
前記オペランドは、第 1 のオペランド及び第 2 のオペランドを含み、
前記複数のデータレイアウトの中から一つのオペランド（1 O P）データレイアウトを

決定した場合、

前記第 1 のオペランドは、前記 D R A M バンクの外部から前記 I M C モジュールに供給され、

前記第 2 のオペランドは、前記 D R A M セルから前記オペランドレジスタを介して前記 I M C モジュールに供給されることを特徴とするメモリモジュール。

【請求項 2】

前記第 1 のオペランドは、複数の第 1 のタイルに分割され、

前記第 2 のオペランドは、複数の第 2 のタイルに分割され、

各タイルは、複数の値を含み、

前記メモリコントローラは、

前記オペランドレジスタに前記第 1 のオペランドの第 1 のタイルを格納し、

前記オペランドレジスタに格納された前記第 1 のオペランドの第 1 のタイル、及び前記第 2 のオペランドの複数の第 2 のタイルの各々に対して、前記 A L U による演算を実行するように更に構成されることを特徴とする請求項 1 に記載のメモリモジュール。

【請求項 3】

前記第 1 のオペランドは、複数の第 1 のタイルに分割され、

前記第 2 のオペランドは、複数の第 2 のタイルに分割され、

各タイルは、複数の値を含み、

前記複数のデータレイアウトの中から同じページ (S R) データレイアウトを決定した場合、

前記メモリコントローラは、前記 D R A M セルの同じページに、一つ以上の前記第 1 のタイル及び一つ以上の前記第 2 のタイルを格納するように構成されることを特徴とする請求項 1 に記載のメモリモジュール。

【請求項 4】

前記メモリコントローラは、

前記オペランドレジスタに前記一つ以上の前記第 1 のタイルの中の一つの第 1 のタイルを格納し、

前記オペランドレジスタに格納された前記第 1 のタイル、及び前記 D R A M セルのアレイの前記第 1 のタイルと同じページに格納された前記一つ以上の前記第 2 のタイルのそれぞれに対して、前記 A L U による演算を実行するように更に構成されることを特徴とする請求項 3 に記載のメモリモジュール。

【請求項 5】

前記 D R A M バンクの前記 I M C モジュールは、アキュムレータを更に含み、

前記アキュムレータは、累積値を格納するように構成されたアキュムレータレジスタを含み、

前記アキュムレータは、

前記 A L U によって演算された出力を受信し、

前記累積値と前記出力との合計で前記アキュムレータレジスタを更新するように構成され、

前記命令は、前記第 1 のオペランドと前記第 2 のオペランドとの内積を計算することを含み、

前記第 1 のタイルの中の一つの第 1 のタイルは、行データを格納し、

前記第 2 のタイルの中の一つの第 2 のタイルは、列データを含むことを特徴とする請求項 4 に記載のメモリモジュール。

【請求項 6】

前記第 1 のタイルは、第 1 の数の値を有し、

前記第 2 のタイルは、第 2 の数の値を有し、

前記 D R A M バンクの前記 I M C モジュールは、出力バッファを含み、

前記出力バッファは、前記第 1 の数の値と前記第 2 の数の値との積以上を格納するためのサイズを有し、

10

20

30

40

50

前記命令は、前記第 1 のオペランドと前記第 2 のオペランドとの外積を計算することを含み、

前記第 1 のタイルの中の一つの第 1 のタイルは、行データを格納し、

前記第 2 のタイルの中の一つの第 2 のタイルは、列データを含むことを特徴とする請求項 4 に記載のメモリモジュール。

【請求項 7】

前記第 1 のタイルは、第 1 の数の値を有し、

前記第 2 のタイルは、第 2 の数の値を有し、

前記 D R A M バンクの前記 I M C モジュールは、出力バッファを含み、

前記出力バッファは、前記第 1 の数の値及び前記第 2 の数の値の中の大きい方の値以上を格納するためのサイズを有し、

10

前記命令は、前記第 1 のオペランドと前記第 2 のオペランドとのテンソル積を計算することを含み、

前記第 1 のタイルの中の一つの第 1 のタイルは、行データを格納し、

前記第 2 のタイルの中の一つの第 2 のタイルは、列データを含むことを特徴とする請求項 4 に記載のメモリモジュール。

【請求項 8】

前記第 1 のオペランドは、複数の第 1 のタイルに分割され、

前記第 2 のオペランドは、複数の第 2 のタイルに分割され、

各タイルは、複数の値を含み、

20

前記複数のデータレイアウトの中から異なるページ (D R) データレイアウトを決定した場合、

前記メモリコントローラは、

前記 D R A M セルのアレイの第 1 のページに前記第 1 のタイルのサブセットを格納し、

前記 D R A M セルのアレイの第 2 のページに前記第 2 のタイルのサブセットを格納するように構成されることを特徴とする請求項 1 に記載のメモリモジュール。

【請求項 9】

前記メモリコントローラは、

前記オペランドレジスタに前記第 1 のページから前記第 1 のオペランドの第 1 のタイルを格納し、

30

前記オペランドレジスタに格納された前記第 1 のオペランドの第 1 のタイル、及び前記第 2 のページからの前記第 2 のオペランドの複数の第 2 のタイルのそれぞれに対して、前記 A L U による演算を実行するように更に構成されることを特徴とする請求項 8 に記載のメモリモジュール。

【請求項 10】

前記 D R A M バンクの前記 I M C モジュールは、前記 A L U によって演算された出力をバッファリングするように構成されたハードウェアのバッファを更に含むことを特徴とする請求項 1 に記載のメモリモジュール。

【請求項 11】

前記ハードウェアのバッファは、前記 I M C モジュールの結果レジスタのサイズの 4 倍以上のサイズを有することを特徴とする請求項 10 に記載のメモリモジュール。

40

【請求項 12】

前記 D R A M バンクの前記 I M C モジュールは、アキュムレータを更に含み、

前記アキュムレータは、累積値を格納するように構成されたアキュムレータレジスタを含み、

前記アキュムレータは、

前記 A L U によって演算された出力を受信し、

前記累積値と前記出力との合計で前記アキュムレータレジスタを更新するように構成されることを特徴とする請求項 1 に記載のメモリモジュール。

【請求項 13】

50

前記メモリモジュールは、シリコン貫通電極によって接続されたメモリダイのスタックを含む高帯域幅メモリ（HBM）モジュールであり、

前記メモリダイのスタックは、前記メモリダイを含むことを特徴とする請求項 1 に記載のメモリモジュール。

【請求項 1 4】

インメモリ（in-memory）計算を実行する方法であって、

メモリモジュールのメモリコントローラによって、ホストプロセッサからオペランド及び命令を受信する段階と、

前記メモリコントローラによって、前記命令に基づいて、複数のデータレイアウトの中からデータレイアウトを決定する段階と、

前記決定されたデータレイアウトに従って、前記メモリモジュールのダイナミックランダムアクセスメモリ（DRAM）バンクに前記オペランドを供給する段階と、

前記命令に従って、前記オペランドに対して算術論理ユニット（ALU）による演算を実行するように、前記 DRAM バンクのインメモリコンピューティング（IMC）モジュールを制御する段階と、を有し、

前記 DRAM バンクは、

複数のページに配列された DRAM セルのアレイと、

前記 ALU を含む前記 IMC モジュールと、を含み、

前記データレイアウトは、前記 DRAM セルのアレイの前記複数のページから前記 IMC モジュールへの前記オペランドの配置を指定し、

前記 IMC モジュールは、オペランドレジスタを更に含み、

前記オペランドは、第 1 のオペランド及び第 2 のオペランドを含み、

前記複数のデータレイアウトの中から一つのオペランド（OP）データレイアウトを決定した場合、

前記第 1 のオペランドは、前記 DRAM バンクの外部から前記 IMC モジュールに供給され、

前記第 2 のオペランドは、前記 DRAM セルから前記オペランドレジスタを介して前記 IMC モジュールに供給されることを特徴とする方法。

【請求項 1 5】

前記第 1 のオペランドは、複数の第 1 のタイルに分割され、

前記第 2 のオペランドは、複数の第 2 のタイルに分割され、

各タイルは、複数の値を含み、

前記メモリコントローラは、

前記オペランドレジスタに前記第 1 のオペランドの第 1 のタイルを格納し、

前記オペランドレジスタに格納された前記第 1 のオペランドの第 1 のタイル、及び前記第 2 のオペランドの複数の第 2 のタイルの各々に対して、前記 ALU による演算を実行するように更に構成されることを特徴とする請求項 1 4 に記載の方法。

【請求項 1 6】

前記第 1 のオペランドは、複数の第 1 のタイルに分割され、

前記第 2 のオペランドは、複数の第 2 のタイルに分割され、

各タイルは、複数の値を含み、

前記複数のデータレイアウトの中から同じページ（SR）データレイアウトを決定した場合、

前記メモリコントローラは、前記 DRAM セルの同じページに、一つ以上の前記第 1 のタイル及び一つ以上の前記第 2 のタイルを格納するように構成されることを特徴とする請求項 1 4 に記載の方法。

【請求項 1 7】

前記メモリコントローラは、

前記オペランドレジスタに前記一つ以上の前記第 1 のタイルの中の一つの第 1 のタイルを格納し、

10

20

30

40

50

前記オペランドレジスタに格納された前記第 1 のタイル、及び前記 D R A M セルのアレイの前記第 1 のタイルと同じページに格納された前記一つ以上の前記第 2 のタイルのそれぞれに対して、前記 A L U による演算を実行するように更に構成されることを特徴とする請求項 1 6 に記載の方法。

【請求項 1 8】

前記 D R A M バンクの前記 I M C モジュールは、アキュムレータを更に含み、
前記アキュムレータは、累積値を格納するように構成されたアキュムレータレジスタを含み、
前記アキュムレータは、
前記 A L U によって演算された出力を受信し、
前記累積値と前記出力との合計で前記アキュムレータレジスタを更新するように構成され、
前記命令は、前記第 1 のオペランドと前記第 2 のオペランドとの内積を計算することを含み、
前記第 1 のタイルの中の一つの第 1 のタイルは、行データを格納し、
前記第 2 のタイルの中の一つの第 2 のタイルは、列データを含むことを特徴とする請求項 1 7 に記載の方法。

10

【請求項 1 9】

前記第 1 のタイルは、第 1 の数の値を有し、
前記第 2 のタイルは、第 2 の数の値を有し、
前記 D R A M バンクの前記 I M C モジュールは、出力バッファを含み、
前記出力バッファは、前記第 1 の数の値と前記第 2 の数の値との積以上を格納するためのサイズを有し、
前記命令は、前記第 1 のオペランドと前記第 2 のオペランドとの外積を計算することを含み、
前記第 1 のタイルの中の一つの第 1 のタイルは、行データを格納し、
前記第 2 のタイルの中の一つの第 2 のタイルは、列データを含むことを特徴とする請求項 1 7 に記載の方法。

20

【請求項 2 0】

前記第 1 のタイルは、第 1 の数の値を有し、
前記第 2 のタイルは、第 2 の数の値を有し、
前記 D R A M バンクの前記 I M C モジュールは、出力バッファを含み、
前記出力バッファは、前記第 1 の数の値及び前記第 2 の数の値の中の大きい方の値以上を格納するためのサイズを有し、
前記命令は、前記第 1 のオペランドと前記第 2 のオペランドとのテンソル積を計算することを含み、
前記第 1 のタイルの中の一つの第 1 のタイルは、行データを格納し、
前記第 2 のタイルの中の一つの第 2 のタイルは、列データを含むことを特徴とする請求項 1 7 に記載の方法。

30

【請求項 2 1】

前記第 1 のオペランドは、複数の第 1 のタイルに分割され、
前記第 2 のオペランドは、複数の第 2 のタイルに分割され、
各タイルは、複数の値を含み、
前記複数のデータレイアウトの中から異なるページ (D R) データレイアウトを決定した場合、
前記メモリコントローラは、
前記 D R A M セルのアレイの第 1 のページに前記第 1 のタイルのサブセットを格納し、
前記 D R A M セルのアレイの第 2 のページに前記第 2 のタイルのサブセットを格納するように構成されることを特徴とする請求項 1 4 に記載の方法。

40

【請求項 2 2】

50

前記メモリコントローラは、

前記オペランドレジスタに前記第 1 のページから前記第 1 のオペランドの第 1 のタイルを格納し、

前記オペランドレジスタに格納された前記第 1 のオペランドの第 1 のタイル、及び前記第 2 のページからの前記第 2 のオペランドの複数の第 2 のタイルのそれぞれに対して、前記 A L U による演算を実行するように更に構成されることを特徴とする請求項 2 1 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、インメモリコンピューティングに対するデータ配置のためのシステム及び方法に関する。

【背景技術】

【0002】

高帯域幅メモリ (High Bandwidth Memory ; HBM) は、グラフィックス処理装置 (GPU) 用高性能メモリとしてしばしば使用される。HBM は、一般的な DRAM 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 アーキテクチャは、シリコン貫通電極 (through silicon via ; TSV) を使用して接続される複数のスタック DRAM ダイ (例えば、ダイス (dice)) と、HBM のバッファ及び GPU の HBM メモリコントローラとして機能するロジックダイとを含む。メモリ内のプロセス (例えば、インメモリ処理 (in-memory processing)) 機能をメモリシステムに追加することによって、さらに性能が向上する。

【0003】

上述の内容は、本発明の実施形態の背景に対する理解を深めるためだけのものであり、従来技術を構成しない情報を含む。

【先行技術文献】

【特許文献】

【0004】

【文献】特開 2019 - 075101 号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

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、メモリモジュールの性能を向上させることができる方法、及び性能が向上したメモリモジュールを提供することにある。

【課題を解決するための手段】

【0006】

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メモリモジュールは、複数のダイナミックランダムアクセスメモリ (DRAM) バンクを含むメモリダイと、メモリコントローラと、を含み、前記 DRAM バンクの各々は、複数のページに配列された DRAM セルのアレイと、前記複数のページの内の開かれたページの値を格納する行バッファと、入出力 (IO) モジュールと、インメモリコンピューティング (IMC) モジュールと、を含み、前記ページの各々は複数の前記 DRAM セルを含み、前記 DRAM セルの各々はビット値を格納し、前記 IMC モジュールは、前記行バッファ又は前記 IO モジュールからオペランドを受信し、前記オペランド及び複数の算術論理演算から選択された一つの算術論理演算に基づいて出力を計算するように構成された算術論理ユニット (ALU) と、前記 ALU によって計算された前記出力を格納するように構成された結果レジスタと、を含み、前記メモリコントローラは、ホストプロセッサから、第 1 のオペランド、第 2 のオペランド、及び命令を受信し、前記命令に基づいて、複数のデータレイアウトから一つのデータレイアウトを決定し、前記一つのデータレイアウトに従って、前記 DRAM バン

10

20

30

40

50

クに前記第 1 のオペランド及び前記第 2 のオペランドを供給し、前記命令に従って、前記第 1 のオペランド及び第 2 のオペランドに対して前記複数の算術論理演算のうち前記一つの算術論理演算を実行するように、前記 D R A M バンクの前記 I M C モジュールを制御するように構成されたことを特徴とする。

【 0 0 0 7 】

前記複数のデータレイアウトは、一つのオペランド (1 O P) データレイアウトを含み、前記第 1 のオペランドは、前記 D R A M セルに書き込まれ、前記第 2 のオペランドは、前記ホストプロセッサから前記 D R A M バンクの前記 I M C モジュールに直接供給され得る。

【 0 0 0 8 】

前記 I M C モジュールは、オペランドレジスタをさらに含み、前記メモリコントローラは、さらに、前記オペランドレジスタに前記第 1 のオペランドの第 1 のタイルを格納し、前記オペランドレジスタに格納された前記第 1 のオペランドの第 1 のタイル、及び前記第 2 のオペランドの複数の第 2 のタイルの各々に対して、算術論理演算を実行するように構成され得る。

【 0 0 0 9 】

前記第 1 のオペランドは複数の第 1 のタイルに分割され、前記第 2 のオペランドは複数の第 2 のタイルに分割され、前記タイルの各々は複数の値を含み、前記複数のデータレイアウトは同じ行 (S R) データレイアウトを含み、前記メモリコントローラは、前記 D R A M セルのアレイの同じページに一つ以上の前記第 1 のタイル及び一つ以上の前記第 2 のタイルを格納し得る。

【 0 0 1 0 】

前記 I M C モジュールは、オペランドレジスタをさらに含み、前記メモリコントローラは、さらに、前記オペランドレジスタに前記一つ以上の第 1 のタイルの内の一つの第 1 のタイルを格納し、前記オペランドレジスタに格納された前記第 1 のオペランド、及び前記 D R A M セルのアレイの前記第 1 のタイルと同じページに格納された一つ以上の前記第 2 のタイルのそれぞれに対して、算術論理演算を実行するように構成され得る。

【 0 0 1 1 】

前記 D R A M バンクの前記 I M C モジュールは、アキュムレータ (a c c u m u l a t o r) をさらに含み、前記アキュムレータは、累積値を格納するように構成されたアキュムレータレジスタを含み、前記アキュムレータは、前記 A L U によって計算された出力を受信し、累積値と出力との合計で前記アキュムレータレジスタを更新するように構成され、前記命令は、前記第 1 のオペランドと前記第 2 のオペランドの内積を計算することを含み、前記第 1 のタイルの内の前記一つの第 1 のタイルは行データを格納し、前記第 2 のタイルの内の一つの第 2 のタイルは列データを含み得る。

【 0 0 1 2 】

前記第 1 のタイルは第 1 の数の値を有し、前記第 2 のタイルは第 2 の数の値を有し、前記 D R A M バンクの前記 I M C モジュールは、出力バッファを含み、前記出力バッファは、前記第 1 の数の値と前記第 2 の数の値の積以上を格納するためのサイズを有し、前記命令は、前記第 1 のオペランドと前記第 2 のオペランドの外積を計算することを含み、前記第 1 のタイルの内の前記一つの第 1 のタイルは行データを格納し、前記第 2 のタイルの内の一つ第 2 のタイルは列データを含み得る。

【 0 0 1 3 】

前記第 1 のタイルは第 1 の数の値を有し、前記第 2 のタイルは第 2 の数の値を有し、前記 D R A M バンクの前記 I M C モジュールは、出力バッファを含み、前記出力バッファは、前記第 1 の数の値と前記第 2 の数の値の内の大きい方の値以上を格納するためのサイズを有し、前記命令は、前記第 1 のオペランドと前記第 2 のオペランドのテンソル積を計算することを含み、前記第 1 のタイルの内の前記一つの第 1 のタイルは行データを格納し、前記第 2 のタイルの内の一つの第 2 のタイルは列データを含み得る。

【 0 0 1 4 】

前記第 1 のオペランドは複数の第 1 のタイルに分割され、前記第 2 のオペランドは複数の第 2 のタイルに分割され、前記タイルの各々は複数の値を含み、前記複数のデータレイアウトは、異なる行 (D R) データレイアウトを含み、前記メモリコントローラは、前記 D R A M セルのアレイの第 1 のページに前記第 1 のタイルのサブセットを格納し、前記 D R A M セルのアレイの第 2 のページに前記第 2 のタイルのサブセットを格納し得る。

【 0 0 1 5 】

前記 I M C モジュールは、オペランドレジスタをさらに含み、前記メモリコントローラは、さらに、前記オペランドレジスタの前記第 1 のページから前記第 1 のオペランドの第 1 のタイルを格納し、前記オペランドレジスタに格納された前記第 1 のオペランド、及び前記第 2 のページから前記第 2 のオペランドの複数の第 2 のタイルのそれぞれに算術論理演算を実行するように構成されることが好ましい。

10

【 0 0 1 6 】

前記 D R A M バンクの各々の前記 I M C モジュールは、前記 A L U によって計算された前記出力をバッファリングするように構成されたバッファをさらに含むことが好ましい。

【 0 0 1 7 】

前記バッファは、前記結果レジスタのサイズの 4 倍以上のサイズを有し得る。

【 0 0 1 8 】

前記 D R A M バンクの各々の前記 I M C モジュールは、アキュムレータをさらに含み、前記アキュムレータは、累積値を格納するように構成されたアキュムレータレジスタを含み、前記アキュムレータは、前記 A L U によって計算された前記出力を受信し、前記累積値と前記出力との合計で前記アキュムレータレジスタを更新するように構成されることが好ましい。

20

【 0 0 1 9 】

前記メモリモジュールは、シリコン貫通電極によって接続されたメモリダイのスタックを含む高帯域幅メモリ (H B M) モジュールであり、前記メモリダイのスタックは、前記メモリダイを含み得る。

【 0 0 2 0 】

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方法は、インメモリ計算を実行する方法において、メモリモジュールのメモリコントローラによって、第 1 のオペランド、第 2 のオペランド、命令を受信する段階と、前記メモリコントローラによって、前記命令に基づいて複数のデータレイアウトから一つのデータレイアウトを決定する段階と、前記データレイアウトに従って、前記メモリモジュールの少なくとも一つの D R A M バンクに前記第 1 のオペランド及び前記第 2 のオペランドを供給する段階と、前記命令に従って、前記第 1 のオペランド及び前記第 2 のオペランドに対して複数の算術論理演算の内の一つの算術論理演算を実行するように、前記 D R A M バンクの I M C モジュールを制御する段階と、を含み、前記 D R A M バンクの各々は、複数のページに配列された D R A M セルのアレイと、前記複数のページの内の開かれたページの値を格納する行バッファと、 I O モジュールと、前記 I M C モジュールと、を含み、前記ページの各々は前記 D R A M セルを含み、前記 D R A M セルの各々はビット値を格納し、前記 I M C モジュールは、前記行バッファ又は前記 I O モジュールからオペランドを受信し、前記オペランド及び複数の算術論理演算から選択された一つの算術論理演算に基づいて出力を計算するように構成された A L U と、前記 A L U によって計算された前記出力を格納するように構成された結果レジスタと、を含むことを特徴とする。

30

【 0 0 2 1 】

前記複数のデータレイアウトは、一つのオペランド (1 O P) データレイアウトを含み、前記第 1 のオペランドは、前記 D R A M セルに書き込まれ、前記第 2 のオペランドは、前記ホストプロセッサから前記 D R A M バンクの前記 I M C モジュールに直接供給され得る。

【 0 0 2 2 】

前記 I M C モジュールは、オペランドレジスタをさらに含み、前記メモリコントローラ

50

は、さらに、前記オペランドレジスタに前記第 1 のオペランドの第 1 のタイルを格納し、前記オペランドレジスタに格納された前記第 1 のオペランドの第 1 タイル、及び前記第 2 のオペランドの複数の第 2 のタイルの各々に対して、算術論理演算を実行するように構成され得る。

【 0 0 2 3 】

前記第 1 のオペランドは複数の第 1 のタイルに分割され、前記第 2 のオペランドは複数の第 2 のタイルに分割され、前記タイルの各々は複数の値を含み、前記複数のデータレイアウトは、同じ行 (S R) データレイアウトを含み、前記メモリコントローラは、前記 D R A M セルのアレイの同じページに一つ以上の前記第 1 のタイル及び一つ以上の前記第 2 のタイルを格納し得る。

10

【 0 0 2 4 】

前記 I M C モジュールは、オペランドレジスタをさらに含み、前記メモリコントローラは、さらに、前記オペランドレジスタに前記一つ以上の第 1 のタイルの内の一つの第 1 のタイルを格納し、前記オペランドレジスタに格納された前記第 1 のオペランドと、前記 D R A M セルのアレイの前記第 1 のタイルと同じページに格納された前記一つ以上の第 2 のタイルのそれぞれに対して、算術論理演算を実行するように構成され得る。

【 0 0 2 5 】

前記少なくとも一つの D R A M バンクの前記 I M C モジュールは、アキュムレータをさらに含み、前記アキュムレータは、累積値を格納するように構成されたアキュムレータレジスタを含み、前記アキュムレータは、前記 A L U によって計算された前記出力を受信し、前記累積値と前記出力との合計で前記アキュムレータレジスタを更新するように構成され、前記命令は、前記第 1 のオペランドと前記第 2 のオペランドの内積を計算する段階を含み、前記第 1 のタイルの内の前記一つの第 1 のタイルは行データを格納し、前記第 2 のタイルの内の一つの第 2 のタイルは列データを含み得る。

20

【 0 0 2 6 】

前記第 1 のタイルは第 1 の数の値を有し、前記第 2 のタイルは第 2 の数の値を有し、前記少なくとも一つの D R A M バンクの前記 I M C モジュールは、出力バッファを含み、前記出力バッファは、前記第 1 の数の値と前記第 2 の数の値の積以上を格納するためのサイズを有し、前記命令は、前記第 1 のオペランドと前記第 2 のオペランドの外積を計算する段階を含み、前記第 1 のタイルの内の前記一つの第 1 のタイルは行データを格納し、前記第 2 のタイルの内の一つの第 2 のタイルは列データを含み得る。

30

【 0 0 2 7 】

前記第 1 のタイルは第 1 の数の値を有し、前記第 2 のタイルは第 2 の数の値を有し、前記少なくとも一つの D R A M バンクの前記 I M C モジュールは、出力バッファを含み、前記出力バッファは、前記第 1 の数の値と前記第 2 の数の値の内の大きい方の値以上を格納するためのサイズを有し、前記命令は、前記第 1 のオペランドと前記第 2 のオペランドのテンソル積を計算する段階を含み、前記第 1 のタイルの前記一つの第 1 のタイルは行データを格納し、前記第 2 のタイルの内の一つの第 2 のタイルは列データを含み得る。

【 0 0 2 8 】

前記第 1 のオペランドは複数の第 1 のタイルに分割され、前記第 2 のオペランドは複数の第 2 のタイルに分割され、前記タイルの各々は複数の値を含み、前記複数のデータレイアウトは、異なる行 (D R) データレイアウトを含み、前記メモリコントローラは、前記 D R A M セルのアレイの第 1 のページに前記第 1 のタイルのサブセットを格納し、前記 D R A M セルのアレイの第 2 のページに前記第 2 のタイルのサブセットを格納し得る。

40

【 0 0 2 9 】

前記 I M C モジュールは、オペランドレジスタをさらに含み、前記メモリコントローラは、さらに、前記オペランドレジスタの前記第 1 のページから前記第 1 のオペランドの第 1 のタイルを格納し、前記オペランドレジスタに格納された前記第 1 のオペランド、及び前記第 2 のページから前記第 2 のオペランドの複数の第 2 のタイルのそれぞれに対して、算術論理演算を実行するように構成されることが好ましい。

50

【発明の効果】

【 0 0 3 0 】

本発明によるインメモリコンピューティングに対するデータ配置のための方法及びその方法が適用されたメモリモジュールによれば、メモリモジュールの性能を向上させることができる。

【図面の簡単な説明】

【 0 0 3 1 】

【図 1】本発明の一実施形態によるメモリ（例えば、HBM）システムのアーキテクチャを示すブロック図である。

【図 2 A】本発明の一実施形態による、埋め込み算術論理ユニット（e m b e d d e d A L U）を有するメモリバンクの概略ブロック図である。

10

【図 2 B】本発明の一実施形態によるDRAMセルの例を示す回路図である。

【図 2 C】DRAMセルのアレイ、行デコーダ、IOSA、IMCモジュール、及び列デコーダを含む、本発明の一実施形態によるDRAMバンクの概略図である。

【図 3】本発明の一実施形態によるDRAMブロックのアレイを示す概略図である。

【図 4 A】第 1 のオペランド（行列 A）がDRAMに格納され、第 2 のオペランド（行列 B）が内蔵されたメモリモジュールの外部からブロードキャストされる場合のGEMMに対するデータの配置を、本発明の一実施形態による統合されたインメモリコンピューティング（IMC）と共に概略的に示す図である。

【図 4 B】双方のオペランド（行列 A 及び行列 B）がメモリモジュールのDRAMの同じページに格納されている場合のGEMMに対するデータの配置を、本発明の一実施形態による統合されたIMCと共に概略的に示す図である。

20

【図 4 C】双方のオペランド（行列 A 及び行列 B）がメモリモジュールのDRAMの異なるページに格納されている場合のGEMMに対するデータの配置を、本発明の一実施形態による統合されたIMCと共に概略的に示す図である。

【図 5 A】一つのオペランド（1OP）データレイアウトにおける行列 A の第 1 の行と行列 B の第 1 の列との乗算の概略図である。

【図 5 B】1OPデータレイアウトにおける行列 A の第 1 の行の第 1 の値と行列 B の各列の第 1 の値との乗算をデータの再利用と共に示した概略図である。

【図 6 A】同じ行（SR）データレイアウトにおける行列 A の第 1 の行と行列 B の第 1 の列との乗算の概略図である。

30

【図 6 B】SRデータレイアウトにおける行列 A の第 1 の行の第 1 の値と行列 B の各列の第 1 の値に対するデータの再利用による乗算の概略図である。

【図 7 A】異なる行（DR）データレイアウトにおける行列 A の第 1 の行と行列 B の第 1 列との乗算の概略図である。

【図 7 B】DRデータレイアウトにおける行列 B の各列の第 1 の値と行列 A の第 1 の行の第 2 の値により、行列 B の各列の第 2 の値による行列 A の第 1 の行の第 1 の値のデータ再利用による乗算の概略図である。

【図 8】DRAMバンクのIMCモジュールの概略ブロック図である。

【図 9】本発明の一部の実施形態による、同じ行（SR）レイアウトを使用する計算を説明するためのGEMMの一例の概略図である。

40

【図 10】本発明の実施形態による、IMCモジュールを有するDRAMバンク内のデータの配置を制御するための方法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 3 2 】

本発明の特徴及びそれを達成するための方法は、以下の実施形態の詳細な説明及び図面を参照することにより、容易に理解される。以下、図面を参照しながら実施形態をより詳細に説明するが、図面全体にわたって同じ参照番号は同じ要素を指す。しかし、本発明は、様々な異なる形態で実施することができ、本明細書で説明する実施形態に限定されない。むしろ、これらの実施形態は、本発明が徹底且つ完全なものとなり、本発明の態様及び

50

特徴が当業者に十分に伝わるように例として提供するものである。したがって、本発明の態様及び特徴の完全な理解のために、当業者に不必要なプロセス、要素、及び技術は説明しない。特に記載がない限り、同様の参照番号は、図面及び明細書に記載された説明全体を通じて同様の要素を指すので、その説明は繰り返さない。図面において、要素、層、及び領域の相対的なサイズは、明確にするために誇張されている場合がある。

【0033】

以下の説明では、説明を目的として、様々な実施形態の完全な理解のために、多数の特定の詳細な説明を提示する。しかし、様々な実施形態は、これらの特定の詳細な説明がなくても実施でき、又は一つ以上の同等の構成で実施してもよい。他方、公知の構造及び装置は、様々な実施形態を曖昧にすることを避けるために、ブロック図の形で示す。

10

【0034】

本明細書で使用される用語は、特定の実施形態を説明するためのものであり、本発明を限定することを意図するものではない。単数形の表現は、文脈がそうでないことを明確に示さない限り、複数形の表現も含む。本明細書で使用される「含む」及び「有する」などの用語は、開示する特徴、数字、ステップ、演算、要素、及び/又は構成要素の存在を特定するが、一つ以上の他の特徴、数字、ステップ、演算、要素、構成要素、及び/又はそれらの組み合わせの存在又は追加を排除するものではない。本明細書で使用される用語「及び/又は」は、一つ以上の関連して挙げられた項目の任意の且つ全ての組み合わせを含む。

【0035】

20

本明細書に記載する本発明の実施形態による電子又は電気装置及び/又は任意の他の関連装置又は構成要素は、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（Application-Specific Integrated Circuit））、ソフトウェア、又はそれらの組み合わせを利用して実装される。例えば、一部の環境では、これらの装置の様々な構成要素は、一つの集積回路（IC）チップ上に、又は別のICチップ上に形成される。また、これらの装置の様々な構成要素は、フレキシブルプリント回路フィルム、テープキャリアパッケージ（TCP）、プリント回路基板（PCB）上に実装されてもよく、又は一つの基板上に形成されてもよい。さらに、これらの装置の様々な構成要素は、コンピュータプログラム命令を実行し、本明細書で説明する様々な機能を行うために他のシステムの構成要素と相互作用する一つ以上のコンピューティングデバイス内における一つ以上のプロセッサ上で実行されるプロセス又はスレッドである。コンピュータプログラム命令は、例えば、ランダムアクセスメモリ（RAM）等の標準的なメモリデバイスを使用してコンピューティングデバイス内に実装されるメモリに格納される。コンピュータプログラム命令はまた、例えば、CD-ROM、フラッシュドライブなどの他の非一時的なコンピュータ読取可能媒体に格納される。また、当業者は、本発明の例示的な実施形態の思想及び範囲内で様々なコンピューティングデバイスの機能が単一のコンピューティングデバイス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デバイスの機能が一つ以上の他のコンピューティングデバイスにわたって分散されることを認識すべきである。

30

【0036】

40

特に定義しない限り、本明細書で使用する技術用語又は科学用語を含むすべての用語は、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に理解するものと同じ意味を有する。また、通常使用される辞書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、理想的又は過度に公式的な意味に解釈されない。

【0037】

本発明の一部の態様は、一般に、インメモリコンピューティング（in-memory compute）の文脈におけるメモリ内のデータ配置の管理に関する。インメモリコンピューティングに関する内容の一例は、メモリを含むDRAM（Dynamic Random Access Memory）ダイと、HBMロジックダイ上のALU及びメモリ

50

コントローラとを含むHBMシステムであり、ここで、HBMロジックダイ上のALUは、インメモリコンピューティングを実行する。HBMロジックダイのメモリコントローラは、DRAMダイのメモリへのデータの格納、及びDRAMダイからのデータの読み取りを制御する。

【0038】

明確化のため、本明細書で使用する用語「インメモリコンピューティング」は、DRAMダイに格納されたデータを使用して、外部データバスを通過することなく、高帯域幅メモリモジュール等のメモリモジュール内で計算を実行することを意味する。比較コンピュータシステムでは、プロセッサは外部DRAMデータバスを介してメインメモリ（例えば、DRAM）に接続され、メインメモリからのデータへのアクセスは、プロセッサ内のレジスタファイル内のデータ及び／又はプロセッサにより近いハードウェアキャッシュ（例えば、L1キャッシュ、L2キャッシュ）内のデータへのアクセスよりも大幅に遅い（例えば、数十倍遅い）。メモリ又はその近くにさらなるプロセッサ（例えば、「インメモリプロセッサ（in-memory processor）」）を配置することで、外部バスを通過することによって惹起される遅延を回避することができ、これにより、高性能の計算が達成される。

【0039】

本発明の実施形態の態様は、ALU等の計算回路（computational circuitry）をDRAMバンクと同じダイ上に配置する（例えば、それぞれのDRAMバンクのセンス増幅器又は行バッファに直接接続される）ことに関する。

【0040】

DRAMモジュールの設計及び性能特性により、メモリ内のデータの特定の配列は、インメモリコンピューティングの性能に影響を与え得る。したがって、本発明の実施形態の一部の態様は、メモリモジュール（例えば、HBMメモリモジュール）のDRAMモジュール内にデータを配置するためのシステム及び方法に関し、ここでデータの配置は、IMCモジュールによって実行される計算の特性に基づいて制御される。

【0041】

図1は、本発明の一実施形態によるメモリ（例えば、HBM）システムのアーキテクチャを示すブロック図である。

【0042】

図1を参照すると、本発明の実施形態は、FIM（Function-In-Memory）メモリシステム100のためのシステムを提供する。メモリシステム100（又は、HBMシステム）は、メモリモジュール110（又は、HBMモジュール）に統合されるさらなる（additional）コンピューティングリソースをサポートする。例えば、様々な実施形態において、メモリシステム100は、一部のデータコンピューティング及び移動がメモリ内で実行されるようにし、且つ大容量メモリスクラッチパッド（high-capacity memory scratchpad）を提供する。メモリシステム100は、グラフィックス処理装置（GPU）又は中央処理装置（CPU）等のホストプロセッサ170に接続された少なくとも一つのメモリモジュール110を含む。様々な実施形態において、メモリモジュール110は、内部メモリバス130を介してメモリコントローラ140（例えば、ロジックダイ上）に接続された一つ以上のDRAMダイ120を含む。様々な実施形態において、ホストプロセッサ170は、メモリモジュール110とインタフェースするためのホストメモリコントローラ180（又は、ホストコントローラ）を含む。しかし、本発明の実施形態は、これに限定されるものではない。たとえば、ホストメモリコントローラ180は、ホストプロセッサ170から分離されてもよい（例えば、ホストプロセッサ170とは別のダイ又は同じダイとして）。

【0043】

様々な実施形態によると、メモリコントローラ140は、ホストプロセッサ170からの命令の実行を調整するように構成される。命令は、通常の命令とFIM命令の両方を含む。例えば、通常の命令（例えば、メモリ内機能（function-in-memory

10

20

30

40

50

y ; F I M) 命令ではなく、伝統的なロード及びストア機能)は、ホストメモリコントローラ180によって送信され、メモリコントローラ140によって受信されて、通常の方法で実行される。例えば、通常の命令は、外部バス190を介して受信したデータをDRAMダイ120に格納する命令、及びDRAMダイ120からデータを検索し、外部バス190を介してホストプロセッサ170にデータを送信する命令を含む。一部の実施形態において、通常の命令及びFIM命令は、DRAMダイの特定の位置(例えば、特定のバンクの特定のページ)にデータを格納する動作を含む。これらのデータは2つの異なるオペランドを含み、ここで、それぞれのオペランドは複数の値(例えば、浮動小数点又は整数値)を含み、以下でより詳細に説明するように、これらのオペランドの値は、様々なデータ配置戦略に従って、DRAMダイの異なる位置に分散して配置(d i s t r i b u t e)される。

10

【0044】

本発明の実施形態の態様は、IMC(in-memory compute)の使用に関する。いくつかの比較HBMシステムは、DRAMダイ120の外部(例えば、メモリコントローラ140に配置)にあるALUを含み、ALUが外部バス190を通過することなく、1つ以上のDRAMダイ120に格納されたデータに対する演算(operation)(例えば、算術演算(arithmetic operations))を実行できるように、DRAMダイ120のメモリバンクによって共有される。例えば、メモリコントローラ140は、DRAMダイ120の異なる部分の間でデータを移動又はコピーするために、データ移動演算(例えば、ロード/ストア ペア命令(load/store pair instructions))を実行する。例えば、メモリコントローラは、ALUを利用する計算FIM命令(例えば、アトミック命令及びALU命令)の実行を調整することで、元々は複数の通常命令であったFIM命令を実行する。別の例として、一部の場合、ホストプロセッサ170から受信したFIM命令は、IMCが統合されたメモリモジュールに、学習済みの機械学習モデル(例えば、ニューラルネットワーク(neural network))を使用して、ニューラルネットワークの訓練中に逆伝播(backpropagation)を実行させるか、又は2つの行列を乗算するために、供給された入力に基づいて推論を計算する等の特定の計算を実行させる。このような場合、メモリコントローラ140は、受信したデータ(例えば、命令のオペランド)をDRAMバンクの特定のページに格納し、FIM命令に関連するデータを格納するDRAMバンクに特定のALU演算を提供することによって、これらの命令の実行を調整する。メモリコントローラ140は、特定のFIM命令を実行する際の計算性能を向上させる方法で、受信したデータ(オペランド)をDRAMバンクの特定のページに配置する。その結果は、DRAMダイ120に格納されるか、又は外部バス190を介してホストプロセッサ170に戻される。

20

30

【0045】

本発明の実施形態の一部の態様は、IMCモジュールをDRAMダイ120のメモリバンクに統合することによってメモリ境界演算を加速することに関し、これにより、DRAMダイ120とメモリコントローラ140との間で内部メモリバス130を通過すること(traversal)を避ける。例えば、IMCモジュールは、データを保持するDRAMバンクと同じ物理的半導体ダイ上にある。それぞれのDRAMバンクは、関連するIMCモジュールを有し、それによって、データが内部メモリバス130をメモリコントローラに通過することなく(例えば、DRAMバンクからデータを送信することなく)、DRAMバンクに格納されたデータに対して計算が実行できる。また、上記計算は、DRAMバンクのIMCモジュール間で並列化(parallelized across)できる。

40

【0046】

本発明の実施形態の態様は高帯域幅メモリに関して説明しているが、実施形態はこれに限定されず、他のタイプのDRAMシステムにおいて、DRAMダイにIMCモジュールを統合することにも適用される。

50

【 0 0 4 7 】

図 2 A は、本発明の一実施形態による、埋め込み算術論理ユニット (e m b e d d e d A L U) を有するメモリバンクの概略ブロック図である。図 2 A に示すように、D R A M バンク 2 0 0 は、行及び列 (又はページ及び列) に配列された D R A M セル 2 1 0 のアレイを含む。図 2 A に例示するように、D R A M バンク 2 0 0 は、 n 行 (又はページ) 及び m 列に配列された D R A M セル 2 1 0 を含む。複数のビットライン $B_1 \sim B_m$ は列方向に沿って延在し、複数の行イネーブルライン (*row enable lines*) $R_1 \sim R_n$ はアレイの行方向に沿って延在し、ビットラインと交差する。各ビットラインは、対応する列のすべてのセル (D R A M セル 2 1 0) に接続される (例えば、アレイの i 番目の列のすべてのセルは、ビットライン B_i に接続される)。同様に、各行イネーブルライン $R_1 \sim R_n$ は、対応する行の各 D R A M セル 2 1 0 に接続される (例えば、アレイの j 番目の行又はページのすべてのセルは、行イネーブルライン R_j に接続される)。また、D R A M バンク 2 0 0 の D R A M セル 2 1 0 の行は、D R A M ページと呼ばれる。

10

【 0 0 4 8 】

図 2 B は、本発明の一実施形態による D R A M セルの例を示す回路図である。それぞれの D R A M セル 2 1 0 は、一般に、データ電圧を格納するためのコンデンサ 2 1 2 (例えば、ビット値、ここで各コンデンサは 0 ビットを表す電圧又は 1 ビットを表す電圧を格納する) と、コンデンサ 2 1 2 にデータ電圧を送信するためのスイッチ 2 1 4 とを含むものとしてモデル化される。図 2 B に示す特定の D R A M セル 2 1 0 は、アレイの i 番目の行及び j 番目の列にある。したがって、図 2 B に示す D R A M セル 2 1 0 のスイッチ 2 1 4 は、 i 番目のビットライン B_i とコンデンサ 2 1 2 の一方の端子との間に接続され、コンデンサ 2 1 2 の他方の端子は接地に接続される。図 2 B に示すように、D R A M セル 2 1 0 のスイッチ 2 1 4 のゲート電極は、 j 番目の行イネーブルライン R_j に接続され、スイッチ 2 1 4 がオンになると、コンデンサ 2 1 2 がビットライン B_i に接続される。

20

【 0 0 4 9 】

図 2 A を再び参照すると、D R A M バンク 2 0 0 は、行イネーブルライン $R_1 \sim R_n$ に接続された行デコーダ 2 2 0 を含み、行デコーダ 2 2 0 は、例えば、メモリコントローラ 1 4 0 から供給される行アドレスに対応する行イネーブルラインの特定の一つに、行イネーブル信号を供給するように構成される。D R A M セル 2 1 0 のアレイの特定の行 r (又はページ) にデータを書き込むか又は読み取る際、行デコーダ 2 2 0 は、特定の行 (又はページ) に対応する行イネーブルラインに行イネーブル信号を供給する。データを書き込む際には、特定の行 (又はページ) がイネーブルされている間、書き込まれるデータに対応する電圧がビットライン $B_1 \sim B_m$ に供給される。同様に、D R A M セル 2 1 0 のアレイの特定の行 (又はページ) からデータを読み取る際、コンデンサ 2 1 2 に格納された電圧に対応する電圧は、ビットライン $B_1 \sim B_m$ を沿って送信され、センス増幅器 2 3 2 を含む入出力センス増幅器層 2 3 0 (又は I O S A) によって読み取られる。センス増幅器 2 3 2 の各センス増幅器は、ビットラインの内の対応する一つに接続される (例えば、センス増幅器 2 3 2 は、 m 個のセンス増幅器を含む)。例えば、一部の実施形態において、D R A M セル 2 1 0 のアレイは、 $8,192$ 個の列、及び $8,192$ 個の対応するセンス増幅器 2 3 2 に接続された $8,192$ 個の対応するビットライン (例えば、ビットライン $B_1 \sim B_{8,192}$) を含む (例えば、各ページには $8,192$ ビット又は 8 キロビットのデータを格納できる)。センス増幅器 2 3 2 は、「プリチャージ (*precharge*)」コマンドによって消去されるまで、現在の行 (又はページ) から読み取られたデータを格納するので、センス増幅器 2 3 2 は「行バッファ (*row buffer*)」と呼ばれる。

30

40

【 0 0 5 0 】

列デコーダ 2 4 0 は、マルチプレクサ 2 3 4 を使用してデータ列のサブセット (*subset*) を選択するために使用され、読み取られたデータは、そのデータの計算を実行するために、グローバル I O 層 2 3 6 を介して I M C モジュール 2 5 0 に供給される。例えば、一部の実施形態において、列デコーダ 2 4 0 及びマルチプレクサ 2 3 4 は、D R A

50

Mセル210の8, 192個の列から256ビット(256b)のデータの選択を可能にする。

【0051】

センス増幅器232に現在格納されているページとは異なるDRAMバンク200のページからデータをロードする場合、「プリチャージ」(PRE)コマンドは現在のページを閉じ、次のアクセスのためにDRAMバンク200を準備するために使用される。次に、「活性化(activate)」(ACT)コマンドを使用して、DRAMバンクの特定の行又はページを開き、その開かれたページのデータをセンス増幅器232に格納する。その後、データがセンス増幅器232から読み取られて(READ)、IMCモジュール250に送信される。

10

【0052】

一方、すでに開いているページからIMCモジュール250にデータをロードする場合には、例えば、列デコーダ240を使用してセンス増幅器232に既に格納されているデータの適切なサブセットを選択することにより(PRE及びACTコマンドは省略してもよい)、READコマンドでデータをロードするのに十分である。

【0053】

様々な実施形態によれば、IMCモジュール250(又はALU&Reg)は、ALU252及び1つ以上のレジスタを含む。図2Aに示す実施形態において、IMCモジュール250は、オペランドレジスタ(Rop)254(又は入力バッファ)及び結果レジスタ(Rz)256を含む。マルチプレクサ257及び258は、(例えば、ALU252に対する第1のオペランド及び第2のオペランドとして)ALU252の2つの入力へのデータの流れを制御するために使用される。例えば、図2Aに示す実施形態において、オペランドレジスタ(Rop)254は、ALU252の第1のオペランド入力に接続され、第1のマルチプレクサ257は、センス増幅器232からグローバルIO層236を介して、又は入出力(IO)モジュール260(又は、書き込み入力/出力及び読み取り入力/出力、又はWIO及びRIO)を介して、外部ソース(例えば、ホストプロセッサ)からのデータを書き込むために、オペランドレジスタ(Rop)254に接続される。図2Aに示す実施形態のように、第2のマルチプレクサ258もまた、グローバルIO層236を介してセンス増幅器232から、又はIOモジュール260を介して外部ソースから、ALU252の第2のオペランド入力に直接データを供給するように構成される。ALU252は、その計算を結果レジスタ(Rz)256に出力し、データは、結果レジスタ(Rz)256からグローバルIO層236を介してDRAMセル210に書き戻されるか、又はIOモジュール260(又は、WIO及びRIO)を介してホストプロセッサ170に送信される。

20

30

【0054】

一部の実施形態によると、ALU252は、様々な計算演算(例えば、簡単な計算コマンド)を実行するように構成される。例えば、ALU252は、算術演算、ビット単位(bitwise)、シフト演算(shift operations)等を実行するように構成された16ビットALU、32ビットALU、又は64ビットALUである。様々な実施形態において、ALU252は、整数演算(integer operations)、浮動小数点演算(floating point operations)、又はその両方を実行する回路を含む。例えば、ALU252は、ADD(+), SUBTRACT(-), MULTIPLY(x), 及びDIVIDE(÷)等の算術演算、AND(&), OR(|), XOR(^), 及びNOT(~)演算、並びにテンソル演算(tensor operations)等のビット演算を実行するように構成される。また、一部の実施形態において、ALU252は、単一命令、複数データ(SIMD)、又はデータのベクトルに対する演算を並列に実行するためのベクトル命令を実装する。本発明の実施形態によるALU252によって実装されるベクトル演算の例は、内積('), 外積(⊗)

40

50

、整流線形ユニット (R e L U)、平方 (v s S q r)、及び平方根 (v s S q r t) を含む。 A L U 2 5 2 は、アトミック及び非アトミック演算に利用される。以下の表 1 は、本発明の一部の実施形態による A L U 2 5 2 によってサポートされる演算を挙げている。

【 0 0 5 5 】

【表 1】

I D	演算	説明
0	R o p = G I O	R o p に格納された列読み取りデータ
1	R o p = W I O	R o p に格納された列書き込みデータ
2	R o p = R z	A L U 出力 R z を R o p にコピー
3	G I O = R z	バンクに書き戻す
4	R I O = G I O	D Q 出力への通常の読み取り
5	R I O = R z	R z を D Q 出力に駆動
6	R z = R o p (o p) G I O	R o p 及びバンクからのデータを使用した演算
7	R z = R o p (o p) W I O	R o p 及びブロードキャストデータを使用した演算
8	R z = W I O (o p) G I O	ブロードキャストデータ及びバンクを使用した演算
9	G I O = W I O	D Q 入力からの通常の書き込み

10

20

【 0 0 5 6 】

図 2 C は、上述した D R A M セル 2 1 0 のアレイ、行デコーダ 2 2 0、入出力センス増幅器層 (I O S A) 2 3 0、I M C モジュール (A L U & R e g) 2 5 0、及び列デコーダ 2 4 0 を含む、本発明の一実施形態による D R A M バンク 2 0 0 の概略図である。

【 0 0 5 7 】

図 3 は、本発明の一実施形態による D R A M ブロックのアレイを示す概略図である。図 3 に示す実施形態のように、16 個の D R A M バンク 2 0 0 は、4 × 4 アレイに配列され、D R A M バンク A ~ P としてラベル付けされ、B G 0 (D R A M バンク A、B、C、及び D を含む)、B G 1 (D R A M バンク E、F、G、及び H を含む)、B G 2 (D R A M バンク I、J、K、及び L を含む)、並びに B G 3 (D R A M バンク M、N、O、及び P を含む) としてラベル付けされた 4 つのバンクグループに配列される。図 2 C に関して上述したように、図 3 に示す各 D R A M バンク 2 0 0 は、D R A M ダイ 1 2 0 内で (例えば、外部バスを通過することなく) 計算を実行するための I M C モジュール 2 5 0 を含む。また、図 3 に示すように、D R A M ダイ 1 2 0 は、外部ソースからのデータを (例えば、D R A M ダイ 1 2 0 を複数の他のスタック型 D R A M ダイ 1 2 0 及びメモリコントローラ 1 4 0 に接続するシリコン貫通電極又は T S V を介して) ブロックの 4 つの列に多重化するように構成されたマルチプレクサ 3 0 0 (例えば、4 : 1 マルチプレクサ) をさらに含む。例えば、マルチプレクサ 3 0 0 は、すべての D R A M バンク 2 0 0 に、2 5 6 ビット (2 5 6 b) データベクトルをブロードキャスト (b r o a d c a s t) するか、又は D R A M バンク 2 0 0 の特定の列 (B G 0、B G 1、B G 2、又は B G 3) にデータベクトル

30

40

【 0 0 5 8 】

D R A M ダイ 1 2 0 に統合された I M C モジュール 2 5 0 等の I M C を含むメモリシステム 1 0 0 は、演算を実行するために、データが外部バス (例えば、バス 1 9 0) のボトルネックを通過する必要がないので、メモリ境界のホスト演算の性能を加速させる。ただし、I M C は、依然として、A L U パイプライン処理 (p i p e l i n i n g) 及び D R A M プロセスの形態のコンピューティングオーバーヘッド、並びにデータ配置及び D R A M タイミングの形態のメモリオーバーヘッドに遭遇する。

【 0 0 5 9 】

したがって、本発明の実施形態の態様は、統合された I M C モジュール 2 5 0 によって

50

インメモリコンピューティングを実行する際に、DRAMタイミングオーバーヘッドの影響を回避又は減少させるために、DRAM内にデータを配置するシステム及び方法に関する。本発明の実施形態の一部の態様は、改善された性能を達成するためのソフトウェア及びハードウェアの共同設計に関する。

【0060】

様々な実施形態において、メモリコントローラ140は、DRAMダイ120への演算及びDRAMダイ120からの演算を提供し、データの入力及び出力を管理する。したがって、本発明の実施形態の一部の態様は、ホストプロセッサ170によってメモリシステム100のメモリコントローラ140に提供された命令に従って、DRAMダイ120内にデータを配置するように構成されたメモリコントローラ140に関する。例えば、本発明の実施形態の一部の態様は、インメモリコンピューティングと、APIを使用してプログラムのソースコードをコンパイル又は解釈する際に、APIへの呼び出しに従ってデータを配置するようにメモリコントローラ140を制御するためのコマンドを生成するように構成されたコンパイラ（例えば、データコンパイラ）と、を有するHBMと相互作用するアプリケーション・プログラミング・インターフェース（API）を提供することに関する。例えば、APIは、GEMM（General Matrix-Matrix Multiplication）を実行するための関数呼び出し（function call）を提供し、コンパイラは、データに対して実行される演算（例えば、内積、外積、行列乗算など）を含む要素、及びデータのサイズ（例えば、データがメモリのページに適合するか否か）に基づくデータのナイーブな（naive）配置よりも性能が改善する方法で、オペランド行列を表すデータをDRAMダイ120に配置するようにメモリコントローラ140を制御するコマンドのシーケンスを生成する。APIを介してIMCでDRAMを使用するようにソフトウェアを作成する際に、本発明の一部の実施形態によるコンパイラ又はデータコンパイラは、DRAMダイ120の特定の位置にデータを配置し、DRAMダイ120の計算を実行するためのIMCモジュールを制御し、結果を格納するために、ソフトウェアのソースコードの少なくとも一部をメモリコントローラ140によって実行されるコマンドに変換する。

【0061】

一実施形態によると、16レーン（lane）のALUは、半精度浮動小数点（half-precision floating-point）（FP-16）計算で8 GFLOPS（ギガ浮動小数点演算/秒）のピーク性能を達成する。（本発明の実施形態によるIMCモジュールの性能をFLOPSに関して本明細書に説明しているが、本発明の実施形態は、浮動小数点演算を実行することに限定されず、様々なデータレイアウトの相対的な性能は、例えば、整数演算の実行する際と似ている。）したがって、第2世代高帯域幅メモリ標準（HBM2）を使用する4つのダイ（4H又は4-Hi）のスタックを使用する、本発明の一実施形態によるIMCの実装は、FP-16計算（ダイあたり256バンク×4ダイのスタック=1,024バンク、各バンクは対応する16レーンのALUを有する）で8 TFLOPS（テラ浮動小数点演算/秒）を達成する。

【0062】

TFLOPSで測定されるピーク計算性能は、異なるデータレイアウトシナリオの下で変わる。

一番目の場合（2OPと表記）、メモリモジュール110の外部からの2つのオペランドがインメモリコンピューティングALUに完全に供給され、その結果はバッファリングされて（buffered）完全に累積され、これにより、上述の8 TFLOPSのピーク計算性能が得られる。

【0063】

二番目の場合（1OPと表記）、第1のオペランドはHBMの外部からIMCに完全に供給されるが、第2のオペランドはDRAMダイ120の任意の位置から読み取られる。これは、約6.5 TFLOPSにピーク計算性能を低下させる。

【0064】

10

20

30

40

50

三番目の場合（D R と表記）、双方のオペランドはD R A Mの異なるページに配置され、結果がD R A Mに書き戻される。このシナリオでは、約0.8 T F L O P S が測定された計算性能である（例えば、2 O P の場合よりも1桁遅い）。

【0065】

四番目の場合（S R と表記）で、双方のオペランドはD R A Mブロックの同じ行又はページに位置する。これはD R の場合よりも性能が大幅に向上し、約3.3 T F L O P S になる。

【0066】

したがって、データの考慮に基づいて、様々なデータレイアウトのトレードオフが行われる。例えば、1 O P の場合、第1のオペランドがD R A Mにあり、第2のオペランドがH B Mの外部からブロードキャストされると、性能が高くなる（上記のように約6.5 T F L O P S ）が、ホストプロセッサ170にオーバーヘッドが発生し、第2のオペランドをH B Mに提供しなければならない。

【0067】

図4Aは、第1のオペランド（行列A）がD R A Mに格納され、第2のオペランド（行列B）がメモリモジュールの外部からブロードキャストされる場合のG E M Mに対するデータの配置を、本発明の一実施形態による統合されたI M Cと共に概略的に示す図である。例示のために、D R A Mバンク200-Oをより詳細に示している。図4Aに示すように、第1のオペランド行列Aに関連するデータは、D R A Mバンク200-Oの一つのページ401（例えば、第1の行又は第1のページ）に配置され、第2のオペランド行列Bに関連するデータは、（例えば、マルチプレクサ300を介して）D R A Mダイ120の外部からブロードキャストされる。計算結果Cは、D R A Mバンク200-Oの異なるページ402に配置される。

【0068】

双方のオペランドがD R A M内にある場合、双方のオペランドを同じページ又は同じ行（S R ）に配置することは、以下で詳細に説明するように、実行する必要があるP R E 及びA C T 演算の数を部分的に減らすことによって、計算性能を向上させる（例えば、約3.3 T F L O P S ）が、D R A Mの正しい部分にデータを配置することに関してより多くの制約を課す。

【0069】

図4Bは、双方のオペランド（行列A及び行列B）が、メモリモジュールのD R A Mの同じページに格納されている場合のG E M Mに対するデータの配置を、本発明の一実施形態による統合されたI M Cと共に概略的に示す図である。例示のために、D R A Mバンク200-Oをより詳細に示している。図4Bに示すように、第1のオペランド行列A及び第2のオペランドの行列Bに関連するデータは、D R A Mバンク200-Oのページ411（例えば、第1の行又は第1のページ）に配置される。より詳細には、ページ411の前半部は第1のオペランド行列Aからのデータで満たされ、ページ411の後半部は第2のオペランド行列Bからのデータで満たされる。行列A及び行列Bが図4Aに関して上述したものと同じサイズと仮定すると、残りのデータを格納するためにさらにページが必要になる。このように、第1のオペランド行列A及び第2のオペランドの行列Bの双方に関連するデータもページ412に配置される。行列乗算の結果Cは、D R A Mバンク200のページ413に格納される。

【0070】

一方、オペランドを異なるページに配置することは、より柔軟で、レイアウトに対する制約を減らす（例えば、固定サイズのメモリのページにきちんと適合しないサイズを有するデータに適合する）が、一般に計算性能が低下する。

【0071】

図4Cは、双方のオペランド（行列A及び行列B）が、メモリモジュールのD R A Mの異なるページに格納されている場合のG E M Mに対するデータの配置を、本発明の一実施形態による統合されたI M Cと共に概略的に示す図である。例示のために、D R A Mバン

ク 2 0 0 - 0 をより詳細に示している。図 4 C に示すように、第 1 のオペランド行列 A に関連するデータは、ページ 4 2 1 に配置され、第 2 のオペランド行列 B に関連するデータは、ページ 4 2 2 に配置され、結果 C はページ 4 2 3 に配置される。

【 0 0 7 2 】

図 4 A、図 4 B、及び図 4 C に示す異なるデータ配置戦略 (1 O P、S R、及び D R) の様々な性能への影響は、以下の図 5 A、図 5 B、図 6 A、図 6 B、図 7 A、及び図 7 B を参照して、より詳細に説明する。例えば、図 5 A に示すように、行列 A は $M \times K$ 行列であり、行列 B は $K \times N$ 行列であるので、行列 A と行列 B の積である行列 C のサイズは、 $M \times N$ である。説明のために、以下の例では $K = 5$ の場合を説明しているが、本発明の実施形態はそれに限定されない。標準行列乗算に従って、結果行列 C の左上の値 C 0 0 は、行列 A (各位置 A_{ij} は、例えば 1 6 個の半精度浮動小数点値のベクトル又は「タイル」を水平順に示し、図 5 A、図 5 B、図 6 A、図 6 B、図 7 A、及び図 7 B は、A 0 0、A 0 1、A 0 2、A 0 3、及び A 0 4 を示す) の第 1 の行に、行列 B (各位置 B_{ij} は、例えば、1 6 個の半精度浮動小数点値のベクトル又は「タイル」を垂直順に示す) の第 1 の列をペアとして乗算して計算される。つまり、C 0 0 は $A_{00} \cdot B_{00} + A_{01} \cdot B_{10} + A_{02} \cdot B_{20} + A_{03} \cdot B_{30} + A_{04} \cdot B_{40}$ を格納する。図 5 A、図 5 B、図 6 A、図 6 B、図 7 A、及び図 7 B において、シェーディング (s h a d i n g) 処理は共に乗算されるオペランドを識別するために使用する。より詳細には、同じパターンを使用してシェーディング処理された 2 つのオペランドが、図示する計算の一部として共に乗算される。本発明の実施形態の態様は、浮動小数点オペランドに対して浮動小数点演算を実行するように構成された I M C に関して説明しているが、本発明の実施形態はそれに限定されず、例えば、整数オペランドに対する整数演算を実行するように構成された I M C に適用してもよい。

【 0 0 7 3 】

図 5 A は、一つのオペランド (1 O P) データレイアウトにおける行列 A の第 1 の行と行列 B の第 1 の列との乗算の概略図である。ここで、本発明の一実施形態による一つのオペランドが外部から供給され、一つのオペランドがインメモリコンピューティングで D R A M バンクに格納される。図 5 A に示すように、行列 A の第 1 の行のタイル A 0 0、A 0 1、A 0 2、A 0 3、及び A 0 4 は、D R A M バンク 2 0 0 の同じページ 4 0 1 (行) に格納され、行列 B のベクトル又はタイル B 0 0、B 1 0、B 2 0、B 3 0、及び B 4 0 は、外部から供給され、結果 (例えば、C 0 0) は、D R A M バンク 2 0 0 の別のページ 4 0 2 に格納される。

【 0 0 7 4 】

乗算を計算するプロセスは、D R A M バンク 2 0 0 から値 A 0 0 を読み取り、 $A_{00} \cdot B_{00}$ を計算することによって開始する。これには、ページ 4 0 1 を開くことが含まれ、したがって、センス増幅器 2 3 2 を準備するためにプリチャージ (P R E) コマンドが必要となり、その後にページ 4 0 1 をセンス増幅器 2 3 2 にロードするための活性化 (A C T) コマンド、及びセンス増幅器 2 3 2 から I M C モジュール 2 5 0 に A 0 0 の値をロードするための読み取り (R E A D) コマンドが続く。上述のように、B 0 0 は外部から入力として提供されるため、この値の取得に D R A M 演算は必要でない。A L U 2 5 2 は、その後、乗算 $A_{00} \cdot B_{00}$ を計算し、出力バッファ (例えば、出力レジスタ R z) に一時的な結果を格納する。

【 0 0 7 5 】

次に、A L U 2 5 2 は、D R A M から A 0 1 を読み取ることによって乗算 $A_{01} \cdot B_{10}$ を計算する。これはまた、P R E コマンド、A C T コマンド、及び R E A D コマンドが必要である。ベクトル又はタイル B 1 0 は、外部から入力として提供されるため、 $A_{01} \cdot B_{10}$ が計算され、 $A_{00} \cdot B_{00}$ を格納するバッファ (例えば、出力レジスタ R z) の一時的な結果に追加される。このプロセスは、行列 A 及び行列 B の残りの値に対して繰り返される。その結果、それぞれの計算 (例えば、タイル A 0 0 及び B 0 0 等の 2 つのオペランドの乗算) は、計算ごとに 1 つの P R E、1 つの A C T、及び 1 つの R E A D を必

要とする。図 8 に関して以下でより詳細に説明するように、一部の実施形態では、IMC は値を格納し、先に格納された値と新たに受信した値の合計 (sum) で格納された値を更新するように構成されるアキュムレータをさらに含む。

【0076】

図 5 B は、1OP データレイアウトにおける行列 A の第 1 の行の第 1 の値と行列 B の各列の第 1 の値との積をデータ再利用と共に示した概略図である。ここで、本明細書の一実施形態によれば、一つのオペランドが外部から供給され、1つのオペランドがインメモリコンピューティングでDRAMバンクに格納される。図 5 B は、図 5 A に示すものとは異なり、DRAMからロードされたデータは、行列 B の異なる列に対して再利用される（例えば、オペランドレジスタ (Rop) 254 に格納される）。特に、2つの行列を乗算する際に、行列 A のすべての行の j 番目の要素は、行列 B の j 番目の行のすべての要素と乗算される。したがって、行列 A の各要素を一度ロードし、それを行列 B のすべての列 (N 列) に乗算することにより、DRAMバンク200からのデータロードのコストは、N 列にわたって償却 (amortized) される。より詳細には、行列 A 及び行列 B を乗算する場合、従来どおり、タイル A00 は、PRE コマンド、ACT コマンド、及び READ コマンドを使用して DRAM バンク 200 から読み取られ、タイル B00 は、外部から入力として受信される。ALU252 は、C00 の合計の一部を計算するために $A00 \cdot B00$ を計算する。しかし、DRAMからタイル A01 をロードする代わりに（例えば、別の PRE、ACT、及び READ シーケンスを使用）、C01 に対して計算する乗算の内の一つである $A00 \cdot B01$ を計算するために、タイル A00 を再利用して（外部から受信された）B01 を乗算する。その結果、各計算には、 $1/N$ PRE、 $1/N$ ACT が必要であり、計算ごとに 1 つの READ が必要である（行列 A の別の部分をロードするための PRE 及び ACT コマンドが、行列 B の N 列に対して償却されるため）。

【0077】

図 6 A は、同じ行 (SR) データレイアウトにおける行列 A の第 1 の行と行列 B の第 1 の列との乗算の概略図である。ここで、本発明の一実施形態によれば、双方のオペランドは、インメモリコンピューティングを有する DRAM バンクの同じページに格納される。図 6 A に示すように、行列 A の第 1 の行のタイル A00、A01、及び A02、並びに行列 B の第 1 の行のタイル B00、B10、及び B20 は、DRAMバンク200の同じページ411に格納され、行列 A のタイル A03 及び A04、並びに行列 B のタイル B30 及び B40 は、DRAMバンク200のページ412に格納される。その結果（例えば、C00）は、DRAMバンク200のページ413に格納される。

【0078】

結果（例えば、内積）を計算するプロセスは、DRAMバンク200からタイル A00 を読み取り、 $A00 \cdot B00$ を計算することによって開始する。これには、ページ411を開くことが含まれ、したがって、プリチャージ (PRE) コマンドが必要となり、その後には活性化 (ACT) コマンド、及びセンス増幅器 232 から IMC モジュール 250 にタイル A00 をロードするための読み取り (READ) コマンドが続く。タイル B00 は DRAM から読み取る。しかし、タイル B00 が A00 と同じページ411にあり、その値は既にセンス増幅器 232 に格納されているので、READ コマンドで十分である（タイル B00 を IMC モジュール 250 に読み取るために、さらに PRE 及び ACT を実行する必要はない）。したがって、タイル A00 及び B00 が読み取られると、ALU252 は $A00 \cdot B00$ を計算し、一時的な結果をバッファに格納する。同様に、タイル A01 とタイル B10 は、いずれもページ411にもあり、したがって、ページ411が A00 を読み取るために最初に開かれた際に、センス増幅器 232 に予め格納されているため、PRE と ACT なしで $A01 \cdot B10$ を計算するためのタイル A01 及び B10 の読み取りも同様に、READ コマンドを使用して実行する。そのため、各計算は、計算ごとに $1/r$ PRE、 $1/r$ ACT、及び 2 READ 演算を実行する。ここで、r は、DRAM バンク 200 の同じページに格納されている一致する値のペアの数である。例えば、上述のように、図 6 A は、行列 A の第 1 の行のタイル A00、A01、及び A02 と、行列 B

の第1の行のタイルB 0 0、B 1 0、及びB 2 0とが、D R A Mバンクのページ4 1 1に格納される場合を示す。したがって、計算A 0 0・B 0 0、A 0 1・B 1 0、及びA 0 2・B 2 0は、ページ4 1 1に3組のタイル（例えば、 $r = 3$ ）が含まれているため、それぞれ $1/3$ P R E、 $1/3$ A C T、及び2 R E A Dコマンドを償却する。A 0 3・B 3 0及びA 0 4・B 4 0を計算する際、ページ4 1 2に2組の値（例えば、 $r = 2$ ）が含まれているため、このような計算は、それぞれ $1/2$ P R E、 $1/2$ A C T、及び2 R E A Dコマンドをそれぞれ償却する。計算ごとに必要なP R E及びA C Tコマンドの数が減ることにより、全体の計算性能が向上する。

【0079】

図6 Bは、S Rデータレイアウトにおける行列Aの第1の行の第1の値と行列Bの各列の第1の値に対するデータの再利用による乗算の概略図である。ここで、本発明の一実施形態による双方のオペランドは、インメモリコンピューティングを有するD R A Mバンクの同じページに格納される。図5 A及び図5 Bの配列間の比較と同様の方法で、図6 Bに示す計算プロセスは、D R A Mからロードされた値が再利用されるという点で図6 Aに示すものと異なる。より詳細には、行列の乗算は、第1のオペランドの指定された行のi番目の要素と第2のオペランドの各列のi番目の要素との乗算を含むということに基づいて、行列Aと行列Bの要素は、指定された行列Aの行の各要素のデータが、乗算される行列Bの行の値と同じページに配置されるように、D R A Mバンク2 0 0に配列される。

【0080】

例えば、図6 Bに示すように、行列Aに行列Bを乗算すると、行列AのタイルA 0 0は、結果行列Cの第1の行の一部（例えば、C 0 0、C 0 1、C 0 2、C 0 3、C 0 4、C 0 5等の項のいずれか）を計算するプロセスにおいて、行列Bのすべての列の第1のタイル（行列Bの第1の行のすべての要素、例えば、B 0 0、B 0 1、B 0 2、B 0 3、B 0 4、B 0 5、...）に乗算される。同様に、行列Aの値A 0 1に行列Bのすべての列の第2のタイルが乗算される（図6 Bに示すように、行列Bの第2の行のすべての要素、例えば、B 1 0、B 1 1、B 1 2、B 1 3、B 1 4、B 1 5、...）。

【0081】

このように、行列Aの少なくとも一つの値は、行列Bの対応値と同じページに格納される。図6 Bに示す特定の例において、ページ4 1 1は、行列AからタイルA 0 0、行列BからタイルB 0 0、B 0 1、B 0 2、B 0 3、B 0 4、B 0 5、...を格納し、ページ4 1 2は、行列AからタイルA 0 1、行列BからタイルB 1 0、B 1 1、B 1 2、B 1 3、B 1 4、B 1 5、...を格納する。D R A M内のデータのこのような配列により、D R A Mのページが少なくとも $N + 1$ エントリを格納すると仮定すると、各計算は、 $1/N$ P R Eコマンド、 $1/N$ A C Tコマンド、及び $(N + 1)/N$ R E A Dコマンド（P R Eコマンド及びA C Tコマンドは、行列BのN列で償却されるため）を実行する。

【0082】

図7 Aは、異なる行（D R）データレイアウトにおける行列Aの第1の行と行列Bの第1の列との乗算の概略図である。ここで、オペランドは、本発明の一実施形態によるインメモリコンピューティングでD R A Mバンクの異なるページに格納される。図7 Aに示すように、行列Aの第1の行のタイルA 0 0、A 0 1、A 0 2、A 0 3、及びA 0 4は、ページ4 2 1に格納される一方、行列Bの第1の列のタイルB 0 0、B 1 0、B 2 0、B 3 0、及びB 4 0は、別のページ4 2 2に格納される。

【0083】

タイルB 0 0がタイルA 0 0とは異なるページ4 2 2（行）にあるため、C 0 0を計算するプロセスは、P R E、A C T、及びR E A Dコマンドのシーケンスを使用して、タイルB 0 0を読み取り、その後に、P R E、A C T、及びR E A Dコマンドのシーケンスを使用して、D R A Mバンク2 0 0のページ4 2 1（行）からA 0 0を読み取ることによって、A 0 0・B 0 0を計算することから始まる。A L U 2 5 2は、A 0 0・B 0 0を計算し、その結果を一時的なバッファに格納する。C 0 0計算を継続するために、タイルA 0 1は、ページ4 2 1（行）から読み取られ、ページ4 2 2（行）からB 1 0が読み取られ

10

20

30

40

50

、ここで、それぞれの値は、D R A Mバンク 2 0 0 に対する P R E、A C T、及び R E A Dシーケンスの実行を含む。

【 0 0 8 4 】

図 7 B は、D R データレイアウトにおける行列 B の各列の第 1 の値と行列 A の第 1 の行の第 2 の値とにより、行列 B の各列の第 2 の値による行列 A の第 1 の行の第 1 の値のデータ再利用による乗算の概略図である。ここで、オペランドは、本発明の一実施例によるインメモリコンピューティングで D R A Mバンクの異なるページに格納される。

【 0 0 8 5 】

図 5 B における説明及び 1 O P におけるデータ再利用と同様の方法により、行列 B で乗算を実行する際に行列 A から取り出したデータを再利用することで、メモリ演算の数が減る。図 7 A における説明のように、指定された行列 A の行の各 i 番目の値は、行列 B の i 番目の行の各値と乗算される。したがって、指定された行列 B の行のすべての値が D R A Mバンク 2 0 0 の同じページに格納されると、P R E 及び A C T コマンドの数が減少し、それによって性能が向上する。

【 0 0 8 6 】

例えば、行列 A に行列 B を乗算するプロセスは、 N 個の部分和（例えば、 $C 0 0$ 、 $C 0 1$ 、 $C 0 2$ 、 \dots 、 $C 0 N$ の一部）を計算するために、行列 A のタイル $A 0 0$ に行列 B の第 1 の行の N 個のタイル（ $B 0 0$ 、 $B 0 1$ 、 $B 0 2$ 、 \dots 、 $B 0 N$ ）のそれぞれを乗算することから始まる。このプロセスは、D R A Mバンク 2 0 0 からタイル $A 0 0$ をロードすることにより始まる。これには、ページ 4 2 1（行）を開き、行列 A のタイル $A 0 0$ をオペランドレジスタ（R o p）2 5 4 にロードするため、P R E コマンド、A C T コマンド、及び R E A D コマンドが含まれる。 $A 0 0$ をロードした後、行列 B の第 1 の行のタイル $B 0 0$ 、 $B 0 1$ 、 $B 0 2$ 、 \dots 、 $B 0 N$ がタイル $A 0 0$ に乗算されるようロードされる。図 7 B に示すように、行列 B のこれらの値がすべて同じページにある場合（例えば、 $B 0 0$ 、 $B 0 1$ 、 $B 0 2$ 、 $B 0 3$ 、 $B 0 4$ 、 \dots は、図 7 B に示す D R A Mバンク 2 0 0 のページ 4 2 2 にある）、このページにアクセスするために使用された P R E 及び A C T コマンドが、そのページに格納された N 個の値で償却される。したがって、行列 A の一つの値に行列 B の一つの行のすべての値を乗算するには、2 つの P R E コマンド、2 つの A C T コマンド、及び $N + 1$ R E A D コマンドが必要となり、行の N 個の値を償却すると、計算ごとに $2 / N$ P R E コマンド、 $2 / N$ A C T コマンド、及び $(N + 1) / N$ R E A D コマンドが生成される。上述の例と同様の方法で、D R A M コマンドの数の減少は、全体の計算性能（例えば、行列乗算演算）を向上させる。

【 0 0 8 7 】

上述のデータ配置オプションは、サイクルレベル（c y c l e - l e v e l）の高帯域幅メモリ - メモリ内機能（H B M - F I M）シミュレータを使用して実験的にテストし、これは G E M M 記録と共に I M C 用にカスタマイズされた。実験アーキテクチャには 4 つの 4 H H B M 2 モジュールを含む H B M が含まれており、ここで D R A M のアレイサイズは 1 6 , 3 8 4（1 6 K i b）行 \times 8 , 1 9 2（8 K i b）列（例えば、各ページのサイズは 8 , 1 9 2 b）であり、行バッファのサイズは、8 , 1 9 2 ビット（8 K i b）である。I M C モジュールは、待ち時間周期を備えた 1 6 レーンの F P - 1 6 ベクトルユニットと 7 6 8 ビットバッファを含み、パイプライン動作が可能であった。様々な実施形態において、バッファはより小さくてもよく（例えば、2 5 6 ビットの入力バッファ及び 2 5 6 ビットの出力バッファを含む 5 1 2 ビットのバッファ）、又はより大きくてもよい（例えば、2 5 6 ビットの入力バッファ及び 1 , 0 2 4 ビットの出力バッファを備えた 1 , 2 8 0 ビットのバッファ）。一部の実施形態において、入力バッファは、2 5 6 ビットよりも大きい（例えば、5 1 2 ビット）。

【 0 0 8 8 】

本発明の実施形態の一部の態様は、D R A Mバンク 2 0 0 のレベルでさらなるバッファ及びアキュムレータ（例えば、D R A Mバンクでバンクごとに提供されるさらなるハードウェア）を含むことに関する。

10

20

30

40

50

【 0 0 8 9 】

図 8 は、D R A Mバンクの I M Cモジュールの概略ブロック図である。ここで、I M Cモジュールは、本発明の一実施形態による結果バッファ、アキュムレータ、及びバッファにさらに接続される。

【 0 0 9 0 】

図 8 に示すように、A L U 2 5 2 は、入力オペランド A (図 8 では 2 5 6 ビットのオペランド A [0 : 2 5 5] と表示) 及び入力オペランド B (図 8 では 2 5 6 ビットのオペランド B [0 : 2 5 5] と表示) を受信する。A L U は、2 つの入力オペランド (例えば、加算、乗算、内積、外積など) に対する演算を実行し、結果 C (図 8 では 2 5 6 ビットの結果 C [0 : 2 5 5] と表示) を計算する。

【 0 0 9 1 】

図 8 を参照すると、I M Cモジュール 2 5 0 は、A L U 2 5 2 の出力に接続されたアキュムレータ 8 0 2 をさらに含む。例えば、A L U 2 5 2 の出力は、結果レジスタ (R z) 2 5 6 及びアキュムレータ 8 0 2 の双方に接続され、結果 C が結果レジスタ (R z) 2 5 6 に格納され、アキュムレータ 8 0 2 に供給される。アキュムレータ 8 0 2 は、累積値 (例えば、2 5 6 ビット値) を格納するアキュムレータレジスタを含む。アキュムレータ 8 0 2 が A L U 2 5 2 から新たな結果を受信すると、新たな結果は、アキュムレータレジスタに既に格納されている累積値に追加 (例えば、累算) される (例えば、アキュムレータ 8 0 2 のアキュムレータレジスタに格納されている値が更新されるか、又は新たな結果と先にアキュムレータ 8 0 2 に格納されている値の合計として設定される)。一部の実施形態において、アキュムレータ 8 0 2 は、リセットコマンドに応答して、アキュムレータレジスタをリセットする (例えば、アキュムレータレジスタに格納された累積値をゼロに設定する) ように構成される。アキュムレータ 8 0 2 は、行列の乗算を計算する場合 (例えば、結果行列の各値が第 1 のオペランドの行と第 2 のオペランドの列の内積である場合) 等、内積 (inner products 又は dot products) の計算の際に特に役立つ。図 8 に示す実施形態において、I M Cモジュール 2 5 0 は、5 1 2 ビットの第 1 のバッファ 8 1 2 及び 5 1 2 ビットの第 2 のバッファ 8 1 4 として示される 1 , 0 2 4 ビットの出力バッファをさらに含む。多数の結果値を格納するさらに大きな出力バッファは、メモリコントローラ 1 4 0 が D R A Mバンクの別のページを開いてその結果を格納する前に、一度に多数の結果を計算するために、D R A Mバンクを制御できるようにする。例えば、図 6 B に示すように S R データレイアウトでデータを再利用する場合、結果を格納するためにページを開くには、P R E 及び A C T コマンドをさらに必要とするオペランドを格納するページ (例えば、図 6 A 及び図 6 B に示すページ 4 1 1) とは異なる D R A Mバンク 2 0 0 の別のページ (例えば、図 6 A 及び図 6 B に示すページ 4 1 3) に書き込まれた部分和を計算するために、第 1 のオペランド行列 A の一つの値に、第 2 のオペランド行列 B の行に対応する異なる値を乗算する。しかし、出力バッファが大きいほど、演算を完了するのに必要なページ切り替え量が減り (各ページの切り替えには P R E 及び A C T が必要なため)、そのため、計算性能が向上する。図 8 は、アキュムレータ 8 0 2 及びより大きな出力バッファ (8 1 2、8 1 4) の双方を有する I M Cモジュール 2 5 0 を示すが、本発明の実施形態はそれに限定されず、I M Cモジュール 2 5 0 がさらなる出力バッファ 8 1 2 及び 8 1 4 なしでアキュムレータ 8 0 2 を含む実施形態、並びに I M Cモジュール 2 5 0 がアキュムレータ 8 0 2 なしでさらなる出力バッファ 8 1 2 及び 8 1 4 を含む実施形態を含む。本発明の一部の実施形態において、I M Cモジュール 2 5 0 は、(例えば、多数の値が同時に蓄積されるように) 並列に配列された多数のアキュムレータ 8 0 2 を含む。本発明の実施形態はさらに、2 つの 5 1 2 ビットの出力バッファ 8 1 2 及び 8 1 4 を有する場合に限定されず、5 1 2 ビットよりも大きい又は小さい、及び / 又は 2 つ以上の出力バッファ又は 2 つ未満の出力バッファを含む。I M Cモジュール 2 5 0 の様々な部分は、異なる量の待ち時間を有する。例えば、乗算演算を実行する A L U 2 5 2 が 4 サイクルの待ち時間を導入し、累積演算は 1 サイクルの待ち時間を含む。

【 0 0 9 2 】

本発明の実施形態の態様は、GEMMを実装することに関して、以下でより詳細に説明する。

【0093】

図9は、本発明の一部の実施形態による、同じ行(SR)レイアウトを使用する計算を説明するためのGEMMの一例の概略図である。より具体的には、図9は、行列A及び行列Bの積を示し、行列A及び行列Bの積をアキュムレータ802に現在格納されている値に追加することにより、アキュムレータ(例えば、アキュムレータ802)に格納された結果Cを更新することを示す(結果C += 行列A × 行列B)。図9に示す計算では、行列Aの16タイルと行列Bの16タイルが共に乗算され、ここで、各タイルは、16個のFP-16要素(256ビット)を有する。特定の内部タイルの構成レイアウトは、以下でより詳細に説明するが、GEMMの全体的な複雑さ(complexity)は、別のレイアウトでも同様である。図9に示すように、同様のパターンでシェーディング処理されたタイルは、共に乗算されたタイルである。上述のように、図9は、同じ行(SR)データレイアウトを示し、ここで、行列A及び行列Bの32個の値すべてが、DRAMバンク200の同じページ414に格納され、各計算のオペランド(例えば、タイルA00、並びにタイルB00、B01、B02、及びB03)は、すべて同じページ414に格納される。個々の計算の結果は、結果レジスタ(Rz)256に蓄積され、結果Cを計算する。

【0094】

タイル-レベル(tile-level)乗算の一つの特定例として、DRAMバンク200の同じページからロードされた2つのタイルA00及びB00の内積(inner products又はdot products)を計算する際、タイルA00は、行列Bの点線タイルB00である。一部の実施形態において、上述のように、各タイルは、16個のFP-16要素を含む。例えば、タイルA00は要素a00、a01、...、a15を含み、タイルB00は要素b00、b01、...、b15を含み、タイルA00及びB00の内積は積のペアの合計(= a00 × b00 + a01 × b01 + ... + a15 × b15)である。したがって、2つのタイルの内積(inner products又はdot products)は、単一の要素又は単一の値を生成することから、本発明の実施形態によるIMCモジュールは、結果を格納するために、より少ないバッファレジスタ(例えば、アキュムレータレジスタ802)で内積を計算する。そのため、アキュムレータレジスタを含む本発明の一部の実施形態によるIMCモジュールは、内積の計算等、値の累算を含む計算を実行するのに非常に適している。一部の実施形態において、内積の計算に適したアキュムレータを含むIMCモジュールは、学習済みのニューラルネットワークを使用して推論(又は順方向伝搬)を実行するために使用され、学習済みのニューラルネットワークで第1のオペランドが入力(例えば、以前の層からの活性化)を示し、第2のオペランドは、学習済みのニューラルネットワークの層のニューロンに関連する重みを示す。

【0095】

タイル-レベル乗算の別の例として、DRAMバンク200の同じページからロードされた2つのタイルの外積を計算する場合、第1のタイルのすべての値は、第2のタイルのすべての値に乗算され、2つのタイルを同時にロードする。例えば、上述のように、各タイルが16個の値を含む場合、2つのタイルの外積は16 × 16 = 256の出力値を有する。ツリー-加算器(tree-adder)の累積が不要であり、256個の出力値を並列に計算するため、外積は内積よりも簡単なハードウェアで計算することができる。しかし、外積の計算には、(例えば、外積の結果のすべての値を格納するために)本発明の実施形態によるIMCモジュール250において多くの出力バッファレジスタ(例えば、出力バッファ812及び814)を必要とする。例えば、各タイルが最大16個のFP-16値を含むと仮定すると、外積を計算するように構成されたIMCモジュール250は、256個の値を格納するのに十分な大きさの出力バッファ(例えば、256 × 16ビット = 4,096ビットのバッファ)を含む。したがって、出力バッファを含む本発明の一部の実施形態によるIMCモジュールは、第1および第2のオペランド(上記の例では、

10

20

30

40

50

オペランドの各タイルは16個の値を含むため、出力バッファには少なくとも256個の値が格納される)のタイル内の値の数の積以上のサイズを格納する大きさを有するオペランドの外積を含む計算を実行するのに適する。

【0096】

タイル・レベル乗算の第3の例として、DRAMバンク200の同じページからロードされた2つのタイルのテンソル積を計算する場合、2つのタイルは、行列に配列された値を含む。例えば、各タイルが16個のFP-16値を含む場合、各タイルは、 4×4 行列値に配列される。このタイルのテンソル積を計算すると、 4×4 の結果行列が生成される。例えば、各タイルが最大16個のFP-16値を含むと仮定すると、16値の内の2つの 4×4 タイルのテンソル積を計算するように構成されたIMCモジュール250は、16個の値を格納するのに十分な大きさの出力バッファ(例えば、 16×16 ビット=256ビットのバッファ)を有する。したがって、本発明の実施形態によるIMCモジュールは、第1および第2のオペランドの内の大きい方の値の数と同じ値の数を格納するのに十分な大きさの出力バッファを含む計算(オペランドのテンソル積を含む)を実行するのに適する。本発明の実施形態によるテンソル積を実装するIMCモジュールは、内積の計算及び外積の計算に適したIMCモジュールハードウェアに比べて、より簡単なハードウェア及び中間数のバッファレジスタを備えた中間グラウンド(middle ground)を示す。

10

【0097】

より詳細には、2つの 4×4 行列AとBとの間でテンソル積を実装することは、64の乗算演算、48の加算演算、及び行列A/行列Bでの転置(transpose)を含む。本発明の一実施形態による16レーンのe-ALU252を使用する場合、16個の演算が並列に実行される。これらの $64 + 48 = 112$ の演算は、7サイクルのALU(一回のサイクルで乗算を実行すると仮定)で実行される。本発明の別の実施形態において、性能を向上させるために、64レーンのALU252は、2~3サイクル又は4~6ナノ秒(nanoseconds)で2つの 4×4 行列を表す2つのタイルのテンソル積の計算に使用される。したがって、ALU252でレーンの数を増加させることは、計算間で可能な並列化の量を増加させ、性能を向上させる。

20

【0098】

したがって、本発明の実施形態の態様は、インメモリコンピューティングを備えたDRAMシステムにおけるデータ配置のためのシステム及び方法に関する。SR(単一ページ)データレイアウトに関して、内積、外積、及びテンソル積を計算する上記の3つの例において、単一の行は、(行列の)列データを有する第2のタイルと同じページに(行列の)行データを有する第1のタイルを含む。

30

【0099】

本発明の一部の実施形態によると、メモリモジュール110のメモリコントローラ140(又はクライアント側メモリコントローラ)は、ホストプロセッサ170から受信した命令に従って、ホストプロセッサ170から受信したデータ(オペランド)の配置を制御する。例えば、命令が内積、外積、又はテンソル積を計算すべきか否かに基づいて、及びオペランドのサイズ(例えば、行列の寸法(dimensions))に基づいて、メモリコントローラ140は、1OP、SR又はDRデータレイアウトを使用し、それに応じてデータを格納するようにDRAMバンクを制御する。

40

【0100】

図10は、本発明の実施形態による、IMCモジュールを有するDRAMバンク内のデータの配置を制御するための方法を示すフローチャートである。ステップS1110において、メモリコントローラ140は、オペランド(例えば、第1のオペランド及び第2のオペランド)と、オペランドに適用される命令(例えば、関数)とを受信する。例えば、これらの命令は、内積、外積の計算、学習済みのニューラルネットワーク及び入力ベクトルを使用した推論の実行、ニューラルネットワークの学習のための逆伝搬{でんぱん}アルゴリズム(backpropagation algorithm)のステップの実行

50

等を含む。

【 0 1 0 1 】

ステップ S 1 1 3 0 において、メモリコントローラ 1 4 0 は、命令に基づいて、D R A Mダイ上の D R A Mバンクに一方又は両方のオペランドを格納するために使用するデータレイアウト（例えば、1 O P、S R、又は D R）を決定する。一部の実施形態において、（例えば、ホストメモリコントローラ 1 8 0 からの）命令は、どのデータレイアウトを使用するかを明示して指定する。一部の実施形態において、実行する計算のタイプ、及び D R A Mバンクにおけるそのような計算を加速するためのメモリ又はハードウェア要件に基づいて、レイアウトが選択される（例えば、内積を計算するための命令は、アキュムレータを有する D R A Mバンクにデータが配置されるのに対し、外積を計算する命令は、より多くの出力バッファを有する D R A Mバンクにデータが配置される）。

10

【 0 1 0 2 】

ステップ S 1 1 5 0 において、メモリコントローラ 1 4 0 は、選択されたデータレイアウトに基づいて、第 1 のオペランド及び第 2 のオペランドを D R A Mバンクに供給する。一例として、1 O P の場合、メモリコントローラ 1 4 0 は、第 1 のオペランドの少なくとも第 1 のタイルを格納し、第 2 のオペランドの第 2 のタイルを D R A Mバンクの I M Cモジュールに直接供給するように、D R A Mバンク 2 0 0 を制御する。他の例として、S R の場合、メモリコントローラ 1 4 0 は、D R A Mバンクの同じ行又は同じページに、第 1 および第 2 のオペランドに対応するタイルを格納するように、D R A Mバンク 2 0 0 を制御する。

20

【 0 1 0 3 】

ステップ S 1 1 7 0 において、メモリコントローラ 1 4 0 は、入力命令に基づいた演算を実行するように、D R A Mバンク 2 0 0 の I M Cモジュールを制御する。例えば、学習済みのモデルを使用して推論を実行する命令の場合、演算は、一つのオペランドに基づいて入力ベクトルを用意し、第 2 のオペランドに格納されたパラメータに基づいて入力ベクトルの値に重みを付けることが含まれる。

【 0 1 0 4 】

本明細書では、特定の例示的な実施形態を提示しているが、本発明は、開示された実施形態に限定されず、本発明の思想や技術範囲を逸脱しない範囲内で多様に変更実施することが可能である。

30

【 符号の説明 】

【 0 1 0 5 】

- 1 0 0 : メモリシステム
- 1 1 0 : メモリモジュール
- 1 2 0 : D R A Mダイ
- 1 3 0 : 内部メモリバス
- 1 4 0 : メモリコントローラ
- 1 7 0 : ホストプロセッサ
- 1 8 0 : ホストメモリコントローラ
- 1 9 0 : 外部バス
- 2 0 0 : D R A Mバンク
- 2 1 0 : D R A Mセル
- 2 1 2 : コンデンサ
- 2 1 4 : スイッチ
- 2 2 0 : 行デコーダ
- 2 3 0 : 入出力センス増幅器層 (I O S A)
- 2 3 2 : センス増幅器
- 2 3 4 : マルチプレクサ
- 2 3 6 : グローバル I O 層
- 2 4 0 : 列デコーダ

40

50

250 : IMC モジュール (ALU & Reg)

252 : ALU

254 : オペランドレジスタ (Rop)

256 : 結果レジスタ (Rz)

257 : 第1のマルチプレクサ

258 : 第2のマルチプレクサ

260 : 入出力 (IO) モジュール

300 : マルチプレクサ

401、402、411、412、413、421、422、423 ページ

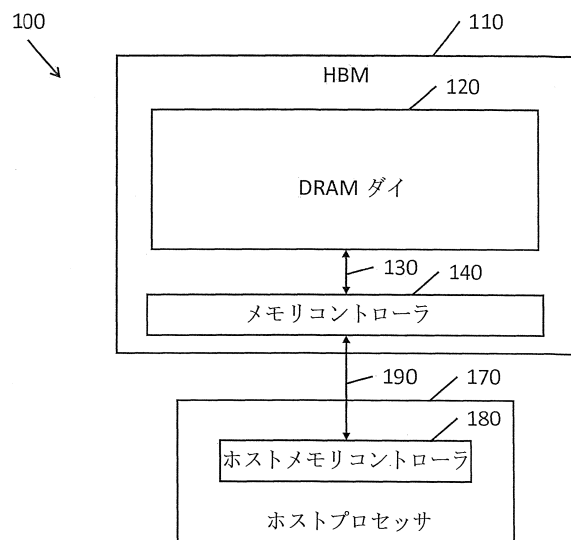
802 アキュムレータ

812 第1のバッファ

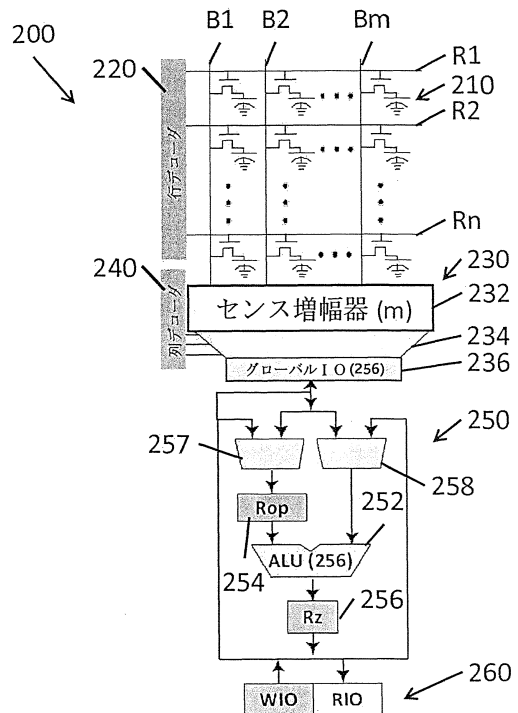
814 第2のバッファ

【図面】

【図1】



【図2A】



10

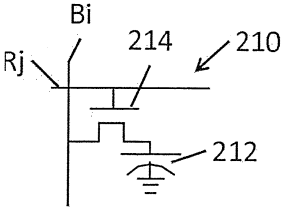
20

30

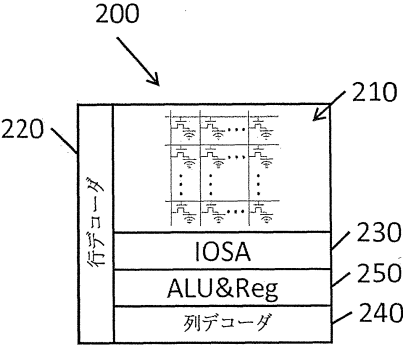
40

50

【図 2 B】

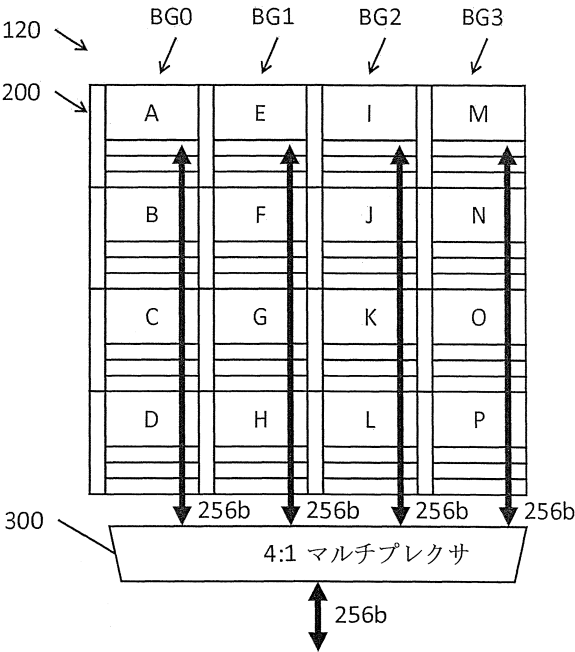


【図 2 C】

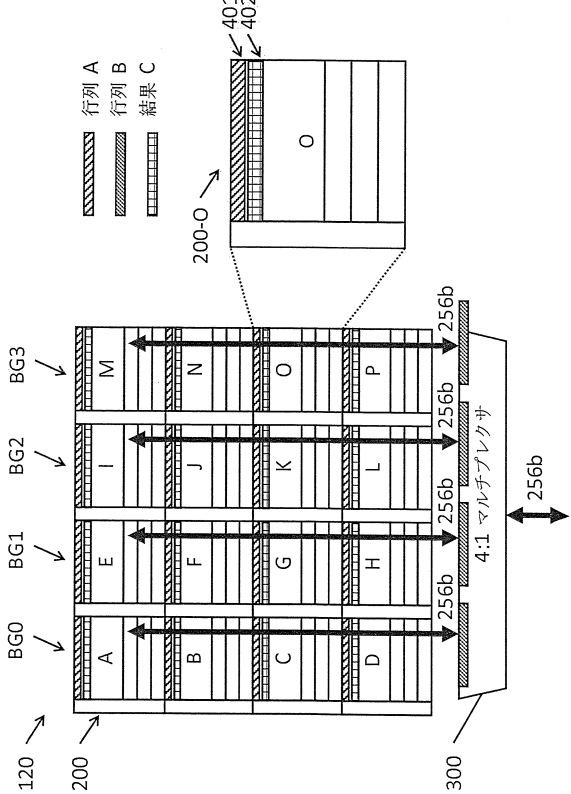


10

【図 3】



【図 4 A】



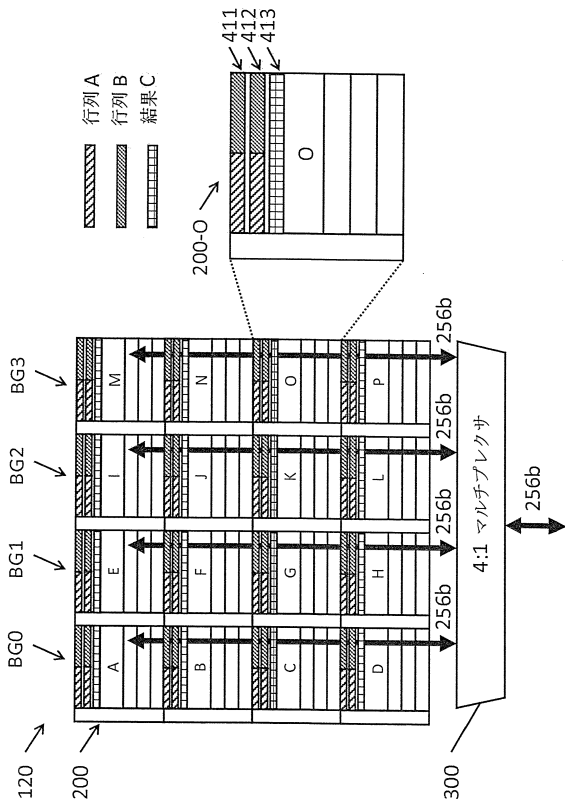
20

30

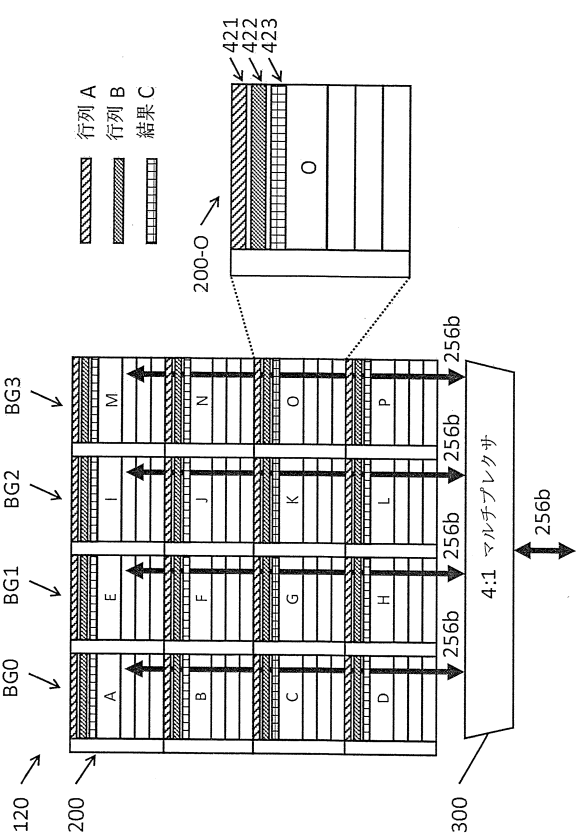
40

50

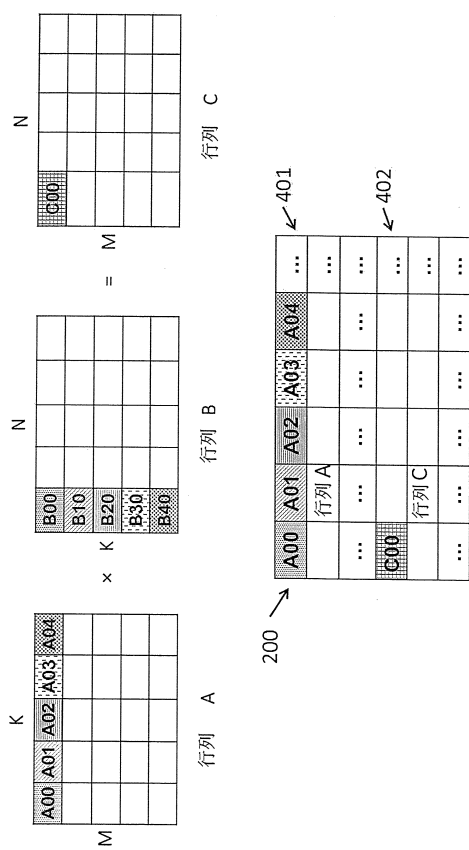
【図 4 B】



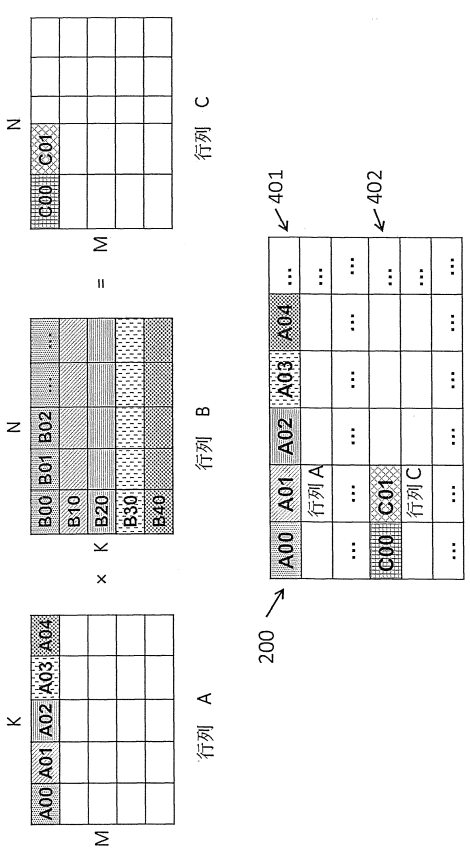
【図 4 C】



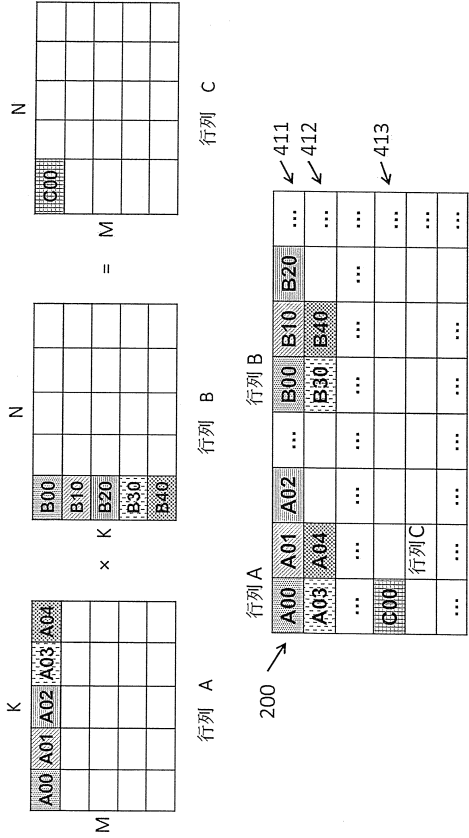
【図 5 A】



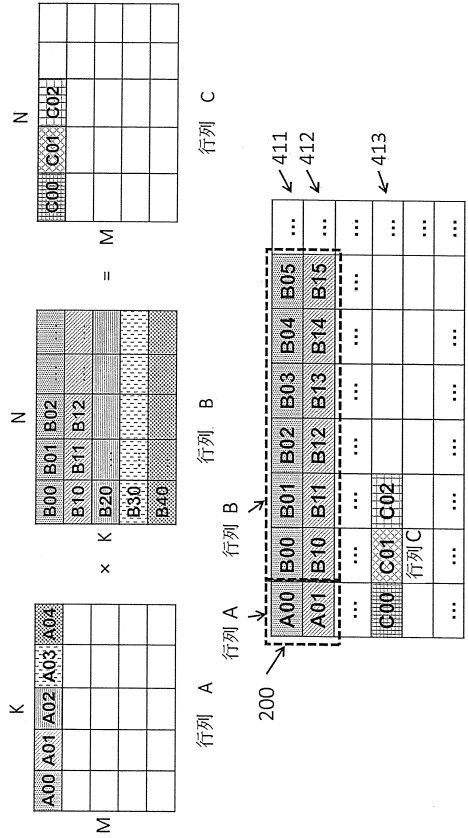
【図 5 B】



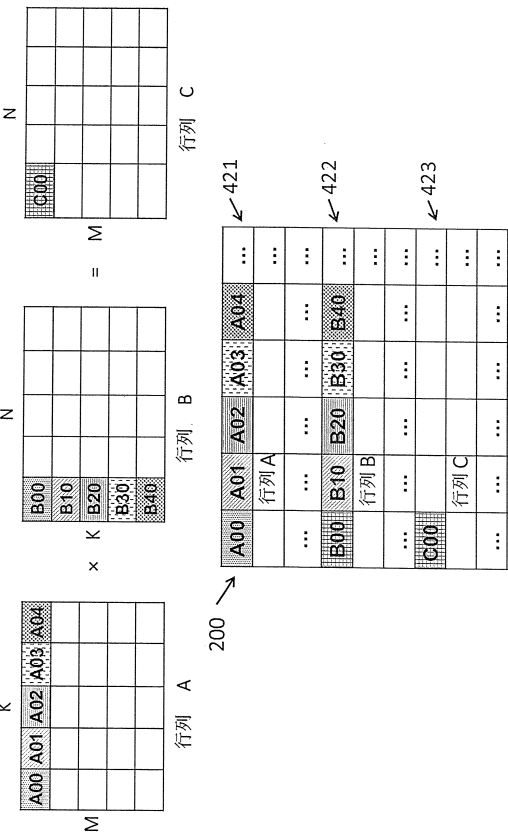
【図 6 A】



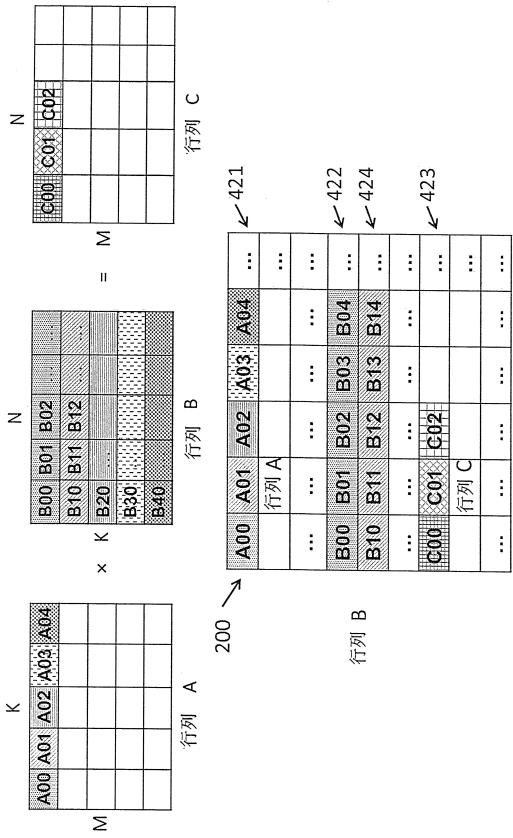
【図 6 B】



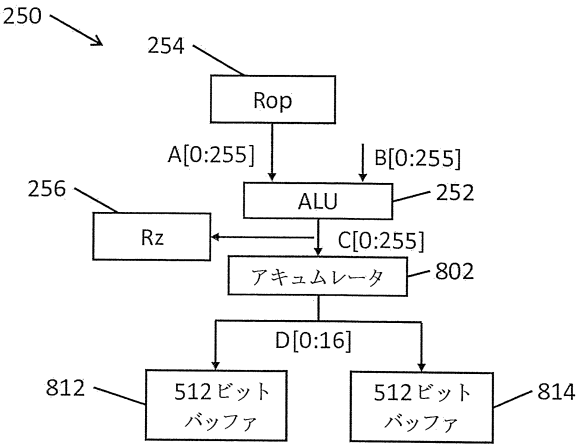
【図 7 A】



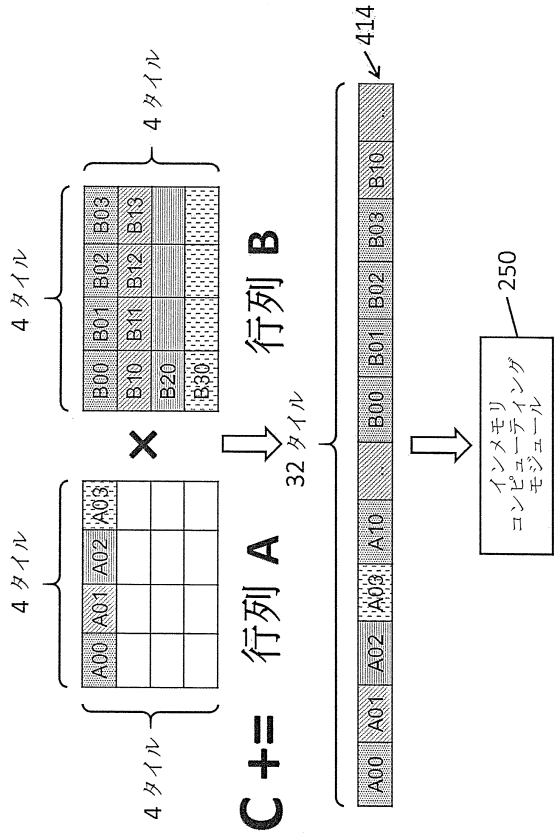
【図 7 B】



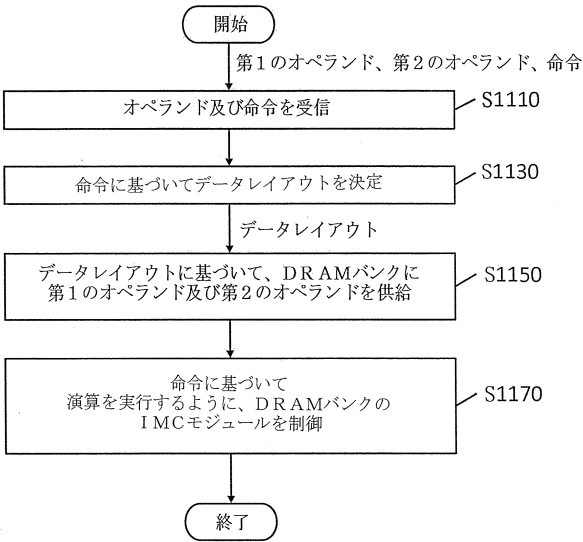
【図 8】



【図 9】



【図 10】



10

20

30

40

50

フロントページの続き

(51)国際特許分類

H 1 0 B

99/00 (2023.01)

F I

H 1 0 B

12/00

6 8 1 F

H 1 0 B

99/00

4 9 5

ライブ 4 1 9 6

(72)発明者

皇甫文沁

アメリカ合衆国, 9 5 1 3 1 カリフォルニア州, サンノゼ, ベンダー プレイス 1 6 1 0

審査官 北村 学

(56)参考文献

特開平 0 6 - 2 1 5 1 6 0 (J P , A)

特開 2 0 0 3 - 2 7 2 3 8 4 (J P , A)

特開 2 0 0 8 - 1 2 3 4 7 9 (J P , A)

特開 2 0 1 9 - 0 2 8 5 7 2 (J P , A)

特開 2 0 1 9 - 0 7 5 1 0 1 (J P , A)

米国特許第 0 5 9 5 3 7 3 8 (U S , A)

米国特許出願公開第 2 0 1 7 / 0 2 5 5 3 9 0 (U S , A 1)

(58)調査した分野 (Int.Cl., D B 名)

G 0 6 F 1 2 / 0 0

G 0 6 F 9 / 3 4

G 0 6 F 1 2 / 0 6

G 1 1 C 5 / 0 4

G 1 1 C 1 1 / 4 0 9 3

H 1 0 B 1 2 / 0 0

H 1 0 B 9 9 / 0 0